# Learning Neuro-symbolic Programs for Language Guided Robot Manipulation

Namasivayam K[*1], Himanshu Singh[*1], Vishal Bindal[*1], Arnav Tuli[1], Vishwajeet Agrawal[#2], Rahul Jain[#2],
Parag Singla[1] and Rohan Paul[1]

[1]Affiliated with IIT Delhi. [2]Work done when at IIT Delhi. [*] and [#] denote equal contribution.

*Abstract*— Given a natural language instruction and an input scene, our goal is to train a model to output a *manipulation program* that can be executed by the robot. Prior approaches for this task possess one of the following limitations: (i) rely on hand-coded symbols for concepts limiting generalization beyond those seen during training [1] (ii) infer action sequences from instructions but require dense sub-goal supervision [2] or (iii) lack semantics required for deeper object-centric reasoning inherent in interpreting complex instructions [3]. In contrast, our approach can handle linguistic as well as perceptual variations, end-to-end trainable and requires no intermediate supervision. The proposed model uses symbolic reasoning constructs that operate on a latent neural object-centric representation, allowing for deeper reasoning over the input scene. Central to our approach is a modular structure consisting of a *hierarchical instruction parser* and an *action simulator* to learn disentangled action representations. Our experiments on a simulated environment with a 7-DOF manipulator, consisting of instructions with varying number of steps and scenes with different number of objects, demonstrate that our model is robust to such variations and significantly outperforms baselines, particularly in the generalization settings. The code, dataset and experiment videos are available at `https://nsrmp.github.io`

## I. INTRODUCTION

As robots enter human-centric environments, they must learn to act and achieve the intended goal based on natural language instructions from a human partner. We address the problem of learning to translate high level language instructions into executable programs grounded in the robot's state and action space. We focus on multi-step manipulation tasks that involve object interactions such as stacking and assembling objects referred to by their attributes and spatial relations. Figure 1 provides an example where the robot is instructed to *"put the block which is behind the green dice to the right of the red cube"*. We assume the presence of natural supervision from a human teacher in the form of input and output scenes, along with linguistic description for a high-level manipulation task. Our goal is to learn representations for visual and action concepts that can be composed to achieve the task. The learning problem is hard since (1) object attributes and actions have to be parsed from the underlying sentence (2) object references need to be grounded given the image and (3) the effect of executing the specified actions has to be deciphered in the image space, requiring complex natural language as well as image level reasoning. Further, the model needs to be trained end-to-end, learning representations for concepts via distant supervision.

Prior efforts can be broadly categorized as (i) Traditional methods which learn a mapping between phrases in the natural language to symbols representing robot state and actions in a pre-annotated dataset [1], [4]–[9]; they lack the flexibility to learn the semantics of concepts and actions on their own, an important aspect required for generalizability (ii) Approaches that model an instruction as a sequence of action labels to be executed, without any deeper semantics, and requiring intermediate supervision for sub-goals, which may not be always be available [2], [10]–[16] (iii) Recent end-to-end learning approaches [3], [17]–[23]; they have limited reasoning capability both at the level of instruction parsing, and their ability to learn varied action semantics.

In response, we introduce a neuro-symbolic approach for jointly learning representations for concepts that can be composed to form *manipulation programs*: symbolic programs that explain how the world scene is likely to affected by the input instruction. Our approach makes use of a Domain Specific Language (DSL) which specifies various concepts whose semantics are learned by the model. We build on the concept learning framework by [24] for grounding of concepts in a natural language instruction in the image. The *manipulation program* is grounded into the scene to get a sequence of actions that are fed to the low-level motion planner of the robot, which is assumed to be given.

The contributions of this work are: (i) A novel neuro-symbolic model that learns to perform complex object manipulation tasks requiring reasoning over scenes for a natural language instruction, given initial and final world scenes. (ii) A demonstration of how dense representations for robot manipulation actions can be acquired using only the initial and final world states (scenes) as supervision without the need for any intermediate supervision. (iii) An RL-based learning approach that facilitates parsing of the complex instruction into symbolic concepts, and allows for learning of disentangled actions in the robot's space. (iv) Evaluation in instruction following experiments with a simulated 7-DOF robot manipulator, demonstrating robust generalization to novel settings improving on the state-of-the-art.

