

VoxPoser: Composable 3D Value Maps for Robotic Manipulation with Language Models

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Large language models (LLMs) are shown to possess a wealth of ac-
2 tionable knowledge that can be extracted for robot manipulation in the form of
3 reasoning and planning. Despite the progress, most still rely on pre-defined mo-
4 tion primitives to carry out the physical interactions with the environment, which
5 remains a major bottleneck. In this work, we aim to synthesize robot trajec-
6 tories, i.e., a dense sequence of 6-DoF end-effector waypoints, for a large variety of
7 manipulation tasks given an *open-set of instructions* and an *open-set of objects*.
8 We achieve this by first observing that LLMs excel at inferring affordances and
9 constraints given a free-form language instruction. More importantly, by leverag-
10 ing their code-writing capabilities, they can interact with a vision-language model
11 (VLM) to compose 3D value maps to ground the knowledge into the observation
12 space of the agent. The composed value maps are then used in a model-based plan-
13 ning framework to *zero-shot* synthesize closed-loop robot trajectories with robust-
14 ness to dynamic perturbations. We further demonstrate how the proposed frame-
15 work can benefit from online experiences by efficiently learning a dynamics model
16 for scenes that involve contact-rich interactions. We present a large-scale study of
17 the proposed method in both simulated and real-robot environments, showcasing
18 the ability to perform a large variety of everyday manipulation tasks specified in
19 free-form natural language. Project website: voxposer-anon.github.io.

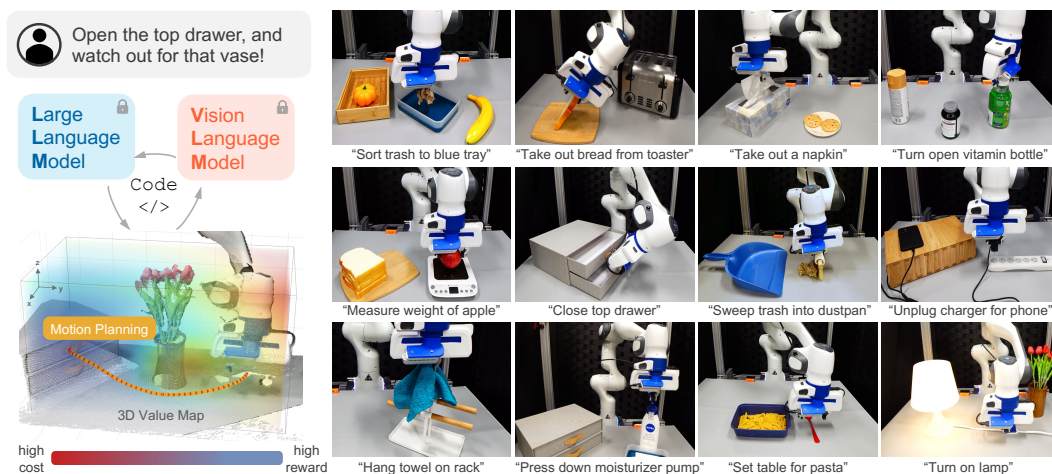


Figure 1: VOXPOSER extracts language-conditioned **affordances** and **constraints** from LLMs and grounds them to the perceptual space using VLMs, using a code interface and without additional training to either component. The composed map is referred to as a 3D value map, which enables **zero-shot** synthesis of trajectories for large varieties of everyday manipulation tasks with an **open-set of instructions** and an **open-set of objects**.

20 1 Introduction

21 Language is a compressed medium through which humans distill and communicate their knowledge
22 and experience of the world. Large language models (LLMs) have emerged as a promising approach
23 to capture this abstraction, learning to represent the world through projection into language space [1–
24 4]. While these models are believed to internalize generalizable knowledge as text, it remains a
25 question about how to use it to enable embodied agents to physically *act* in the real world.

26 We look at the problem of grounding abstract language instructions (e.g., “set up the table”) in robot
27 actions [5]. Prior works have leveraged lexical analysis to parse the instructions [6–8], while more
28 recently language models have been used to decompose the instructions into a textual sequence of
29 steps [9–11]. However, to enable physical interactions with the environment, existing approaches
30 typically rely on a repertoire of pre-defined motion primitives (i.e., skills) that may be invoked by
31 an LLM or a planner, and this reliance on individual skill acquisition is often considered a major
32 bottleneck of the system due to the lack of large-scale robotic data. The question then arises: how can
33 we leverage the wealth of internalized knowledge of LLMs at the even fine-grained action level for
34 robots, without requiring laborious data collection or manual designs for each individual primitive?

35 In addressing this challenge, we first note that it is impractical for LLMs to directly output control
36 actions in text, which are typically driven by high-frequency control signals in high-dimensional
37 space. However, we find that LLMs excel at inferring language-conditioned *affordances* and *con-*
38 *straints*, and by leveraging their code-writing capabilities, they can compose dense 3D voxel maps
39 that ground them in the visual space by orchestrating perception calls (e.g., via CLIP [12] or open-
40 vocabulary detectors [13–15]) and array operations (e.g., via NumPy [16]). For example, given an
41 instruction “open the top drawer and watch out for the vase”, LLMs can be prompted to infer: 1)
42 the top drawer handle should be grasped, 2) the handle needs to be translated outwards, and 3) the
43 robot should stay away from the vase. By generating Python code to invoke perception APIs, LLMs
44 can obtain spatial-geometric information of relevant objects or parts and then manipulate the 3D
45 voxels to prescribe reward or cost at relevant locations in observation space (e.g., the handle region
46 is assigned high values while the surrounding of the vase is assigned low values). Finally, the com-
47 posed value maps can serve as objective functions for motion planners to directly synthesize robot
48 trajectories that achieve the given instruction¹, without requiring additional training data for each
49 task or for the LLM. An illustration diagram and a subset of tasks we considered are shown in Fig. 1.

50 We term this approach **VOXPOSER**, a formulation that extracts affordances and constraints from
51 LLMs to compose 3D value maps in observation space for guiding robotic interactions. Rather than
52 relying on robotic data that are often of limited amount or variability, the method leverages LLMs
53 for *open-world reasoning* and VLMs for *generalizable visual grounding* in a model-based planning
54 framework that directly enables *physical* robot actions. We demonstrate its zero-shot generalization
55 for *open-set* instructions with *open-set* objects for various everyday manipulation tasks. We fur-
56 ther showcase how VoxPoser can also benefit from limited online interactions to efficiently learn a
57 dynamics model that involves contact-rich interactions.

58 2 Related Works

59 **Grounding Language Instructions.** Language grounding has been studied extensively both in
60 terms of intelligent agents [19–22] and of robotics [23, 6, 24, 25, 5, 7, 26], where language can be
61 used as a tool for compositional goal specification [5, 27–33], semantic anchor for training multi-
62 modal representation [12, 34, 35], or as an intermediate substrate for planning and reasoning [36–
63 38, 9, 10, 39, 40]. Prior works have looked at using classical tools such as lexical analysis, formal
64 logic, and graphical models to interpret language instructions [27, 7, 6, 26]. More recently, end-
65 to-end approaches, popularized by successful applications to offline domains [41–43, 1], have been
66 applied to directly ground language instructions in robot interactions by learning from data with

¹The approach also bears resemblance and connections to potential field methods in path planning [17] and constrained optimization methods in manipulation planning [18].

67 language annotations, spanning from model learning [44], imitation learning [45, 46, 30, 47–54], to
68 reinforcement learning [55–57]. Most closely related to our work is Sharma et al. [50], where an
69 end-to-end cost predictor is optimized via supervised learning to map language instructions to 2D
70 costmaps, which are used to steer a motion planner to generate preferred trajectories in a collision-
71 free manner. In contrast, we rely on pre-trained language models for their open-world knowledge
72 and tackle the more challenging robotic manipulation in 3D.

