Sketch-Plan-Generalize: LEARNING INDUCTIVE REPRE SENTATIONS FOR GROUNDED SPATIAL CONCEPTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Our goal is to enable embodied agents to learn inductive representations for grounded spatial concepts, e.g., learning staircase as an inductive composition of towers of increasing height. Given few human demonstrations, we seek a learning architecture that infers a succinct inductive *program* representation that *explains* the observed instances. The approach should generalize to learning of novel structures of different size or complexity expressed as a hierarchical composition of previously learned concepts. Existing approaches that use code generation capabilities of pre-trained large (visual) language models as well as purely neural models show poor generalization to *a-priori* unseen complex concepts. Our key insight is to factor inductive concept learning as: (i) Sketch: detecting and inferring a coarse signature of a new concept (ii) Plan: performing MCTS search over grounded action sequences (iii) Generalize: abstracting out grounded plans as inductive programs. Our pipeline facilitates generalization and modular re-use enabling continual concept learning. Our approach combines the benefits of code generation ability of large language models (LLMs) along with grounded neural representations, resulting in neuro-symbolic programs that show stronger inductive generalization on the task of constructing complex structures vis-á-vis LLM-only and purely neural approaches. Further, we demonstrate reasoning and planning capabilities with learned concepts for embodied instruction following.

028 029 030 031

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

1 INTRODUCTION

032 The ability to learn *inductive* representation for novel grounded concepts is one of the hallmarks 033 of human intelligence (Tenenbaum et al., 2011). Humans are highly data efficient – observing a 034 few instances of *towers* of a certain heights, we can generalize to constructing *towers* of *any* height. Further, we interpret increasingly complex concepts as hierarchical composition over simpler ones, 035 e.g., a tower as a sequence of blocks placed on top of each other, or a staircase composed of towers of increasing height. This paper considers the problem of learning a program representation, from a few 037 demonstrations, that models the *inductive* realization of grounded spatial concepts. Learning of such concepts is a challenging task due to an expansive space of programs and the need to reason about their physical plausibility. Further, the representation must support inductive generalization over 040 learned concepts as well as express complex hierarchical concepts via modular re-use of concepts 041 learnt previously. 042

Prior efforts such as Liang et al. (2023) uses a LLM (Large Language Model) to generate control 043 program for a given task specification but fail to generalize to complex spatial concepts which are 044 difficult to tokenize. Extension of this to VLMs (Vision Language Models), Achiam et al. (2023), 045 also fail to generalize when presented with linguistically novel concepts, due to over reliance on 046 prior knowledge and their inability to effectively learn novel concepts from given demonstrations. 047 On the other hand, neural approaches such as Liu et al. (2023), learn from given demonstrations but 048 generalize poorly due to their inability to (a) explicitly model symbolic concept of induction and (b) modularize as well as re-use previously acquired concepts. Approaches such as Li et al. (2019) train RL-based policies to attain spatial assembly by encoding an inductive spatial prior using GNNs 051 within the policy architecture. However, generalization is still limited and assumes an elicitation of the goal as per object positions, resulting in lack of ability to take in goal description such as 052 "construct a tower of size three". In essence, we attribute poor generalization of such approaches to an *implicit entangling* of the following objectives: (i) postulating a high-level program for new concepts,

Training Demonstrations Gen. over visual attributes Gen. over size Composition Gen. Construct a row of size 3 using cyan cubes Construct a row of size 5 using green die Construct a pyramid of height 3 using cyan die Construct a row of size 3 using vellow cubes a a a E E Construct a tower of height 3 using green cubes Construct a white die of height 3 using Construct a to cyan die ver of height 5 using Construct a staircase of pink cubes having 4 steps 633

Figure 1: Problem Overview. Our goal is to enable an embodied agent to learn grounded and generalizeable representations for spatial abstractions possessing a notion of induction (e.g., constructing a tower, row or their combinations such as staircases, boundary etc.). Learning is enabled by querying prior knowledge from large pretrained models, performing search in the action space guided by observations of a human demonstration for few examples and finally generalizing as compact programs. (Left) A human demonstrates the construction of a row and tower of size three. (Right) The agent learns program representation that enables inductive generalization to novel structures (varied sizes and visual attributes) and expresses complex concepts as hierarchical composition of previously acquired ones. E.g., learning a tower as a sequence of blocks placed one on top of another and a pyramid as rows of decreasing size.

074 075

054

056

060 061

062

063 064 065

066

067

068

069

071

073

(ii) evaluating plausibility of grounded plans to align with human demonstration of concepts to be 076 learned and (iii) abstracting out a grounded program to facilitate inductive generalization and modular 077 re-use in a continual manner.

This paper introduces an approach, termed SPG, that *factorizes* the concept learning task as: (a) Sketch: 079 Given a language-annotated demonstration of a novel concept, an LLM is used to postulate a function signature. (b) *Plan*: Refinement of program sketches via MCTS search, rapidly evaluating actions 081 sequences guided by a reward associated with constructing a concept. The search is accelerated by training a neural action predictor that uses the given demonstrations. (c) Generalize: Leveraging the 083 code generation capability of an LLM to distill grounded plans into a program that is inductively 084 generalize-able. This results in a continually evolving library of concepts which can be used 085 to hierarchically learn complex concepts in future. The modular architecture enables continual learning by providing the ability to decide whether the new concept encountered either as a symbolic 087 composition of existing concepts, or, a neural embedding trained via gradient update. Our experiments 880 demonstrate accurate learning of simple and complex concepts from few demonstrations for a range of spatial structures. Further, our approach shows inductive generalization in out-of-distribution settings, significantly improving over the baselines. We also present deeper insights around the 090 efficiency gains obtained by combining symbolic MCTS with neural action predictor. Finally, we 091 show how learned concepts can be grounded in the visual input, enabling a robot to follow natural 092 language instructions referring to *a-priori* unseen spatial configurations.

094

2 **RELATED WORKS**

096

Concept Learning: The problem of acquiring higher-order programmatic constructs is often modeled 098 as Bayesian inference over a latent symbol space given observed instances. Seminal works have demonstrated efficient inference over latent generative programs to express hand-written digits (Lake 100 et al., 2015), object arrangements (Ellis et al., 2018), motion plans (Mao et al., 2019), goal-directed 101 policies (Silver et al., 2019) or compressed/refactored code (Grand et al., 2023; Ellis et al., 2021). 102 These works are focused on learning abstract programs without considering their grounding in the 103 3D world or the process of constructing them (e.g., via an embodied agent performing stacking). 104 In contrast, this paper focus on learning a representation of specific class of higher-order spatially-105 grounded concepts, namely those possessing a notion of induction resulting in the construction of a structure. While prior efforts have leveraged program synthesis/search methods in learning concepts, 106 we expose such a search to the assessing the physical construction plausibility thereby learning 107 physically grounded concepts.

108 Learning-to-plan Methods: Our work is complementary to efforts that learn symbolic constructs 109 for efficient planning. Works such as Silver et al. (2023; 2024); Liu et al. (2024), infer state-action 110 abstractions for planning by querying large pre-trained models or by optimizing a goal attainability 111 objective. This paper, instead, focuses on learning a representation for complex spatial assemblies 112 as inductive programs leading to the ability to infer complex goal specifications which can then be combined with aforementioned works for synthesize efficient plans to realize complex assemblies. 113 Works such as Li et al. (2019) learn to construct structures by encoding relational knowledge via 114 graph neural network. However, this effort suffers from poor generalization to unseen examples, (e.g., 115 tower of larger size) and do not possess a mechanism to re-use previously acquired concepts. Works 116 such as Wang et al. (2023a;d) shows lifelong learning of skills by learning to plan high-level tasks 117 through composition of simple skills for simulated agents. Others (Wan et al., 2023; Parakh et al., 118 2023) initiate new skill acquisition upon detecting task failure, building a library of skills over time. 119 However, they do not model deep inductive use of learned concepts and initiate skill acquisition only 120 upon failure as opposed to learning continually even from goal-reaching demonstrations. 121

Robot Instruction Following: Instruction following involves grounding symbolic constructs ex-122 pressed in language with aspects of the state-action space such as object assemblies (Paul et al., 2018; 123 Collins et al., 2024; Lachmy et al., 2022), spatial relations (Tellex et al., 2011; Kim et al., 2024), 124 reward functions (Boularias et al., 2015), or motion constraints (Howard et al., 2014). These works 125 assume the presence of grounded representation for symbolic concepts and only learn associations 126 between language and concepts. In contrast our work *jointly* learns higher-order concepts composed 127 of simpler concepts along with their grounding in the robot's state and action space. Others (Singh 128 et al., 2022; Wang et al., 2023b; Ahn et al., 2022; Liang et al., 2023) leverage prior-knowledge em-129 bodied in large vision-language models to directly translate high-level tasks to robot control programs. Our experiments (reported subsequently) demonstrate their limitation in outputting programs for 130 structure assembly-type tasks that require long-range (inductive) spatial reasoning and consideration 131 of physical plausability of construction. Our approach addresses this problem by *coupling* abstract 132 task knowledge from pre-trained models with physical reasoning in the space of executable plans. 133

- 133
- 135

3 PRELIMINARIES AND PROBLEM SETTING

136 137

138 We consider an embodied agent that uses a visual and depth sensor to observe its environment and can 139 grasp and release objects at specified poses. We represent the robot's domain as a goal-conditioned 140 MDP $\langle S, A, T, q, \mathcal{R}, \gamma \rangle$ where S is the state space, A is the action space, T is the transition 141 function, q is the goal, R is the reward model and γ is the discount factor. The agent's objective is to 142 learn a policy that generates a sequence of actions from an initial state s_0 to achieve the goal q in response to an instruction Λ specifying the intended goal from a human. We assume that the agent 143 possesses a model of semantic relations (e.g., left(), right() etc.) as well as semantic actions such 144 as moving an object by grasping and releasing at a target location. Such modular and composable 145 notions can be acquired from demonstrations via approaches outlined in Kalithasan et al. (2023); Mao 146 et al. (2019; 2022). Such notions populate a library of concepts \mathcal{L} available as grounded executable 147 function calls. Following recent efforts (Liang et al., 2023; Huang et al., 2022; Ahn et al., 2022; 148 Singh et al., 2022) in representing robot control directly as executable programs, we represent action 149 sequence corresponding to a plan as a program consisting of function calls to executable actions and 150 grounded spatial reasoning. 151

Our goal is to enable a robot to interpret and learn the concepts in instructions such as "construct a 152 tower with red blocks of height five". Specifically, we aim to learn spatial constructs like a tower that 153 requires sequential actions that repeatedly place a block on top of a previously constructed assembly, 154 a process akin to induction. Given a few demonstrations of constructing a spatial assembly, \mathcal{D} , each 155 consisting of natural language description Λ ("construct a tower of red blocks of size five") and a 156 sequence of key frame states $\{S_1, \dots, S_q\}$ associated with the construction process, we seek to learn 157 a program that models the inductive nature of the concept of tower. This learned representation should 158 enable the agent to generalize inductively to new instructions, such as "construct a tower of blue 159 blocks of height ten." Moreover, the learned concepts should facilitate the learning of more complex structures, which are challenging to represent using primitive actions alone. For example, the concept 160 of a "tower" should assist in learning a "staircase," which can be represented as a sequence of towers 161 of increasing heights.

¹⁶² 4 Representing Inductive Spatial Concepts

We formalize the notion of inductive spatial concepts and formulate the learning objective.

Inductive Spatial Concepts: A spatial structure is an inductive concept if its construction can be 166 described recursively using a similar structure of smaller size or as a composition of other simpler 167 structures. Formally, let $C_1, \dots, C_{|\mathcal{L}|}$ represent the concepts in the concept library \mathcal{L} . We define a 168 partial order on \mathcal{L} where a concept C is "dependent on" \tilde{C} if \tilde{C} is a substructure of C. For example, a staircase is dependent on a tower, and X (cross) is dependent on diagonals, and so on. This partial 170 order is referred to as structural complexity, where a concept C is more structurally complex than C if 171 C is dependent on C. Without loss of generality, assume that $C_1, \dots, C_{|\mathcal{L}|}$ are written in topological 172 order as per their structural complexity. Now, the construction of an inductive spatial concept C_K of 173 size n at a position p, denoted by the function $h(C_k, n, p)$, is defined recursively as: 174

183

185

186

187

188

189 190

191

192

193

194

196

197

199

210 211

163 164

$$h(C_k, n, p) = \underbrace{h^{\lambda}(C_k, n-1, pos(.))}_{\text{Induction (I)}} \circ \underbrace{\prod_{l=1}^{L(C_k)} h(C_{k'_l}, \text{size}(.), pos(.))}_{\text{Composition (C)}} \circ \underbrace{\prod_{l=1}^{L'(C_k)} \eta^l_{\theta}(pos(.))}_{\text{Base (B)}} \quad (1)$$

where, $\lambda \in \{0, 1\}$, k' < k, $0 \le L(C_k)$, $L'(C_k) \le o(|\mathcal{L}|)$, $pos(.) = pos(C_k, l, n, p)$ and size(.) = size(C_k, l, n) are functions that predict the size and position of the structure to be constructed.

- 1. *Induction term:* The first term $h^{\lambda}(C_k, n-1, \text{pos}(.))$ is referred to as the induction term because it represents the possibility of constructing C_k of size n using C_k of size n-1. Here, λ is an integer exponent, either 0 or 1, where $\lambda = 0$ indicates the absence of the induction term, and $\lambda = 1$ indicates its presence.
- 2. Composition term: The second term $\prod_{l=1}^{L(C_k)} h(C_{k'_l}, \mathtt{size}(.), \mathtt{pos}(.))$, called the composition term, allows us to express the construction of C_k as a composition of previously known concepts in the library. The number of required compositions depends on the concept C_k and the size of the library \mathcal{L} .
- 3. *Base term:* The third term $\prod_{l=1}^{L'(C_k)} \eta_{\theta}^l(\text{pos}(.))$ defines the base case where the construction of concept C_k may include L' number of primitive actions. For example, the construction of a tower of size n can be written as a construction of a tower of size n 1 followed by a primitive action of moving a block on top.