## II. RELATED WORKS

**Grounding Instructions to Robot Control.** Traditional approaches for grounding instructions assume a symbolic description of the environment state and robot actions. These approaches learn to associate phrases in an instruction with symbols conveying their meaning from annotated datasets. For example, [1], [4], [5] map language to discrete motion constraints that can be provided to a motion planner for
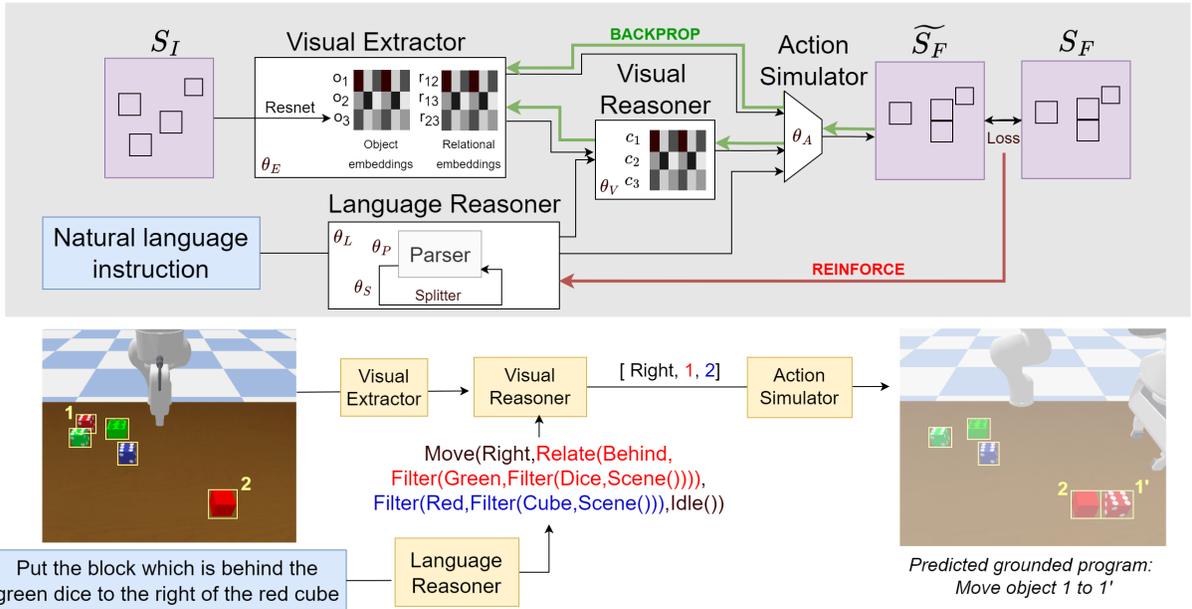
**Fig. 1: Model architecture.** The *Visual Extractor* forms dense object representations from the scene image using pre-trained object detector and feature extractor. The *Language Reasoner* auto-regressively induces a symbolic program from the instruction that represents rich symbolic reasoning over spatial and action constructs inherent in the instruction. The *Visual Reasoner* determines which objects are affected by actions in the plan using symbolic and spatial reasoning. The *Action Simulator* predicts final location of the moved object. The model is trained end-to-end with a loss on the bounding boxes, backpropagated to action and visual modules. REINFORCE is used to train [25] the language reasoner from which symbolic programs are sampled.

trajectory generation, [6], [7] parse language to a logical parse that includes robot control actions and [8], [9] use sequence to sequence modeling to ground language to LTL formulae capturing high-level task specifications. In contrast to such approaches, the proposed work doesn't assume a symbolic model of the world and learns spatial and visual concepts and action semantics from data.

**Inferring Plans from Instructions.** Another set of approaches focus on multi-stage tasks (e.g., cooking, assembly etc.) and propose learners that can infer action sequences by observing human task demonstrations. The work in [2] proposes a recurrent neural architecture that converts a natural language instruction to a sequence of intermediate goals for robot execution. Efforts such as [10]–[12] learn LTL task specifications from humans demonstrations and introduce a model for effectively searching in the space of programs. In another effort, [13], authors introduce a cognitively-inspired approach that infers a likely program from a symbolic generative grammar for a robot manipulator to form and re-arrange block patterns in a table top setting. Misra et al. [16] present a similar approach but additionally assess plan feasibility for inferred plan candidates using a symbolic planner in the training loop. Despite successful demonstrations, these approaches require dense intermediate supervision and treat robot actions in a purely symbolic manner without learning their deep grounded semantics. However, our method learns a latent space of composable symbolic programs which can be executed in the latent space of object representations. The composability of the symbolic programs enables incremental learning through a specified curriculum from only the initial and final state data, forgoing the need for dense intermediate supervision.

**Learning Action Models.** This work builds on the rich literature on acquiring action models for planning and captures *when* actions are applicable and *how* they affect the world. Efforts such as [17] and [18] learn initiation sets for actions implicitly learning to classify the robot's configuration space where an action can be initiated. [20]–[22] present neuro-symbolic approaches for learning a latent object-centric world model for tasks such as stacking, re-arrangement etc. In [3], authors introduce an end-to-end model that jointly predicts object affordances and object displacements specifying metric goals for robot re-arrangement and sorting tasks. Recent parallel work by Wang et al. [26] builds on [3] and represents the task as an executable program parsed from the instruction using a CCG parser, restricting the length and variation of the instructions it can handle. Even though these approaches successfully learn modular action models that are amenable to planning, the scope of reasoning is limited to one-hop reasoning over spatial relations and attributes. In contrast, this work focuses on interpreting task instructions that require compositional/hierarchical spatial reasoning over an extended planning horizon.

## III. PROBLEM FORMULATION

The robot perceives the world state comprising a set of rigid objects placed on a table via a depth sensor that outputs a depth image $S \in \mathbb{R}^{H \times W \times C}$, where $H, W, C$ respectively denote the height, width and the number of channels (including depth) of the imaging sensor. The workspace is co-habited by a human partner who provides language instructions to the robot to perform assembly tasks. We assume a closed world setting, i.e. changes to the world state are caused only by the robot manipulator.

