

# Parametric Visual Program Induction with Function Modularization

Xuguang Duan<sup>1</sup> Xin Wang<sup>1</sup> Ziwei Zhang<sup>1</sup> Wenwu Zhu<sup>1</sup>

## Abstract

Generating programs to describe visual observations has gained much research attention recently. However, most of the existing approaches are based on non-parametric primitive functions, making them unable to handle complex visual scenes involving many attributes and details. In this paper, we propose the concept of parametric visual program induction. Learning to generate parametric programs for visual scenes is challenging due to the huge number of function variants and the complex function correlations. To solve these challenges, we propose the method of function modularization, capable of dealing with numerous function variants and complex correlations. Specifically, we model each parametric function as a multi-head self-contained neural module to cover different function variants. Moreover, to eliminate the complex correlations between functions, we propose the hierarchical heterogeneous Monto-Carlo tree search (H2MCTS) algorithm which can provide high-quality uncorrelated supervision during training, and serve as an efficient searching technique during testing. We demonstrate the superiority of the proposed method on three visual program induction datasets involving parametric primitive functions. Experimental results show that our proposed model is able to significantly outperform the state-of-the-art baseline methods in terms of generating accurate programs.

## 1. Introduction

Studying how to generate computer-executable programs is one of the core interests of the AI community (Waldinger & Lee, 1969; Manna & Waldinger, 1975), and has drawn

<sup>1</sup>Department of Computer Science and Technology, Tsinghua University, Beijing, China. Correspondence to: Xin Wang <xin\_wang@tsinghua.edu.cn>, Wenwu Zhu <wwzhu@tsinghua.edu.cn>.

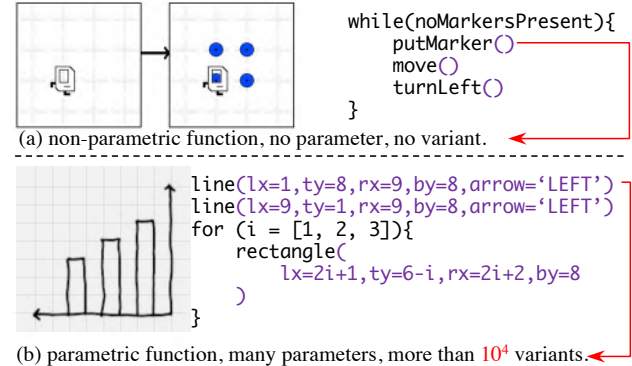


Figure 1. (a): An example of the visual program induction task that only generates non-parametric programs, within which, each function has only one variant, and could be modeled as a symbolic token. (b): An example of the parametric visual program induction task studied in this paper, where parametric primitive functions with many more variants are needed to describe the complex visual scene. However, it is hard to tackle such many function variants.

lots of recent interests in the visual domain thanks to deep learning (Ellis et al., 2020). By leveraging powerful deep models, these works can successfully describe the logic behind visual games (Sun et al., 2018), learn spatial patterns hidden in images (Young et al., 2019), or conduct neural-symbolic reasoning (Yi et al., 2018).

Despite their enormous success, most of the existing approaches are based on non-parametric primitive functions, failing to meet the requirement of the increasing complexity of visual observations, as well as the increasing elaboration of programs. In this paper, we are the first to propose the concept of **Parametric Visual Program Induction**, *i.e.*, *generating programs with parametric primitive functions for complex visual observations*, to the best of our knowledge. By leveraging parametric primitive functions, we can generate much more detailed programs to describe both the hidden logic and visual details.

However, the challenges for solving parametric program induction are two folds. First, the action space for a single function can be huge. Compared with basic non-parametric primitive functions, the parametric primitive functions always have several heterogeneous parameters, resulting in a huge number of function variants. For example, in Figure 1(a), a basic visual program induction task may contain simple primitive functions such as `move()`,

`turnLeft()`; while in Figure 1(b), a parametric function studied in this work tend to have more than  $10^4$  variants due to different parameter combinations.

Second, the function space for the whole program is also very huge. Given that parametric functions may contain multiple parameters, and these parameters and functions are correlated together, it becomes very challenging to model the long-range function transitions within a program. This problem is also known as *program aliasing* (Bunel et al., 2018) in the non-parametric scenario, and becomes more severe for parametric functions.

These two challenges make non-parametric visual program induction methods hard to extend to the parametric domain. To address these challenges, we propose the concept and method of **Function Modularization**, which can model numerous and complex parametric functions. In particular, we treat each function along with its parameters as a self-contained module and learn the module to predict the correct parameters given visual contexts, which is able to solve the challenge of the huge action space. Furthermore, based on the modularized functions, we propose a Hierarchical Heterogeneous Monto-Carlo Tree-Search (H2MCTS) algorithm that can traversal all the program aliases, thus providing uncorrelated training data during training and serving as a powerful search method during inference. To verify the superiority of the concept of **function modularization** and the efficiency of the H2MCTS algorithm, we conduct extensive experiments on a small hand-craft dataset and two well-known datasets (Ellis et al., 2018; Dong et al., 2019). Experimental results show that a modularized function is easier to learn and has higher accuracy compared with vanilla baselines. Also, the proposed H2MCTS algorithm is able to efficiently search over different function combinations and reduce the inference time significantly. In summary, we make the following contributions:

- To the best of our knowledge, we are the first to investigate the problem of parametric visual program induction by proposing the concept and method of **Function Modularization**, which decouples the learning of function parameters and function transitions, resulting in accurate and efficient learning of the parametric programs.
- We propose the H2MCTS algorithm to assist the learning of modularized functions. Our proposed algorithm can provide uncorrelated data to train modularized functions and serve as an efficient search method during inference.
- We conduct extensive experiments to demonstrate that our proposed model can significantly outperform state-of-the-art baselines on all three datasets.

## 2. Related Work

Learning to generate programs has a long history in AI (Waldinger & Lee, 1969; Manna & Waldinger, 1975; 1980). Traditionally, the process of generating programs is based on search-based induction, and one of the most famous works is the Excel *FlashFill* system (Gulwani, 2011). These methods rely on syntax-based pruning (Feser et al., 2015), or use satisfiability modulo theories-based solvers (Lezama, 2008; Feser et al., 2015). With the development of deep learning, this area has gained new attention as learning to generate a program from data directly (Parisotto et al., 2017; Devlin et al., 2017; Ling et al., 2017; Chollet, 2019), including previously unsolvable visual domain tasks (Bunel et al., 2018; Sun et al., 2018; Shin et al., 2019). Besides, the combination of search and learning is also appealing by leveraging advantages from both sides by combining learning and searching (Balog et al., 2016; Irving et al., 2016; Ellis et al., 2020). Balog et al. (2016) and Irving et al. (2016) propose to use neural networks to predict the probability of the next word, and lead into a guided-search schema; Ellis et al. (2020) propose the EC2 algorithm to iteratively learn and search over a Domain-Specific Language. Despite the success of these methods, most existing approaches work with non-parametric or few-parameter primitive functions and solve the task by treating programs as “a sequence of tokens” to learn the token transition dynamics, which cannot effectively handle parametric programs. Besides, Nye et al. (2019) had also tried to solve the problem of generating complex programs, which focus on generating longer programs with complex control flows by proposing a series of control-flow sketches and learning to fill the “sketch-hole”.

Besides, as visual scenes become prevalent, researchers start to work with much more complex visual scenes like  $\LaTeX$  drawings and computer-aided design objects (Eslami et al., 2016; Ellis et al., 2018; Young et al., 2019; Tian et al., 2019; Zhou et al., 2021). Most of these tasks are based on parametric functions and thus make the traditional view of “*treating the program as a sequence of tokens*” collapse due to a large number of variants of parametric functions. Ellis et al. (2018) uses STN (Jaderberg et al., 2015) to model multiple parameters, Tian et al. (2019) aligns all the function parameters such that they could be modeled with the same neural network, while Zhou et al. (2021) uses a grammar-encoded LSTM model. Though obtained remarkable results, these methods are not easy to generalize.

Compared with existing methods, we follow the combination of learning and searching, while, at the same time, tackling those parametric primitive functions. We propose to model each parametric function along with its parameters as a module and propose the H2MCTS algorithm that could benefit both training and inference.

### 3. Problem Formulation

#### 3.1. Notations and Problem Formulation

Following Piantadosi (2011), we define a program as a logical collection of primitive functions. Specifically, given a set of primitive functions  $\mathbb{F}$ , a program  $\mathcal{P} = (f_1^{\Theta_1}, f_2^{\Theta_2}, \dots, f_T^{\Theta_T})$ , where  $f_i^{\Theta_i} \in \mathbb{F}$  is a primitive function  $f$  with parameters  $\Theta_i = (\Theta_{i,0}, \Theta_{i,1}, \dots, \Theta_{i,n_f})$ ,  $n_f$  is the number of parameters for  $f$ , and  $T$  is a program-dependent parameter that indicates the length of the program  $\mathcal{P}$ . Besides, in the main text of this paper, we focus on the parametric functions and simplify our program syntax as context-free grammar (CFG) (Zhou et al., 2021), *i.e.*, programs without loops and other control commands, and we show in the experiments (Section 6.3) and Appendix B that our method could be easily extended to context-based scenarios.