73 **Language Models for Robotics.** Leveraging pre-trained language models for embodied applica-
74 tions is an active area of research, where a large body of works focus on planning and reasoning
75 with language models [9–11, 58, 31, 39, 59–72, 36, 73, 74]. To allow language models to perceive
76 the physical environments, textual descriptions of the scene [39, 11, 59] or perception APIs [75] can
77 be given, vision can be used during decoding [67] or can be directly taken as input by multi-modal
78 language models [68, 2]. In addition to perception, to truly bridge the perception-action loop, an
79 embodied language model must also know how to *act*, which typically is achieved by a library of
80 pre-defined primitives. Liang et al. [75] showed that LLMs exhibit behavioral commonsense that
81 can be useful for low-level control. Despite the promising signs, hand-designed motion primitives
82 are still required, and while LLMs are shown to be capable of composing *sequential* policy logic, it
83 remains unclear whether composition can happen at *spatial* level. A related line of works has also
84 explored using LLMs for reward specification in the context of reward design [76] and exploration
85 in reinforcement learning [77–80], and human preference learning [81]. In contrast, we focus exclu-
86 sively on grounding the reward generated by LLMs in the *3D observation space* of the robot, which
87 we identify as most useful for manipulation tasks.

88 **Learning-based Trajectory Optimization.** Many works have explored leveraging learning-based
89 approaches for trajectory optimization. While the literature is vast, they can be broadly categorized
90 into those that learn the models [82–90] and those that learn the cost/reward or constraints [91–
91 94, 50, 95], where data are typically collected from in-domain interactions. To enable generaliza-
92 tion in the wild, a parallel line of works has explored learning task specification from large-scale
93 offline data [96–98, 35, 34, 44, 99, 100, 54], particularly egocentric videos [101, 102], or leverag-
94 ing pre-trained foundation models [103–105, 33, 106, 107]. The learned cost functions are then
95 used by reinforcement learning [103, 100, 108], imitation learning [98, 97], or trajectory optimiza-
96 tion [96, 35] to generate robot actions. In this work, we leverage LLMs for *zero-shot in-the-wild*
97 cost specification with superior generalization. Compared to prior works that leverage foundation
98 models, we ground the cost directly in *3D observation space* with *real-time* visual feedback, which
99 makes VoxPoser amenable to closed-loop MPC that’s robust in execution.

100 3 Method

101 We first provide the formulation of VoxPoser as an optimization problem (Sec. 3.1). Then we de-
102 scribe how VoxPoser can be used as a general zero-shot framework to map language instructions
103 to 3D value maps (Sec. 3.2). We subsequently demonstrate how trajectories can be synthesized
104 in closed-loop for robotic manipulation (Sec. 3.3). While zero-shot in nature, we demonstrate
105 how VoxPoser can learn from online interactions to efficiently solve contact-rich tasks (Sec. 3.4).

106 3.1 Problem Formulation

107 Consider a manipulation problem given as a *free-form* language instruction \mathcal{L} (e.g., “open the top
108 drawer”). Generating robot trajectories according to \mathcal{L} can be very challenging because \mathcal{L} may
109 be arbitrarily long-horizon or under-specified (i.e., requires contextual understanding). Instead, we
110 focus on individual phases (sub-tasks) of the problem ℓ_i that distinctively specify a manipulation
111 task (e.g., “grasp the drawer handle”, “pull open the drawer”), where the decomposition $\mathcal{T} \rightarrow$
112 $(\ell_1, \ell_2, \dots, \ell_n)$ is given by a high-level planner (e.g., an LLM or a search-based planner)². The
113 central problem investigated in this work is to generate a motion trajectory τ_i^r for robot r and each

²Note that the decomposition and sequencing of these sub-tasks are also done by LLMs in this work, though we do not investigate this aspect extensively as it is not the focus of our contributions.

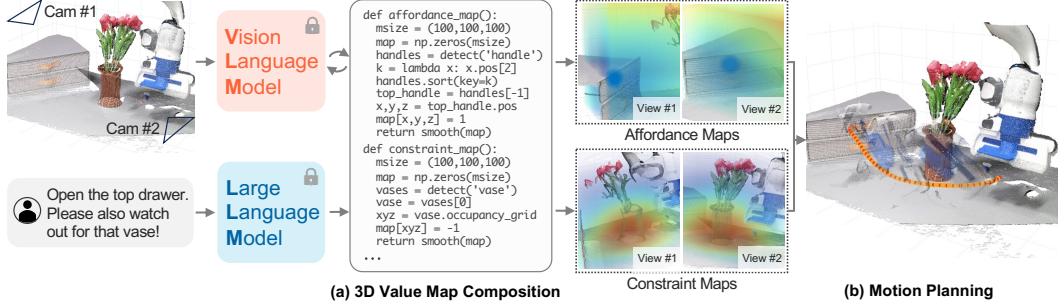


Figure 2: Overview of VOXPOSER. Given the RGB-D observation of the environment and a language instruction, LLMs generate code, which interacts with VLMs, to produce a sequence of 3D affordance maps and constraint maps (collectively referred to as value maps) grounded in the observation space of the robot (a). The composed value maps then serve as objective functions for motion planners to synthesize trajectories for robot manipulation (b). The entire process does not involve any additional training.

114 manipulation phase described by instruction ℓ_i . We represent τ_i^r as a sequence of dense end-effector
 115 waypoints to be executed by an Operational Space Controller [109], where each waypoint consists
 116 of a desired 6-DoF end-effector pose, end-effector velocity, and gripper action. However, it is worth
 117 noting that other representations of trajectories, such as joint space trajectories, can also be used.
 118 Given each sub-task ℓ_i , we formulate this as an optimization problem defined as follows:

$$\min_{\tau_i^r} \{ \mathcal{F}_{task}(\mathbf{T}_i, \ell_i) + \mathcal{F}_{control}(\tau_i^r) \} \quad \text{subject to} \quad \mathcal{C}(\mathbf{T}_i) \quad (1)$$

119 where \mathbf{T}_i is the evolution of environment state, and $\tau_i^r \subseteq \mathbf{T}_i$ is the robot trajectory. \mathcal{F}_{task} scores the
 120 extent of \mathbf{T}_i completes the instruction ℓ_i while $\mathcal{F}_{control}$ specifies the control costs, e.g., to encourage
 121 τ_i^r to minimize total control effort or total time. $\mathcal{C}(\mathbf{T}_i)$ denotes the dynamics and kinematics con-
 122 straints, which are enforced by the known model of the robot and a physics-based or learning-based
 123 model of the environment. By solving this optimization for each sub-task ℓ_i , we obtain a sequence
 124 of robot trajectories that collectively achieve the overall task specified by the instruction \mathcal{L} .
 125

126 3.2 Grounding Language Instruction via VoxPoser

127 Calculating \mathcal{F}_{task} with respect to free-form language instructions is extremely challenging, not only
 128 because of the rich space of semantics language can convey but also because of the lack of robot data
 129 labeled with \mathbf{T} and ℓ . However, we provide a critical observation that a large number of tasks can
 130 be characterized by a voxel value map $\mathbf{V} \in \mathbb{R}^{w \times h \times d}$ in robot’s observation space, which guides the
 131 motion of an “entity of interest” in the scene, such as the robot end-effector, an object, or an object
 132 part. For example, consider the task “open the top drawer” and its first sub-task “grasp the top
 133 drawer handle” (inferred by LLMs) in Fig. 2. The “entity of interest” is the robot end-effector, and
 134 the voxel value map should reflect the attraction toward the drawer handle. By further commanding
 135 “watch out for the vase”, the map can also be updated to reflect the repulsion from the vase. We
 136 denote the “entity of interest” as \mathbf{e} and its trajectory as τ^e . Using this voxel value map for a given
 137 instruction ℓ_i , \mathcal{F}_{task} can be approximated by accumulating the values of \mathbf{e} traversing through \mathbf{V}_i ,
 138 formally calculated as $\mathcal{F}_{task} = - \sum_{j=1}^{|\tau_i^e|} \mathbf{V}(p_j^e)$, where $p_j^e \in \mathbb{N}^3$ is the discretized (x, y, z) position
 139 of \mathbf{e} at step j .