The robot's goal is to interpret the human's instruction $\Lambda$ in the context of the initial world state $S_I$ and determine a sequence of low-level motions that result in the final world state $S_F$ conforming to the human's intention. Following [27], [28], planning for a complex task is factorized into (i) high-level task planning to determine a sequence of sub-goals, and (ii) the generation of low-level motions to attain each sub-goal. Formally, a semantic model for a manipulation task denoted as $ManipulationProgram(.)$ takes the initial scene $S_I$ and the instruction $\Lambda$ as input and determines a sequence of sub-goals as $(g_0, g_1, \ldots, g_n) = ManipulationProgram(S_I, \Lambda)$. Each sub-goal $g_i$ aggregates the knowledge of the object, $o_i$, to be manipulated and its target Cartesian $SE(3)$ pose $p_i$. This is provided to the low-level motion planner to synthesize the end-effector trajectory for the robot to execute. The robot's motion planner, which is assumed to be given, includes grasping an object and synthesizing a collision-free trajectory for positioning it at a target pose.

Training data consists of demonstrations, corresponding to task instructions $\Lambda$, in the form of initial and final world states $(S_I, S_F)$. Given data $D=\{S_I^i, S_F^i, \Lambda^i\}_{i=1}^M$, the model parameters are trained by optimizing a loss $\mathcal{L}(\widetilde{S}_F^i, S_F^i)$ where $\widetilde{S}_F^i = Simulate(ManipulationProgram(S_I^i, \Lambda; \Theta))$ is the final state estimated by simulating the plan inferred by $ManipulationProgram(S_I^i, \Lambda)$ on initial state $S_I^i$. We seek strong generalization on novel scenes, instructions and plan lengths beyond those encountered during training, along with interpretability in sub-goals.

## IV. TECHNICAL APPROACH

We propose a neuro-symbolic architecture to solve the task planning problem described in section III. Our architecture is inspired by [24] and is trained end-to-end with no intermediate supervision. We assume that the reasoning required to infer the sub-goals can be represented as a program determined by a domain specific language (DSL). Table I lists the keywords and operators in our DSL, along with the implementation details of the operators. We assume a lexer that identifies all the keywords that are referred to in the instruction $\Lambda$. We do not assume prior knowledge of the semantics of the DSL constructs, and they are learned purely from the data. Our architecture (ref. Figure 1) consists of the following key modules (a) Language Reasoner (LR) to parse the instruction into a hierarchical plan consisting of references to visual and action concepts and operators (b) Visual Extractor (VE) to obtain the initial bounding boxes from the input image (c) Visual Reasoner (VR) to ground the object related concepts present in the parsed instruction with respect to the input image (d) Action Simulator (AS) to learn the semantics of actions to be executed to produce a sequence of sub-goals. Figure 2 illustrates our approach.

### A. Language Reasoner (LR)

The language reasoner (LR) deduces a hierarchical symbolic program that corresponds to the manipulation task implied by the human's utterance to the robot. The deduced

| Keywords and their classes | | |
|---|---|---|
| **Object-level concepts** | | **Other concepts** |
| **Color**: {Red, Blue, Cyan,...} | | **RelCpt**: {Left, Behind, Front,...} |
| **Type**: {Cube, Lego, Dice} | | **ActCpt**: {MovRight, MovTop,...} |
| **Operators: ( Input $\rightarrow$ Output)** | | |
| Scene : None $\rightarrow$ ObjSet | | Unique: ObjSet $\rightarrow$ Obj |
| Filter : (ObjSet, ObjCpt)$\rightarrow$ ObjSet | | Relate : (Obj, RelCpt) $\rightarrow$ ObjSet |
| Move : (ActCpt, World) $\rightarrow$ World | | Idle : World $\rightarrow$ World |
| ObjectSet $\in \mathbb{R}^N$, N = Num objects, Object = one-hot ObjectSet | | |
| World = $\{(b_i, d)\}_{i=1}^N \in \mathbb{R}^5$, bounding boxes and depth for all objects | | |

TABLE I: Domain Specific Language.

program consists of symbolic reasoning constructs that operate on neural concepts grounded in the state space of the scene and the action space of the robot. Since a high-level task instruction may imply a sequence of actions, we adopt a hierarchical model where an LSTM [29] network first splits the instruction into a sequence of sub-instructions, each corresponding to one grounded action. The semantic parse of each sub-instruction is then inferred using a hierarchical *seq2tree* architecture similar to [24], [30]. The sequence of programs thus generated are composed to produce the symbolic program for the original language instruction.

### B. Visual Extractor (VE)

The visual extractor forms an object-centric view of the world by extracting object proposals in the form of bounding boxes for the objects (and class labels) using a pre-trained object detector [31]. Following [24], dense object representations are obtained by passing the bounding boxes (after non-maximal suppression) through a feature extractor as [32]. The data association between object proposals is estimated greedily based on cosine similarity between the dense object features in the initial/final scenes.

### C. Visual Reasoner (VR)

The visual reasoner performs visuo-spatial, object-centric reasoning to compute a sequence of sub-goals grounded in the scene from the symbolic program. E.g., the instruction *"put the block which is behind the green dice"* requires the robot to perform reasoning to resolve the specific world object that conforms to the relations/attributes as being behind a dice with colour attribute of green.