The task of parametric visual program induction is defined as: given an input-output observation pair  $(\mathcal{O}_I, \mathcal{O}_O)$ , find a parametric program  $\mathcal{P}$  to transform the input to the output:

$$\mathcal{P}(\mathcal{O}_I) \rightarrow \mathcal{O}_O. \quad (1)$$

Moreover, based on CFG, Eq. (1) can be rewritten as

$$f_T^{\Theta_T} \circ f_{T-1}^{\Theta_{T-1}} \cdots \circ f_1^{\Theta_1}(\mathcal{O}_I) \rightarrow \mathcal{O}_O, \quad (2)$$

where  $f_i^{\Theta_i} \circ f_j^{\Theta_j}$  is the composition of two functions, *i.e.*,  $f_i^{\Theta_i} \circ f_j^{\Theta_j}(\mathcal{O}_{in}) \doteq f_i^{\Theta_i}(f_j^{\Theta_j}(\mathcal{O}_{in}))$ .

#### 3.2. The Existing Methods

To generate the desired program  $\mathcal{P}$  in Eq. (1), most of the existing works adopt the method of tokenization, *i.e.*, transforming  $(f_1^{\Theta_1}, f_2^{\Theta_2}, \dots, f_T^{\Theta_T})$  into  $(t_1, t_2, \dots, t_N)$ , where  $t_i$  is a token and  $N$  is the number of tokens. The probability of program  $\mathcal{P}$  is calculated by assuming the Markov property:

$$\Pr[\mathcal{P}|\mathcal{O}_I, \mathcal{O}_O] = \prod_{i=1}^N P(t_i|t_{<i}; \mathcal{O}_I, \mathcal{O}_O), \quad (3)$$

where  $P$  is a conditional Markov transition probability. To tackle the problem, traditional rule-based search methods (Manna & Waldinger, 1980) adopt syntax-pruned search strategies, while recent neural program synthesis methods (Bunel et al., 2018) learn the probability function with language models (Figure 2 (a) and (b)). Though this Markov chain modeling works well with zero-parameter functions by treating each function as a token, such approaches encounter great difficulties in modeling parametric functions. Considering an example function of “dot ( $x, y, \text{col}$ )”. If each function variant is treated as a token, *e.g.*, “dot ( $x=1, y=2, \text{col}=\text{red}$ )” is a token, the set of tokens becomes too large considering different parameter combination. On the other hand, if each function

fragment is treated as a token, *e.g.*, “dot”, “(”, “ $x=1$ ”, “ $y=2$ ”, “ $\text{col}=\text{red}$ ”, “)” as 6 tokens, the generated program may be very long ( $N \gg T$ ) and suffer syntax-error.

To train  $P$  in Eq. (3), synthetic data is generated and utilized (Shin et al., 2019). Specifically, a random program  $\mathcal{P}$  is first generated. Then, by inputting a random visual observation  $\mathcal{O}_I$ , the corresponding  $\mathcal{O}_O$  is obtained as  $\mathcal{O}_O = \mathcal{P}(\mathcal{O}_I)$ .  $\mathcal{O}_I, \mathcal{O}_O, \mathcal{P}$  are used as ground-truths to train the model by maximizing the posterior probability. However, training from such synthetic data is biased due to the distribution mismatch between the random program and true programs (Shin et al., 2019). Besides, due to the program aliasing problem (Bunel et al., 2018), naively using the generated data will lead to inefficient training.

### 4. Function Modularization

In this section, we tackle the problem of learning parametric functions by function modularization. Specifically, we transform the parametric program induction problem as learning inter-function transition and intra-function parameter prediction. The former focuses on selecting which function should be used, and the latter models each function along with its parameters as a self-contained module to obtain the most suitable parameters for that function.

#### 4.1. Function Transition and Parameter Prediction

The goal of function modularization is to separate the learning of the program into two parts: **inter-function transition** and **intra-function parameter prediction** as:

$$\Pr[\mathcal{P}|\mathcal{O}_I, \mathcal{O}_O] = \prod_{i=1}^T P(f_i|\hat{\mathcal{O}}_{i-1}) \cdot Q_{f_i}(\Theta_i|\hat{\mathcal{O}}_{i-1}), \quad (4)$$

where  $\hat{\mathcal{O}}_i = (\mathcal{O}_i, \mathcal{O}_O)$  is the pair of the observation at the  $i$ -th step and the target output,  $\mathcal{O}_i = f_i^{\Theta_i}(\mathcal{O}_{i-1})$ , and  $P$  and  $Q$  are two learnable probabilities. By this transformation, we could decouple the learning objective into:

$$\begin{aligned} \text{(Function Transition)} P &: (\mathcal{O}_t, \mathcal{O}_O) \rightarrow f, \\ \text{(Parameter Prediction)} Q_f &: (\mathcal{O}_t, \mathcal{O}_O) \rightarrow \Theta, \end{aligned} \quad (5)$$

*i.e.*,  $P$  proposes the next function  $f_i$ , and the corresponding  $Q_{f_i}$  is used to determine the parameters for the function  $f_i$ . This process iterates until we reach some terminal states or a preset maximum step.

**Function Transition.**  $P$  controls the function transition dynamics. In the literature, to simplify the learning process, most works use a predefined order (*e.g.*, canonical orders as from left to right, top to down (Ellis et al., 2018)) to constrain the execution of functions. However, this strategy is inefficient and problematic because a context-free program should have the flexibility to execute freely as needed. Thus,

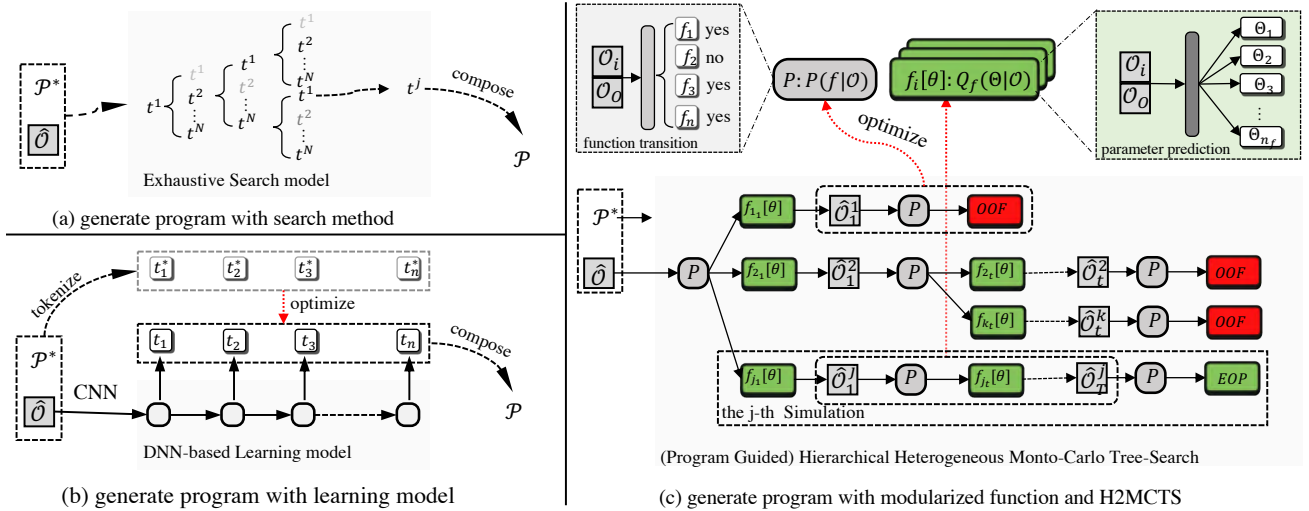


Figure 2. An illustration of (a) the search method; (b) the learning model; and (c) the proposed Function Modularization and H2MCTS. Basically, our proposed method learns the **function transition** model to propose possible new functions, and the **parameter prediction** to generate parameters for the function. Moreover, we use our proposed H2MCTS algorithm to effectively generate supervisions. Failed search processes (end with red OOF, “out of function”) as well as successful search processes (end with green EOP, “end of program”) can help to train the model. (Best view in color.)

our function transition model aims to capture the probability that  $f$  is a suitable function for the  $i$ -th step as follows:

$$P(f|\hat{\mathcal{O}}_i) \doteq \Pr [\exists \Theta : d(f^\Theta(\mathcal{O}_i), \mathcal{O}_O) < d(\mathcal{O}_i, \mathcal{O}_O)], \quad (6)$$

where  $d(\cdot, \cdot)$  is a distance metric (e.g., program distance (Ellis et al., 2018) or image IoU (Tian et al., 2019)). Intuitively, we aim to find  $f$  such that applying  $f$  could make our observation more similar to the target with some parameters  $\Theta$ . Besides modeling primitive functions in  $\mathbb{F}$ , we have two extra functions to determine the terminal states of programs: “End of Program”(EOP) and “Out of Function”(OOF) where EOP means that the program has successfully reached the target  $\mathcal{O}_O$ , while OOF means that it is impossible to generate the target observation  $\mathcal{O}_O$  with the current program.