200 **Learning Objective:** The functional space of inductive concepts (*h*) leads to a hypothesis space \mathcal{H} of 201 associated neuro-symbolic programs. Each goal-reaching demonstration corresponds to a particular instantiation of a given inductive concept, i.e. $h(C_k, n, p)$, where the p comes from the sequence of 202 frames, and n, C_k comes from Λ . We aim to learn a generic representation $H = h(C_k, \cdot, \cdot) \in \mathcal{H}$ 203 for the given concept, which is general for all n and p. Given (few) demonstrations of a human 204 constructing a spatial structure, concept learning can be formulated as the Bayesian posterior (Lake 205 et al., 2015; Shah et al., 2018; Silver et al., 2019), $P_{\mathcal{H}}(H|\Lambda, S_1..S_q) \propto P(S_1..S_q|\Lambda, H) \cdot P(H|\Lambda)$. 206 Here, the likelihood term associates a candidate program, and the prior term regularizes the program 207 space. The maximum *a-posteriori* estimate, representing the learnt program, is obtained by optimizing 208 the following objective: 209

$$H^* = \arg\min_{H \in \mathcal{H}} \left[\text{Loss}(\{S_1 .. S_g\}, Exec(H, \Lambda, S_1)) - \log P(H|\Lambda) \right]$$
(2)

212 Since exact inference is intractable, approximate inference is performed via search in the program 213 space. Note that learning inductive spatial concepts given demonstration considers programs that 214 represent plans that attain physically grounded/feasible structures, an object we model during the 215 search. Additionally, we seek strong generalization from a few instances of an inductive structure to 216 structures with arbitrary sizes, in effect favouring programs with iterative looping constructs.

219

220

221

222

224

225

226

227 228

229

230

231

232 233

234 235

236

237

240

241 242

243

244

245

246

247

248

249

250

252 253



Figure 2: Method Overview. We learn a neuro-symbolic program for inductive spatial concepts factored as (a) Sketch (b) Plan (c) Generalize. The example above shows the progressive realization of a program for the concept of a staircase acquired by observing a single demonstration of building a staircase of size four and its corresponding natural language instruction, "construct a staircase of size four using magenta legos."

5 LEARNING INDUCTIVE CONCEPTS FROM DEMONSTRATIONS

We address the problem of estimating a succinct generalized program, Eq. 2, modeling an inductive spatial concept modeling structures whose construction is observed in a human demonstration. Direct symbolic search in the space of programs is intractable (particularly due to looping constructs needed 238 for modeling induction) but can explicitly reason over previously acquired concepts. Alternatively, 239 neural methods attempting to predict action sequence to attain the assembly are challenged by continual setting where concepts can increase over time building on previously learnt ones but are resilient to noise. Our approach blends both approaches and factors the concept learning task as:

- Sketch: From the natural language instruction (Λ), we extract a task sketch (H_s^{α}) using an LLM that provides the signature (concept name and instantiated arguments) of the concept to be learned. When grounded in the initial scene of the demonstration, the task sketch provides a particular instance of the concept demonstrated in the given demonstration.
- Plan: MCTS-based search using the already learnt concepts that outputs the sequence of grounded actions, best explaining the given demonstration.
- **Generalize:** The grounded plan H_P^* and task sketch H_S^* are provided to an LLM to obtain a general Python program whose execution on the given scene matches the searched plan.

Formally, the factored exploration of the program space for a demonstration is performed as:

$$H_S^* \leftarrow Sketch(\Lambda; \theta_S); H_P^* \leftarrow Plan(S_1..S_g, H_S^*; \theta_P); H_G^* \leftarrow Generalize(H_P^*, H_S^*; \theta_G)$$
(3)

254 Here, θ_S , θ_P and θ_G are the learnable parameters (including hyperparameters) of the Sketch, Plan 255 and Generalize functions, respectively. The concept library \mathcal{L} is initialized with primitive visual and action concepts. Upon acquiring a new inductive concept $H_G^* = H^*$, we update our library 256 accordingly: $\mathcal{L} \leftarrow \mathcal{L} \cup H^*$. An example is provided in Appendix Sec. A.5. Fig. 2 illustrates an 257 example of progressive prog. estimation. Next, we detail each of the three steps mentioned above. 258

259 260

5.1 (SKETCH) GROUNDED TASK SKETCH GENERATION

261 An LLM driven by in-context learning is used to get a program signature (a sketch) for a concept 262 from the natural language instruction. The task sketch is a tree of nested function calls that outlines 263 the function header (name and the parameters) of the inductive concept/program to be learned. A 264 detailed exposition on prompting appears in the Appendix C.1. The task sketch thus obtained is 265 then grounded on the input scene using a quasi-symbolic visual grounding module akin to Mao 266 et al. (2019); Kalithasan et al. (2023); Wang et al. (2023c). This module has three key components: (1) a visual extractor (ResNet-34 based) that extracts the features of all objects in the scene, (2) a 267 concept embedding module that learns disentangled representations for visual concepts like green 268 and *dice*, and (3) a quasi-symbolic executor equipped with pre-defined behaviours such as "filter" 269 to select/ground the objects of interest. For example, grounding the task sketch "Tower (height =3, objects = filter(green, dice))" results in an instantiated function call "Tower(height = 3, objects = [1, 2, 3])" where [1, 2, 3] are the green coloured dice.

272 273

274

5.2 (PLAN) PHYSICAL REWARD GUIDED PLAN SEARCH

The plan search involves finding a generalizable plan that effectively explains the demonstration S_1, \dots, S_g . Specifically, this involves determining the concepts, their respective grounded parameters, and the order of composition as specified in the Equation 1.

278 **Primitive Actions.** Constructing complex structures involves two steps: (1) identifying or imag-279 ining the placement location of an object/structure and (2) picking and placing the object at the 280 imagined location. The position $pos_{\theta}(.)$ for placement is determined using a head, which rep-281 resents a cuboidal enclosure in 3D space. Conceptually, moving the head is akin to the robot's 282 cognitive exploration of potential placements to achieve the desired spatial configuration. We define a set of primitive functions, A_p , to guide the movement and placement of objects in 283 two stages: (1) move_head(direction): This primitive moves the abstract head to a de-284 sired relative position and (2) keep_at_head(objects): This primitive places the target 285 object from the list objects at the current location of the head. It is important to note that 286 move head (direction) is a neural operator, which takes the head's current position and pre-287 dicts its new location based on the specified direction. This operator is trained on a corpus of 288 pick-and-place instructions, such as "move the green object to the right of the red cube," similar to 289 the approach in Kalithasan et al. (2023). 290

MCTS Search. We use an object-centric state representation defined by bounding boxes (including 291 the depth of the center) and visual attributes of all the objects that are present on the table. For each 292 learned inductive concept <cpt>, we define a macro-action Make_<cpt>(size) that executes 293 the corresponding program with the given size argument, resulting in the construction of the desired 294 concept. Thus, the action space A is the union of primitive actions A_p and compound/macro-295 actions \mathcal{A}_c . Intersection over Union (IoU) between the attained state and the expected state in the 296 demonstration is provided as a reward for all macro actions and keep_at_head (objects); all 297 other actions yield zero reward. An MCTS procedure similar to Khandelwal et al. (2016) is performed 298 to find a plan that maximizes the reward. The node expansion process and reward calculation for the 299 MCTS procedure is detailed in Appendix A.3. The search outputs a sequence of grounded actions for an instantiation of the given inductive concept by the task arguments. 300

Modularity and Scalability. MCTS that searches for a plan only in terms of primitive actions may not be generalizable due to lack of modularity C.3. The use of macro-actions in the search ensures that the plan H_p^* for a given demonstration is concise, modular, and easily generalizable. This can be seen as a form of regularization in terms of the length of concept description by making the prior $P(H) \propto |H|^{-\alpha}$ (where $\alpha > 0$) in equation (2)

$$H^* = \arg\min_{H \in \mathcal{H}} \left[\texttt{Loss}(\{S_1..S_n\}, H(\Lambda)) + \alpha \log |H| \right]$$

308 However, as the action space expands with the learning of more concepts, the search becomes 309 slower, necessitating the pruning of the search space. To avoid searching over the size parameter in macro-actions, we greedily select the smallest size that achieves the maximum average reward 310 from the current state. Additionally, to prune primitive actions, we train a reactive policy π_{neural} 311 which, given the current state \tilde{s}_t and the next expected state s_{t+1} (from the demonstration), outputs 312 one of the primitive actions $a_t^* \in \mathcal{A}_p$. Consequently, the effective branching factor of the search is 313 reduced from $|\mathcal{A}_c| + |\mathcal{A}_p|$ to $|\mathcal{A}_c| + 1$. Thus, our MCTS algorithm is modular through the hierarchical 314 composition of learned concepts and efficient through action space pruning, and is referred to as 315 MCTS+L+P. Further details regarding modifications in MCTS are given in the Appendix A.3.

316 317

318

306 307

5.3 (GENERALIZE) PLAN TO PROGRAM ABSTRACTION

Leveraging the code generation and pattern matching abilities of LLMs (Mirchandani et al., 2023), we use GPT-4 to distil out a general Python program from the sequence of grounded actions as determined by MCTS+*L*+*P*. The learnt program is incorporated in the concept library, \mathcal{L} , for modular reuse in subsequent learning tasks. Additional details, prompting mechanism and use of learnt programs in the search step of future learning tasks are described in Appendix C.2, A.3. We take a curriculum learning approach beginning from learning of primitive actions and visual attributes, followed by structures of increasing complexity. Appendix C.2, Fig. 2, A.4 details the curriculum used for concept learning, architecture details, and learning from multiple demonstrations.

6 EVALUATION SETUP

324

325

326 327 328

Corpus. A corpus is created using a simulated Robot Manipulator assembling spatial structures on a table-top viewed by a visual-depth sensor. Demonstration data (3 demonstrations per structure, with up to 20 objects present in the scene) includes observations (via a visual-depth camera) of the action sequence (picking and placing of blocks) resulting in the construction of the final assembly using varied block instances and types (e.g., cubes, dice, lego etc.).

The scope of concepts and associated evaluation tasks 335 are adapted from closely related works. The stair-336 case and enclosure construction tasks are inspired from 337 from Silver et al. (2019), adapted to 3D from the original 338 2D grid world setting. Structures such as boundaries in-339 volving repetitive use of columns and rows (w/o explicit 340 joint fastening) are inspired by a robotic assembly data 341 set (Collins et al., 2024). Finally, the arc-bridge and 342 x-shaped patterns are inspired from concept learning 343 works as Lake et al. (2015). A total of 15 structures 344 types are incorporated and are additionally modulated in size/spatial arrangement for generalization evaluation. 345



Figure 3: Illustrative examples of spatial structures from corpus showing inductive composition over simpler structures. Details and visualizations in Appendix B

- Three evaluation data sets are formed each with *simple* structures and *complex* structures composed of simpler concepts (e.g., staircase consists of towers as substructure). *Dataset I* and *II* contain demonstrations constructing structures with size(.) \in [3, 5], where size is defined in 4. *Dataset II* reverses the linguistic labels used (e.g., the "tower" in I becomes "rewot" in II) to assess model reliance on pre-training knowledge in presence of new labels for concepts. *Dataset III* includes concepts of larger size than those in training to test generalization.
- **Baselines.** Four baselines are formed from two alternative approaches as follows.

(1) *Purely-Neural*: An end-to-end neural model inspired by StructDiffusion (Liu et al., 2023) that
treats structure construction as a rearrangement problem. We consider two variations of the model:
(1.1) Struct-Diff (SD): End-to-end approach without any additional supervision regarding which
objects need to be moved. (1.2) With-Grounder (SD+G): Similar to (1.1) except that we assume a
perfect object selector/grounder that identifies the relevant set of objects which are to be moved.

358 (2) Pre-trained models that directly output symbolic programs: (2.1) LLMs for Scene-Graph 359 Reasoning: This approach uses a Pre-trained Language Model (GPT-4) to generate Python programs 360 from instructions which describe the given demonstration. To help the LLM understand the underlying 361 structure, it is provided with the symbolic spatial relationships (e.g., left(a,b)) between objects in the 362 demonstration. For this baseline we further assume absence of distractor objects in the scene. (2.2) 363 Vision Language Model (GPT-4V): Similar to (2.1) but has the ability to take input demonstration as images. For learning the program of a new inductive concept, we give the demonstration to the VLM 364 in the form of Λ , $(S_1...S_q)$. Additional details on prompting method in Appendix C.6. Experiments were also conducted with open source code-generation LLMs such as CodeLlama (70Bq), Due to 366 significantly poorer performances w.r.t. GPT-4, GPT-4 was retained as the primary LLM baseline. 367

Model variants. We implement three variants of the MCTS search to perform a grounded plan search over the action space A: (i) MCTS+*P*+*L*: Our approach as described in section 5.2, that uses the learnt concepts from \mathcal{L} as macro actions in subsequent searches (e.g., Make_Tower (3, objects) $\in A_c$) (*L*). Further it performs pruning of A_p using π_{neural} (*P*), (ii) MCTS-*P*+*L*: Our approach without neural pruning and (iii) MCTS+*P*-*L*: No access to library of concepts during continual learning and thereby lacks ability to use macro actions. This method greedily selects the action from A_p as given by π_{neural} . We provide more details about the 3 methods in Appendix D.4.

375 Metrics. We adopt the following metrics to evaluate our models: (i) *Program Accuracy:* A binary
 376 score obtained through human evaluation. 1 for constructing the structure fully, 0 otherwise. (ii)
 377 *Target Construction IoU:* Intersection over Union (2D-IoU) between bounding boxes. (iii) *Target Construction Loss:* Mean Squared Error (MSE) loss over the bounding boxes + depth of the center.

378 7 RESULTS

379 380 381

Our experiments evaluate the following questions. **Q1:** How does our model perform when compared to baselines in terms of concept learning and execution ability (In-Distribution)? Q2: How does our 382 model generalize to concept instances not seen (larger) during training (Out-of-Distribution)? Q3: How robust and efficient is our concept learning pipeline? **Q4:** How can the acquired concepts be 384 used in for embodied instruction following tasks?

385 386

387

388

389

390

391

392

393

394

395

396

397

398

399

419

430

Q1: CONCEPT LEARNING ACCURACY

We compare the program accuracy (Table 1) and the IoU/MSE values (Table 2) of the final states attained by SPG and the baselines w.r.t. the gold states in the in-distribution setting and find that SPG significantly outperforms other approaches. Values for Purely neural approaches are marked NA because Neural Outputs are not physically grounded. We make the following observations: (i) For *complex* compositional structures, the accuracy of the pre-trained models is poor (zero), indicating their inability to reason over the numerous and complex spatial relations present in these structures. (ii) While program inference via the LLM is better than the VLM for learning *simple* structures, it is worse for *complex* structures. This indicates the inherent weakness of the textual descriptions of complex spatial relations present in complex structures. (iii) While the data-intensive purely neural approaches perform much better on *complex* structures when compared to the pre-trained foundation models, they are still weaker than SPG.