The visual reasoner parameterizes object-level and relational concepts in the DSL (Table I) with neural embeddings. Similar to [24], a differentiable framework is defined for the execution of operators such as *Filter, Unique, Relate, Idle*. As in [24], intermediate results of execution are represented in a probabilistic manner. For e.g., a set of objects is represented as a vector of probabilities where the $i$-th element represents the probability of the $i$-th object being an element of the set. *Filter, Unique, Relate* etc. sequentially modify this vector to ultimately yield the referred object/set of objects. This execution framework is used to execute the symbolic program on top of the visual features (extracted by the visual extractor) resulting in the grounding of the program, $\Pi$, on the world state.
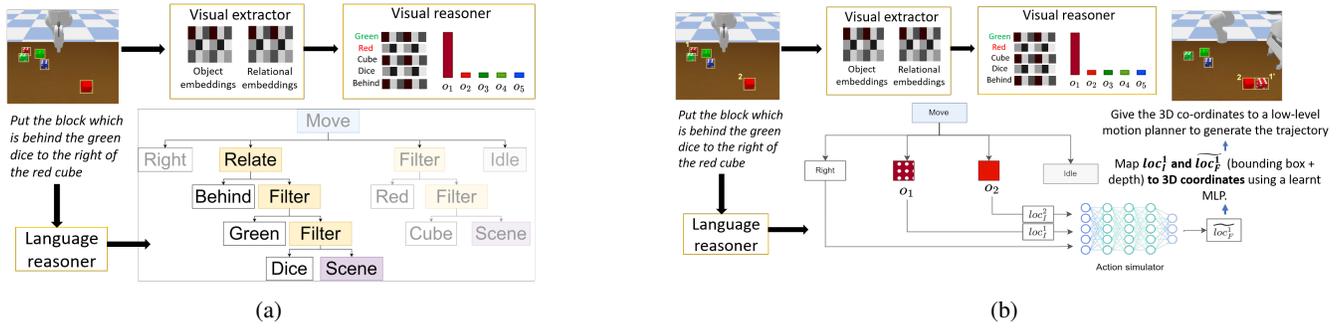
Fig. 2: **Quasi-symbolic program execution.** We postulate a latent space of hierarchical, symbolic programs that performs explicit reasoning over action, spatial and visual concepts. a) The language reasoner infers a program belonging to this space, capturing the semantics of the language instruction. The program is executed over the latent object representations extracted from the initial scene to get a grounded program. b) The grounded program is used by the action simulator to compute the final location of the moved object. This is fed into a low-level motion planner for computing the low-level trajectory.

### D. Action Simulator (AS)

The action simulator learns the semantics of action concepts of the DSL. It takes as input, the action concept, the initial locations of the object to be manipulated and the reference object respectively and outputs the target location of the former. The outputted target location of the object being manipulated serves as a sub-goal for the low-level motion planner. We determine the location of an object by the corners $b = (x_1, y_1, x_2, y_2)$ of the enclosing bounding box and the depth of the object from the camera face. Overall, the model can be summarized as follows:

- $\texttt{P} \leftarrow \texttt{LR}(\Lambda; \theta_P, \theta_S)$, where $\texttt{P}$ is a symbolic program composed of DSL constructs.
- $\Pi \leftarrow \texttt{VR}(\texttt{P}, \texttt{VE}(S_I; \theta_E); \theta_V)$. The visual reasoner then grounds $\texttt{P}$ and outputs $\Pi$, the grounded program. For example, in Figure 1, red and blue subprograms are grounded to object 1 and 2 respectively.
- $G \leftarrow \texttt{AS}(\Pi; \theta_A)$. The action simulator then takes $\Pi$ and returns a sequence of sub-goals, $G = (g_0, g_1, ..., g_{n-1})$, for the motion planner.

Here, $\theta$'s are the parameters of the corresponding module.

### E. Loss Function and Model Training

Given a single-step instruction $\Lambda$, the parser predicts a symbolic program $\texttt{P}$. The visual reasoner grounds $\texttt{P}$ to predict a sequence of sub-goals. The action simulator computes the low-level action corresponding to each sub-goal. The sequence of low-level actions is executed on the initial scene $S_I$ to get the predicted final locations of the objects $\{\widetilde{loc}_F^i\}_{i=1}^N$. Let $\{loc_F^i\}_{i=1}^N$ be the true locations in the gold final state, $S_F$. As mentioned above $loc = (b, d)$, where $b$ is the corners of bounding box and $d$ is the depth. The loss function $L_{act} := \sum_i^N \|\widetilde{loc}_F^i - loc_F^i\| + \beta(1 - \text{IoU}(\widetilde{b}_F^i, b_F^i))$ is used to train the action simulator and the visual modules. Since there is no explicit supervision to the parser, we train the parser using the policy gradient algorithm REINFORCE with the reward set to $-L_{act}$. During initial training, an explicit expectation (subtracting the mean action loss as the baseline) is computed over all programs to inform the loss, ameliorating the variance issue arising in REINFORCE. The

Language Reasoner is trained using REINFORCE and the other modules using backpropagation.