**Parameter prediction.** After obtaining the function  $f$  from Equation 6, we model each function with its parameters  $\Theta$  as a self-contained module  $Q_f$ . The goal of  $Q_f$  is to predict the best  $\Theta$  given the context  $\hat{\mathcal{O}}_i$ :

$$Q_f(\Theta|\hat{\mathcal{O}}_i) \doteq \arg \min_{\Theta} d(f^\Theta(\mathcal{O}_i), \mathcal{O}_O). \quad (7)$$

Next, we introduce the instantiation of  $P$  and  $Q_f$ .

#### 4.2. Instantiation

In our considered setting, the raw observations come from the visual domain. Therefore, we first encode the visual observations with a convolutional-based encoder  $\mathcal{E}$ , which transforms the visual observations into a hidden state. Then, the function transition and parameter prediction are performed in the hidden state.

To instantiate the transition model, we set  $P$  as the following function:

$$P(f|\hat{\mathcal{O}}_i) = \sigma(\mathcal{E}(\mathcal{O}_i, \mathcal{O}_O)^\top \cdot \mathbf{e}_f), \quad (8)$$

where  $\mathbf{e}_f$  is a vector representation of the function  $f$ ,  $\sigma$  is the normalization function.

To instantiate the parameter prediction model, we set  $Q_f$  as a multi-head self-contained deep neural network where each head corresponds to one parameter:

$$Q_{f,j}(\Theta_j|\hat{\mathcal{O}}_i) = \text{MLP}_{f,j}(\mathcal{E}(\mathcal{O}_i, \mathcal{O}_O)), j = 1, 2 \dots, n_f. \quad (9)$$

$\text{MLP}_{f,j}$  corresponds to the parameter  $\Theta_{f,j}$  for function  $f$ 's  $j$ -th parameter, with an appropriate activation function, e.g., sigmoid for numerical parameters and softmax for categorical parameters. Finally, for each  $f \in \mathbb{F}$ , the parameter prediction function is as follows:

$$Q_f(\Theta|\hat{\mathcal{O}}_i) = \prod_{j=1}^{n_f} Q_{f,j}(\Theta_j|\hat{\mathcal{O}}_i). \quad (10)$$

### 5. Learning and Inference with H2MCTS

In this section, we propose the Hierarchical-Heterogeneous Monto-Carlo tree search (H2MCTS) algorithm in conjunction with the modularized function. Compared to using the naive synthetic data introduced in Section 3.2, our proposed H2MCTS can provide high-quality uncorrelated training data during training, and serves as an efficient search technique during inference.



## 5.1. H2MCTS

Following the standard Monto-Carlo tree search, we start from the initial observation  $\mathcal{O}_I$  and aim to find the target observation  $\mathcal{O}_O$  by hierarchically running the search algorithm on two types of tree node: `function` node for function transition and `parameter` node for parameter prediction. The probability of visiting node  $v$  depends on the following score function:

$$\text{score}(v) = \frac{1}{\sqrt{1 + \beta \cdot \alpha(v)}} p(v), \quad (11)$$

where  $\alpha(v)$  is the visit count of node  $v$ ,  $\beta$  is a scaling hyper-parameter, and  $p(\cdot)$  is  $P(\cdot)$  for `function` node and  $Q_f(\cdot)$  for `parameter` node. In essence, the learnable distribution encourages searching most potential candidates, while penalizing frequently visiting the same nodes within a tree and thus helps exploration. Concretely, in the  $l$ -th round of search, starting from  $\hat{\mathcal{O}}_0 = (\mathcal{O}_I, \mathcal{O}_O)$ , H2MCTS includes the following steps:

- **Select function.** Given the current observation  $\hat{\mathcal{O}}_i$ , we choose the next function  $f_i$  based on the score of each `function` node.
- **Select parameter.** After selecting the function  $f_i$ , we choose  $\Theta_i$  from the corresponding `parameter` nodes.
- **Expand Node.** We obtain  $\mathcal{O}_{i+1} = f_i^{\Theta_i}(\mathcal{O}_i)$ , and add  $\hat{\mathcal{O}}_{i+1} = (\mathcal{O}_{i+1}, \mathcal{O}_O)$  into the search tree.
- **BackUp.** Repeat the above steps until we reach the terminal state `EOF`, `EOF` or the maximum step. If we reach `EOF` then return the observation and the program from the current iteration as the output. If we reach `EOF` or the maximum step, which indicates that this search process fails, we update the visit count of each node and start the next search process.

More details of H2MCTS are provided in Appendix 1.

## 5.2. Learning with H2MCTS

In this subsection, we introduce how to transform the synthetic data as introduced in Section 3.2 using H2MCTS and adopt the transformed data to more effectively train the model. We mainly need to train the transition model  $P$  and the parameter prediction function  $\{Q_f\}_{f \in \mathbb{F}}$ .

Specifically, consider a randomly generated training data  $(\mathcal{O}_I, \mathcal{O}_O, \mathcal{P})$  so that  $\mathcal{P}(\mathcal{O}_I) = \mathcal{O}_O$ . We show how  $\mathcal{P}$  could be used to guide H2MCTS. In the  $i$ -th search step, we define the **program-guided** probability  $P^*$  and  $Q_f^*$  with  $\mathcal{P}_0 = \mathcal{P}$ :

$$P^*(f|\hat{\mathcal{O}}_i) = \mathbf{1}(f \in \mathcal{P}_i),$$

$$Q_f^*(\Theta|\hat{\mathcal{O}}_i) = \mathbf{1}(\Theta = \text{Top}(\mathcal{P}_i, f)),$$

where  $\mathbf{1}(\cdot)$  is the indicator function and  $\text{Top}(\cdot)$  finds the first parameter  $\Theta$  of the function  $f$  in  $\mathcal{P}_i$ . After selecting a function  $f$  and parameter  $\Theta$  in the  $i$ -th step, we update the program as  $\mathcal{P}_{i+1} = \mathcal{P}_i \setminus \{f^\Theta\}$ . We emit the `EOF` token when  $\mathcal{P}_i$  is empty. Using  $P^*$  and  $Q_f^*$  in the H2MCTS algorithm, we actually obtain an alias of program  $\mathcal{P}$  which could transform  $\mathcal{O}_I$  into  $\mathcal{O}_O$  as  $\mathcal{P}$  does. We then form the training data by collecting the pair along the  $(\mathcal{O}_i, \mathcal{O}_O, f_i, \Theta_i), \forall 1 \leq i \leq T$ . The same procedure can be performed for other search traces and the overall training set is  $\mathcal{T} = \bigcup_{l, \mathcal{P}} \{(\mathcal{O}_i, \mathcal{O}_O, f_i, \Theta_i)\}_{i=1}^T$ .

Finally, we optimize the following objective functions

$$\begin{aligned} \max_P \mathbb{E}_{(\mathcal{O}_i, \mathcal{O}_O, f_i, \Theta_i) \in \mathcal{T}} [P(f_i|\mathcal{O}_i, \mathcal{O}_O)] \\ \max_{Q_f} \mathbb{E}_{(\mathcal{O}_i, \mathcal{O}_O, f_i, \Theta_i) \in \mathcal{T}} [Q_f(\Theta_i|\mathcal{O}_i, \mathcal{O}_O)], f \in \mathbb{F} \end{aligned} \quad (12)$$

Notice that in the above learning process, instead of only using functions with the generated order, we enumerate and use all possible functions variants using H2MCTS as long as we can approach the target observation, which greatly enriches our training examples.

## 5.3. Inference with H2MCTS

With learned  $P(f|\hat{\mathcal{O}})$  and  $\{Q_f(\Theta|\hat{\mathcal{O}})|f \in \mathbb{F}\}$  (Sec. 5.2), the H2MCTS algorithm could be used as a search technique directly by using the score function to calculate visiting probabilities of nodes. Therefore, we also adopt H2MCTS during inference. In experiments, we show that H2MCTS empirically outperforms other search methods in both search efficiency and accuracy.