Table 1. Program Accuracy

Table 2: In-distribution Performance (Mean \pm Std-error)

			Model	Sin	ıple	Complex	
Model	Simple	Complex	Widder	IoU	MSE (1e-3)	IoU	MSE (1e-3)
SPG(Ours)	1.00	0.83		0.01 0.00	0.01 + 0.00	0.05 0.02	2 0 C 1 02
GPT-4V	0.33	0.00	SPG(Ours)	0.96 ± 0.00	0.01 ± 0.00	0.85 ± 0.02	2.06 ± 1.02
GPT-4	0.78	0.00	GPT-4V	0.75 ± 0.01	4.33 ± 0.41	0.50 ± 0.02	7.29 ± 1.10
SD+G	NA	NA	GPT-4	0.89 ± 0.01	1.36 ± 0.26	0.28 ± 0.02	13.5 ± 1.65
SD	NA	NA	SD+G	0.74 ± 0.01	1.42 ± 0.29	0.61 ± 0.02	2.43 ± 0.48
			SD	0.49 ± 0.01	1.48 ± 0.24	0.46 ± 0.02	3.71 ± 1.53

Q2: GENERALIZATION PERFORMANCE

410 Table 3, compares the generalization performance on *Dataset III* for models trained on *Dataset I* 411 (full table in Appendix, 8). We see that SPG outperforms other approaches. We further consider 412 the relative decrease (R.D.) in performance (2D-IoU) on going from the in-distribution to the out-of-413 distribution (OOD) setting. We make the following observations: (i) SPG suffers a relative decrease 414 of 7.27% for simple and 5.74% for complex structures. (ii) In contrast, the SD+G baseline shows 415 a large R.D. of 63.25% on simple structures and 74.72% on complex structures; highlighting the 416 inability of Purely Neural Models to generalize inductively. (iii) Pre-trained models also have a large 417 R.D. in perf. for complex structures (GPT-4 : 53.87% & GPT4V : 41.64%), which is attributed to 418 their inability to generate the correct program that can generalize inductively to unseen data.

420	Table 3: OOD Performance. R.D% is the relative
421	decrease in IoU from Table 2. MSE is in 1e-3 units

Table 4: Perf. on *Dataset II* with Reversed Names. Acc. is Prog. Accuracy, MSE in 1e-3 units

Model		Simple			Complex	x			<u> </u>			<u>a 1</u>	
Model	IoU	R.D%	MSE	IoU	R.D%	MSE	Model	Simple		e		Complex	
SPG(Ours)	0.89	7.27	0.43	0.80	5.74	1.49		Acc.	IoU	MSE	Acc.	IoU	MSE
GPT-4V	0.58	23.33	13.2	0.29	41.64	10.9	SPG(Ours)	0.88	0.86	1.74	0.78	0.78	3.93
PT-4	0.78	12.61	5.51	0.13	53.87	19.1	GPT-4V	0.23	0.71	3.92	0.00	0.09	21.29
D+G	0.27	63.25	6.21	0.15	74.72	14.2	GPT-4	0.67	0.78	3.16	0.00	0.00	22.73
)	0.24	51.84	6.86	0.15	67.67	11.6				-			

Q3: ROBUSTNESS AND EFFICIENCY ANALYSIS

Reliance on pre-trained Knowledge vs. Demonstration. Next, we evaluate the degree to which 431 concept learning relies on prior knowledge vs. the action sequences observed in demonstrations. We

432 compare pre-trained models against our approach by learning programs on Dataset II (6), which 433 uses arbitrary names for concepts. This forces all models to rely on demonstration data because there 434 is no real-world knowledge associated with the name of the concept, say "rewot" instead of "tower". 435 Table 4 indicates the corresponding performances, with our approach outperforming others. For 436 the IoU/MSE values along with standard errors refer to Appendix Table 9. The relative decrease in performance (program accuracy w.r.t. Table 1) for simple structures for our approach (12%) is 437 lower than GPT-4 (14%) and much lower than GPT-4V (30%). The poorer generalization of pre-438 trained models can be attributed to their over-reliance on prior knowledge and failure to effectively 439 incorporate the data from demonstrations. In contrast, SPG better captures the semantics of a novel 440 concept, especially ones whose knowledge may not be available for the LLMs/VLMs at training time. 441

442 MCTS Variants for Concept Learning.

466

467

468

443 Figure 4 compares the program accuracy for the three meth-444 ods of plan search. For the MCTS-L+P method, the program 445 accuracy is expected to be independent of expansion steps as 446 it greedily chooses the action for which π_{neural} gives the highest 447 probability. For the MCTS+L based methods the accuracy in-448 creases beyond 0.6 with time, which demonstrates that having 449 a composable library of concepts allows us to learn a much richer class of inductive concepts. MCTS+P+L saturates to a 450 program accuracy of 0.933 in just 4000 expansion steps com-451 pared to MCTS+L-P taking 512000 expansion steps, which 452 demonstrates a significant increase in learning efficiency via 453 use of the neural pruner. For very low number of expansions 454



Figure 4: **MCTS Variants.** Num. of expansion steps in search (log scale) (X-axis) vs Program accuracy (Y-axis).

steps (<40) accuracy of MCTS+L based methods is lower than MCTS-L as the former expends expansion steps on UCB exploration (instead of greedy actions).

Significance of MCTS in SPG. To assess the necessity and importance of MCTS, we carry out an 457 ablation study where we replace it with an LLM planner during the planning stage, referred to as 458 SPG-M+LMP. In the "plan" stage of our pipeline, GPT-4V is prompted to output a plan given the 459 concept library and RGB keyframes from the demonstration. Our experiment shows that GPT-4V 460 struggles to generate correct plans, particularly for complex structures like pyramids, arch bridge 461 and boundaries, resulting in significantly lower performance than SPG, see 5. Additionally, some 462 plans generated by GPT-4V are not physically grounded, leading to errors in both the planning 463 and generalization stages, which compounds the inaccuracies. This demonstrates that combining 464 symbolic search with LLMs offers a substantial advantage over using only LLMs. 465

Figure 5: Ablation studies and Disentanglement. Left: Ablations with SPG-M+LMP and GPT-4V+VR. MSE values are in 1e-3. Right: The acquisition of new visual concepts. Plot shows an increase in the likelihood of correct grounding of an object referenced with a new neural concept (*chocolate* color) with training iterations.

Model		Simple	e	(Comple	ex	
viouei	Acc.	IoU	MSE	Acc.	IoU	MSE	у тэ _{«7-}
SPG(Ours)	1.0	0.96	0.01	0.83	0.85	2.06	op)
SPG-M+LMP	0.55	0.68	11.1	0.16	0.19	20.0	
GPT-4V+VRF	0.66	0.75	6.8	0.16	0.46	12.0	2 ^{es}

477 Effectiveness of Continual Learning of Visual Concepts Having a disentangled representations 478 allows us to (i) intersperse learning of new visual attributes with learning of new inductive concepts 479 (ii) avoid catastrophic forgetting of already learnt attributes. For example, the model can learn the 480 chocolate from an instruction "construct a tower using chocolate blocks of size 4", even if it has not 481 seen the color in the pretraining phase. Because of our modular architecture, we can learn the color 482 as a new embedding in the space of visual attributes. The plot in Fig. 5 demonstrates the benefit of 483 having such disentangled representations. As the training proceeds the probability of being able to select the chocolate blocks when required increases with time, while keeping the ability of selecting 484 a magenta colored object (when required) remains the same. Additional details for continual learning 485 of visual concepts appear in Appendix D.2.

 Ablating with Pre-trained Models + Visual Reward Filter In line with program synthesis techniques using LLMs (Li et al., 2022; Chen et al., 2021), we sample five programs from GPT-4V and rank them according to the visual reward obtained from their execution. Furthermore, we provide the ground-truth programs of the concepts that are needed to learn the given new concept, thus employing a form of teacher forcing in program generation. Even with these measures, it performs significantly worse than SPG, see Table 5 (GPT-4V+VRF). While this performance is better than that of GPT-4V, it still unable to generate correct programs, especially for complex structures.

493 494

520

521

522

523

524 525 526

527

Q4: APPLICATION OF LEARNT CONCEPTS FOR ROBOT INSTRUCTION FOLLOWING

495 **Complex instruction execution via LLM.** We demonstrate our ability to use the learnt program 496 representations to perform complex guided robot manipulation tasks. We instruct the robot to perform 497 tasks like : "Construct a tower of green die having the same height as the existing tower of white 498 die." and "Construct a tower of total 6 blocks using alternating blue and red blocks". For both the 499 above tasks we prompt GPT-4 by providing it with the set of learnt inductive concepts, the set of 500 primitive actions, and some pre-defined helper functions by using Python import statements in a 501 manner similar to Liang et al. (2023). GPT-4 generates an executable Python code in terms of these 502 functions, which, on running, generates the resultant and required action sequence. Figure 6 (top) 503 illustrates task execution (also see Appendix D.3).

Grounding learnt concepts into visual input for plan synthesis. We further demonstrate that the concepts we have acquired can help us to perform goal conditioned planning. Fig. 6 (bottom) demonstrates the results of our approach for the tasks of constructing a staircase beginning from adversarial and assistive initial states. Note, the planner that we learn is a grounded neuro-symbolic planner, as a PDDL based planner cannot be hand-coded easily (D.5), and LLMs/VLMs are unable to perform such complex reasoning tasks (see Appendix D.5).



Figure 6: Application of learnt concepts. Top: Using LLM to generate the executable code for a novel tasks, given the concept definitions. Bottom: Integrating a neuro-symbolic planner over the concepts. Bottom-left: The planner is able to optimally replace the green cube from the adversarial initial state by unstacking and re-stacking the faulty tower. Bottom-right: The planner is able to complete a staircase from an initially constructed row by layering rows upon rows, a method of construction it has not seen while learning staircase.

8 CONCLUSION

528 This paper introduces a novel approach for learning inductive representation of grounded spatial 529 concepts as neuro-symbolic *programs* via language-guided demonstrations. Our approach factors pro-530 gram learning as: *Sketch:* generating the high-level program signature via an LLM, *Plan:* searching 531 for a grounded plan that maximises the total discounted reward with the respect to the demonstration, 532 and Generalize: abstracting the grounded plan into an inductively generalize-able abstract plan via 533 an LLM. Continual learning is achieved via learning of modular programs by giving preference to 534 shorter programs through composition of learnt ones. Extensive evaluation demonstrates accurate program learning and stronger generalization in relation to purely LLM based as well as purely 536 neural baselines. Grounding of learned concepts in visual data facilitates reasoning and planning for embodied instruction following. Limitations include reliance on perfect demonstrations, assumption of full observability of all objects and experiments confined to simulation. Incorporating noisy 538 demonstrations, reasoning with beliefs and interleaving planning and execution remains part of future work.

540 REFERENCES

562

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea
 Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, et al. Do as i can, not as i say:
 Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- Abdeslam Boularias, Felix Duvallet, Jean Oh, and Anthony Stentz. Grounding spatial relations for outdoor robot navigation. In 2015 IEEE International Conference on Robotics and Automation (ICRA), pp. 1976–1982. IEEE, 2015.
- 552 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Ka-553 plan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, 554 Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, 556 Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher 558 Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, 559 Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language 561 models trained on code, 2021. URL https://arxiv.org/abs/2107.03374.
- Jack Collins, Mark Robson, Jun Yamada, Mohan Sridharan, Karol Janik, and Ingmar Posner. Ramp:
 A benchmark for evaluating robotic assembly manipulation and planning. *IEEE Robotics and Automation Letters*, 9(1):9–16, January 2024. ISSN 2377-3774. doi: 10.1109/lra.2023.3330611.
 URL http://dx.doi.org/10.1109/LRA.2023.3330611.
- Kevin Ellis, Daniel Ritchie, Armando Solar-Lezama, and Josh Tenenbaum. Learning to infer graphics
 programs from hand-drawn images. *Advances in neural information processing systems*, 31, 2018.
- Kevin Ellis, Catherine Wong, Maxwell Nye, Mathias Sablé-Meyer, Lucas Morales, Luke Hewitt, Luc Cary, Armando Solar-Lezama, and Joshua B Tenenbaum. Dreamcoder: Bootstrapping inductive program synthesis with wake-sleep library learning. In *Proceedings of the 42nd acm sigplan international conference on programming language design and implementation*, pp. 835–850, 2021.
- Gabriel Grand, Lionel Wong, Matthew Bowers, Theo X Olausson, Muxin Liu, Joshua B Tenenbaum, and Jacob Andreas. Lilo: Learning interpretable libraries by compressing and documenting code. *arXiv preprint arXiv:2310.19791*, 2023.
- Thomas M Howard, Stefanie Tellex, and Nicholas Roy. A natural language planner interface for
 mobile manipulators. In 2014 IEEE International Conference on Robotics and Automation (ICRA),
 pp. 6652–6659. IEEE, 2014.
- Wenlong Huang, F. Xia, Ted Xiao, Harris Chan, Jacky Liang, Peter R. Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, Pierre Sermanet, Noah Brown, Tomas Jackson, Linda Luu, Sergey Levine, Karol Hausman, and Brian Ichter. Inner monologue: Embodied reasoning through planning with language models. In *Conference on Robot Learning*, 2022.
- 586 Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax, 2017.
- Namasivayam Kalithasan, Himanshu Singh, Vishal Bindal, Arnav Tuli, Vishwajeet Agrawal, Rahul
 Jain, Parag Singla, and Rohan Paul. Learning neuro-symbolic programs for language guided robot
 manipulation. 2023.
- Piyush Khandelwal, Elad Liebman, Scott Niekum, and Peter Stone. On the analysis of complex backup strategies in monte carlo tree search. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning Volume 48*, ICML'16, pp. 1319–1328. JMLR.org, 2016.