The following training curriculum is adopted: (i) training on single step commands with reasoning involving individual object features only, (ii) the trained action simulator module is frozen and additional training on single step commands involving reasoning over spatial relations between objects is carried out. (iii) the sentence splitter in the Language Reasoner module is fine-tuned on multiple-step instructions.

### F. Scene Reconstruction

We additionally train a neural model to synthesize the scene corresponding to each sub-goal, given only the initial scene and predicted object locations. This enables us to visualise the scene modification without the need of execution by a robot manipulator along with providing interpretability to the model's latent program space. The reconstruction architecture is adapted from [33], where the scene graph is constructed with nodes having object features and bounding boxes. This graph is updated with predicted bounding boxes at each step, yielding generated scenes. Presence of initial and final scenes in our data means we can train in a supervised manner, unlike [33].

## V. EVALUATION AND RESULTS

### A. Experimental Setup

**Data collection.** The dataset is collected in a PyBullet tabletop environment using a simulated Franka Emika Panda robot arm. The workspace consists of a tabletop with blocks of different shapes and colors. Each datapoint consists of an initial scene paired with a language instruction and the expected resulting final scene. For training, close to 5000 synthetic scenes and language instructions(with a 80:20 train-test split) are sampled with 3-5 blocks of varying colors and shapes and placed at randomized orientations on the table.

The dataset consists of both single-step(e.g. *"Put the red block on the green block"*) and double step commands(e.g. *"Put the red block which is behind the white dice to right of the blue cube, then put the lego block on top of red block"*). On the basis of the level of reasoning involved in the given instruction, we can classify the dataset into (i) *simple* that involve reasoning over individual object features only and

TABLE II: Comparison between Proposed Model and NMN+ (BB: Bounding Boxes)

| Model | Overall | | | Single-step | | Double step | | Simple | | Complex | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | IOU | IOU-M | Program (Action/Subj/Pred) | IOU | IOU-M | IOU | IOU-M | IOU | IOU-M | IOU | IOU-M |
| NMN+ (gold BB) | 0.77 | 0.55 | –/0.81/0.76 | 0.80 | 0.56 | 0.71 | 0.52 | 0.90 | 0.71 | 0.64 | 0.31 |
| Ours (gold BB) | 0.87 | 0.72 | 0.99/0.99/0.94 | 0.91 | 0.64 | 0.87 | 0.62 | 0.92 | 0.73 | 0.87 | 0.64 |
| NMN+ (extracted BB) | 0.69 | 0.32 | –/0.81/0.55 | 0.78 | 0.49 | 0.58 | 0.22 | 0.79 | 0.55 | 0.60 | 0.20 |
| Ours (extracted BB) | 0.76 | 0.41 | 0.99/0.74/0.76 | 0.80 | 0.43 | 0.69 | 0.54 | 0.83 | 0.62 | 0.69 | 0.24 |

TABLE III: Comparison b/w Proposed Model and CLIPort

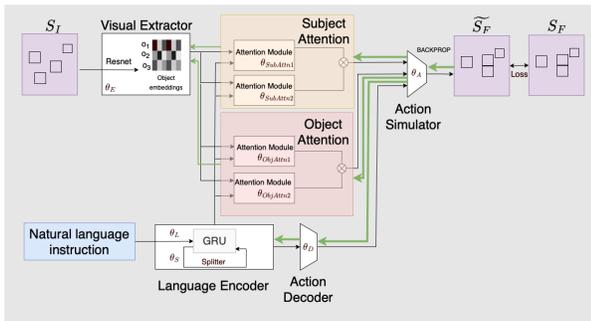| Model | Overall | | Simple | | Complex | |
|---|---|---|---|---|---|---|
| | Id. Acc. | Pl. Acc. | Id. Acc. | Pl. Acc. | Id. Acc. | Pl. Acc. |
| CLIPort | 0.46 | 0.39 | 0.55 | 0.47 | 0.32 | 0.26 |
| CLIPort(5x) | 0.76 | 0.61 | 0.90 | 0.74 | 0.55 | 0.39 |
| Ours (gold BB) | 0.94 | 0.93 | 1.0 | 1.0 | 0.83 | 0.81 |



Fig. 3: **Object-centric baseline NMN+.** For a fair comparison, our model and NMN+ share the same language encoder (with LSTM-based splitter) and the visual extractor. Attention blocks compute language-guided attention over object embeddings to get the *subject* and *predicate* for the manipulation action. This is fed into the action simulator along with the action embedding from the action decoder to get the predicted final location of the object.

*complex* that additionally involve reasoning over inter-object relationships.

### B. Baselines, Metrics and Comparisons

The proposed model (i) performs visuo-linguistic reasoning in a symbolic latent space and (ii) assumes an object-centric state representation. We provide comparisons with baselines that forego these two assumptions.

**NMN+**, inspired from Neural Module Networks [34], assumes an object-centric state representation. However, in contrast to the proposed model, visuo-linguistic reasoning required for manipulation is performed using language-guided attention over object embeddings instead of using symbolic reasoning constructs. Figure 3 provides the detailed architecture. The entire architecture is neural and is trained end-to-end with the loss same as that of the proposed model.