140 Notably, we observe large language models, by being pre-trained on Internet-scale data, exhibit ca-
 141 pabilities not only to identify the “entity of interest” but also to compose value maps that accurately
 142 reflect the task instruction by writing Python programs. Specifically, when an instruction is given
 143 as a comment in the code, LLMs can be prompted to 1) call perception APIs (which invoke vision-
 144 language models (VLM) such as an open-vocabulary detector [13–15]) to obtain spatial-geometrical
 145 information of relevant objects, 2) generate NumPy operations to manipulate 3D arrays, and 3) pre-
 146 scribe precise values at relevant locations. We term this approach as **VOXPOSER**. Concretely, we
 147 aim to obtain a voxel value map $\mathbf{V}_i^t = \text{VoxPoser}(\mathbf{o}^t, \ell_i)$ by prompting an LLM and executing the
 148 code via a Python interpreter, where \mathbf{o}^t is the RGB-D observation at time t and ℓ_i is the current

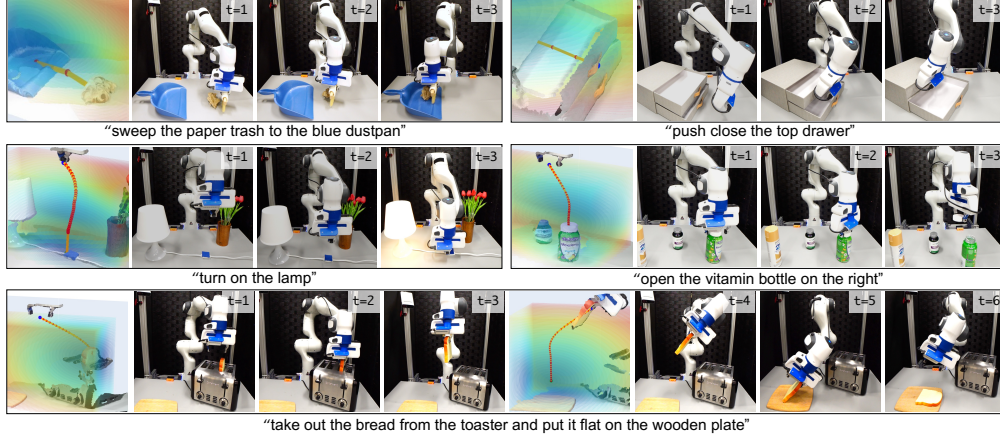


Figure 3: Visualization of composed 3D value maps and rollouts in real-world environments. The top row demonstrates where “entity of interest” is an object or part, and the value maps guide them toward target positions. The bottom two rows showcase tasks where “entity of interest” is the robot end-effector. The bottom-most task involves two phases, which are also orchestrated by LLMs.

149 instruction. Additionally, because \mathbf{V} is often sparse, we densify the voxel maps via smoothing
 150 operations, as they encourage smoother trajectories optimized by motion planners.

151 **Additional Trajectory Parametrization.** The above formulation of VoxPoser uses LLMs to com-
 152 pose $\mathbf{V} : \mathbb{N}^3 \rightarrow \mathbb{R}$ to map from discretized coordinates in voxel space to a real-valued “cost”, which
 153 we can use to optimize a path consisting only of the positional terms. To extend to $SE(3)$ poses,
 154 we can also use LLMs to compose rotation maps $\mathbf{V}_r : \mathbb{N}^3 \rightarrow SO(3)$ at coordinates relevant to the
 155 task objectives (e.g., “end-effector should face the support normal of the handle”). Similarly, we
 156 further compose gripper maps $\mathbf{V}_g : \mathbb{N}^3 \rightarrow \{0, 1\}$ to control gripper open/close and velocity maps
 157 $\mathbf{V}_v : \mathbb{N}^3 \rightarrow \mathbb{R}$ to specify target velocities. Note that while these additional trajectory parametriza-
 158 tions are not mapped to a real-valued “cost”, they can also be factored in the optimization procedure
 159 (Equation 1) to parametrize the trajectories.

160 3.3 Zero-Shot Trajectory Synthesis with VoxPoser

161 After obtaining the task cost \mathcal{F}_{task} , we can now approach the full problem defined in Equation 1
 162 to plan a motion trajectory. We use simple zeroth-order optimization by randomly sampling trajec-
 163 tories and scoring them with the proposed objective. The optimization is implemented in a model
 164 predictive control framework that iteratively replans the trajectory at every step using the current
 165 observation to robustly execute the trajectories even under dynamic disturbances³, where either
 166 a learned or physics-based model can be used. However, because VoxPoser effectively provides
 167 “dense rewards” in the observation space and we are able to replan at every step, we surprisingly
 168 find that the overall system can already achieve a large variety of manipulation tasks considered in
 169 this work even with simple heuristics-based models. Since some value maps are defined over “entity
 170 of interest”, which may not necessarily be the robot, we also use the dynamics model to find the
 171 needed robot trajectory to minimize the task cost (i.e., what interactions between the robot and the
 172 environment achieve the desired object motions).

173 3.4 Efficient Dynamics Learning with Online Experiences

174 While Sec. 3.3 presents a zero-shot framework for synthesizing trajectories for robot manipula-
 175 tion, VoxPoser can also benefit from online experiences by efficiently learning a dynamics model.
 176 Consider the standard setup where a robot interleaves between 1) collecting environment transition
 177 data $(\mathbf{o}_t, \mathbf{a}_t, \mathbf{o}_{t+1})$, where \mathbf{o}_t is the environment observation at time t and $\mathbf{a}_t = \text{MPC}(\mathbf{o}_t)$, and 2)
 178 training a dynamics model \mathbf{g}_θ parametrized by θ by minimizing the L2 loss between predicted next

³Although involving an LLM in the loop, closed-loop execution is possible because the generated code remains the same throughout task ℓ_i , which allows us to cache its output for the current task.

179 observation $\hat{\mathbf{o}}_{t+1}$ and \mathbf{o}_{t+1} . A critical component that determines the learning efficiency is the ac-
 180 tion sampling distribution $P(\mathbf{a}_t|\mathbf{o}_t)$ in MPC, which typically is a random distribution over the full
 181 action space \mathbf{A} . This is often inefficient when the goal is to solve a particular task, such as opening
 182 a door, because most actions do not interact with the relevant objects in the scene (i.e., the door
 183 handle) nor do they necessarily interact with the objects in a meaningful way (i.e., pressing down
 184 the door handle). Since VoxPoser synthesizes robot trajectories with LLMs, which have a wealth of
 185 commonsense knowledge, the zero-shot synthesized trajectory τ_0^r can serve as a useful prior to bias
 186 the action sampling distribution $P(\mathbf{a}_t|\mathbf{o}_t, \tau_0^r)$, which can significantly speed up the learning process.
 187 In practice, this can be implemented by only sampling actions in the vicinity of τ_0^r by adding small
 188 noise ε to encourage local exploration instead of exploring in the full action space \mathbf{A} .

189 4 Experiments and Analysis

190 We first discuss our implementation details. Then we validate VoxPoser for real-world everyday ma-
 191 nipulation (Sec. 4.1). We also study its generalization in simulation (Sec. 4.2). We further demon-
 192 strate how VoxPoser enables efficient learning of more challenging tasks (Sec. 4.3). Finally, we
 193 analyze its source of errors and discuss how improvement can be made (Sec. 4.4).

194 **LLMs and Prompting.** We follow prompting structure by Liang et al. [75], which recursively calls
 195 LLMs using their own generated code, where each language model program (LMP) is responsible
 196 for a unique functionality (e.g., processing perception calls). We use GPT-4 [2] from [OpenAI API](#).
 197 For each LMP, we include 5-20 example queries and corresponding responses as part of the prompt.
 198 An example can be found in Fig. 2 (simplified for clarity). Full prompts are in Appendix.

199 **VLMs and Perception.** Given an object/part query from LLMs, we first invoke open-vocab detector
 200 OWL-ViT [15] to obtain a bounding box, then feed it into Segment Anything [110] to obtain a mask,
 201 and finally track the mask using video tracker XMEM [111]. The tracked mask is used with RGB-D
 202 observation to reconstruct the object/part point cloud.

203 **Value Map Composition.** We define the following types of value maps: affordance, avoidance, end-
 204 effector velocity, end-effector rotation, and gripper action. Each type uses a different LMP, which
 205 takes in an instruction and outputs a voxel map of shape $(100, 100, 100, k)$, where k differs for each
 206 value map (e.g., $k = 1$ for affordance and avoidance as it specifies cost, and $k = 4$ for rotation as
 207 it specifies $SO(3)$). We apply Euclidean distance transform to affordance maps and Gaussian filters
 208 for avoidance maps. On top of value map LMPs, we define two high-level LMPs to orchestrate their
 209 behaviors: planner takes user instruction \mathcal{L} as input (e.g., “open drawer”) and outputs a sequence
 210 of sub-tasks $\ell_{1:N}$, and composer takes in sub-task ℓ_i and invokes relevant value map LMPs with
 211 detailed language parameterization.