## 6. Experiments

In this section, we demonstrate the superiority of our method through three experiments: firstly, we compare our modularized function with the token-based models on a small hand-crafted dataset, which is designed as small as possible such that token-based methods would still work; secondly, we compare our method on the  $\text{\LaTeX}$  2D drawing dataset (Ellis et al., 2018), which only includes control-free commands; lastly, we examine our model on the complex 3D shape-synthesis dataset (Tian et al., 2019), which is more difficult with control flows and more primitive functions.

### 6.1. $5 \times 5$ Pixel Grid: Token vs. Modular

**Dataset.** In the first experiment, we aim to compare our modularized function and the token-based model. We create a small dataset on the  $5 \times 5$  pixel grid with ten colors (one background and nine foregrounds) with relatively simple parametric functions such that those token-based methods still work. We consider five primitive function `DOT`, `VLINE`, `HLINE`, `BLOCK`, and `BORDER`, corresponding to drawing a

Table 1. Primitive functions for the Pixel Grid dataset

| Primitive Functions         | Descriptions  | Overall Action Space  |
|-----------------------------|---|---|
| DOT[X,Y,COL]                | draw a <b>dot</b> at position (X, Y) with color COL                     | $C(5,1) \times C(5,1) \times 9 = 5 \times 5 \times 9 = 225$   |
| VLINE[TX,BX,Y,COL]          | draw a <b>vertical line</b> from (TX, Y) to (BX, Y) with color COL      | $C(5,1) \times C(5,2) \times 9 = 5 \times 10 \times 9 = 450$  |
| HLINE[LY,RY,X,COL]          | draw a <b>horizontal line</b> from (X, LY) to (X, RY) with color COL    | $C(5,2) \times C(5,1) \times 9 = 10 \times 5 \times 9 = 450$  |
| BLOCK[TX, LY, BX, RY, COL]  | draw a <b>block</b> from (TX, LY) to (BX, RY) with color COL            | $C(5,2) \times C(5,2) \times 9 = 10 \times 10 \times 9 = 900$ |
| BORDER[TX, LY, BX, RY, COL] | draw a <b>rectangle border</b> from (TX, LY) to (BX, RY) with color COL | $(C(5,2) - 4)^2 \times 9 = 6 \times 6 \times 9 = 324$         |

single pixel dot, a vertical line, a horizontal line, a rectangle block, and a rectangle border, respectively. For example, the definition of the primitive function DOT is:

DOT(X, Y, COL): draw a DOT at position (X, Y) with color COL. On the  $5 \times 5$  grid, this function has  $5 \times 5 \times 9 = 225$  function variants.

We summarize all the primitive functions of this dataset in Table 1. Considering all the function variants, there are 2349 actions in total, which is much larger than the common program learning tasks used tasks (Pattis et al., 1981; Balog et al., 2016).

An example is shown in the left of Figure 4. Since this dataset does not contain control-flow, we only need to learn a sequence of parametric primitive functions to transform  $\mathcal{O}_I$  to  $\mathcal{O}_O$ .

**Baselines.** We compare our proposed function modularization with two token-based baselines:

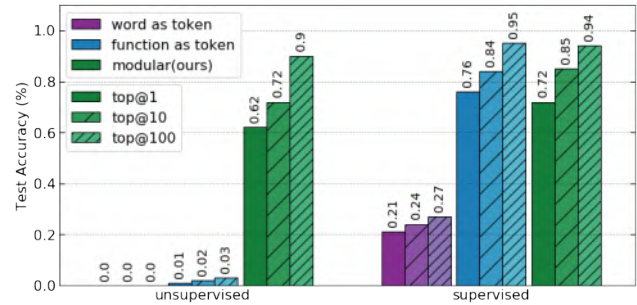
- **Word-as-Token (Balog et al., 2016):** The method treats each word as a token, e.g., DOT(X, Y, COL) in transformed into six tokens: ['DOT', '(', 'X', 'Y', 'COL', ')']. This is similar to the sentence tokenization preprocessing in the NLP community.
- **Function-as-Token:** The method treats each function-parameter variant as a token, e.g., 'DOT(1, 2, RED)' and 'DOT(2, 2, RED)' are considered as two tokens. This is similar to Reinforcement Learning (RL) which learns when to conduct which symbolic action.

We conceptually compare different methods in Table 2. Word-as-Token results in a relatively small number of tokens, but these tokens do not contain syntax information. Function-as-Token contains syntax by itself, but results in a large number of tokens. As for our proposed modularized function, the total number of tokens equals the number of primitive functions, while ensuring the syntax.

We compare different methods in both the unsupervised and the supervised settings. For the unsupervised setting, we train the Word-as-Token model as in Balog et al. (2016), which is a state-of-the-art model on program induction. The

Table 2. A comparison of different methods on the token setting.

| TOKENIZING METHOD | # TOKENS | SYNTAX |
|-------------------|----------|--------|
| WORD-AS-TOKEN     | 29       | NO     |
| FUNCTION-AS-TOKEN | 2,349    | YES    |
| MODULARIZED(OUR)  | 7        | YES    |


 Figure 3. The testing accuracy (%) on the  $5 \times 5$  Pixel Grid dataset.

Function-as-Token model is trained with Reinforcement Learning (Sutton et al., 2000). For the supervised settings, following (Ellis et al., 2018), we provide ground-truth programs to all the models.

**Results.** As shown in Figure 3, for the unsupervised setting, our proposed method reports impressive results by reporting accuracy of 0.9 for the top100 metric, while Word-as-Token (Balog et al., 2016) and Function-as-Token both fail due to the huge search space. The results show that our proposed function modularization can effectively reduce the number of tokens and relieve the burden of supervision signals in the unsupervised setting.

For the supervised setting, Word-as-Token also fails to capture the complex parameters and syntax and performs poorly, while our method and Function-as-Token show satisfactory performance. In particular, both methods could find 95% of the valid program within 100 searches. The results show that Function-as-Token can work reasonably well if enough supervised data is provided, which, however, is expensive or infeasible in practice. In contrast, our proposed function modularization can work with both supervised and unsupervised settings.

Finally, we provide a showcase of different models in Fig. 4. The failure of Word-as-Token is mainly due to the syntax

Table 3. The results on the 2D  $\LaTeX$  Drawing dataset. All the models are trained on the synthesized dataset. The testing dataset includes both synthesized data and another 100 real hand-drawn images.

| MODEL                         | MLP OUTPUT DIMENSION                              | SYNTHESIZED |        | REAL HAND-DRAWN |        |         |
|-------------------------------|---|-------------|--------|-----------------|--------|---------|
|                               |   | TRAIN       | TEST   | TOP@1           | TOP@10 | TOP@100 |
| SEQ2SEQ(F) + CE LOSS          | $\sum_f \prod_j^{n_f}  \Theta_j  > 30,000$        | 66.94%      | 6.30%  | 0%              | 0%     | 0%      |
| SEQ2SEQ(F) + RL LOSS          |   | 44.20%      | 2.10%  | 0%              | 0%     | 0%      |
| SEQ2SEQ(F) + CE-RL LOSS       |   | 76.44%      | 15.41% | 0%              | 0%     | 0%      |
| STN + SMC(ELLIS ET AL., 2018) | $\sum_f n_f +  \mathbb{F}  + 1 = 16$              | —           | —      | 63%             | 70%    | 70%     |
| STN + SMC(OUR IMPLEMENT)      |   | 90.74%      | 71.21% | 64%             | 72%    | 74%     |
| FM + CANONICAL                | $\sum_f \mathcal{O}(n_f) +  \mathbb{F}  + 2 = 20$ | 91.81%      | 74.45% | 67%             | 72%    | 72%     |
| FM + H2MCTS                   |   | 93.91%      | 83.13% | 76%             | 84%    | 86%     |

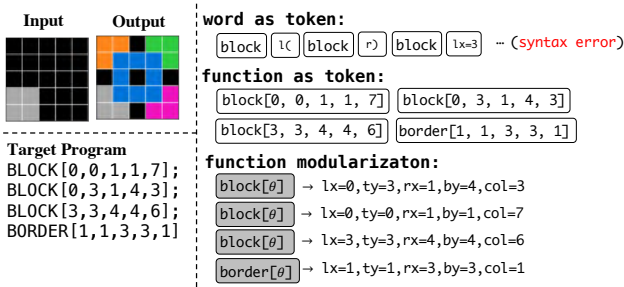


Figure 4. An illustration of the  $5 \times 5$  Pixel Grid dataset. We show the input-output observation and the target program in the left. In the right, we show the generated programs of three methods under the supervised setting.

problem. Function-as-Token could correctly predict the program, with the cost of having to select from the 2,349 candidates tokens (Table 2). In contrast, our proposed method can easily generate the program by function modularization.