594 595 596	Dohyun Kim, Nayoung Oh, Deokmin Hwang, and Daehyung Park. Lingo-space: Language- conditioned incremental grounding for space. <i>arXiv preprint arXiv:2402.01183</i> , 2024.
597 598 599	Royi Lachmy, Valentina Pyatkin, Avshalom Manevich, and Reut Tsarfaty. Draw me a flower: Processing and grounding abstraction in natural language. <i>Transactions of the Association for Computational Linguistics</i> , 10:1341–1356, 2022.
600 601 602	Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept learning through probabilistic program induction. <i>Science</i> , 350(6266):1332–1338, 2015.
603 604	Richard Li, Allan Jabri, Trevor Darrell, and Pulkit Agrawal. Towards practical multi-object manipu- lation using relational reinforcement learning, 2019.
605 606 607 608 609 610 611	Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien de Masson d'Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level code generation with alphacode. <i>Science</i> , 378(6624):1092–1097, December 2022. ISSN 1095-9203. doi: 10.1126/science.abq1158. URL http://dx.doi.org/10.1126/science.abq1158.
612 613 614	Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control, 2023.
615 616	Weiyu Liu, Yilun Du, Tucker Hermans, Sonia Chernova, and Chris Paxton. Structdiffusion: Language- guided creation of physically-valid structures using unseen objects, 2023.
617 618 619	Weiyu Liu, Geng Chen, Joy Hsu, Jiayuan Mao, and Jiajun Wu. Learning planning abstractions from language. <i>arXiv preprint arXiv:2405.03864</i> , 2024.
620 621 622	Jiayuan Mao, Chuang Gan, Pushmeet Kohli, Joshua B. Tenenbaum, and Jiajun Wu. The neuro- symbolic concept learner: Interpreting scenes, words, and sentences from natural supervision, 2019.
623 624 625 626	Jiayuan Mao, Tomás Lozano-Pérez, Josh Tenenbaum, and Leslie Kaelbling. Pdsketch: Integrated domain programming, learning, and planning. Advances in Neural Information Processing Systems, 35:36972–36984, 2022.
627 628 629	Suvir Mirchandani, Fei Xia, Pete Florence, Brian Ichter, Danny Driess, Montserrat Gonzalez Arenas, Kanishka Rao, Dorsa Sadigh, and Andy Zeng. Large language models as general pattern machines, 2023.
630 631 632	Meenal Parakh, Alisha Fong, Anthony Simeonov, Tao Chen, Abhishek Gupta, and Pulkit Agrawal. Lifelong robot learning with human assisted language planners, 2023.
633 634 635	Rohan Paul, Jacob Arkin, Derya Aksaray, Nicholas Roy, and Thomas M Howard. Efficient grounding of abstract spatial concepts for natural language interaction with robot platforms. <i>The International Journal of Robotics Research</i> , 37(10):1269–1299, 2018.
636 637 638 639	Ankit Shah, Pritish Kamath, Julie A Shah, and Shen Li. Bayesian inference of temporal task specifications from demonstrations. <i>Advances in Neural Information Processing Systems</i> , 31, 2018.
640 641	Tom Silver, Kelsey R. Allen, Alex K. Lew, Leslie Pack Kaelbling, and Josh Tenenbaum. Few-shot bayesian imitation learning with logical program policies, 2019.
642 643 644 645	Tom Silver, Rohan Chitnis, Nishanth Kumar, Willie McClinton, Tomás Lozano-Pérez, Leslie Kael- bling, and Joshua B Tenenbaum. Predicate invention for bilevel planning. In <i>Proceedings of the</i> <i>AAAI Conference on Artificial Intelligence</i> , volume 37, pp. 12120–12129, 2023.
646 647	Tom Silver, Soham Dan, Kavitha Srinivas, Joshua B Tenenbaum, Leslie Kaelbling, and Michael Katz. Generalized planning in pddl domains with pretrained large language models. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pp. 20256–20264, 2024.

648 649 650	Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. Progprompt: Generating situated robot task plans using large language models. <i>arXiv preprint arXiv:2209.11302</i> , 2022.
652	Richard S Sutton and Andrew G Barto. Reinforcement learning: An introduction. MIT press, 2018.
653 654 655	Stefanie Tellex, Thomas Kollar, Steven Dickerson, Matthew R Walter, Ashis Gopal Banerjee, Seth Teller, and Nicholas Roy. Approaching the symbol grounding problem with probabilistic graphical models. <i>AI magazine</i> , 32(4):64–76, 2011.
657 658 659 660	Joshua B. Tenenbaum, Charles Kemp, Thomas L. Griffiths, and Noah D. Goodman. How to grow a mind: Statistics, structure, and abstraction. <i>Science</i> , 331(6022):1279–1285, 2011. doi: 10.1126/science.1192788. URL https://www.science.org/doi/abs/10.1126/science.1192788.
661 662	Weikang Wan, Yifeng Zhu, Rutav Shah, and Yuke Zhu. Lotus: Continual imitation learning for robot manipulation through unsupervised skill discovery, 2023.
663 664 665 666	Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. <i>arXiv</i> preprint arXiv:2305.16291, 2023a.
667 668 669	Huaxiaoyue Wang, Gonzalo Gonzalez-Pumariega, Yash Sharma, and Sanjiban Choudhury. Demo2code: From summarizing demonstrations to synthesizing code via extended chain-of-thought. <i>arXiv preprint arXiv:2305.16744</i> , 2023b.
670 671 672	Renhao Wang, Jiayuan Mao, Joy Hsu, Hang Zhao, Jiajun Wu, and Yang Gao. Programmatically grounded, compositionally generalizable robotic manipulation. <i>arXiv preprint arXiv:2304.13826</i> , 2023c.
673 674 675 676	Zihao Wang, Shaofei Cai, Anji Liu, Yonggang Jin, Jinbing Hou, Bowei Zhang, Haowei Lin, Zhaofeng He, Zilong Zheng, Yaodong Yang, Xiaojian Ma, and Yitao Liang. Jarvis-1: Open-world multi-task agents with memory-augmented multimodal language models, 2023d.
678 679	
680 681 682	
683 684 685	
686 687	
688 689 690	
691 692 693	
694 695	
697 698	
699 700 701	

A ADDITIONAL DETAILS ON TECHNICAL APPROACH

Figure 7 illustrates the pipeline for online inference to realize to realize construction of novel structures.



Figure 7: **SPG: Inference** First the library of concept \mathcal{L} is loaded with the corresponding set of learnt programs. Then the given instruction is converted into task-sketch H_s , which is grounded in the initial scene. The required program is fetched from the library, and the grounded task-sketch is executed based on the semantics of the learnt program.

A.1 SYMBOLIC CONSTRUCTS AND THEIR SEMANTICS USED IN PROGRAMS

Table 5 defines the types of the signature and semantics of all the operators. Table 6, includes the type definition of various symbols. Standard Python constructs such as for loops, if else \cdots) as assumed in addition to the constructs defined here.

Table 5: **Symbols and Semantics** Signature and semantics for the primitive concepts and operations that are used in the construction of the programs used to express inductive spatial concepts.

Function	Signature	Semantics
filter	(VisualConcept, ObjSet) \rightarrow ObjSet	Returns the objects that con-
		tain the VisualConcept
move_head	$(\text{Head, Dir}) \rightarrow \text{Head}$	Moves the head to the given
		direction (May or may not
		take input/return the head,
		based on a flag)
assign_head a.k.a	$(\text{Head, ObjIdx}) \rightarrow \text{Head}$	Given the index/one-hot rep-
move_head(overloaded)		resentation for an object, it
		moves the head to the posi-
		tion corresponding to that ob-
		ject
keep_at_head	$(ObjSet, Head) \rightarrow None$	Keeps the argmax of ObjSet
	-	at the head
reset_head	None \rightarrow Head	Sets the head to the top posi-
		tion of stack and pops this po-
		sition from the stack as well
store_head	Head \rightarrow None	Pushes the current position
		of head into the stack

A.2 CURRICULUM LEARNING

751 We follow a curriculum approach where the visual concepts are trained first from simpler linguistically-752 described demonstrations (t_0 in figure 8). This is followed by learning of action concepts through 753 sequentially composed pick and place tasks. Its essential to use such long range sequential instructions 754 in order to ensure that the semantics of action concepts are learnt for placement of objects at a height 755 much above the tabletop (t_1 in figure 8). After the pre-training phase, the agent can continually learn 768 new inductive concepts and visual attributes. (t_2, t_3 in figure 8)

Defined Types	Python Type	Usage
IntArg	int	Argument for the structures that de-
-		fines the size (height, length etc)
Obj	torch.Tensor	One-hot vector whose non-zero in-
		dex represents the selected objects
ObjSet	torch.tensor	Probability mask over the selected
		objects
Dir	string	Primitive directions like left, right,
	-	front, top, etc
ConceptName	string	name of the visual, action or induc-
		tive concept
Head	torch.Tensor	Bounding box with depth. 3D
		cuboidal space.

756	able 6: Symbolic representation. The table lists the type definitions used in the implementation	of
757	PG programs.	



Figure 8: **Continual learning through curriculum:** Using simple pick and place demonstrations we learn visual attributes such as *blue cube*, *yellow lego* (t_0). Using long range instructions which are sequentially concatenated descriptions of pick and place tasks we train our action concepts such as *left, right, top* (t_1). After pre-training the agent can perform continual learning of concepts such as learning generalized representation for *tower* (t_2). Because of disentangled representation of neural and symbolic concepts, interspersed learning of new visual attributes such as *chocolate color* are also possible through few demonstrations of structure creation(t_3).

A.3 DETAILS FOR PLAN-SEARCH AND GENERALIZATION

Modifications to the Simulation and Reward Back propagation Steps: Next, we outline the modifications in the simulation and the reward back propagation steps of the standard MCTS algorithm for our setting. During program search we assume access of intermediate scenes in the demonstration. This allows us to provide intermediate rewards that can guide the search well. We observed that making the following changes in simulation and back propagation step increased the efficiency of our search procedure. Fig 9 illustrates the possible states explored by MCTS and the reward calculation.

• Simulation: Rather than performing Monte Carlo simulations at each newly expanded leaf node (to estimate its value) we completely avoid these simulation steps. This was motivated by the fact that our reward is not completely sparse and the intermediate IoU rewards for each object we place allow us to guide the search effectively.





850 Improving Modularity and Scalability of MCTS procedure: We present additional details on the 851 MCTS procedure for searching for a plan conditioned on a program signature and guided by the 852 demonstration.

849

853 In order to incorporate the objective of searching for physically realizeable plans and facilitating 854 generation via re-use of concepts, the following conceptual changes are incorporated in the standard 855 MCTS procedure Sutton & Barto (2018). 856

Modularity (MCTS + L): We want to allow learning of novel inductive concepts in terms of existing 857 ones. This would ensure that the plan H_P^* corresponding to a given demonstration is concise and 858 can be easily generalized to the H_{α}^{*} . C.3 in appendix give example of two plans for the structure 859 *Pyramid* one which is modular and can be successfully generalized by GPT-4, other for which GPT-4 860 fails in generalization due to lack of modularity. This can be seen as a form of regularization in terms 861 of the length of concept description, by making the prior $P(H) \propto |H|^{-\alpha}$ (where $\alpha > 0$) in equation 862 (2)863

$$H^* = \arg\min_{H \in \mathcal{H}} \left[\text{Loss}(\{S_1 ... S_g\}, H(\Lambda)) + \alpha \log |H| \right]$$
(7)

864 In order to allow modular learning of programs, for every inductive concept already stored in the library we define corresponding action instantiations which can be a potential candidate ac-866 tions during our search. As an example, for the concept Tower we have one of the action in-867 stantiation as Make_Tower(3, objects) which would be the action of constructing desired 868 tower. This can be visualized as a compound/macro-action which is composed of primitive actions keep_at_head(objects), move_head(`top'). We define \mathcal{A}_c as the space of such compound actions, and A_p as the space of primitive action/function consisting reset_head(), 870 move_head(direction), keep_at_head(objects), store_head(). Realization 871 of equation (6) (increased preference of macro-actions over primitive ones) is done through discounted 872 IoU rewards during our search (Make_Tower (3) would have a reward of 1+1+1, as compared to 873 $1 + \gamma^2(1) + \gamma^4(1)$ for a sequence of 3 (keep_at_head(objects), move_head('top')0). 874 We refer to MCTS using concept instantiations from \mathcal{L} as macro actions in search as MCTS+L. 875 **Scalability** (MCTS+P): as more concepts are added to the library \mathcal{L} the action space of our search 876 \mathcal{A} (specifically \mathcal{A}_c , the space of compound concepts) increases, therefore we want to prune the 877 search space effectively. For this during the pre-training phase we train a reactive policy π_{neural} 878 which given the current state, \tilde{s}_t and the next expected state s_{t+1} (part of the demonstration) would 879 output one of the primitive action, $a_t^* \in \mathcal{A}_p$ where \mathcal{A}_p is the primitive action space, i.e. $a_t^* =$ 880 $\pi_{neural}(a_t|\tilde{s}_t, s_{t+1}), a_t \in \mathcal{A}_p$ Note that while expanding our search tree we only search among the space of compound actions \mathcal{A}_c and the action a_t^* , thereby reducing the branching factor of search 882 from $|\mathcal{A}_c \cup \mathcal{A}_p|$ to $|\mathcal{A}_c| + 1$. We refer to MCTS using neural pruning as MCTS+P. Therefore our 883 MCTS algorithm is modular through hierarchical composition of learnt concepts and efficient through pruning of action space and is referred to as MCTS+L+P. 885 **Generalization**: Given multiple equal length plans for a given demonstration, we seek to recover a plan one that can be easily generalized by the LLM. 5 shows a plan which could be correctly 887 abstracted out into a generic program by GPT-4. Whereas 15 shows another plan with similar 888 semantics, for which GPT-4 is unable to correctly find the generalized program (Note that row or 889 column of size 1 is equivalent to keep_at_head). We tackle this problem in the following manner. 890 1. Rather than getting a single plan from the plan search we get the top k plans $\{H_{P,i}\}_{i=1}^{i=k}$. In 891 892

- 1. Rather than getting a single plan from the plan search we get the top k plans $\{H_{P,i}\}_{i=1}^{k-1}$. In order to get these top k plans we expand the complete tree (based on UCB criteria) starting from the root node corresponding to the initial state, till a predefined budget of expansions. Then we select the top k paths(potential plans) from the root node to all the leaf nodes (where the top k ones are those that give the highest accumulated IoU reward with respect to the given demonstration).
 - 2. Later we abstract out each of these k plans into corresponding generalized programs, $\{H_{G,i}\}_{i=1}^{i=k}$ using GPT-4. We again run each of these programs on the given demonstration and then choose the one which gives the highest IoU reward (resolving ties based on predefined order). Note that some program $H_{G,i}$ upon execution may result in a plan $\tilde{H}_{P,i}$ different from the original plan $H_{P,i}$ using which it was generalized. This can be attributed to potential errors in GPT-4s program generalization process.