The following metrics are used for comparing the proposed model with *NMN+*: (i) *Intersection over Union (IOU)* of the predicted and groundtruth bounding box is calculated in the 2D image space assuming a static camera viewing the scene. Average IOU over all objects in the scene and mean IOU for objects moved during execution, termed IOU-M, are reported. (ii) *Program Accuracy*: The grounded program inferred for an (instruction, scene)-pair using the proposed model is compared with the manually annotated ground truth program. We separately report the grounding accuracy for the subject and predicate of our action (assumed binary) and the accuracy of the predicted action inferred from the instruction. Since, there is no explicit notion of grounded actions in the baseline, we do not report this metric for the baseline.



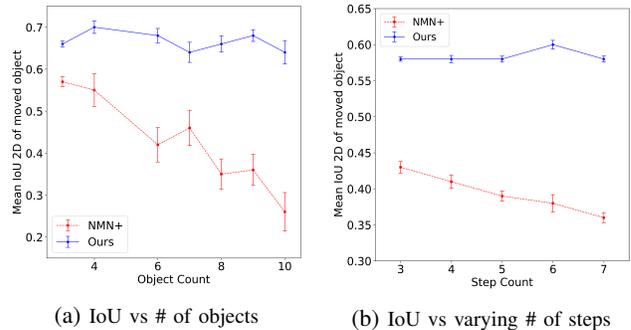(a) IoU vs # of objects     (b) IoU vs varying # of steps

Fig. 4: Performance in generalization settings

Table II reports the performance of our model and the baseline on the test set. We report both numbers: using gold set of bounding boxes and using bounding boxes extracted [1] using the approach described in Section IV-B. Our model outperforms the baseline overall. For instructions with complex reasoning (resolution of binary spatial relations) involved, the proposed model outperforms the baseline by 33 points in the IOU-M metric. Disentangled representations of relations and concepts in *Visual Reasoning* module allow the proposed model to reason over complex instructions.

**CLIPort** [3] takes as input a dense, image-based state representation, in contrast to our object-centric state representation. It proposes a two-stream architecture that computes language-guided attention masks on the input scene for both pick and place locations. It lacks explicit symbolic constructs for reasoning about the scene. Instead, the reasoning required for manipulation is performed by leveraging the visuo-linguistic understanding of the large pre-trained model CLIP [35]. Table III compares *CLIPort* with our model. Since *CLIPort* performs only single-step pick-and-place operations, the numbers are calculated only on single-step examples. The following metrics are used (i) *Placement Accuracy*: fraction of examples in which the predicted place location lies within the groundtruth bounding box of the moved object (ii) *Identification Accuracy*: fraction of examples in which the predicted pick location lies within the groundtruth bounding box of the moved object. For the proposed model, the centre of the predicted bounding box location is used as the predicted place location.

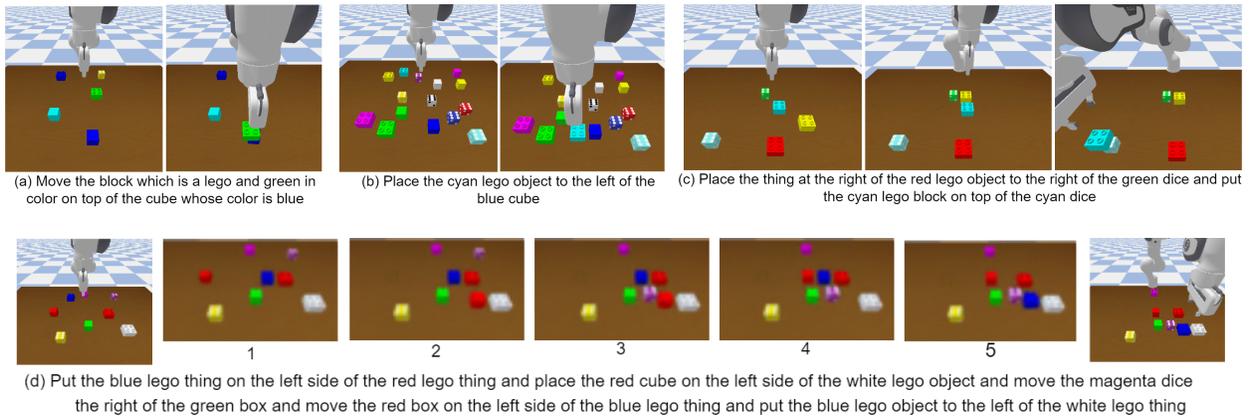---

[1] Excluded about 10% no detection cases for both models.

(a) Move the block which is a lego and green in color on top of the cube whose color is blue

(b) Place the cyan lego object to the left of the blue cube

(c) Place the thing at the right of the red lego object to the right of the green dice and put the cyan lego block on top of the cyan dice



1     2     3     4     5

(d) Put the blue lego thing on the left side of the red lego thing and place the red cube on the left side of the white lego object and move the magenta dice the right of the green box and move the red box on the left side of the blue lego thing and put the blue lego object to the left of the white lego thing

Fig. 5: Execution of robot manipulator on (a) compound instructions, (b) scene with 15 objects, (c) double step instruction with relational attributes, (d) 5-step instruction. (d) also shows reconstruction of the predicted scene before each step of the simulation

We observe that CLIPort on our dataset shows low accuracy in both predicting the object to be moved as well as the place location. Further, its performance deteriorates in examples with relational concepts. We note that the performance and sample-efficiency of CLIPort critically depends upon spatial consistency of input images i.e. objects do not scale or distort depending on the viewing angle as in [3]. To provide additional leverage to the baseline, we train CLIPort on 5 times more data than the original, results of which are reported under the title CLIPort-5x. We observe an improvement in performance on the larger dataset, although still worse than that of the proposed model.