## 6.2. 2D $\LaTeX$ Drawing: Control-free Program Learning

**Dataset.** In the second experiment, we adopt the  $\LaTeX$  2D drawing dataset (Ellis et al., 2018). The goal is to learn  $\LaTeX$  executable programs with visual observations. Following (Ellis et al., 2018), the training set includes 95,000 synthesized data, and the testing set includes 5,000 synthesized data as well as 100 real hand-drawn images.

Specifically,  $\mathcal{O}_I$  is an empty canvas, and  $\mathcal{O}_O$  is an image with size  $256 \times 256$ . This dataset contains 3 primitive functions: Circle, Line, and Rectangle, each of which draws on a discrete  $16 \times 16$  grid coordinates. The synthesized data contains randomly generated programs, while the real hand-drawn images aim to show certain structures. Figure 5 shows some examples of the dataset.

**Baselines.** We compare the following methods. (1) A LSTM-based Seq2Seq language model, which has achieved successes in language translation and image captioning tasks (Vinyals et al., 2015). Based on the results from  $5 \times 5$

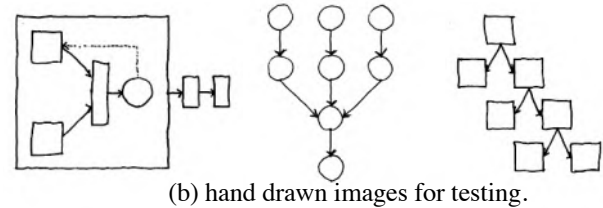
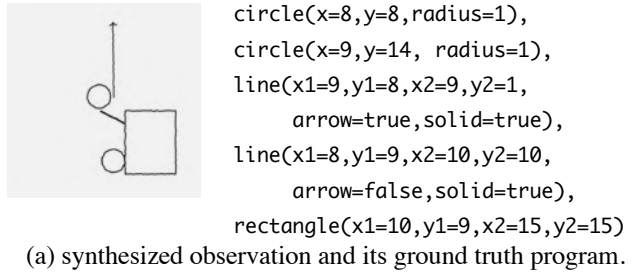


Figure 5. Top: an example of the synthesized image and the program that generated it. Bottom: examples of the true hand-drawn images which correspond to model diagram, flowing chart, and tree structures. The learned programs should create legible figures by rendering in  $\LaTeX$ .

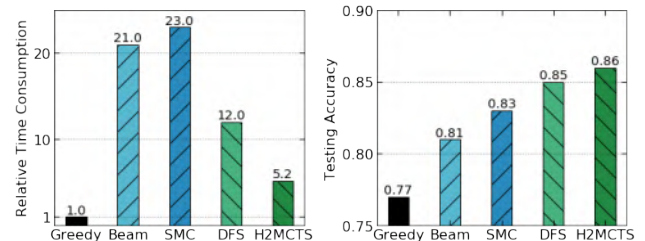


Figure 6. A comparison between different search methods. Left: the relative time consumption, from which we could find that our H2MCTS algorithm consumes much less time than other methods; Right: the testing accuracy, which demonstrate that our H2MCTS obtain the best performance.

Pixel Grid in Section 6.1, we only consider the Function-as-Token tokenizing method and denote it as Seq2seq(F). We compare three versions where the first two use Cross-Entropy (CE) loss, Reinforcement Learning (RL) loss during training respectively, and the third one is pretrained

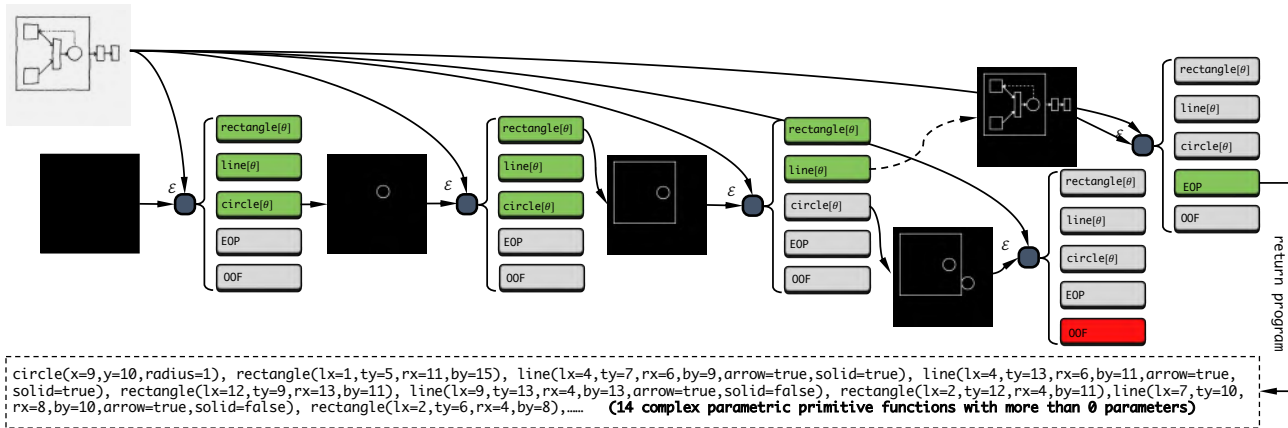


Figure 7. An example of the 2D  $\text{\LaTeX}$  Drawing dataset. In the illustration, the function transition model proposes green function modules as executable ones, and the H2MCTS algorithm selects and executes one function. This process is repeated until we reach the OOF terminal state, marking the failure of this search process, or the EOP terminal state, which marks the success of this search. Then, the programs along with their parameters are returned as the final program. The color of the rendered images is inverted for better visualization.

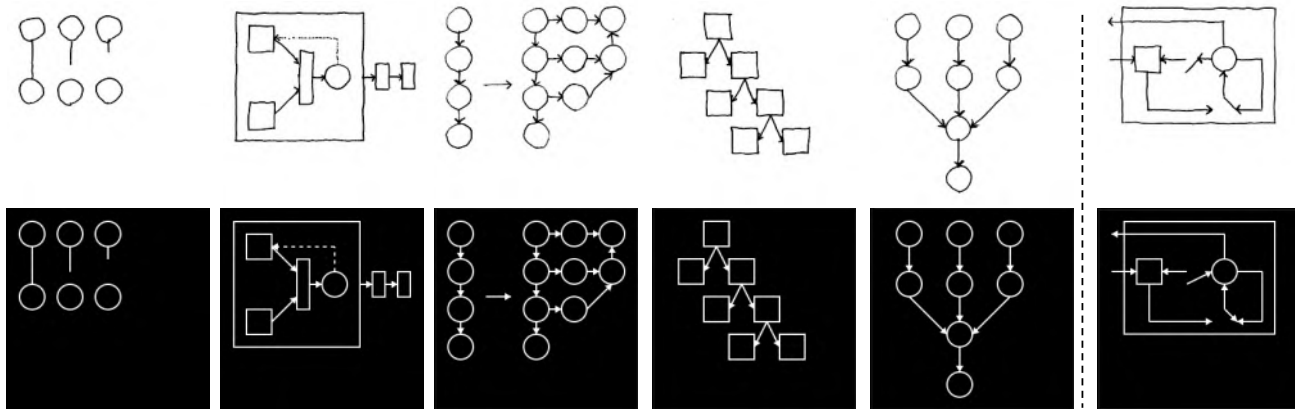


Figure 8. More showcases for the  $\text{\LaTeX}$  2D drawing datasets. Top: the hand drawings; Bottom: the  $\text{\LaTeX}$  rendered output with our generated program.

with the CE loss and further fine-tuned with the RL loss. (2) Spatial Transformer Network model with Sequential Monte-Carlo search (STN+SMC)(Ellis et al., 2018), which achieves the state-of-the-art result. We present two versions: the original results reported in the paper and our own implementation. (3) Our proposed function modularization and H2MCTS model (FM+H2MCTS). We also include an ablation study, which uses the standard canonical function order to replace the H2MCTS training, denoted as FM+Canonical. For all methods except the original result from (Ellis et al., 2018), we use ResNet-18 as the encoder  $\mathcal{E}$  to ensure a fair comparison.

**Results.** The results are shown in Table 3. We make the following observations. Overall, our proposed FM+H2MCTS model reports the best results, consistently and greatly outperforming the most competitive baseline by more than 10% in both the synthesized test set and real hand-drawn images. The results demonstrate the effectiveness of our pro-

posed method in handling the 2D  $\text{\LaTeX}$  dataset. We attribute the reasons into two folds. First, our proposed multi-head self-contained neural module is more flexible to handle different parameters adaptively. For example, we could model coordinates with a regression head and model arrow state with a binary classification head. On the other hand, compared to using the canonical approach, our H2MCTS algorithm also contributes to the performance by getting rid of the predefined function execution order during both training and inference which is extremely inflexible.