A.4 ADDITIONAL DETAILS: LEARNING WITH INCREASING NUMBER OF DEMONSTRATIONS

Given k demonstrations for a novel inductive concept, we independently find k task sketch $\{H_{S,i}^*\}_{i=1}^{i=k}$ and grounded plans $\{H_{P,i}^*\}_{i=1}^{i=k}$. During the generalization phase we give these k pair of task-sketch and corresponding plans to GPT-4 and ask into infer a single abstraction over them. C.5 in appendix gives a concrete example. Equation for generalize step (getting H_G^* from H_P^*, H_S^*) can be modified as follows.

$$H_G^* \leftarrow Generalize(H_G \mid \{H_{P,i}^*, H_{S,i}^*\}_{i=1}^{i=k}; \theta_G), \ H_G \in \mathcal{H}^{\mathcal{G}}$$

$$(8)$$

912 913 A.5 DETAILED EXPERIMENTAL METHODOLOGY

893

894

895

896

897

899

900

901

902 903

904 905

906

907

908

909

910 911

Input to the Method: The input consists of a language instruction and a human demonstration represented as a sequence of RGBD keyframes.

917 Output/Aim of the Model: The goal is to learn a representation of the unknown concept in the instruction, assuming there is only one unknown concept. If the unknown concept is inductive (e.g.,

"tower"), the model learns a program definition deftower() and stores it in the program library. If the unknown concept is a primitive concept, a concept embedding is learned through backpropagation.

Evaluation: The learned model is evaluated based on the correctness of the program representation
 for inductive concepts and the correctness of object placements, measured through the Intersection
 over Union (IoU) metric (see Metrics, line 262). d. Examples: Suppose the current library contains
 the concepts "red", "tower". Given the instruction "construct a staircase of height 3 using red blocks,"
 the process is as follows:

Parsing: The instruction is parsed into a sketch: Staircase(height=3, objects=filter(red, blocks)).

Grounding: The "objects" parameter is grounded using the visual grounder, which identifies the indices of the red blocks, e.g., [1, 2, 3, 4, 5, 6]. That is, filter(red, blocks) = [1,2,3,4,5,6].

Planning: The planning step uses the demonstrations (sequence of keyframes) to identify the sequence of actions that best explains the demonstrations. In this case, the plan might be: Tower(height=1, objects=[1,2,3,4,5,6]), move_head(right), Tower(height=2, objects=[2,3,4,5,6]), move_head(right), Tower(height=3, objects=[4,5,6]

Generalization: The generalization step abstracts the plan obtained from three such demonstrations into a program. The resulting program would be:

```
1 def staircase(height, objects):
2   for i in range(height):
3      tower(height=i, objects)
4      move_head(right)
```

Note: Whenever an object is placed, the objects list is modified in place, and the index of the placed object is removed. primitive actions: The movement of any object is achieved by first determining the placement position by moving the head (an imaginary bounding box) in specific directions and then placing the object to be moved at the head. The head is implemented as a 3D bounding box defined by coordinates (x1, y1, x2, y2, d), where x1, y1, x2, and y2 are the 2D corners of the bounding box, and d is the depth at the center. The primitive action move_head(direction) shifts the bounding box in the required direction. The action keep_at_head(object_list) picks the first object in the list and places it at the center of the bounding box. Two other primitives, store_head and reset_head, are used to save the current position of the head, allowing the search to return to useful positions later if needed.

B ADDITIONAL DETAILS REGARDING DATASETS

Figure 10 demonstrates the kind of inductive concepts for which we want to learn generic (i.e instance agnostic) representations.

Dataset for Pre-training: We use 5k examples of constructing *twin-towers* i.e. 2 towers adjacent to each other, for learning semantics of move_head(dir), a basic set of visual attributes, reactive policy π_{neural} , and neural modules required for grounded planning. The twin towers allow us to learn various action semantics for all possible configurations of blocks in 3D-space (and not being limited to blocks placed directly on table top surface). Since we are not aware of the underlying semantics of *tower* during pre-training phase the corresponding natural language instruction consists of step by step pick and place actions. 11 gives example demonstrations from this dataset.

Dataset for Inductive Structures: We learn a variety of structures which we have divided into
 Simple and Complex structures. A structure is considered complex if it can be expressed as an
 inductive composition of simpler structures. As an example, we can express a staircase to be a
 composition of towers of increasing height. The structures are listed in the Table 7. Figure 12 shows
 the hierarchical relationship among these structures in the form of a DAG (directed acyclic graph).
 F.3 gives the ground truth program representations for each structure.

968

934

935 936

937

938

939 940

941

942

943

944

945

946

947

948 949 950

951

C PROMPTING STRATEGY AND EXAMPLES

969 970

971 This section gives various prompting examples for our approach and baselines, along with examples motivating particular design decisions in our approach.





Simple Structures	Complex Structures			
Row, Column, Tower	X (cross-shape), Staircase			
Inverted-Row, Inverted-Column	Inverted-Staircase, Pyramid			
Diagonal-45, Diagonal-135	Arch-Bridge, Boundary			
Diagonal-225, Diagonal-315				

C.1 PROMPT EXAMPLE FOR TASK SKETCH GENERATION STAGE (Sketch)

In order to get a program representation (high level task sketch) of the given natural language instruction, we prompt GPT-4 with few shot examples in a manner similar to Liang et al. (2023).
Code segment 1 gives an example of getting the task sketch given the demonstration for constructing a *staircase*. We first import the available primitive operators and functions and also give examples in order to demonstrate the signature of the available primitives(line 1-8). Then we give incontext



```
1080
            for i in range(n):
1081
                tower(i+1, objects)
     0
1082
                move_head('right')
     10
1083
                           Listing 2: Plan to Program using GPT-4 (Generalize)
1084
1085
       C.3 BENEFIT OF ESTIMATING MODULAR/SMALLER PLANS
1086
1087
       The below examples demonstrate the benefit of learning new inductive concepts in terms of already
1088
       acquired inductive concepts (more modular representation). The program 3 is obtained through
1089
       generalization of plan H_p that represents pyramid in terms of rows of decreasing size. The generated
1090
       program generalizes to pyramid of different height.
1091
     1 # input: n = 3, objects = ObjSet
1092
     2 # output: row(5, ObjSet), move_head('right'), move_head('top'), row(3,
1093
           ObjSet), move_head('right'), move_head('top'), row(1, ObjSet)
1094
     3
1095
     4 # Program (GPT-4s output)
     5 def pyramid(n, objects)
1096
     6 for i in range (n, 0, -1):
1097
                # Calculate the number of objects in the current row
1098
                row\_count = 2 * i - 1
     8
1099
     9
1100 10
                # Perform actions for the current layer
                row(row_count, objects)
1101 11
                if i > 1: # No need to move head after the last row
1102<sup>12</sup>
                    move_head('right')
    13
1103
                    move_head('top')
     14
1104
                                  Listing 3: Pyramid in terms of rows
1105
1106
       On the other hand program 4 which is obtained in an attempt to generalize the plan expressed in
1107
       terms of primitives is incorrect.
1108
     1 # input: n = 3, objects = ObjSet
1109
     2 # output: store_head(), keep_at_head(ObjSet), move_head('right'),
1110
           keep_at_head(ObjSet), move_head('right'), keep_at_head(ObjSet),
1111
           move_head('right'), keep_at_head(ObjSet), move_head('right'),
1112
           keep_at_head(ObjSet), reset_head(), move_head('top'), move_head('
           right'), store_head(), keep_at_head(ObjSet), move_head('right'),
1113
           keep_at_head(ObjSet), move_head('right'), keep_at_head(ObjSet),
1114
           reset_head(), move_head('right'), move_head('top), keep_at_head(
1115
           ObjSet)
1116
1117
     4 # Program (GPT-4s output)
     5 def pyramid(n, objects)
1118
     6 # Store initial position
1119
            store_head()
1120
1121
            # Initial sequence of placing and moving right
     9
1122 10
            for _ in range(n):
                keep_at_head(objects)
1123 11
                move_head('right')
1124 <sup>12</sup>
1125<sup>13</sup>
            # Reset and move to top right position
     14
1126
    15
            reset_head()
1127 16
           move_head('top')
            move_head('right')
1128 17
1129 <sup>18</sup>
1130<sup>19</sup>
            # Store again and repeat placing and moving right
            store_head()
     20
1131
            for _ in range(n):
    21
1132 22
                keep_at_head(objects)
1133 23
                move_head('right')
     24
```

```
1134
           # Reset and adjust final position
     25
1135
    26
            reset_head()
1136 27
            move_head('right')
1137 28
            move head ('top')
1138 <sup>29</sup> keep_at_head(objects)
                              Listing 4: Pyramid in terms of primitive actions
1139
1140
1141
            COMPARING SAME-LENGTH PLANS IN TERMS OF GENERALIZABILITY
       C.4
1142
1143
       Below we illustrate plan to program generalization for 2 different plans which are correct and have
1144
       same length, but LLMs find it difficult to generalize the second one as compared to first. (Note that
1145
       row and column of size 1 are equivalent to keep_at_head)
1146
     1 # input: n = 3, objects = ObjSet
1147
     2 # output: keep_at_head(ObjSet), move_head('top'), keep_at_head(ObjSet),
1148
           move_head('top'), keep_at_head(ObjSet)
1149
1150
1151 5 # Program
1152 6 def tower(n, objects):
          for _ in range(n):
     7
1153
     8
                     keep_at_head(objects)
1154
                     move_head('top')
     9
1155
                Listing 5: Plan for tower that can be easily generalized (correct generalization)
1156
1157
     1 # input: n = 3, objects = ObjSet
1158
     2 # output: row(1, ObjSet), move_head('top'), keep_at_head(1, ObjSet),
1159
           move_head('top'), column(1, ObjSet)
1160
1161
     4
1162 5 # Program
1163 6 def tower(n, objects):
            for i in range(1, n + 1):
     7
1164
                     row(i, objects)
1165
                     move_head('top')
     9
1166 10
                     keep_at_head(objects)
1167 11
                     move_head('top')
1168 <sup>12</sup>
                     if i < n:
1169<sup>13</sup>
                          column(i, objects)
     14
                          move_head('top')
1170
                Listing 6: Plan for tower that is difficult to generalize (Incorrect generalization)
1171
1172
1173
       C.5 GENERALIZING VIA MULTIPLE DEMONSTRATIONS
1174
1175
       Given multiple demonstrations we independently find task sketch and corresponding grounded plans
1176
       for each demonstration. These are further given to GPT-4 for generalization. Code segment 7 gives
1177
       an example of getting a single Python program from multiple demonstrations. Note that we explicitly
       prompt the LLM that some of the grounded plans might be incorrect (which may lead to more robust
1178
       generalization in case of noisy demonstrations).
1179
1180
1181
     2 # Function Call: wor(height = 3, objects = ObjSet_1)
1182 3 # Execution: keep_at_head(obj = ObjSet_1), move_head(dir = right),
           keep_at_head(obj = ObjSet_1), move_head(dir = right), keep_at_head(
1183
           obj = ObjSet_1),
1184
     4 # Function Call: wor(height = 3, objects = ObjSet_1)
1185
     5 # Execution: keep_at_head(obj = ObjSet_1), move_head(dir = right),
1186
           keep_at_head(obj = ObjSet_1), move_head(dir = right), keep_at_head(
           obj = ObjSet_1),
1187
     6 # Function Call: wor(height = 3, objects = ObjSet_1)
```

```
1188
     7 # Execution: column(size=1, obj = ObjSet_1), move_head(dir = right),
1189
           keep_at_head(obj = ObjSet_1), move_head(dir = right), keep_at_head(
1190
           obj = ObjSet_1),
1191
    8
     9 #Write the function definition, which generalizes the above executions.
1192
           Note that some of the executions can be partially wrong.
1193
     10 '''python
1194
     11 def wor(height, objects):
1195
    12
1196
    13
    14 GPT-4s Output .....
1197
                              Listing 7: Generalizing through multiple plans
1198
1199
1200
       C.6 PROMPT EXAMPLES FOR LEARNING PROGRAMS USING LLM/VLM MODELS
1201
1202
       Below we describe the prompting methodologies for learning programs through LLM/VLM models.
1203
       Note that although the prompt examples described below are for the case of learning novel structure
1204
       from 1 demonstration, we use 3 demonstration per novel structure in our main results (for both our
       approach and LLM/VLM baseline).
1205
1206
       LLM/GPT-4 Code segment 8 depicts our prompting methodology given a demonstration for a new
1207
       concept tower. For this baseline we aim to check demonstration following and spatial reasoning
1208
       abilities of LLMs (GPT-4). We provide supervision of the intermediate scenes by using tokenized
1209
       spatial relations between objects in the scene (Given in the form of Scene = [right(1, 0) ...]). We
1210
       further assume that only those objects that are required to perform the task are present in the scene
1211
       (no distractor objects). For every structure (that needs to be learned at time t) we give LLM a prompt
       providing in-context example on how to generalize (line 19-35), the set of primitive operators (line
1212
       4) available and the set of structures learnt/present in library (till time t-1) (line 5-18). Finally we
1213
       append to this prompt the expected declaration (arguments and keywords arguments) of the inductive
1214
       concept that is to be learnt along with the spatial relations for each scene of the given demonstration
1215
       (36-51). Note that we assume absence of distractor objects for this baseline.
1216
     1 # Consider a block world domain ..... Given a structure creation task
1217
           along with intermediate scnes complete a general Python function for
1218
            it. The function should be in terms of primitive operators and
1219
           already learnt structures that are present in the program library.
1220
           Enclose the function within backtick (''')
1221
     2
     3 primitive_operators = [keep_at_head, move_head ..]
1222
     4 # this would be our program library
1223
     5 learnt_structures = {
1224
            "row": {
     6
1225
     7
                     "program_tree":
                     ...
1226
     8
                     def row(size, objects):
     9
1227
                          for i in range(size):
     10
1228
                              keep_at_head(obj = objects)
     11
1229
                               move_head(dir = 'right')
     12
1230 13
                     ···,
1231 14
                 },
1232 <sup>15</sup>
    16
            . . . . . .
1233
    17
1234
    18 # the example task
1235 19 Example task:- Place all the objects to the right of each other.
1236 20 Final state :- [right(1, 0), right(2, 1), right(3, 2), right(4, 3)]
1237 21 Intermediate scenes :-
1238 22 Scene 0 = []
    23 Scene 1 = []
1239
     24 Scene 2 = [right(1, 0)]
1240 25 Scene 3 = [right(1, 0), right(2, 1)]
1241 26 Scene 4 = [right(1, 0), right(2, 1), right(3, 2)]
     27 Scene 5 = [right(1, 0), right(2, 1), right(3, 2), right(4, 3)]
```

```
1242
     28 Python function :-
1243
     29 '''python
1244 30 def placing_all_right (objects):
1245 31
            for i in range(len(objects)):
                 keep_at_head(objects) # select one object from the objects set
1246 <sup>32</sup>
            and keep the head at this location
1247
     33
             move_head(dir = 'right') # move the head to the right of the
1248
            previous position
1249 34
1250 35 # The current task for which program needs to be found
1251 36 Current task:- Construct a tower of size 6.
1252 37 Final state :- [top(1, 0), top(2, 1), top(3, 2), top(4, 3), top(5, 4)]
     38 Intermediate scenes :-
1253 <sub>39</sub> Scene 0 = []
1254 40 Scene 1 = []
1255 41 Scene 2 = [top(1, 0)]
1256 42 Scene 3 = [top(1, 0), top(2, 1)]

      43
      Scene 4 = [top(1, 0), top(2, 1), top(3, 2)]

      44
      Scene 5 = [top(1, 0), top(2, 1), top(3, 2), top(4, 3)]

1258 _{45} Scene 6 = [top(1, 0), top(2, 1), top(3, 2), top(4, 3), top(5, 4)]
1259 46 Python function :-
1260 47 '`'python
1261 <sup>48</sup> def tower(size, objects):
     49 ??
1262
     50 111
1263
                            Listing 8: Prompting Strategy for LLM baselines (GPT-4)
1264
1265
        VLM/GPT-4-V Unlike LLM, VLMs have the abilities to process the demonstration as a sequence of
1266
        visual frames. Therefore rather than providing the symbolic spatial relations between every scene we
1267
        instead directly provide all the intermediate scenes for the given demonstration. Further we also relax
1268
        the assumption that there are no distractor objects. As shown in figure 13 We first give information
1269
        about the set of primitive operators and the structures that we have already learnt (library of concepts).
1270
        In order to visually ground the semantics of our primitive actions we give 3 example tasks (natural
1271
        language instruction and intermediate scenes) that do not directly refer to any structure, along with
1272
        corresponding sequence of actions taken (# Demonstration for visual grounding). We further provide
1273
        another example (without scenes) demonstrating how to write generalizable Python function for a
1274
        given task using our operators (# Example for generalization). Finally we give the natural language
        instruction and corresponding scenes for the current task along with signature of the program to be
1275
        learnt (# Current task description).
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
```