### C. Combinatorial Generalization

To evaluate the approaches in an out-of-distribution generalization setting, we collect dataset consisting of scenes with more objects than seen during training, and with instructions that involve larger number of steps for carrying out the task.

Figure 4a illustrates the model generalizing combinatorially to scenes with more number of objects. The models were first trained on scenes having up to 5 objects only, and then tested on scenes having up to 10 objects.
The superior generalization demonstrated by the model can be attributed to reliance on an object-centric state representation and the ability to learn dense disentangled representations of spatial and action concepts, facilitating modular and structured reasoning that scales gracefully.

Figure 4b illustrates model generalization to longer horizon manipulation. We observe that the proposed model is able to perform multiple scene manipulation and reasoning steps with considerable accuracy up to 7 steps after being trained with instructions translating to plans up to 2 steps. The performance of the object-centric baseline (NMN+) is worse and the model struggles to generalize to plans extending to longer horizons. We attribute this to the modular structure of our approach compared to the baseline.

### D. Demonstration on a Simulated Robot

We demonstrate the learned model for interpreting instructions provided to a simulated 7-DOF Franka Emika manipulator in a table top setting. The robot is provided language instructions and uses the model to predict a program that once executed transitions the world state to the intended one. The 2-D bounding boxes predicted by the action simulator are translated to 3-D coordinates in the world space via a learned MLP using simulated data. The predicted positions are provided to a low-level motion planner for trajectory generation with crane grasping for picking/placing. Figure 5 shows execution by the robot manipulator on complex instructions, scenes having multiple objects, double step relational instructions, and multi-step instructions. We also visualise reconstruction of the moved objects before each step of the actual execution. The structural similarity index (SSIM) [36] for the reconstruction model is 0.935.

## VI. CONCLUSIONS

We present a neuro-symbolic architecture that learns grounded manipulation programs via visual-linguistic reasoning for instruction understanding over a given scene, to achieve a desired goal. Unlike previous work, we do not assume any sub-goal supervision, and demonstrate how our model can be trained end-to-end. Our experiments show strong generalization to novel scenes and instructions compared to a neural-only baseline. Directions for future work include dealing with richer instruction space including looping constructs, real-time recovery from errors caused by faulty execution, and working with real workspace data.

## VII. ACKNOWLEDGEMENTS

REFERENCES

[1] R. Paul, J. Arkin, N. Roy, and T. M Howard, "Efficient grounding of abstract spatial concepts for natural language interaction with robot manipulators," in *Robotics: Science and Systems Foundation*, 2016.

[2] C. Paxton, Y. Bisk, J. Thomason, A. Byravan, and D. Fox, "Prospection: Interpretable plans from language by predicting the future," in *2019 International Conference on Robotics and Automation (ICRA)*, IEEE, 2019, pp. 6942–6948.

[3] M. Shridhar, L. Manuelli, and D. Fox, "Cliport: What and where pathways for robotic manipulation," in *Conference on Robot Learning*, PMLR, 2022, pp. 894–906.

[4] T. M. Howard, S. Tellex, and N. Roy, "A natural language planner interface for mobile manipulators," in *2014 IEEE International Conference on Robotics and Automation (ICRA)*, IEEE, 2014, pp. 6652–6659.

[5] S. Tellex, T. Kollar, S. Dickerson, M. R. Walter, A. G. Banerjee, S. Teller, and N. Roy, "Approaching the symbol grounding problem with probabilistic graphical models," *AI magazine*, vol. 32, no. 4, pp. 64–76, 2011.

[6] C. Matuszek, E. Herbst, L. Zettlemoyer, and D. Fox, "Learning to parse natural language commands to a robot control system," in *Experimental robotics*, Springer, 2013, pp. 403–415.

[7] R. A. Knepper, T. Layton, J. Romanishin, and D. Rus, "Ikeabot: An autonomous multi-robot coordinated furniture assembly system," in *2013 IEEE International conference on robotics and automation*, IEEE, 2013, pp. 855–862.

[8] N. Gopalan, D. Arumugam, L. L. Wong, and S. Tellex, "Sequence-to-sequence language grounding of non-markovian task specifications.," in *Robotics: Science and Systems*, vol. 2018, 2018.

[9] E. C. Williams, N. Gopalan, M. Rhee, and S. Tellex, "Learning to parse natural language to grounded reward functions with weak supervision," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, IEEE, 2018, pp. 4430–4436.

[10] A. J. Shah, P. Kamath, S. Li, and J. A. Shah, "Bayesian inference of temporal task specifications from demonstrations," 2018.

[11] C. Wang, C. Ross, Y.-L. Kuo, B. Katz, and A. Barbu, "Learning a natural-language to ltl executable semantic parser for grounded robotics," *arXiv preprint arXiv:2008.03277*, 2020.