The Seq2Seq model shows poor results even in the training dataset, not to mention handling real hand-drawn images. A plausible reason is that this baseline is difficult to converge due to the huge number of parameters in MLP (more than 30,000 dimension outputs). The results are consistent with Ellis et al. (2018) that pure DNN-based approaches could not tackle the complicated visual program induction tasks.

We also compare H2MCTS and SMC with other search



Table 4. The results on the 3D Shape dataset. Following (Tian et al., 2019), we compute the intersection over union (IoU) between the target observation and the rendered output of the generated program as the evaluation metric.

| MODEL              | TRAINING SET (IoU $\uparrow$ ) |       | TESTING SET (IoU $\uparrow$ ) |       |         |       |
|--------------------|--------------------------------|-------|-------------------------------|-------|---------|-------|
|                    | TABLE                          | CHAIR | BED                           | SOFA  | CABINET | BENCH |
| TIAN ET AL. (2019) | 0.492                          | 0.469 | 0.283                         | 0.365 | 0.345   | 0.248 |
| FM + H2MCTS (OURS) | 0.641                          | 0.642 | 0.501                         | 0.670 | 0.661   | 0.602 |

strategies including greedy search, beam search, and depth-first search (DFS) regarding their search accuracy and efficiency during inference time based on our trained model. The results are shown in Figure 6. We observe that H2MCTS outperforms other search methods with respect to both search time and testing accuracy. Among baselines, DFS is most competitive based on our well-trained function transition model and parameter prediction model. In comparison, beam search and SMC perform slightly worse. The differences mainly come from the searching strategy: Beam Search and SMC are optimized Breadth-First-Search that keep a number of partials (bandwidth) at every iteration which consumes much more computation; while H2MCTS is an optimized Depth-First-Search method that stores the search statics at each iteration for next round. With a well-trained model, a small bandwidth could be sufficient to tackle most of the problem and thus leads to the high efficiency of H2MCTS and DFS.

The results demonstrate that H2MCTS serves as an efficient inference technique, as discussed in Section 5.3.

Finally, to provide a more intuitive understanding, we present a showcase of our method in Figure 7 and Figure 8, including the raw observation, the learned programs, and the  $\text{\LaTeX}$  rendered images. The figure clearly shows the workflow and the effectiveness of our proposed method in learning parametric programs from real visual scenes. We include more examples in Figure 8, where the first five drawings could be solved perfectly.

### 6.3. 3D Shape: Control-based Program Learning

**Datasets.** In the last experiment, we adopt the 3D-Shape dataset (Tian et al., 2019) containing 18 primitive functions and `for` loops, *i.e.*, control-based programs. A showcase of the dataset is provided in Figure 9.

**Settings and Baselines.** To enable our proposed method to handle the controls, we add one extra type of **Node** to determine the control flow, and extend the bi-level modeling in Eq. (5)) into a tri-level modeling as follows:

$$\begin{aligned}
 \text{(Control Transition)} \quad C &: (\mathcal{O}_i, \mathcal{O}_O) \rightarrow C_i. \\
 \text{(Function Transition)} \quad P &: (\mathcal{O}_i, \mathcal{O}_O) \rightarrow f_i, \\
 \text{(Parameter Prediction)} \quad Q_f &: (\mathcal{O}_i, \mathcal{O}_O, C_i) \rightarrow \Theta.
 \end{aligned} \tag{13}$$

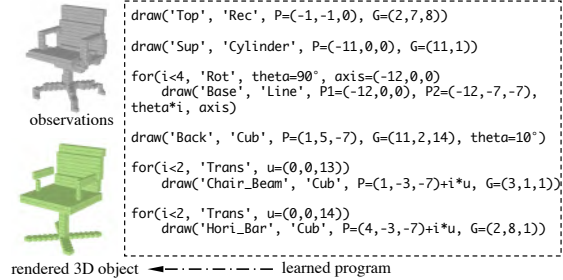


Figure 9. An illustration of the 3D Shape dataset.

With this tri-level modeling, functions within different control blocks are still context-free and therefore our proposed H2MCTS algorithm still applies. See Appendix B for more details. We mainly compare our proposed method with Tian et al. (2019), a state-of-the-art baseline for this dataset. Notice that guided adaptation used in (Tian et al., 2019) is not available in our considered setting.

**Results** As shown in Table 4, the program generated by our method achieves better performance for all categories of objects, demonstrating the general applicability of our method on control-based programs. In Figure 9, we provide a showcase of our method, which successfully generates a program to describe the visual observation.

## 7. Conclusion

In this paper, we investigate the parametric visual program induction task by decoupling the learning of parametric function as learning function transition and function parameter prediction. We propose the concept of function modularization and the H2MCTS algorithm. Our method outperforms state-of-the-art baselines with higher accuracy and efficiency. Future works include exploring more visual program induction scenarios using our proposed method.

## Acknowledgement

This work is supported by the National Key Research and Development Program of China No. 2020AAA0106300 and National Natural Science Foundation of China No. 62102222.

## References

- Balog, M., Gaunt, A. L., Brockschmidt, M., Nowozin, S., and Tarlow, D. Deepcoder: Learning to write programs. In *International Conference on Learning Representations*, 2016.
- Bunel, R., Hausknecht, M., Devlin, J., Singh, R., and Kohli, P. Leveraging grammar and reinforcement learning for neural program synthesis. *Proceedings of the 4th International Conference on Learning Representations (ICLR)*, 2018.
- Chollet, F. On the measure of intelligence. *arXiv preprint arXiv:1911.01547*, 2019.
- Devlin, J., Bunel, R., Singh, R., Hausknecht, M., and Kohli, P. Neural program meta-induction. *Neural Information Processing Systems (NIPS)*, 2017.
- Dong, H., Mao, J., Lin, T., Wang, C., Li, L., and Zhou, D. Neural logic machines. *Proceedings of the 4th International Conference on Learning Representations (ICLR)*, 2019.
- Ellis, K., Ritchie, D., Solar-Lezama, A., and Tenenbaum, J. B. Learning to infer graphics programs from hand-drawn images. *Neural Information Processing Systems (NIPS)*, 2018.
- Ellis, K., Wong, C., Nye, M., Sable-Meyer, M., Cary, L., Morales, L., Hewitt, L., Solar-Lezama, A., and Tenenbaum, J. B. Dreamcoder: Growing generalizable, interpretable knowledge with wake-sleep bayesian program learning. *arXiv preprint arXiv:2006.08381*, 2020.
- Eslami, S., Heess, N., Weber, T., Tassa, Y., Szepesvari, D., Hinton, G. E., et al. Attend, infer, repeat: Fast scene understanding with generative models. *Advances in Neural Information Processing Systems (NIPS)*, 29:3225–3233, 2016.
- Feser, J. K., Chaudhuri, S., and Dillig, I. Synthesizing data structure transformations from input-output examples. *ACM SIGPLAN Notices*, 50(6):229–239, 2015.
- Gulwani, S. Automating string processing in spreadsheets using input-output examples. *ACM Sigplan Notices*, 46(1):317–330, 2011.
- Irving, G., Szegedy, C., Alemi, A. A., Eén, N., Chollet, F., and Urban, J. Deepmath-deep sequence models for premise selection. *Advances in Neural Information Processing Systems*, 29:2235–2243, 2016.
- Jaderberg, M., Simonyan, K., Zisserman, A., et al. Spatial transformer networks. *Advances in neural information processing systems (NIPS)*, 28:2017–2025, 2015.
- Lezama, A. S. *Program synthesis by sketching*. PhD thesis, PhD thesis, EECS Department, University of California, Berkeley, 2008.
- Ling, W., Yogatama, D., Dyer, C., and Blunsom, P. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, pp. 158–167, 2017.
- Manna, Z. and Waldinger, R. Knowledge and reasoning in program synthesis. *Artificial intelligence*, 6(2):175–208, 1975.
- Manna, Z. and Waldinger, R. A deductive approach to program synthesis. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 2(1):90–121, 1980.
- Nye, M., Hewitt, L., Tenenbaum, J., and Solar-Lezama, A. Learning to infer program sketches. In *International Conference on Machine Learning*, pp. 4861–4870. PMLR, 2019.
- Parisotto, E., Mohamed, A.-r., Singh, R., Li, L., Zhou, D., and Kohli, P. Neuro-symbolic program synthesis. *Proceedings of the 4th International Conference on Learning Representations (ICLR)*, 2017.
- Pattis, R., Roberts, J., and Stehlik, M. Karel the robot. *A gentele introduction to the Art of Programming*, 1981.
- Piantadosi, S. T. *Learning and the language of thought*. PhD thesis, Massachusetts Institute of Technology, 2011.
- Shin, R., Kant, N., Gupta, K., Bender, C., Trabucco, B., Singh, R., and Song, D. Synthetic datasets for neural program synthesis. *Proceedings of the 4th International Conference on Learning Representations (ICLR)*, 2019.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- Sun, S.-H., Noh, H., Somasundaram, S., and Lim, J. Neural program synthesis from diverse demonstration videos. In *International Conference on Machine Learning*, pp. 4790–4799. PMLR, 2018.
- Sutton, R. S., McAllester, D. A., Singh, S. P., and Mansour, Y. Policy gradient methods for reinforcement learning with function approximation. In *Advances in neural information processing systems (NIPS)*, pp. 1057–1063, 2000.
- Tian, Y., Luo, A., Sun, X., Ellis, K., Freeman, W. T., Tenenbaum, J. B., and Wu, J. Learning to infer and execute 3d shape programs. In *International Conference on Learning Representations*, 2019.