212 **Motion Planner.** We consider only affordance and avoidance maps in the planner optimization,
 213 which finds a sequence of collision-free end-effector positions $p_{1:N} \in \mathbb{R}^3$ using greedy search. Then
 214 we enforce other parametrization at each p by the remaining value maps (e.g., rotation map, velocity
 215 map). The cost map used by the motion planner is computed as the negative of the weighted sum
 216 of normalized affordance and avoidance maps with weights 2 and 1. After a 6-DoF trajectory is
 217 synthesized, the first waypoint is executed, and then a new trajectory is re-planned at 5 Hz.

218 **Dynamics Model.** We use the known robot dynamics model in all tasks, where it is used in motion
 219 planning for the end-effector to follow the waypoints. For the majority of our considered tasks where
 220 the “entity of interest” is the robot, no environment dynamics model is used (i.e., scene is assumed
 221 to be static), but we replan at every step to account for the latest observation. For tasks in which
 222 the “entity of interest” is an object, we study only a planar pushing model parametrized by contact
 223 point, push direction, and push distance. We use a heuristic-based dynamics model that translates
 224 an input point cloud along the push direction by the push distance. We use MPC with random
 225 shooting to optimize for the action parameters. Then a pre-defined pushing primitive is executed
 226 based on the action parameters. However, we note that a primitive is not necessary when action
 227 parameters are defined over the end-effector or joint space of the robot, which would likely yield

Task	LLM + Prim. [75]		VoxPoser		Train/Test	Category	U-Net	Language Models	
	Static	Dist.	Static	Dist.			MP [50]	Prim. [75]	MP (Ours)
Move & Avoid	0/10	0/10	9/10	8/10	SI SA	Object Int.	21.0%	41.0%	64.0%
Set Up Table	7/10	0/10	9/10	7/10	SI SA	Composition	53.8%	43.8%	77.5%
Close Drawer	0/10	0/10	10/10	7/10	SI UA	Object Int.	3.0%	46.0%	60.0%
Open Bottle	5/10	0/10	7/10	5/10	SI UA	Composition	3.8%	25.0%	58.8%
Sweep Trash	0/10	0/10	9/10	8/10	UI UA	Object Int.	0.0%	17.5%	65.0%
Total	24.0%	0.0%	88.0%	70.0%	UI UA	Composition	0.0%	25.0%	76.7%

Table 1: Success rate in real-world domain. VoxPoser performs everyday manipulation tasks with high success and is more robust to disturbances than the baseline using action primitives.

Table 2: Success rate in simulated domain. “SI” and “UI” are seen and unseen instructions. “SA” and “UA” are seen and unseen attributes. VoxPoser outperforms both baselines across 13 tasks from two categories on both seen and unseen tasks and maintains similar success rates.

228 smoother trajectories but takes more time for optimization. We also explore the use of a learning-
 229 based dynamics model in Section 4.3, which enables VoxPoser to benefit from online experiences.

230 4.1 VoxPoser for Everyday Manipulation Tasks

231 We study whether VoxPoser can zero-shot synthesize robot trajectories to perform everyday manip-
 232 ulation tasks in the real world. Details of the environment setup can be found in Appendix A.3.
 233 While the proposed method can generalize to an open-set of instructions and an open-set of objects
 234 as shown in Fig. 1, we pick 5 representative tasks to provide quantitative evaluations in Table 1.
 235 Qualitative results including environment rollouts and value map visualizations are shown in Fig. 3.
 236 We find that VoxPoser can effectively synthesize robot trajectories for everyday manipulation tasks
 237 with a high average success rate. Due to fast replanning capabilities, it is also robust to external dis-
 238 turbances, such as moving targets/obstacles and pulling the drawer open after it has been closed by
 239 the robot. We further compare to a variant of Code as Policies [75] that uses LLMs to parameterize a
 240 pre-defined list of simple primitives (e.g., `move_to_pose`, `open_gripper`). We find that compared to
 241 chaining sequential policy logic, the ability to *compose spatially* while considering other constraints
 242 under a joint optimization scheme is a more flexible formulation, unlocking the possibility for more
 243 manipulation tasks and leading to more robust execution.

244 4.2 Generalization to Unseen Instructions and Attributes

245 To provide rigorous quantitative evaluations on generalization, we set up a simulated environment
 246 that mirrors our real-world setup [112] but features 13 highly-randomizable tasks with 2766 unique
 247 instructions. Each task comes with a templated instruction (e.g., “push [obj] to [pos]”) that con-
 248 tains randomizable attributes chosen from a pre-defined list. Details are in Appendix A.4. Seen
 249 instructions/attributes may appear in the prompt (or in the training data for supervised baselines).
 250 The tasks are grouped into 2 categories, where “Object Interactions” are tasks that require interac-
 251 tions with objects, and “Spatial Composition” are tasks involving spatial constraints (e.g., moving
 252 slower near a particular object). For baselines, we ablate the two components of VoxPoser, LLM
 253 and motion planner, by comparing to a variant of [75] that combines an LLM with primitives and
 254 to a variant of [50] that learns a U-Net [113] to synthesize costmaps for motion planning. Table 2
 255 shows the success rates averaged across 20 episodes per task. We find VoxPoser exhibits superior
 256 generalization in all scenarios. Compared to learned cost specification, LLMs generalize better by
 257 explicitly reasoning about affordances and constraints. On the other hand, grounding LLM knowl-
 258 edge in robot perception through *value map composition* rather than directly specifying primitive
 259 parameters offers more flexibility that generalizes beyond the prompt examples.

260 4.3 Efficient Dynamics Learning with Online Experiences

261 As discussed in Sec. 3.4, we investigate how VoxPoser can optionally benefit from online experi-
 262 ences for tasks that involve more intricacies of contact, such as opening doors, fridges, and windows,
 263 in a simulated environment. Specifically, we first synthesize k zero-shot trajectories using VoxPoser,

Task	Zero-Shot		No Prior		w/ Prior	
	Success	Time(s)	Success	Time(s)	Success	Time(s)
Door	6.7% \pm 4.4%	58.3 \pm 4.4%	TLE	88.3% \pm 1.67%	142.3 \pm 22.4	
Window	3.3% \pm 3.3%	36.7% \pm 1.7%	TLE	80.0% \pm 2.9%	137.0 \pm 7.5	
Fridge	18.3% \pm 3.3%	70.0% \pm 2.9%	TLE	91.7% \pm 4.4%	71.0 \pm 4.4	

Table 3: VoxPoser enables efficient dynamics learning by using zero-shot synthesized trajectories as prior. TLE (time limit exceeded) means exceeding 12 hours. Results are reported over 3 runs different seeds.

264 each represented as a sequence of end-effector waypoints, that act as priors for exploration (e.g.,
 265 “handle needs to be pressed down first in order to open a door”). Then an MLP dynamics model is
 266 learned through an iterative procedure where the agent alternates between data collection and model
 267 learning. During data collection, we add $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ to each waypoint in τ_0^r to encourage local
 268 exploration. As shown in Tab. 3, we find zero-shot synthesized trajectories are typically meaningful
 269 but insufficient. However, we can learn an effective dynamics model with less than 3 minutes of
 270 online interactions by using these trajectories as exploration prior, leading to high eventual success
 271 rates. In comparison, exploring without prior all exceed the maximum 12-hour limit.

272 4.4 Error Breakdown

273 In this section, we analyze the errors resulting from each component of VoxPoser and how the overall
 274 system can be further improved. We conduct experiments in simulation where we have access to
 275 ground-truth perception and dynamics model (i.e., the simulator). “Dynamics error” refers to errors
 276 made by the dynamics model⁴. “Perception error” refers to errors made by the perception module⁵.
 277 “Specification error” refers to errors made by the module specifying cost or parameters for the low-
 278 level motion planner or primitives. Examples for each method include 1) noisy prediction by the
 279 U-Net, 2) incorrect parameters specified by the LLM, and 3) incorrect value maps specified by the
 280 LLM. As shown in Fig. 4, VoxPoser achieves lowest “specification error” due to its generalization
 281 and flexibility. We also find that having access to a more robust perception pipeline and a physically-
 282 realistic dynamics model can contribute to better overall performance. This observation aligns with
 283 our real-world experiment, where most errors are from perception. For example, we find that the
 284 detector is sensitive to initial poses of objects and is less robust when detecting object parts.