Figure 13: **Prompt example for VLM-baseline.** Figure shows the prompting strategy for GPT-4-V that includes the primitive actions, the example tasks for visual grounding, example of writing generalizable Python functions and the demonstration frames.

D SUPPLEMENTARY RESULTS

Out-of-Distribution Performance:

Table 8: Out-of-Distribution Performance (mean \pm std-error)

Model		Simple	Complex			
1110401	IoU	MSE	IoU	MSE		
SPG(Ours) GPT-4 GPT-4V SD	$\begin{array}{c} \textbf{0.892} \pm 0.065 \\ 0.776 \pm 0.023 \\ 0.575 \pm 0.026 \\ 0.236 \pm 0.005 \\ 0.272 \pm 0.004 \end{array}$	$\begin{array}{c} \textbf{4.386e-4} \pm 2.387e\text{-}4 \\ 0.006 \pm 0.001 \\ 0.013 \pm 0.001 \\ 0.006 \pm 7.495e\text{-}4 \\ 0.006 \pm 7.495e\text{-}4 \end{array}$	$\begin{array}{c} \textbf{0.804} \pm 0.025 \\ 0.131 \pm 0.019 \\ 0.290 \pm 0.016 \\ 0.150 \pm 0.011 \\ 0.154 \pm 0.010 \end{array}$	$0.001 \pm 5.391e-4$ $0.019 \pm 1.498e-3$ $0.011 \pm 1.316e-3$ $0.011 \pm 2.860e-3$		

Performance *Dataset II* (i.e. name reversed evaluation): Table 4 gives the corresponding program accuracies, while Table 9 give the corresponding IoU/MSE metrics along with standard errors.

1343
 D.1 QUALITATIVE COMPARISON BETWEEN PURELY-NEURAL (STRUCT-DIFF+GROUNDER) VS. OURS(SPG)
 1345

Figure 14 gives compares the qualitative results for our approach against Struct-Diffusion with
grounder on both in-distribution, *Dataset I* and out-of-distribution (larger size), *Dataset III*. In
in-distribution setting the our method performs slightly better in terms of structure creation for
both simple and complex structures, but the difference is not significant. However for out-ofdistribution setting structures created by our approach are much better than those created by Struct-

1351						
1352		Sin	nple	Complex		
1353	Niodel	IoI	MSF (1e-3)	IoU	MSF (10-3)	
1354		100	MBE (10-3)	100	MBE (10-3)	
1355	SPG(Ours)	$\textbf{0.86} \pm 0.03$	$\textbf{1.74} \pm 0.45$	$\textbf{0.78} \pm 0.02$	$\textbf{3.93} \pm 1.09$	
1356	GPT-4	0.78 ± 0.03	3.16 ± 0.51	0.00 ± 0.00	22.73 ± 1.48	
1357	GPT-4V	0.71 ± 0.01	3.92 ± 0.45	0.09 ± 0.02	21.29 ± 1.59	
1358						

 Table 9: Performance on Dataset II (Names Reversed)

Diffusion+Grounder. Further for this setting structure creation by Struct-Diffusion seems to be muchworse for complex structures than simple ones.





1401 D.2 CONTINUAL LEARNING OF NEURAL CONCEPTS

Given demonstration for the task "Construct a tower of height 4 using *chocolate* cubes", we would like to learn the neural embedding for the unknown color *chocolate* (where we assume that tower has

1404 already been learnt and stored in the library \mathcal{L}). First the instruction is converted into corresponding 1405 plan sketch \mathcal{H}_S = tower(4, Filter(**chocolate**, cubes)), which is passed to the visual grounder. The 1406 grounder detects the presence of an unknown attribute *chocolate* as an argument to filter, and randomly 1407 initializes a new neural embedding for it. Using this new embedding along with the already present 1408 embedding of cube and the visual features found through ResNet-34, the quasi-symbolic executor 1409 outputs a grounded task-sketch. The executor executes the grounded task-sketch by getting the semantics of the underlying function i.e. tower from the library \mathcal{L} . MSE+IoU loss computed over the 1410 final scene obtained and the expected final scene is backpropogated through the network. Note that 1411 during backpropogation all the neural modules (action semantics, visual attributes, ResNet-34) are 1412 frozen, except for the newly initialized embedding for chocolate. For the purpose of differentiable 1413 sampling during tower construction we use gumbel-softmax Jang et al. (2017) with masking. Figure 1414 15 illustrates our approach. 1415



Figure 15: Continual learning of visual primitives

1428 D.3 DETAILS FOR INFERENCE ON NOVEL TASKS USING AN LLM 1429

1426 1427

1430 Below we show the Liang et al. (2023) inspired prompting methodology that we use to get the 1431 executable code corresponding to a language specified manipulation task. We initially begin by 1432 importing the helper functions, spatial direction, primitive functions, and learnt inductive concept-1433 s/structures (line 5-11). Then we give few examples for how to use and compose the various functions for different tasks (line 16-83). Finally we give the instruction of current task, and expect GPT-4 to 1434 output the corresponding executable code. 1435

```
1436
     1 # Given a task you have to provide Python code for executing the task
1437
    2
    3 # importing available functions
1438
1439
     5 from spatial_directions import top, front, back, left, right
1440
     6
1441
     7 from primitives import assign_head, move_head, keep_at_head
1442
     8 # HEAD is a imaginary pointer keeping track of the current spatial
          location in consideration
1443
     9
1444
    10 from helpers import find_size, filter
1445
    11 from structures import row, column, tower
1446 12
1447 13
1448 14 #function signature of the imported functions
    15 # finds all the objects with the given color and shape, returns a mask
1449
          denoting the probability of object selection
1450
    16 filter(color, shape)
1451 17
1452 18 # finds the size of the structure struct_name that is formed with objects
           of the given type, returns the size of the structure (whose type is
1453
          integer), arguments for this should be provided as kwargs
1454
    19 find_size(struct_name = str_name, objects = ObjSet)
1455
    20
1456
    21 # assigns the head to the location of the object
1457 22 assign_head(at_obj_loc)
    23
```

```
1458
    24 # moves the head in the given dir
1459 25 move_head(dir)
1460 26
1461 27 # keeps the object obj at the head
1462 28 keep_at_head(obj)
1463<sup>29</sup>
    30
1464 31 #Examples:
1465 32 #Instruction: Move the green block to the left of the red dice
1466 33 assign_head(at_obj_loc = filter(red, dice))
1467 34 move_head(left)
1468 <sup>35</sup> keep_at_head(obj = filter(green, cube))
    36
1469 37 # Instruction: Find the size of the tower made of yellow legos
1470 38 find_size(struct_name = tower, objects = filter(yellow, lego))
1471 39
1472 40 #Instruction: Find the size of the row made of orange cubes
1473 <sup>41</sup> find_size(struct_name = row, objects = filter(orange, cube))
    42
1474 \frac{42}{43} # Instruction: Find the size of the column made of cyan cubes
1475 44 find_size(struct_name = column, objects = filter(cyan, cube))
1476 45
1477 46 # Instruction: Move the green block to the left of the red dice and the
     yellow block to the top of the green block
1478 47 assign_head(at_obj_loc = filter(red, dice))
1479 48 move_head(left)
1480 49 keep_at_head(obj = filter(green, cube))
1481 50 assign_head(at_obj_loc = filter(green, cube))
1482 51 move_head(top)
1483 <sup>52</sup> keep_at_head(obj = filter(yellow, cube))
    53
1484 54 # Instruction: Construct a row of green legos of length 3 to the right of
the blue block
1486 55 assign_head(at_obj_loc = filter(blue, block))
1487 56 move_head(right)
1488 57 row(length = 3, objects = filter(green, legos))
    58
1489 59
1490 60 # Instruction: Construct a tower of size 3 using red cubes
1491 61 tower (height = 3, objects = filter (red, cube))
1492 62
1493 63 # Instruction: Construct a row of size 5 using blue legos
    64 row(length = 5, objects = filter(blue, lego))
1494 65
1495 66 # Instruction: Construct a column of size 6 using green die
1496 67 column(length = 6, objects = filter(green, dice))
1497 68
1498 69 # Instruction: Place 3 green blocks so that one block is to the right of
         the other
1499 70 green_blocks = filter(green, block)
1500 71 keep_at_head(green_blocks)
1501 72 move_head(right)
1502 73 keep_at_head(green_blocks)
1503 <sup>74</sup> move_head(right)
    75 keep_at_head(green_blocks)
1504 76
1505 77 # Instruction: Place 3 red legos on top of one another
1506 78 red_legos = filter(red, lego)
1507 79 keep_at_head(red_legos)
1508 <sup>80</sup> move_head(top)
    81 keep_at_head(red_legos)
1509 82 move_head(top)
1510 83 keep_at_head(red_legos)
1511 84
    85
```

```
1512
    86 # CURRENT TASK
1513 87 # Instruction: Construct tower of white cubes to the same height as
1514 88 # existing tower of green die
1515 89
1516 90 # GPT-4s output
1517 91 #First, we have to find the height of the tower of green dice,
    92 #then construct a tower of white cubes of the same size
1518 93
1519 94 tower_size = find_size(struct_name = tower, objects = filter(green, die))
1520 95 tower (height = tower_size , objects = filter (white, cube))
1521
      Listing 9: Prompting method for the task of constructing tower of white cubes to the same height as
1522
      existing tower of green die
1523
1524 | # Given a task you have to provide Python code for executing the task
1525 2
1526 3 # importing available functions
1527 <sup>4</sup> from spatial_directions import top, front, back, left, right
    5 from primitives import assign_head, move_head, keep_at_head
1528
    6 # HEAD is a imaginary pointer keeping track of the current spatial
1529
          location in consideration
1530 7 from helpers import find_size, filter
1531 8 from structures import row, column, tower
1532 9 #function signature of the imported functions
    10 # finds all the objects with the given color and shape, returns a mask
1533
          denoting the probability of object selection
1534 11 filter(color, shape)
1535 12 # finds the size of the structure struct_name that is formed with objects
           of the given type, returns the size of the structure (whose type is
1536
          integer), arguments for this should be provided as kwargs
1537
    13 find_size(struct_name = str_name, objects = ObjSet)
1538 _{14} # assigns the head to the location of the object
1539 15 assign_head(at_obj_loc)
1540 16 # moves the head in the given dir
1541 17 move_head(dir)
1542 <sup>18</sup> # keeps the object obj at the head
1543 20
    19 keep_at_head(obj)
1544 <sub>21</sub> #Examples:
1545 22 #Instruction: Move the green block to the left of the red dice
1546 23 assign_head(at_obj_loc = filter(red, dice))
1547 <sup>24</sup> move_head(left)
    25 keep_at_head(obj = filter(green, cube))
1548
    26
1549 27 # Instruction: Find the size of the tower made of yellow legos
1550 28 find_size(struct_name = tower, objects = filter(yellow, lego))
1551<sup>29</sup>
1552 30 #Instruction: Find the size of the row made of orange cubes
    31 find_size(struct_name = row, objects = filter(orange, cube))
1553
1554 33 # Instruction: Find the size of the column made of cyan cubes
1555 34 find_size(struct_name = column, objects = filter(cyan, cube))
1556 <sup>35</sup>
1557 36 # Instruction: Move the green block to the left of the red dice and the
          yellow block to the top of the green block
1558 37 assign_head(at_obj_loc = filter(red, dice))
1559 38 move_head(left)
1560 39 keep_at_head(obj = filter(green, cube))
1561 40 assign_head(at_obj_loc = filter(green, cube))
1562 <sup>41</sup> move_head(top)
    42 keep_at_head(obj = filter(yellow, cube))
1563 43
1564 44 # Instruction: Construct a row of green legos of length 3 to the right of
1565
         the blue block
    45 assign_head(at_obj_loc = filter(blue, block))
```

```
1566
     46 move_head(right)
1567 47 row(length = 3, objects = filter(green, legos))
1568 48
1569 49
1570 50 # Instruction: Construct a tower of size 3 using red cubes
1571 51 tower (height = 3, objects = filter (red, cube))
     52
1572 53 # Instruction: Construct a row of size 5 using blue legos
1573 54 row(length = 5, objects = filter(blue, lego))
1574 55
1575 56 # Instruction: Construct a column of size 6 using green die
1576 57 column(length = 6, objects = filter(green, dice))
     58
1577 59 # Instruction: Place 3 green blocks so that one block is to the right of
1578
           the other
1579 60 green_blocks = filter(green, block)
1580 61 keep_at_head(green_blocks)
1581 <sup>62</sup> move_head(right)
    63 keep_at_head(green_blocks)
1582 <sub>64</sub> move_head(right)
1583 65 keep_at_head(green_blocks)
1584 66
1585 67 # Instruction: Place 3 red legos on top of one another
1586 68 red_legos = filter(red, lego)
     69 keep_at_head(red_legos)
1587 <sub>70</sub> move_head(top)
1588 71 keep_at_head(red_legos)
1589 72 move_head(top)
1590 <sup>73</sup> keep_at_head(red_legos)
     74
1591
     75
1592 76 # CURRENT TASK
1593 _{77} # Instruction: Construct a tower of height 6 using red and blue blocks
1594
         that are alternating
1595 <sup>78</sup>
     79 # GPT-4s output
1596
     80 # Python code:
1597 81
1598 82 # Define the red and blue blocks
1599 83 red_blocks = filter('red', 'block')
1600 <sup>84</sup> blue_blocks = filter('blue', 'block')
    85
1601
     86 # Start at the bottom and alternate building the tower
1602 <sub>87</sub> for i in range(6):
1603 <sub>88</sub>
            if i%2 == 0: # if the stack position is even
1604 89
                keep_at_head(red_blocks)
1605 <sup>90</sup>
            else: # if the stack position is odd
                keep_at_head(blue_blocks)
    91
1606
            if i != 5: # if not at the top of the tower
     92
1607 <sub>93</sub>
               move_head(top)
1608
        Listing 10: Prompting method for the task of constructing tower of alternating red and blue cubes
1609
1610
1611
       To find the size of a given structure in the given scene we define the function find_size, which
1612
       takes the name of structure, all the objects in the initial scene, mask of the objects (a distribution over
1613
       the objects based on the attributes), and the initial state. (we assume that this function has access to
1614
       the semantics of all the concepts learnt so far, through a transition function). Algorithm 1 gives the
1615
       pseudocode for the function find_structure. Below we provide a brief explanation for it.
1616
```