- Vinyals, O., Toshev, A., Bengio, S., and Erhan, D. Show and tell: A neural image caption generator. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3156–3164, 2015.
- Waldinger, R. J. and Lee, R. C. Prow: A step toward automatic program writing. In *Proceedings of the 1st international joint conference on Artificial intelligence*, pp. 241–252, 1969.
- Yi, K., Wu, J., Gan, C., Torralba, A., Kohli, P., and Tenenbaum, J. B. Neural-symbolic vqa: Disentangling reasoning from vision and language understanding. *Neural Information Processing Systems (NIPS)*, 2018.
- Young, H., Bastani, O., and Naik, M. Learning neurosymbolic generative models via program synthesis. In *International Conference on Machine Learning*, pp. 7144–7153. PMLR, 2019.
- Zhou, C., Li, C.-L., and Póczos, B. Unsupervised program synthesis for images by sampling without replacement. In *Uncertainty in Artificial Intelligence*, pp. 408–418. PMLR, 2021.

## A. Implementation and Training

In this section, we briefly introduce more details about the modularized functions, as well as our training framework.

### A.1. Function Modular

Parametric primitive functions are easy to implement, but how to determine their parameters is hard. This is one of the motivation of this paper. To achieve this, we organize each primitive function and its parameter into a module with a multi-head MLP (each head is a two-layer MLP as ParamNet), leaving each head corresponding to one parameter. An examples of the DOT primitive function and its modularized form Dot are shown as follows:

```

1 import numpy as np
2 import torch
3
4 def DOT(In, X, Y, COL):
5     Out = copy.deepcopy(In)
6     Out[X, Y] = COL
7     return Out
8
9 class Dot(torch.nn.Module):
10     def __init__(self, H):
11         super().__init__()
12         self.fn = DOT
13         # ParamNet is another mlp to predict parameters.
14         self.params['X'] = ParamNet("X", "REG", range=MAXSIZE[0])
15         self.params['Y'] = ParamNet("Y", "REG", range=MAXSIZE[0])
16         self.params['C'] = ParamNet("C", "CLS", range=NUM_OF_COLOR)
17
18     def inference(In, Target):
19         s = encoder(Target) - encoder(In)
20         param_x = self.params['X'](s)
21         param_y = self.params['Y'](s)
22         param_c = self.params['C'](s)
23         return self.fn(In, param_x, param_y, param_c)

```

Listing 1. An example of Function Modular.

### A.2. Overall Network Architecture

With the idea of function modularization, the whole deep model is clear to go as a shared convolutional encoder, followed with **function transition** head to predict the transition dynamics for next function, while several function modular to predict the parameter for each function. As for the encoder, we adopt a 4-layer convolutional networks for  $5 \times 5$  Pixel Grid, use ResNet18 for 2D L<sup>A</sup>T<sub>E</sub>X Drawing, use a shapenet as Tian et al. (2019) for 3D Shape.

Moreover, to accelerate the training speed, we follow the setting of MuZero (Silver et al., 2016) to implement a distributed system based on Ray<sup>1</sup>. The system consists of 60 explorers to continuously explore the function space via the H2MCTS algorithm (Alg.1), and explored traces are used to train the model via an online trainer. The whole framework is launched on a GPU server with two Intel(R) Xeon(R) Gold 6240 CPU @ 2.60GHz CPU processors and two Nvidia GeForce RTX 3090 GPU processors.

### A.3. The H2MCTS algorithm

Our H2MCTS algorithm follows and generalizes the basic MCTS algorithm and (Silver et al., 2016). A basic MCTS algorithm includes four parts as: **Selection, Expansion, Simulation, BackUp**. Both Silver et al. (2016) and our work use the neural network to conduct the **Simulation** step. Moreover, we consider two different types of nodes in the simulation process, and thus our algorithm is called H2MCTS. The detailed algorithm is shown in Algorithm 1.

<sup>1</sup><https://www.ray.io/>



**Algorithm 1** H2MCTS

---

**Input:** The Function Transition model  $P(\cdot|\cdot)$ ; A set of Parameter Prediction models  $\{Q_f(\cdot|\cdot)\}$ ; {Eq. (5)}  
**Input:** A raw observation input-output pair  $\mathcal{O}_I, \mathcal{O}_O$ .  
**init:**  $j \leftarrow 0$ .  
**init:**  $\text{ROOT} \leftarrow (\mathcal{O}_I, \mathcal{O}_I, 0)$  {ROOT node with  $\mathcal{O}_I, \mathcal{O}_O$ , and visit count 0}  
**repeat**  
    $\text{PNode} \leftarrow \text{ROOT}$  {We starts from root at every round}  
   **repeat**  
      $f = \arg \max_f 1/(1 + \beta * \alpha(\text{PNode})) * P(f|\text{PNode.context})$  {**Select Function** node according to Eq. (11)}  
     **if**  $f$  is OOF **then**  
        $\text{node} \leftarrow \text{PNode}$  {**BackUp** node visit count}  
       **repeat**  
          $\alpha(\text{node}) \leftarrow \alpha(\text{node}) + 1$   
          $\text{node} \leftarrow \text{node.parent}$   
       **until**  $\text{node}$  is None  
       break inner loop  
     **end if**  
     **if**  $f$  is EOP **then**  
       break outer loop  
     **end if**  
     **if**  $f \notin \text{PNode.children}$  **then**  
        $\text{FNode} \leftarrow (f, 0)$  {create new function node with  $f$ , and visit count 0}  
        $\text{PNode.children.add}(\text{FNode}(f, 0))$  {Add  $f$  to PNode's children}  
     **end if**  
      $\Theta = \arg \max_{\Theta} 1/(1 + \beta * \alpha(\text{PNode})) * Q(\Theta|\text{PNode.context})$  {**Select Parameter Node** for  $f$  according to Eq. (11)}  
      $\mathcal{O}_{new} \leftarrow f^{\Theta}(\text{PNode.I})$  {Obtain new observations}  
     **if**  $\Theta \notin \text{FNode.children}$  **then**  
        $\mathcal{O}_{new} \leftarrow f^{\Theta}(\text{PNode.Input})$  {render the new observation and **Expand** the search tree}  
        $\text{PNode} \leftarrow (\mathcal{O}_{new}, \text{PNode.Output}, 0)$  {A new PNode}  
        $\text{FNode.children.add}(\text{PNode})$   
     **end if**  
     **until**  $K$  exceed maximum search depth  
   **until**  $j$  exceed maximum number of simulation  
    $\mathcal{P} \leftarrow []$   
   **repeat**  
      $\mathcal{P.add}(\text{PNode.parent.f}, \text{PNode}.\Theta)$   
      $\text{PNode} \leftarrow \text{PNode.parent.parent}$   
   **until** PNode is None  
**Return**  $\mathcal{P}$

---

## B. Extension to Context-based Scenarios

To extend our formulation to context-based scenarios, we firstly rewrite a program as  $\mathcal{P} = (C_1, C_2, \dots, C_{N_P})$ , where  $C_i = (C^{\Theta_C}, f_{i,1}^{\Theta_{i,1}}, f_{i,2}^{\Theta_{i,2}}, \dots)$  is the  $i$ -th logic collection block and  $C^{\Theta_C}$  is the control unit (e.g., loop, if-else), and  $f_{i,j}^{\Theta_{i,j}}$  is primitive function  $f$  with parameters  $\Theta_{i,j} = (\Theta_{i,j,0}, \Theta_{i,j,1}, \dots, \Theta_{i,j,n_f})$ ,  $n_f$  is the number of parameters for  $f$  (e.g.,  $\text{line}(lx, ty, rx, by)$ ). Then we define a tri-level modeling as:

$$\Pr[\mathcal{P}|\mathcal{O}_I, \mathcal{O}_O] = \prod_{i=1}^{N_P} C(c_i|\hat{\mathcal{O}}_{i-1}) \cdot P(f_i|\hat{\mathcal{O}}_{i-1}) \cdot Q_{f_i}(\Theta_i|\hat{\mathcal{O}}_{i-1}, c_i). \quad (14)$$

This means we firstly determine the control unit, then determine the function, and finally the parameter. Moreover, we add an extra control unit as NUL which directly executes its enclosed commands, and the H2MCTS algorithm is correspondingly extended with the extra Control node.