285 5 Conclusion, Limitations, & Future Works

286 In this work, we present **VOXPOSER**, a general framework for extracting affordances and con-
 287 straints, grounded in 3D perceptual space, from LLMs and VLMs for everyday manipulation tasks
 288 in the real world, offering significant generalization advantages for open-set instructions and ob-
 289 jects. Despite compelling results, VoxPoser has several limitations. First, it relies on external per-
 290 ception modules, which is limiting in tasks that require holistic visual reasoning or understand-
 291 ing of fine-grained object geometries. Second, while applicable to efficient dynamics learning, a
 292 general-purpose dynamics model is still required to achieve contact-rich tasks with the same level of
 293 generalization. Third, our motion planner considers only end-effector trajectories while whole-arm
 294 planning is also feasible and likely a better design choice [115–117]. Finally, manual prompt engi-
 295 neering is required for LLMs. We also see several exciting venues for future work. For instance,
 296 recent success of multi-modal LLMs [68, 2, 118] can be directly translated into VoxPoser for direct
 297 visual grounding. Methods developed for alignment [119, 120] and prompting [121–124] may be
 298 used to alleviate prompt engineering effort. Finally, more advanced trajectory optimization methods
 299 can be developed that best interface with value maps synthesized by VoxPoser.

⁴LLM + Primitives [75] does not use model-based planning, thus not having a dynamics module.

⁵U-Net + MP [50] maps RGB-D to costmaps using U-Net [113, 114], thus not having perception module. Errors by which are attributed to “specification error”.



Figure 4: Error breakdown of components. VoxPoser significantly reduces specification error.

References

- [1] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [2] OpenAI. Gpt-4 technical report. *arXiv*, 2023.
- [3] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- [4] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- [5] S. Tellex, N. Gopalan, H. Kress-Gazit, and C. Matuszek. Robots that use language. *Annual Review of Control, Robotics, and Autonomous Systems*, 2020.
- [6] S. Tellex, T. Kollar, S. Dickerson, M. Walter, A. Banerjee, S. Teller, and N. Roy. Understanding natural language commands for robotic navigation and mobile manipulation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 25, pages 1507–1514, 2011.
- [7] T. Kollar, S. Tellex, D. Roy, and N. Roy. Toward understanding natural language directions. In *2010 5th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pages 259–266. IEEE, 2010.
- [8] M. Bollini, S. Tellex, T. Thompson, N. Roy, and D. Rus. Interpreting and executing recipes with a cooking robot. In *Experimental Robotics*, pages 481–495. Springer, 2013.
- [9] W. Huang, P. Abbeel, D. Pathak, and I. Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International Conference on Machine Learning*. PMLR, 2022.
- [10] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, D. Ho, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, E. Jang, R. J. Ruano, K. Jeffrey, S. Jesmonth, N. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, K.-H. Lee, S. Levine, Y. Lu, L. Luu, C. Parada, P. Pastor, J. Quiambao, K. Rao, J. Rettinghouse, D. Reyes, P. Sermanet, N. Sievers, C. Tan, A. Toshev, V. Vanhoucke, F. Xia, T. Xiao, P. Xu, S. Xu, and M. Yan. Do as i can and not as i say: Grounding language in robotic affordances. In *arXiv preprint arXiv:2204.01691*, 2022.
- [11] A. Zeng, A. Wong, S. Welker, K. Choromanski, F. Tombari, A. Purohit, M. Ryoo, V. Sindhwani, J. Lee, V. Vanhoucke, et al. Socratic models: Composing zero-shot multimodal reasoning with language. *arXiv preprint arXiv:2204.00598*, 2022.
- [12] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.
- [13] X. Gu, T.-Y. Lin, W. Kuo, and Y. Cui. Open-vocabulary object detection via vision and language knowledge distillation. *arXiv preprint arXiv:2104.13921*, 2021.
- [14] A. Kamath, M. Singh, Y. LeCun, G. Synnaeve, I. Misra, and N. Carion. Mdetr-modulated detection for end-to-end multi-modal understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1780–1790, 2021.

- 343 [15] M. Minderer, A. Gritsenko, A. Stone, M. Neumann, D. Weissenborn, A. Dosovitskiy, A. Ma-
344 hendran, A. Arnab, M. Dehghani, Z. Shen, et al. Simple open-vocabulary object detection
345 with vision transformers. *arXiv preprint arXiv:2205.06230*, 2022.
- 346 [16] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau,
347 E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk,
348 M. Brett, A. Haldane, J. F. del Río, M. Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard,
349 T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant. Array programming with
350 NumPy. *Nature*, 585(7825):357–362, Sept. 2020. doi:10.1038/s41586-020-2649-2. URL
351 <https://doi.org/10.1038/s41586-020-2649-2>.
- 352 [17] Y. K. Hwang, N. Ahuja, et al. A potential field approach to path planning. *IEEE transactions*
353 *on robotics and automation*, 8(1):23–32, 1992.
- 354 [18] M. Toussaint, J. Harris, J.-S. Ha, D. Driess, and W. Hönig. Sequence-of-constraints mpc:
355 Reactive timing-optimal control of sequential manipulation. In *2022 IEEE/RSJ International*
356 *Conference on Intelligent Robots and Systems (IROS)*, pages 13753–13760. IEEE, 2022.
- 357 [19] J. Andreas, D. Klein, and S. Levine. Learning with latent language. *arXiv preprint*
358 *arXiv:1711.00482*, 2017.
- 359 [20] R. Zellers, A. Holtzman, M. Peters, R. Mottaghi, A. Kembhavi, A. Farhadi, and Y. Choi.
360 Piglet: Language grounding through neuro-symbolic interaction in a 3d world. *arXiv preprint*
361 *arXiv:2106.00188*, 2021.
- 362 [21] R. Zellers, X. Lu, J. Hessel, Y. Yu, J. S. Park, J. Cao, A. Farhadi, and Y. Choi. Merlot:
363 Multimodal neural script knowledge models. *Advances in Neural Information Processing*
364 *Systems*, 2021.
- 365 [22] V. Shwartz, P. West, R. L. Bras, C. Bhagavatula, and Y. Choi. Unsupervised commonsense
366 question answering with self-talk. *arXiv preprint arXiv:2004.05483*, 2020.
- 367 [23] T. Winograd. Procedures as a representation for data in a computer program for understanding
368 natural language. Technical report, MASSACHUSETTS INST OF TECH CAMBRIDGE
369 PROJECT MAC, 1971.
- 370 [24] V. Blukis, R. A. Knepper, and Y. Artzi. Few-shot object grounding and mapping for natural
371 language robot instruction following. *arXiv preprint arXiv:2011.07384*, 2020.
- 372 [25] S. Tellex, R. Knepper, A. Li, D. Rus, and N. Roy. Asking for help using inverse semantics.
373 *Robotics: Science and Systems Foundation*, 2014.
- 374 [26] T. Kollar, S. Tellex, D. Roy, and N. Roy. Grounding verbs of motion in natural language
375 commands to robots. In *Experimental robotics*, pages 31–47. Springer, 2014.
- 376 [27] J. Thomason, S. Zhang, R. J. Mooney, and P. Stone. Learning to interpret natural language
377 commands through human-robot dialog. In *Twenty-Fourth International Joint Conference on*
378 *Artificial Intelligence*, 2015.
- 379 [28] J. Thomason, A. Padmakumar, J. Sinapov, N. Walker, Y. Jiang, H. Yedidsion, J. Hart, P. Stone,
380 and R. Mooney. Jointly improving parsing and perception for natural language commands
381 through human-robot dialog. *Journal of Artificial Intelligence Research*, 67:327–374, 2020.
- 382 [29] E. Jang, A. Irpan, M. Khansari, D. Kappler, F. Ebert, C. Lynch, S. Levine, and C. Finn.
383 Bc-z: Zero-shot task generalization with robotic imitation learning. In *Conference on Robot*
384 *Learning*, pages 991–1002. PMLR, 2021.
- 385 [30] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan,
386 K. Hausman, A. Herzog, J. Hsu, et al. Rt-1: Robotics transformer for real-world control
387 at scale. *arXiv preprint arXiv:2212.06817*, 2022.