1. First we assign our head to every block in the available blocks (line 6)

1619 2. Then we begin constructing/visualizing the corresponding structure from that block beginning with a size of 1. (line 7)

	3. For each structure created/visualized we compare the blocks moved for the structure creation
	with corresponding blocks originally present in the scene, and perform a matching between
	these blocks and a subset of the blocks originally present (line 11-26).
	4. If we are able to find a mapping for each moved block, such that each mapped pair has an
	IoU greater than a threshold, we increase the next potential size to test by 1 (line 26-27).
	5. The final size is the size corresponding to 2nd last iteration, before termination (line 28).
	6. We return the maximum of all the possible structures that are found (line 34)
Algo	rithm 1 Find Size
Requ	ire: name: structure name, objs: object list, state: initial state, mask: object mask
Ensu	re: Size of the maximum sized structure found
1: 1	$ound \leftarrow []$
2: f	or each cand in $objs$ do
3:	$size \leftarrow 1$
4:	$curr \leftarrow state.copy()$
5:	while true do
6:	$new_state \leftarrow assign_head(cand)$
7:	$vis_state, num_mov, rew \leftarrow transition($
8:	$new_state, name, [('size', size), ('objects', mask)])$
9:	$topk \leftarrow torch.topk(mask, num_mov)$
10:	$possible \leftarrow true$
11:	$matches \leftarrow []$
12:	for each $idx \sin topk$ do
13:	$possible \leftarrow (idx in objs)$
14:	match $ok \leftarrow false$
15:	for $m_i dx$, obi in enumerate(state.state) do
16:	$iou \leftarrow iou2d(vis state.state[idx])$
17.	state state $[m \ idr])$
18.	match $ok \leftarrow (iou > 0.75)$ and
10.	$(m \ idr in \ ohis)$
20.	if match ok then
20.	$m match_{ob} a annend(m idx)$
21:	heal
22:	ord if
23:	end for
24:	end for
25:	II not match_ok inen
26:	break
27:	end II
28:	end for
29:	if possible then
30:	$size \leftarrow size + 1$
31:	else
32:	found.append((size-1,matches))
33:	break
34:	end if
35:	end while
36: e	nd for
37: r	eturn max(size for size, match in found)
F :	
D.4	DETAILS ON MCTS VARIANTS FOR PLAN SEARCH
Here	we provide the details for 3 different plan search methods, that search over the space \mathcal{A} =
$\mathcal{A}_c \cup$	\mathcal{A}_p

- 1672
- 1673
- MCTS+L+P: This is the approach that we describe in section 4.2. For every concept say Tower ∈ L we have a corresponding set of macro action say Make_Tower (3,

1698

1700

1701 1702

1708

objects) (L). Further we use a neural pruner π_{neural} that outputs a primitive action a_n^* 1675 (given current state and next expected state). We only consider the actions $\mathcal{A}_c \cup \{a_n^*\}$ during 1676 our search from the given state. This helps to reduce the effective branching factor and 1677 allows to search longer length plans within the same computational budget. 1678 1679 • MCTS+L-P: Here we do not use the reactive policy, therefore the branching factor for every node becomes $\mathcal{A} = \mathcal{A}_c \cup \mathcal{A}_p$. 1683 • MCTS-L+P: We only search among the space of primitive actions i.e A_p . Given a state \tilde{s}_t and corresponding next expected state s_{t+1} we greedily pick the action $a_p^* = \pi_{\text{neural}}(\tilde{s}_t, s_{t+1})$. This method is much more faster than the previous 2 methods as there is no explicit search. 1687 However the corresponding policy is trained only to output an action $a \in \mathcal{A}_p$ and lacks the ability to output modular plans composed of macro actions such as Make_Tower (3, 1689 objects) $\in A_c$ (the action space A_c is increasing with time and the architecture of network needs to be changed accordingly). As a result the plans found are not modular and difficult to generalize. Further, the reactive policy is not trained to output reset_head(), store_head() as additional annotated data is required in order to train a classifier over 1693 them. This further decreases the space of grounded plans (and therefore corresponding generic programs) such a policy can represent. Training a reactive policy that can handle 1695 actions such as reset head () and an action space \mathcal{A}_c that grows with time is part of future work.

D.5 GOAL-CONDITIONED PLANNING WITH LEARNT CONCEPTS

Why is it difficult to hand encode a PDDL for our domain? Most of the PDDL description of blocks world assume actions involving only the spatial relation on Top, which limits their applicability to describing different structures like *row* that need spatial relations like onRight. Further a single action might lead to varied effects/post-conditions based on the initial state. 16 gives 2 example of the same action moveOnTop (A, B) which would end up giving adding different number of spatial relations.

1709 Approach Overview. Given an instruction Λ = "Con-1710 struct a staircase of magenta die having 3 steps", we first 1711 convert it into corresponding grounded task sketch H_S^* = 1712 staircase(3, $[3, 2, 1, \cdots]$). Executing the 1713 corresponding program of staircase (by getting the semantics from the library \mathcal{L}) on the desired objects we get the 1714 expected final scene S'_f in bounding box space. The ini-1715 tial scene S_i and expected final scene S'_f are converted 1716 into scene graph SG'_i and SG'_f (described in D.5). The 1717 1718 relations between the task relevant objects in SG'_{f} act as 1719 propositions/relations for goal check and the initial scene graph act as the initial state. Then a neuro-symbolic planner is used to obtain the optimal plan from the start state



Figure 16: **Difficult to encode postconditions**. Illustration of a domain where encoding a PDDL for direct planning is challenging.

to a state that satisfies the goal. Below we also detail different aspects of the approach.

1723Scene-graph Extraction.0 gives the algorithm used for generating scene graph from a given1724scene (set of bounding boxes). Suppose we need to check whether there exists a relation of the1725form (i, j, direction) i.e. block i is in the direction direction of block j, in a given scene. We1726first initialize the head at the position/bounding-box of block j (line 7). Then we move the head in1727direction (line 9). We claim that the relationship would exists if bounding box for block i1728has IoU > 0.75 with the predicted_head (line 10-12).

Algo	rithm 2 Get Scene Graph				
1: 1	procedure GET_SCENE_GRAPH(bboxes) // bboxes describe the corresponding scene for which				
s	scene graph is to be found				
2:	$spatial_relations \leftarrow []$				
3:	$directions \leftarrow {\text{'right', 'front', 'top'}}$				
4:	for $i, bbox$ in $enumerate(bboxes)$ do				
5:	for $j, other_bbox$ in $enumerate(bboxes)$ do				
6:	$final_head \leftarrow bboxes[i]$				
7:	$init_head \leftarrow bboxes[j]$				
8:	for direction in directions do				
9:	$pred_head \leftarrow move_head(direction, init_head)$				
10:	$iou_score \leftarrow IoU(final_neaa, prea_neaa)$				
11:	$n iou_score > 0.15$ then spatial relations $annend((i \ i \ direction))$				
12.	end if				
14:	end for				
15:	end for				
16:	end for				
17:	return spatial_relations				
18: e	end procedure				
Pre-	conditions. We define the following 2 preconditions (and learn their grounding) in order to				
ensu	re that the generated plans are physically possible.				
	1 is a place (h) where h = 1				
	direction dir. For this we simply move the head in the direction dir with respect to the block				
	blk. If the resulting position of head has 0 overlap with bounding boxes of all the other				
	objects the predicate is True otherwise False.				
	2. will-not-be-floating (pred loc): We need to check whether the resultant/pre-				
	dicted location of an object on taking an action would be dynamically stable or not. The				
	location would be stable if either it is on top of some already placed object or it is on the				
	table surface. The former can be checked through the resultant scene graph itself (that is				
	obtained by applying the algorithm 0), while for the later we train an <i>on-table</i> classifier.				
	This would take as input a bounding box and predict whether the box is on the table or				
	not. For training this classifier we use the dataset of pretraining phase. The blockwise				
	positive and negative sample annotation can be done automatically by giving GP1-4 the				
	aren't 17 gives an example. (Though we have not taken this approach for the complete				
	dataset of 5k samples due to high cost)				
	dataset of 5K samples due to mgn cost).				
Actio	ons. We define the following two types actions				
	1. place-random (blk) : To place a block at a random free position on the table. For this				
	we train a generative model (VAE) which learns the underlying distribution of bounding				
	boxes for all the blocks that lie on the table. Given a scene we would sample position				
	(bounding box) from this until we get a position that is not overlapping with the existing				
	DIOCKS in the scene. For training the VAE we assume that for every demonstration in the				
	preuanning data, the first scene has all the objects randomly placed on the table (We could have also used the positive examples used for training on table classifier)				
	nave also used the positive examples used for training off-table classifier). $2 \mod (rel = b k b k)$. This action corresponds to moving the blk in the di				
	rection rel of blk2 resulting in the addition of a relation (rel blk1 blk2) in the set				
	of spatial relations. This action is defined as a sequential composition of the actions				
	assign head(blk2), move head(rel), keep at head(blk1) (blk1 is a				
	one hot tensor for the corresponding block).				
	· · · · · · · · · · · · · · · · · · ·				
Tech	niques and heuristic for efficient planning. Since the action space for the planner could				
be o	(n^2) , where n is the number of objects we adopt the following techniques to make planning				
scala	ble/efficient:				

1782					
1783		Input Drompt			
1784		<u>Input Prompt</u>			
1785	Suppose the spatial relation (top, 1, 2) means 1 is on top of object 2 and in				
1786		this case object 2 is on table			
1787		Consider the following spatial relations:			
1788		(C, A, 'right')			
1789		(B, C, 'top')			
1790		(H, A, 'top')			
1791		(B, H, 'right')			
1792		(X, B, Top') A C			
1793					
1794		(Y, H, LOP)			
1705		identity which objects are on the horizontal table surface			
1706		CDT Ale Outeut			
1707		<u>GP1-4 S Oulpul</u>			
1700		The objects directly on the horizontal table surface are A and C.			
1790					
1799		D'an an 17. Mathed to and an an at the for the balance of the sheet for			
1001		Figure 17: Method to get annotation for training on-table classifier.			
1801					
1802		1 Heuristic - We define the heuristic value $h(s)$ for a state s as the number of relations that			
1803		are present in the goal but are absent in the scene graph corresponding to the state s. Even			
1804		though this heuristic is not admissible (as it may over overestimate the cost to goal), it was			
1805	found to work optimally in most of the cases				
1806		2. Greedy-pruning - We assume that all the actions resulting in states with higher or same			
1807	heuristic value would be of the form place-random (blk). This means among the				
1808		actions of the form move (rel, blk1, blk2) we only select those that lead to states			
1809	with decreased heuristic value.				
1810	3. Relevant-object-set - Suppose O is the set of objects that are part of atleast one of the				
1811	predicate in goal. We define O' as the transitive closure of O with respect to the relation				
1812		Related in the initial state s_i , where $SG(s_i)$ is scene graph for the initial state			
1813		$Related(a, b, s_i) \iff \exists dir((dir, a, b) \in SG(s_i) \lor (dir, b, a) \in SG(s_i)) $ (9)			
1814		1 = 1 = 1 = 1 = 1 = 1 = 1 = 1 = 1 = 1 =			
1815		We assume O' is the relevant set of object for completing the task and actions that move any			
1816		other object should not be taken.			
1817					
1818	Е	BROADER IMPACT			
1819					
1820	Thi	s work creates foundational knowledge in understanding human-like spatial abstractions. This			
1821	wor	k contributes towards the development of explainable and interpret-able learning architectures			
1822	that	t may eventually contribute towards the development of embodied agents collaborating with and			
1823	assi	sting humans in performing tasks. No negative impact of this work is envisioned.			
1824					
1825	F	HVDEDDADAMETEDS ADCUITECTUDE DETAILS AND GOOIND TOUTU			
1826	T	CONCEPTS			
1827		CUNCEPIS			
1828	г 1				
1829	F. I	ARCHITECTURE FOR NEURAL MODULES			
1830	Act	ion Simulator:			
1831	4 101				
1832 ¹	imp	port torch.nn as nn			
1833		ass ActionSimulatorNetwork(nn Module).			
1834 4	, CIC	def init (self, bbox mode, hidden size = 256):			
1835 5	5	<pre>super(ActionSimulatorNetwork, self)init()</pre>			