## C. More Experiments Results

### C.1. Pixel Grid

#### C.1.1. DATASET

To generate the dataset, we randomly sample a sequence of functions and parameters from the primitive functions  $\mathcal{P} = [f_1^{(\Theta_1)}, f_2^{(\Theta_2)}, \dots, f_T^{(\Theta_T)}]$  with  $T < T_{max}$ , and generate a random input  $\mathcal{O}_I$ , and apply the program to  $\mathcal{O}_I$  to obtain  $\mathcal{O}_O$ . Then  $(\mathcal{O}_I, \mathcal{O}_O)$  is used as an input-output pair. Moreover, as some functions could overwrite previous functions (e.g.,  $f_j = \text{DOT}(X, Y, C2)$  will overwrite  $f_i = \text{DOT}(X, Y, C1)$  if  $j > i$ ), we carefully compare the intermediate results  $[\mathcal{O}_0, \mathcal{O}_1, \dots, \mathcal{O}_T]$  to remove functions which have been overwritten.

#### C.1.2. COMPARISON WITH RULE-BASED METHODS.

In our experiments, the training set contains of 10,000 input-output pairs, which consists of 20% programs with length 1, 20% programs with length 2, and 60% programs with length 3. The testing set contains 1,000 input-output pairs with the same length distribution.

### C.2. ~~La~~TeX 2D drawings

We show the training curve of our model in Figure 10. From the figure, we can see that most of the functions and their parameters could reach an accuracy of more than 95% after 50k iterations.

### C.3. 3D shape-synthesis

We provide more showcases for this dataset in Figure 11.

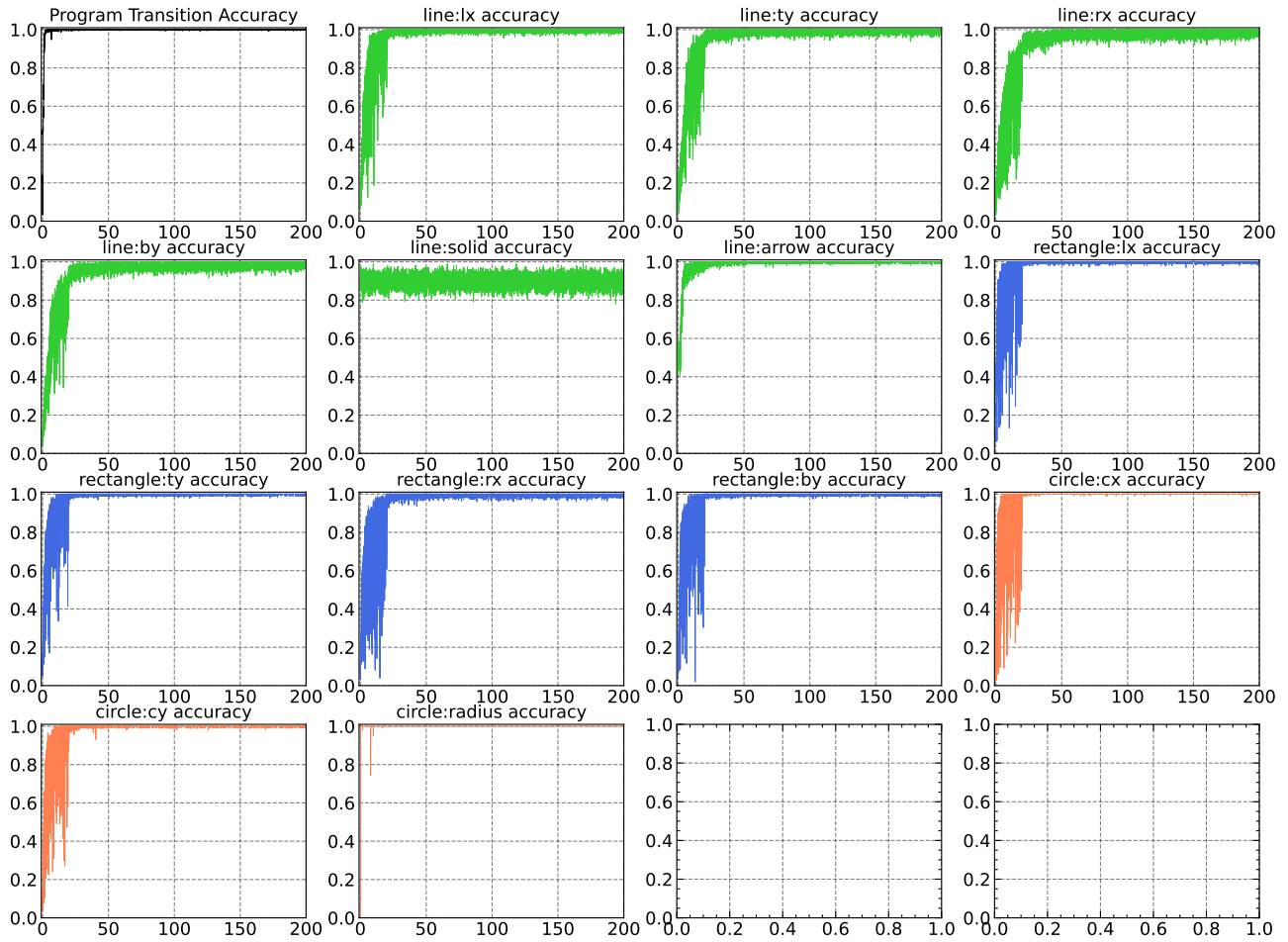


Figure 10. Training Curves for the L<sup>A</sup>T<sub>E</sub>X 2D drawing datasets. We can observe that most of the modules could converge within 50k iterations.

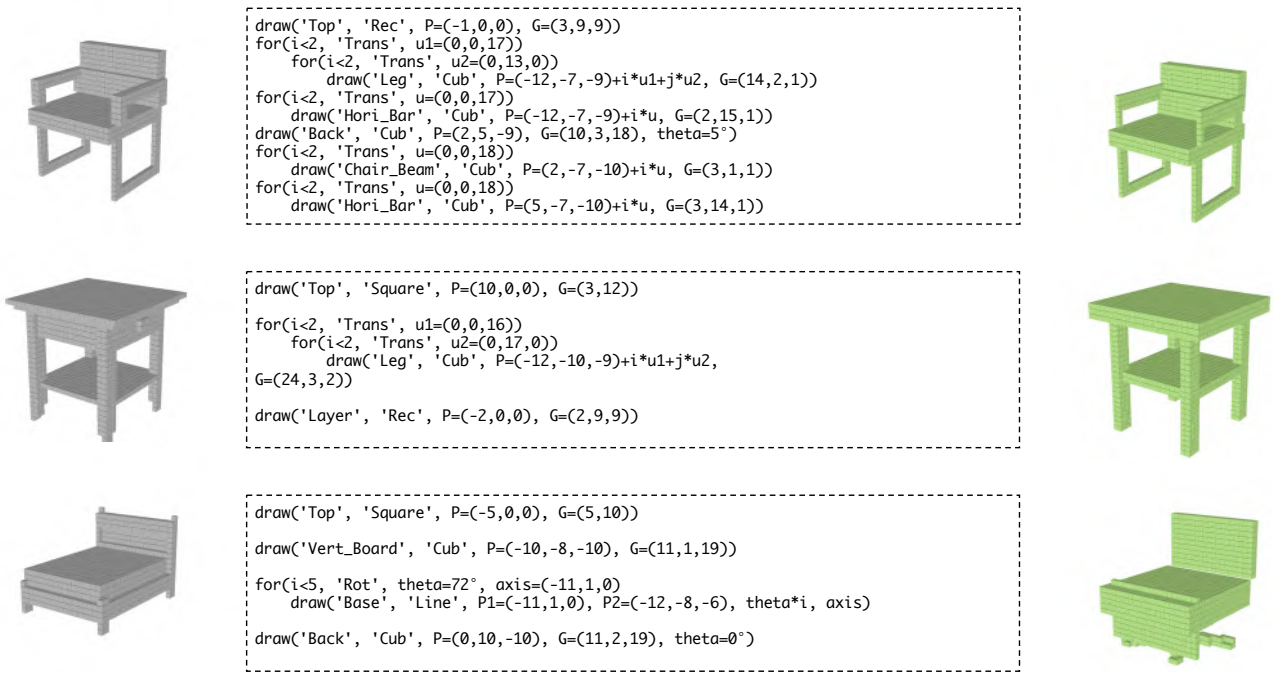


Figure 11. More showcases for the 3D shape dataset. Left: the raw observations. Middle: the generated program. Right: the rendered results.