- 388 [31] D. Shah, B. Osinski, B. Ichter, and S. Levine. Lm-nav: Robotic navigation with large pre-
389 trained models of language, vision, and action. *arXiv preprint arXiv:2207.04429*, 2022.
- 390 [32] Y. Cui, S. Karamcheti, R. Palleti, N. Shivakumar, P. Liang, and D. Sadigh. ”no, to the right”–
391 online language corrections for robotic manipulation via shared autonomy. *arXiv preprint*
392 *arXiv:2301.02555*, 2023.
- 393 [33] A. Stone, T. Xiao, Y. Lu, K. Gopalakrishnan, K.-H. Lee, Q. Vuong, P. Wohlhart, B. Zitkovich,
394 F. Xia, C. Finn, et al. Open-world object manipulation using pre-trained vision-language
395 models. *arXiv preprint arXiv:2303.00905*, 2023.
- 396 [34] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta. R3m: A universal visual represen-
397 tation for robot manipulation. *arXiv preprint arXiv:2203.12601*, 2022.
- 398 [35] Y. J. Ma, V. Kumar, A. Zhang, O. Bastani, and D. Jayaraman. Liv: Language-image repre-
399 sentations and rewards for robotic control. *arXiv e-prints*, 2023.
- 400 [36] P. A. Jansen. Visually-grounded planning without vision: Language models infer detailed
401 plans from high-level instructions. *arXiv preprint arXiv:2009.14259*, 2020.
- 402 [37] V. Micheli and F. Fleuret. Language models are few-shot butlers. *arXiv preprint*
403 *arXiv:2104.07972*, 2021.
- 404 [38] P. Sharma, A. Torralba, and J. Andreas. Skill induction and planning with latent language.
405 *arXiv preprint arXiv:2110.01517*, 2021.
- 406 [39] W. Huang, F. Xia, T. Xiao, H. Chan, J. Liang, P. Florence, A. Zeng, J. Tompson, I. Mor-
407 datch, Y. Chebotar, P. Sermanet, N. Brown, T. Jackson, L. Luu, S. Levine, K. Hausman, and
408 B. Ichter. Inner monologue: Embodied reasoning through planning with language models. In
409 *arXiv preprint arXiv:2207.05608*, 2022.
- 410 [40] B. Z. Li, W. Chen, P. Sharma, and J. Andreas. Lampp: Language models as probabilistic
411 priors for perception and action. *arXiv e-prints*, pages arXiv–2302, 2023.
- 412 [41] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale
413 Hierarchical Image Database. In *CVPR09*, 2009.
- 414 [42] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolu-
415 tional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.
- 416 [43] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. Language models are
417 unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 418 [44] S. Nair, E. Mitchell, K. Chen, S. Savarese, C. Finn, et al. Learning language-conditioned
419 robot behavior from offline data and crowd-sourced annotation. In *Conference on Robot*
420 *Learning*, pages 1303–1315. PMLR, 2022.
- 421 [45] M. Shridhar, L. Manuelli, and D. Fox. Cliport: What and where pathways for robotic manip-
422 ulation. In *Conference on Robot Learning*, pages 894–906. PMLR, 2022.
- 423 [46] M. Shridhar, L. Manuelli, and D. Fox. Perceiver-actor: A multi-task transformer for robotic
424 manipulation. In *Proceedings of the 6th Conference on Robot Learning (CoRL)*, 2022.
- 425 [47] S. Li, X. Puig, Y. Du, C. Wang, E. Akyurek, A. Torralba, J. Andreas, and I. Mordatch. Pre-
426 trained language models for interactive decision-making. *arXiv preprint arXiv:2202.01771*,
427 2022.
- 428 [48] O. Mees, L. Hermann, and W. Burgard. What matters in language conditioned robotic imita-
429 tion learning. *arXiv preprint arXiv:2204.06252*, 2022.

- 430 [49] O. Mees, J. Borja-Diaz, and W. Burgard. Grounding language with visual affordances over
431 unstructured data. *arXiv preprint arXiv:2210.01911*, 2022.
- 432 [50] P. Sharma, B. Sundaralingam, V. Blukis, C. Paxton, T. Hermans, A. Torralba, J. An-
433 dreas, and D. Fox. Correcting robot plans with natural language feedback. *arXiv preprint*
434 *arXiv:2204.05186*, 2022.
- 435 [51] W. Liu, C. Paxton, T. Hermans, and D. Fox. Structformer: Learning spatial structure for
436 language-guided semantic rearrangement of novel objects. In *2022 International Conference*
437 *on Robotics and Automation (ICRA)*, pages 6322–6329. IEEE, 2022.
- 438 [52] C. Lynch and P. Sermanet. Language conditioned imitation learning over unstructured data.
439 *Robotics: Science and Systems*, 2021. URL <https://arxiv.org/abs/2005.07648>.
- 440 [53] C. Lynch, A. Wahid, J. Tompson, T. Ding, J. Betker, R. Baruch, T. Armstrong, and P. Florence.
441 Interactive language: Talking to robots in real time. *arXiv preprint arXiv:2210.06407*, 2022.
- 442 [54] L. Shao, T. Migimatsu, Q. Zhang, K. Yang, and J. Bohg. Concept2robot: Learning manip-
443 ulation concepts from instructions and human demonstrations. *The International Journal of*
444 *Robotics Research*, 40(12-14):1419–1434, 2021.
- 445 [55] J. Luketina, N. Nardelli, G. Farquhar, J. N. Foerster, J. Andreas, E. Grefenstette, S. Whiteson,
446 and T. Rocktäschel. A survey of reinforcement learning informed by natural language. In
447 *IJCAI*, 2019.
- 448 [56] J. Andreas, D. Klein, and S. Levine. Modular multitask reinforcement learning with policy
449 sketches. *ArXiv*, abs/1611.01796, 2017.
- 450 [57] Y. Jiang, S. S. Gu, K. P. Murphy, and C. Finn. Language as an abstraction for hierarchical
451 deep reinforcement learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- 452 [58] B. Chen, F. Xia, B. Ichter, K. Rao, K. Gopalakrishnan, M. S. Ryoo, A. Stone, and D. Kappler.
453 Open-vocabulary queryable scene representations for real world planning. *arXiv preprint*
454 *arXiv:2209.09874*, 2022.
- 455 [59] I. Singh, V. Blukis, A. Mousavian, A. Goyal, D. Xu, J. Tremblay, D. Fox, J. Thomason, and
456 A. Garg. Progprompt: Generating situated robot task plans using large language models.
457 *arXiv preprint arXiv:2209.11302*, 2022.
- 458 [60] C. Huang, O. Mees, A. Zeng, and W. Burgard. Visual language maps for robot navigation.
459 *arXiv preprint arXiv:2210.05714*, 2022.
- 460 [61] S. S. Raman, V. Cohen, E. Rosen, I. Idrees, D. Paulius, and S. Tellex. Planning with large
461 language models via corrective re-prompting. *arXiv preprint arXiv:2211.09935*, 2022.
- 462 [62] C. H. Song, J. Wu, C. Washington, B. M. Sadler, W.-L. Chao, and Y. Su. Llm-planner: Few-
463 shot grounded planning for embodied agents with large language models. *arXiv preprint*
464 *arXiv:2212.04088*, 2022.
- 465 [63] B. Liu, Y. Jiang, X. Zhang, Q. Liu, S. Zhang, J. Biswas, and P. Stone. Llm+ p: Empowering
466 large language models with optimal planning proficiency. *arXiv preprint arXiv:2304.11477*,
467 2023.
- 468 [64] S. Vemprala, R. Bonatti, A. Bucker, and A. Kapoor. Chatgpt for robotics: Design principles
469 and model abilities. *2023*, 2023.
- 470 [65] K. Lin, C. Agia, T. Migimatsu, M. Pavone, and J. Bohg. Text2motion: From natural language
471 instructions to feasible plans. *arXiv preprint arXiv:2303.12153*, 2023.