[12] H. Kress-Gazit, G. E. Fainekos, and G. J. Pappas, "Translating structured english to robot controllers," *Advanced Robotics*, vol. 22, no. 12, pp. 1343–1359, 2008.

[13] M. Lázaro-Gredilla, D. Lin, J. S. Guntupalli, and D. George, "Beyond imitation: Zero-shot task transfer on robots by learning concepts as cognitive programs," *Science Robotics*, vol. 4, no. 26, 2019.

[14] M. Tenorth, D. Nyga, and M. Beetz, "Understanding and executing instructions for everyday manipulation tasks from the world wide web," in *2010 ieee international conference on robotics and automation*, IEEE, 2010, pp. 1486–1491.

[15] G. Lisca, D. Nyga, F. Bálint-Benczédi, H. Langer, and M. Beetz, "Towards robots conducting chemical experiments," in *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, IEEE, 2015, pp. 5202–5208.

[16] D. K. Misra, J. Sung, K. Lee, and A. Saxena, "Tell me dave: Context-sensitive grounding of natural language to manipulation instructions," *The International Journal of Robotics Research*, vol. 35, no. 1-3, pp. 281–300, 2016.

[17] G. Konidaris, L. P. Kaelbling, and T. Lozano-Perez, "From skills to symbols: Learning symbolic representations for abstract high-level planning," *Journal of Artificial Intelligence Research*, vol. 61, pp. 215–289, 2018.

[18] Z. Wang, C. R. Garrett, L. P. Kaelbling, and T. Lozano-Pérez, "Learning compositional models of robot skills for task and motion planning," *The International Journal of Robotics Research*, vol. 40, no. 6-7, pp. 866–894, 2021.

[19] L. S. Zettlemoyer, H. Pasula, and L. P. Kaelbling, "Learning planning rules in noisy stochastic worlds," in *AAAI*, 2005, pp. 911–918.

[20] V. Xia, Z. Wang, K. Allen, T. Silver, and L. P. Kaelbling, "Learning sparse relational transition models," in *International Conference on Learning Representations*, 2018.

[21] T. Silver, K. R. Allen, A. K. Lew, L. P. Kaelbling, and J. Tenenbaum, "Few-shot bayesian imitation learning with logical program policies," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, 2020, pp. 10 251–10 258.

[22] Y. Zhu, J. Tremblay, S. Birchfield, and Y. Zhu, "Hierarchical planning for long-horizon manipulation with geometric and symbolic scene graphs," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, IEEE, 2021, pp. 6541–6548.

[23] A. Zeng, P. Florence, J. Tompson, S. Welker, J. Chien, M. Attarian, T. Armstrong, I. Krasin, D. Duong, V. Sindhwani, *et al.*, "Transporter networks: Rearranging the visual world for robotic manipulation," *arXiv preprint arXiv:2010.14406*, 2020.

[24] J. Mao, C. Gan, P. Kohli, J. B. Tenenbaum, and J. Wu, "The Neuro-Symbolic Concept Learner: Interpreting Scenes, Words, and Sentences From Natural Supervision," in *International Conference on Learning Representations*, 2019. [Online]. Available: https://openreview.net/forum?id=rJgMlhRctm.

[25] R Williams, "A class of gradient-estimation algorithms for reinforcement learning in neural networks," in *Proceedings of the International Conference on Neural Networks*, 1987, pp. II–601.

[26] R. Wang, J. Mao, J. Hsu, H. Zhao, J. Wu, and Y. Gao, "Programmatically grounded, compositionally generalizable robotic manipulation," in *The Eleventh International Conference on Learning Representations*, 2023. [Online]. Available: `https://openreview.net/forum?id=rZ-wylY5VI`.

[27] L. P. Kaelbling and T. Lozano-Pérez, "Hierarchical planning in the now," in *Workshops at the Twenty-Fourth AAAI Conference on Artificial Intelligence*, 2010.

[28] Y. Zhu, J. Tremblay, S. Birchfield, and Y. Zhu, "Hierarchical planning for long-horizon manipulation with geometric and symbolic scene graphs," *arXiv preprint arXiv:2012.07277*, 2020.

[29] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, 1997.

[30] L. Dong and M. Lapata, "Language to logical form with neural attention," in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2016, pp. 33–43.

[31] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.

[32] S. Targ, D. Almeida, and K. Lyman, "Resnet in resnet: Generalizing residual architectures," *arXiv preprint arXiv:1603.08029*, 2016.

[33] H. Dhamo, A. Farshad, I. Laina, N. Navab, G. D. Hager, F. Tombari, and C. Rupprecht, "Semantic image manipulation using scene graphs," in *CVPR*, 2020.

[34] J. Andreas, M. Rohrbach, T. Darrell, and D. Klein, "Neural module networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 39–48.

[35] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, *Learning transferable visual models from natural language supervision*, 2021. DOI: `10.48550/ARXIV.2103.00020`. [Online]. Available: `https://arxiv.org/abs/2103.00020`.

[36] Z. Wang, A. Bovik, H. Sheikh, and E. Simoncelli, "Image quality assessment: From error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004. DOI: `10.1109/TIP.2003.819861`.