6 self.bbox_mode = bbox_mode

```
1836
                self.hidden_size = hidden_size
1837
     8
1838
                self.action_semantics_encoder = nn.Sequential(
    9
1839 10
                    nn.Linear(5, hidden_size),
1840 <sup>11</sup>
                    nn.ReLU(),
1841 <sup>12</sup>
                    nn.Linear(hidden_size, hidden_size),
                    nn.ReLU()
     13
1842
                )
     14
1843 15
                self.argument_encoder = nn.Sequential(
1844 16
                    nn.Linear(5, hidden_size),
1845 <sup>17</sup>
                    nn.ReLU(),
                    nn.Linear(hidden_size, hidden_size),
    18
1846
                    nn.ReLU()
     19
1847
                )
     20
1848 21
                self.decoder = nn.Sequential(
                    nn.Linear(hidden_size, hidden_size),
1849 22
                    nn.ReLU(),
1850 <sup>23</sup>
                    nn.Linear(hidden_size, 5),
    24
1851
     25
                    nn.Tanh()
1852
    26
1853
                            Listing 11: Action Simulator Network in PyTorch
1854
1855
       Reactive Policy(\pi_{neural}):
1856
     1 import torch.nn as nn
1857
     2
1858
     3 class NeuralSearch(nn.Module):
1859
           def __init__(self, action_space=6):
     4
1860 5
                super(NeuralSearch, self).__init__()
1861 6
               self.action_space = action_space
1862 7
               self.fc1 = nn.Linear(10, 256)
                # self.bn1 = nn.BatchNorm1d(256)
1863
     0
                self.bn1 = nn.Identity()
1864
                self.fc2 = nn.Linear(256, 256)
     10
1865
    11
                # self.bn2 = nn.BatchNorm1d(256)
1866
    12
                self.bn2 = nn.Identity()
                self.fc3 = nn.Linear(256, 256)
1867 13
                # self.bn3 = nn.BatchNorm1d(256)
1868 <sup>14</sup>
                self.bn3 = nn.Identity()
     15
1869
               self.fc4 = nn.Linear(256, action_space)
     16
1870
                                 Listing 12: Neural Search in PyTorch
1871
1872
       Random Position predictor (for grounding of place-random(blk)):
1873
1874
     1 import torch.nn as nn
     2
1875
     3 class VAE(nn.Module):
1876
           def __init__(self, input_dim, latent_dim):
1877
                super(VAE, self).__init__()
1878 6
               self.input_dim = input_dim
               self.latent_dim = latent_dim
1879 7
                # Encoder
1880 <sup>8</sup>
                self.fc1 = nn.Linear(input_dim, 512)
     9
1881
                self.bn1 = nn.BatchNorm1d(512)
     10
1882
                self.fc2 = nn.Linear(512, 512)
    11
1883 12
               self.bn2 = nn.BatchNorm1d(512)
               self.fc3 = nn.Linear(512, 512)
1884 13
1885 <sup>14</sup>
               self.bn3 = nn.BatchNorm1d(512)
1886<sup>15</sup>
               self.fc4 = nn.Linear(512, 512)
                self.bn4 = nn.BatchNorm1d(512)
     16
1887
    17
                self.fc51 = nn.Linear(512, latent_dim) # Mean of the latent
1888
           space
1889 18
                self.fc52 = nn.Linear(512, latent_dim) # Log-variance of the
           latent space (log-var for numerical stability)
```

1890 19 1891 # Decoder 20 **1892** 21 self.fc5 = nn.Linear(latent_dim, 512) **1893** 22 self.bn5 = nn.BatchNorm1d(512) **1894** ²³ self.fc6 = nn.Linear(512, 512)**1895**²⁴ self.bn6 = nn.BatchNorm1d(512) 25 self.fc7 = nn.Linear(512, 512)1896 self.bn7 = nn.BatchNorm1d(512) 26 **1897** 27 self.fc8 = nn.Linear(512, 512)**1898** 28 self.bn8 = nn.BatchNorm1d(512) **1899** ²⁹ self.fc9 = nn.Linear(512, input_dim) Listing 13: VAE in PyTorch 1900 1901 1902 **On-table classifier** (for grounding of will-not-be-floating (pred_loc): 1903 1 import torch.nn as nn 1904 1905 3 class TableClassifier(nn.Module): 1906 4 def __init__(self): super(TableClassifier, self).__init__() 1907 5 self.fc1 = nn.Linear(5, 16)6 1908 7 self.bn1 = nn.BatchNorm1d(16) 1909 self.fc2 = nn.Linear(16, 16)8 1910 self.bn2 = nn.BatchNorm1d(16) 0 **1911** 10 self.fc3 = nn.Linear(16, 16)self.bn3 = nn.BatchNorm1d(16) **1912** ¹¹ **1913** ¹² self.fc4 = nn.Linear(16, 1)13 self.bn4 = nn.BatchNorm1d(1) 1914 self.sigmoid = nn.Sigmoid() 14 1915 Listing 14: Table Classifier in PyTorch 1916 1917 1918 F.2 HYPERPARAMETERS USED IN EXPERIMENT 1919 1920 As indicated in A.3 for the purpose of generalization through multiple candidate plans (from 1 1921 demonstration) we chose the top-k plans (as measured by overall IoU achieved). The k chosen for all 1922 our experiments involving MCTS was 5. (The performance of our best approach was found to be the 1923 same for k=5 to 20). For every plan we obtain 3 programs from GPT-4 by re-prompting it 3 times with the same input prompt (with temperature > 0). From the pool of these 3*k programs we chose 1924 the one with highest IoU reward by running each of them on the given demonstration. The discount 1925 factor kept for our search is $\gamma = 0.95$, and unless explicitly specified the number of expansions steps 1926 used = 5000. 1927 1928 F.3 GROUND-TRUTH INDUCTIVE CONCEPTS 1929 1930 1 ###### 1931 2 **#** row 1932 3 def row(length, objects): 1933 4 for i in range(length): keep_at_head(obj = objects) 1934 5 move_head(dir = "right") 6 1935 7 1936 8 ###### 1937 9 # tower 1938 10 def tower(height, objects): **1939** ¹¹ for i in range(height): **1940**¹² keep_at_head(obj = objects) move_head(dir = 'top') 13 1941 14 **1942** 15 **##### 1943** 16 # column 17 def column(size, objects):

```
1944
        for _ in range(size):
    18
1945 19
              keep_at_head(obj = objects)
1946 20
               move_head(dir = 'front')
1947 21
1948 <sup>22</sup> ######
1949 <sup>23</sup> # staircase
    24 def staircase(steps, objects):
1950 25 for step in range (1, steps+1):
1951 26
           tower(height = step, objects = objects)
1952 27
               move_head(dir = 'right')
1953 <sup>28</sup>
1954 <sup>29</sup> #####
    30 # inverted_row
1955 31 def inverted_row(num, objects):
1956 <sub>32</sub> for i in range(num):
           keep_at_head(obj=objects)
1957 33
1958 <sup>34</sup>
               move_head(dir='left')
1959 35
    36 ######
1960 37 # inverted_column
1961 38 def inverted_column(size, objects):
1962 39 for _ in range(size):
           keep_at_head(obj = objects)
1963 <sup>40</sup>
1964 <sup>41</sup>
              move_head(dir = 'back')
          return None
    42
1965 43
1966 44 #####
1967 45 # inverted_staircase
1968 46 def inverted_staircase(steps, objects):
1969 47 for step in range(1, steps+1):
            tower(height = step, objects = objects)
    48
1970 49
              move_head(dir = "left")
1971 50
1972 51 # # # # # #
1973 <sup>52</sup> # diagonal_135
1974 <sup>53</sup> def diagonal_135(length, objects):
        for i in range(length):
    54
1975 55
           keep_at_head(obj = objects)
1976 56
              move_head(dir = 'front')
              move_head(dir = 'left')
1977 57
1978 58
          return
1979 <sup>59</sup>
    60 ######
1980 61 # diagonal_315
1981 62 def diagonal_315(length, objects):
        for i in range(length):
1982 63
            keep_at_head(obj = objects)
1983 <sup>64</sup>
              move_head(dir = 'back')
1984 <sup>65</sup>
            move_head(dir = 'right')
    66
1985 67
          return
1986 68
1987 69 #####
1988 70 # diagonal_225
1989 <sup>71</sup> def diagonal_225(length, objects):
        for _ in range(length):
    72
1990 73
            keep_at_head(obj = objects)
1991 74
              move head(dir = 'back')
              move_head(dir = 'left')
1992 75
1993 <sup>76</sup>
1994 77 #####
    78 # diagonal_45
1995 79 def diagonal_45(length, objects):
        for _ in range(length):
1996 80
1997 81
            keep_at_head(obj = objects)
             move_head(dir = 'front')
    82
```

```
1998
                 move_head(dir = 'right')
     83
1999
     84
2000 85 ######
2001 86 # boundary
2002 87 def boundary(size, objects):
2003 <sup>88</sup>
          row(length=size-1, objects=objects)
           for _ in range(size-1):
2004 90
             move_head(dir = 'right')
2005 91
            move_head(dir = 'front')
2006 92
2007 <sup>93</sup>
            column(length=size-1, objects=objects)
2008 <sup>94</sup>
            for _ in range(size-1):
                move_head(dir = 'front')
     95
2009 96
            move_head(dir = 'left')
2010 97
            inverted_row(length=size-1, objects=objects)
2011 98
            for _ in range(size-1):
2012 99
2013<sup>100</sup><sub>101</sub>
                 move_head(dir = 'left')
            move_head(dir = 'back')
2014 102
2015 103
            inverted_column(length=size-1, objects=objects)
2016 104
           for _ in range(size-1):
2017 <sup>105</sup>
               move_head(dir = 'back')
2018<sup>106</sup>
            move_head(dir = 'right')
2019 108 ######
2020 109 # arch_bridge
2021 110 def arch_bridge(height, objects):
2022<sup>111</sup>
            staircase(steps = height, objects = objects)
            move_head(dir = 'left')
2023<sup>112</sup><sub>113</sub>
            inverted_staircase(steps = height, objects = objects)
2024 114
           return
2025 115
2026 116 ######
2027 <sup>117</sup> # x-shaped structure
2028 118 def x(size, objects):
119 diagonal_45(lengt
            diagonal_45(length = size, objects = objects)
2029 120
            move_head(dir = 'back')
2030 121
          diagonal_315(length = size, objects = objects)
         move_head(dir = 'left')
2031 122
2032 <sup>123</sup>
         diagonal_225(length = size, objects = objects)
2033<sup>124</sup>
          move_head(dir = 'front')
            diagonal_135(length = size, objects = objects)
     125
2034 126
2035 127 # # # # # # #
2036 128 # pyramid
2037 129 def pyramid(height, objects):
         for i in range(height):
2038<sup>130</sup><sub>131</sub>
               row_length = (height * 2) - (i * 2) - 1
2039 132
                row(length = row_length, objects = objects)
2040 133
                if i != height - 1:
2041 134
                     move_head(dir = 'top')
2042 <sup>135</sup>
                      move_head(dir = 'right')
     136 ######
2043
                                Listing 15: Definition of inductive concepts
2044
2045
2046
           COMPUTATIONAL REQUIREMENTS: DETAILS
        G
2047
2048
        All our experiments were run on a server with the following machine specifications.
2049
2050
        CPU Specification:
2051
```

2052		
2053	Specification	Value
2054	Architecture	x86_64
2055	CPU op-mode(s)	32-bit, 64-bit
2056	Address sizes	46 bits physical, 57 bits virtual
2057	Byte Order	Little Endian
2058	CPU(s)	112
2059	On-line CPU(s) list	0-111
2060	Vendor ID	GenuineIntel
2061	Model name	Intel(R) Xeon(R) Gold 6330 CPU @ 2.00GHz
2062	CPU family	6
2063	Model	106
2064	Thread(s) per core	2
2004	Core(s) per socket	28
2005	Socket(s)	2
2066	Stepping	6
2067	CPU max MHz	3100.0000
2068	CPU min MHz	800.0000
2069	BogoMIPS	4000.00
2070		

GPU Specification:

2071

2073	Specification	Value		
2074	Specification			
2075	Description	GPU I		
2076	Description	VGA compatible controller		
2077	Vandar	Metroy Electronics Systems Ltd		
2078	Dhusiaal ID	Matrox Electronics Systems Ltd.		
2079	Physical ID Dug Info			
2080	Logical Nama	pc1@0000.05.00.0		
2081	Logical Maine			
2001	Version			
2002	Width	32 DITS		
2083	Clock	00MHZ		
2084	Capabilities	pm vga_controller bus_master cap_list rom fb		
2085	Configuration	depth=32 driver=mgag200 mingnt=16		
2086	Resources	1rq:16 memory:91000000-91ffffff memory:92808000-9280bfff		
2087		memory:92000000-92/fffff memory:c0000-dffff		
2088	D. I.I.	GPU 2		
2089	Description	3D controller		
2090	Product	GA102GL [A40]		
2091	Vendor	NVIDIA Corporation		
2092	Physical ID	0		
2093	Bus Info	pci@0000:17:00.0		
2094	Version	al		
2005	Width	64 bits		
2005	Clock	33MHz		
2090	Capabilities	pm bus_master cap_list		
2097	Configuration	driver=nvidia latency=0		
2098	Resources	iomemory:21000-20fff iomemory:21200-211ff irq:18		
2099		memory:9c000000-9cffffff memory:210000000000-210fffffffff		
2100		memory:21200000000-212001ffffff memory:9d000000-9d7fffff		
2101		memory:21100000000-211fffffffff memory:21200200000-212041ffffff		
2102		GPU 3		
2103	Description	3D controller		
2104	Product	GA102GL [A40]		
2105	Vendor	NVIDIA Corporation		

Physical ID Bus Info pci@0000:ca:00.0 Version a1 Width 64 bits Clock 33MHz Capabilities pm bus_master cap_list Configuration driver=nvidia latency=0 Resources iomemory:28000-27fff iomemory:28200-281ff irq:18 memory:e7000000-e7ffffff memory:28000000000-280fffffffff memory:28200000000-282001ffffff memory:e8000000-e87fffff memory:28100000000-281ffffffff memory:28200200000-282041ffffff

Time Required: The time required for pretraining phase of all the neural modules is around 36 hours. For learning of inductive concepts the time taken varies from 5 minutes to 1 day depending on the search method used and the specific set of hyperparameters. However for our best approach we get the maximum performance in approx 12 minutes. Time taken for our approach during inference is less than 2 minutes per dataset.