- 472 [66] Y. Ding, X. Zhang, C. Paxton, and S. Zhang. Task and motion planning with large language
473 models for object rearrangement. *arXiv preprint arXiv:2303.06247*, 2023.
- 474 [67] W. Huang, F. Xia, D. Shah, D. Driess, A. Zeng, Y. Lu, P. Florence, I. Mordatch, S. Levine,
475 K. Hausman, et al. Grounded decoding: Guiding text generation with grounded models for
476 robot control. *arXiv preprint arXiv:2303.00855*, 2023.
- 477 [68] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson,
478 Q. Vuong, T. Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint*
479 *arXiv:2303.03378*, 2023.
- 480 [69] H. Yuan, C. Zhang, H. Wang, F. Xie, P. Cai, H. Dong, and Z. Lu. Plan4mc: Skill reinforce-
481 ment learning and planning for open-world minecraft tasks. *arXiv preprint arXiv:2303.16563*,
482 2023.
- 483 [70] Y. Xie, C. Yu, T. Zhu, J. Bai, Z. Gong, and H. Soh. Translating natural language to planning
484 goals with large-language models. *arXiv preprint arXiv:2302.05128*, 2023.
- 485 [71] Y. Lu, P. Lu, Z. Chen, W. Zhu, X. E. Wang, and W. Y. Wang. Multimodal procedural planning
486 via dual text-image prompting. *arXiv preprint arXiv:2305.01795*, 2023.
- 487 [72] D. Patel, H. Eghbalzadeh, N. Kamra, M. L. Iuzzolino, U. Jain, and R. Desai. Pretrained
488 language models as visual planners for human assistance. *arXiv preprint arXiv:2304.09179*,
489 2023.
- 490 [73] G. Wang, Y. Xie, Y. Jiang, A. Mandlekar, C. Xiao, Y. Zhu, L. Fan, and A. Anandku-
491 mar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint*
492 *arXiv:2305.16291*, 2023.
- 493 [74] J. Yang, W. Tan, C. Jin, B. Liu, J. Fu, R. Song, and L. Wang. Pave the way to grasp
494 anything: Transferring foundation models for universal pick-place robots. *arXiv preprint*
495 *arXiv:2306.05716*, 2023.
- 496 [75] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng. Code as
497 policies: Language model programs for embodied control. *arXiv preprint arXiv:2209.07753*,
498 2022.
- 499 [76] M. Kwon, S. M. Xie, K. Bullard, and D. Sadigh. Reward design with language models. *arXiv*
500 *preprint arXiv:2303.00001*, 2023.
- 501 [77] A. Tam, N. Rabinowitz, A. Lampinen, N. A. Roy, S. Chan, D. Strouse, J. Wang, A. Banino,
502 and F. Hill. Semantic exploration from language abstractions and pretrained representations.
503 *Advances in Neural Information Processing Systems*, 35:25377–25389, 2022.
- 504 [78] J. Mu, V. Zhong, R. Raileanu, M. Jiang, N. Goodman, T. Rocktäschel, and E. Grefenstette.
505 Improving intrinsic exploration with language abstractions. *arXiv preprint arXiv:2202.08938*,
506 2022.
- 507 [79] C. Colas, T. Karch, N. Lair, J.-M. Dussoux, C. Moulin-Frier, P. Dominey, and P.-Y. Oudeyer.
508 Language as a cognitive tool to imagine goals in curiosity driven exploration. *Advances in*
509 *Neural Information Processing Systems*, 33:3761–3774, 2020.
- 510 [80] Y. Du, O. Watkins, Z. Wang, C. Colas, T. Darrell, P. Abbeel, A. Gupta, and J. Andreas.
511 Guiding pretraining in reinforcement learning with large language models. *arXiv preprint*
512 *arXiv:2302.06692*, 2023.
- 513 [81] H. Hu and D. Sadigh. Language instructed reinforcement learning for human-ai coordination.
514 *arXiv preprint arXiv:2304.07297*, 2023.

- 515 [82] I. Lenz, R. A. Knepper, and A. Saxena. Deepmpc: Learning deep latent features for model
516 predictive control. In *Robotics: Science and Systems*, volume 10. Rome, Italy, 2015.
- 517 [83] L. Hewing, K. P. Wabersich, M. Menner, and M. N. Zeilinger. Learning-based model pre-
518 dictive control: Toward safe learning in control. *Annual Review of Control, Robotics, and*
519 *Autonomous Systems*, 3:269–296, 2020.
- 520 [84] M. B. Chang, T. Ullman, A. Torralba, and J. B. Tenenbaum. A compositional object-based
521 approach to learning physical dynamics. *arXiv preprint arXiv:1612.00341*, 2016.
- 522 [85] P. Battaglia, R. Pascanu, M. Lai, D. Jimenez Rezende, et al. Interaction networks for learning
523 about objects, relations and physics. *Advances in neural information processing systems*, 29,
524 2016.
- 525 [86] Z. Xu, J. Wu, A. Zeng, J. B. Tenenbaum, and S. Song. Densephysnet: Learning dense physical
526 object representations via multi-step dynamic interactions. *arXiv preprint arXiv:1906.03853*,
527 2019.
- 528 [87] A. Byravan and D. Fox. Se3-nets: Learning rigid body motion using deep neural networks.
529 In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 173–180.
530 IEEE, 2017.
- 531 [88] A. Nagabandi, K. Konolige, S. Levine, and V. Kumar. Deep dynamics models for learning
532 dexterous manipulation. In *Conference on Robot Learning*, pages 1101–1112. PMLR, 2020.
- 533 [89] A. Sanchez-Gonzalez, N. Heess, J. T. Springenberg, J. Merel, M. Riedmiller, R. Hadsell,
534 and P. Battaglia. Graph networks as learnable physics engines for inference and control. In
535 *International Conference on Machine Learning*, pages 4470–4479. PMLR, 2018.
- 536 [90] Y. Li, J. Wu, R. Tedrake, J. B. Tenenbaum, and A. Torralba. Learning particle dynamics for
537 manipulating rigid bodies, deformable objects, and fluids. *arXiv preprint arXiv:1810.01566*,
538 2018.
- 539 [91] C. Finn, S. Levine, and P. Abbeel. Guided cost learning: Deep inverse optimal control via
540 policy optimization. In *International conference on machine learning*, pages 49–58. PMLR,
541 2016.
- 542 [92] J. Fu, K. Luo, and S. Levine. Learning robust rewards with adversarial inverse reinforcement
543 learning. *arXiv preprint arXiv:1710.11248*, 2017.
- 544 [93] D. Driess, O. Oguz, J.-S. Ha, and M. Toussaint. Deep visual heuristics: Learning feasibility of
545 mixed-integer programs for manipulation planning. In *2020 IEEE International Conference*
546 *on Robotics and Automation (ICRA)*, pages 9563–9569. IEEE, 2020.
- 547 [94] B. Amos, I. Jimenez, J. Sacks, B. Boots, and J. Z. Kolter. Differentiable mpc for end-to-end
548 planning and control. *Advances in neural information processing systems*, 31, 2018.
- 549 [95] M. Mittal, D. Hoeller, F. Farshidian, M. Hutter, and A. Garg. Articulated object interaction
550 in unknown scenes with whole-body mobile manipulation. In *2022 IEEE/RSJ International*
551 *Conference on Intelligent Robots and Systems (IROS)*, pages 1647–1654. IEEE, 2022.
- 552 [96] S. Bahl, A. Gupta, and D. Pathak. Human-to-robot imitation in the wild. *arXiv preprint*
553 *arXiv:2207.09450*, 2022.
- 554 [97] S. Bahl, R. Mendonca, L. Chen, U. Jain, and D. Pathak. Affordances from human videos
555 as a versatile representation for robotics. In *Proceedings of the IEEE/CVF Conference on*
556 *Computer Vision and Pattern Recognition*, pages 13778–13790, 2023.

- 557 [98] C. Wang, L. Fan, J. Sun, R. Zhang, L. Fei-Fei, D. Xu, Y. Zhu, and A. Anandkumar.
558 Mimicplay: Long-horizon imitation learning by watching human play. *arXiv preprint*
559 *arXiv:2302.12422*, 2023.
- 560 [99] H. Bharadhwaj, A. Gupta, S. Tulsiani, and V. Kumar. Zero-shot robot manipulation from
561 passive human videos. *arXiv preprint arXiv:2302.02011*, 2023.
- 562 [100] Y. J. Ma, S. Sodhani, D. Jayaraman, O. Bastani, V. Kumar, and A. Zhang. Vip: Towards
563 universal visual reward and representation via value-implicit pre-training. *arXiv preprint*
564 *arXiv:2210.00030*, 2022.
- 565 [101] D. Damen, H. Doughty, G. M. Farinella, S. Fidler, A. Furnari, E. Kazakos, D. Moltisanti,
566 J. Munro, T. Perrett, W. Price, et al. Scaling egocentric vision: The epic-kitchens dataset. In
567 *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 720–736, 2018.
- 568 [102] K. Grauman, A. Westbury, E. Byrne, Z. Chavis, A. Furnari, R. Girdhar, J. Hamburger,
569 H. Jiang, M. Liu, X. Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video.
570 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
571 pages 18995–19012, 2022.
- 572 [103] Y. Cui, S. Niekum, A. Gupta, V. Kumar, and A. Rajeswaran. Can foundation models perform
573 zero-shot task specification for robot manipulation? In *Learning for Dynamics and Control*
574 *Conference*, pages 893–905. PMLR, 2022.
- 575 [104] T. Yu, T. Xiao, A. Stone, J. Tompson, A. Brohan, S. Wang, J. Singh, C. Tan, J. Peralta,
576 B. Ichter, et al. Scaling robot learning with semantically imagined experience. *arXiv preprint*
577 *arXiv:2302.11550*, 2023.
- 578 [105] Z. Mandi, H. Bharadhwaj, V. Moens, S. Song, A. Rajeswaran, and V. Kumar. Cacti: A
579 framework for scalable multi-task multi-scene visual imitation learning. *arXiv preprint*
580 *arXiv:2212.05711*, 2022.
- 581 [106] T. Xiao, H. Chan, P. Sermanet, A. Wahid, A. Brohan, K. Hausman, S. Levine, and J. Tompson.
582 Robotic skill acquisition via instruction augmentation with vision-language models. *arXiv*
583 *preprint arXiv:2211.11736*, 2022.
- 584 [107] C. Wang, D. Xu, and L. Fei-Fei. Generalizable task planning through representation pretrain-
585 ing. *IEEE Robotics and Automation Letters*, 7(3):8299–8306, 2022.
- 586 [108] C. Li, R. Zhang, J. Wong, C. Gokmen, S. Srivastava, R. Martín-Martín, C. Wang, G. Levine,
587 M. Lingelbach, J. Sun, et al. Behavior-1k: A benchmark for embodied ai with 1,000 everyday
588 activities and realistic simulation. In *Conference on Robot Learning*, pages 80–93. PMLR,
589 2023.
- 590 [109] O. Khatib. A unified approach for motion and force control of robot manipulators: The
591 operational space formulation. *IEEE Journal on Robotics and Automation*, 3(1):43–53, 1987.
- 592 [110] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead,
593 A. C. Berg, W.-Y. Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.
- 594 [111] H. K. Cheng and A. G. Schwing. Xmem: Long-term video object segmentation with an
595 atkinson-shiffrin memory model. In *Computer Vision–ECCV 2022: 17th European Con-*
596 *ference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXVIII*, pages 640–658.
597 Springer, 2022.
- 598 [112] F. Xiang, Y. Qin, K. Mo, Y. Xia, H. Zhu, F. Liu, M. Liu, H. Jiang, Y. Yuan, H. Wang, et al.
599 Sapien: A simulated part-based interactive environment. In *Proceedings of the IEEE/CVF*
600 *Conference on Computer Vision and Pattern Recognition*, pages 11097–11107, 2020.

- 601 [113] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image
602 segmentation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI*
603 *2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings,*
604 *Part III 18*, pages 234–241. Springer, 2015.
- 605 [114] Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger. 3d u-net: learn-
606 ing dense volumetric segmentation from sparse annotation. In *Medical Image Computing*
607 *and Computer-Assisted Intervention–MICCAI 2016: 19th International Conference, Athens,*
608 *Greece, October 17-21, 2016, Proceedings, Part II 19*, pages 424–432. Springer, 2016.
- 609 [115] L. E. Kavvaki, P. Svestka, J.-C. Latombe, and M. H. Overmars. Probabilistic roadmaps for
610 path planning in high-dimensional configuration spaces. *IEEE transactions on Robotics and*
611 *Automation*, 12(4):566–580, 1996.
- 612 [116] N. D. Ratliff, J. Issac, D. Kappler, S. Birchfield, and D. Fox. Riemannian motion policies.
613 *arXiv preprint arXiv:1801.02854*, 2018.
- 614 [117] T. Marcucci, M. Petersen, D. von Wrangel, and R. Tedrake. Motion planning around obstacles
615 with convex optimization. *arXiv preprint arXiv:2205.04422*, 2022.
- 616 [118] J. Li, D. Li, S. Savarese, and S. Hoi. Blip-2: Bootstrapping language-image pre-training with
617 frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023.
- 618 [119] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agar-
619 wal, K. Slama, A. Ray, et al. Training language models to follow instructions with human
620 feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- 621 [120] Y. Bai, S. Kadavath, S. Kundu, A. Askell, J. Kernion, A. Jones, A. Chen, A. Goldie, A. Mirho-
622 seini, C. McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint*
623 *arXiv:2212.08073*, 2022.
- 624 [121] J. Wei, X. Wang, D. Schuurmans, M. Bosma, E. Chi, Q. Le, and D. Zhou. Chain of thought
625 prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*,
626 2022.
- 627 [122] Y. Wang, Y. Kordi, S. Mishra, A. Liu, N. A. Smith, D. Khashabi, and H. Hajishirzi.
628 Self-instruct: Aligning language model with self generated instructions. *arXiv preprint*
629 *arXiv:2212.10560*, 2022.
- 630 [123] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa. Large language models are zero-
631 shot reasoners. *arXiv preprint arXiv:2205.11916*, 2022.
- 632 [124] S. Yao, D. Yu, J. Zhao, I. Shafran, T. L. Griffiths, Y. Cao, and K. Narasimhan. Tree
633 of thoughts: Deliberate problem solving with large language models. *arXiv preprint*
634 *arXiv:2305.10601*, 2023.
- 635 [125] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma,
636 D. Zhou, D. Metzler, et al. Emergent abilities of large language models. *arXiv preprint*
637 *arXiv:2206.07682*, 2022.
- 638 [126] T. Gupta and A. Kembhavi. Visual programming: Compositional visual reasoning with-
639 out training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
640 *Recognition*, pages 14953–14962, 2023.
- 641 [127] D. Surís, S. Menon, and C. Vondrick. Vipergpt: Visual inference via python execution for
642 reasoning. *arXiv preprint arXiv:2303.08128*, 2023.
- 643 [128] Y. Zhu, A. Joshi, P. Stone, and Y. Zhu. Viola: Imitation learning for vision-based manipulation
644 with object proposal priors. *6th Annual Conference on Robot Learning*, 2022.

645 **A Appendix**

646 **A.1 Emergent Behavioral Capabilities**

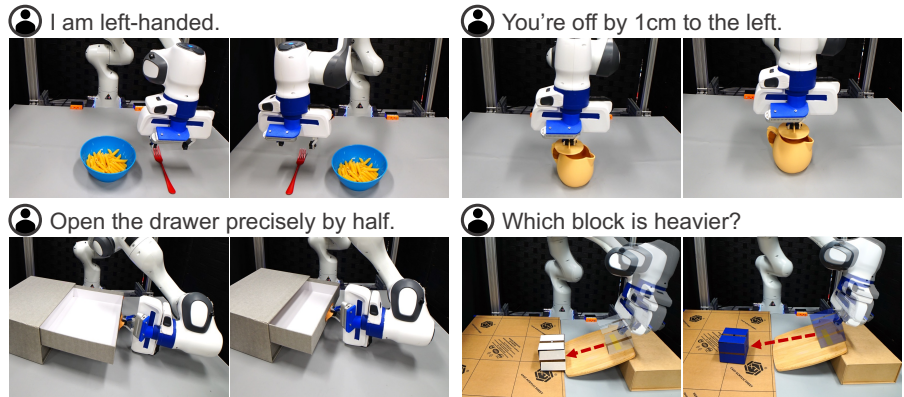


Figure 5: Emergent behavioral capabilities by VoxPoser inherited from the language model, including behavioral commonsense reasoning (**top left**), fine-grained language correction (**top right**), multi-step visual program (**bottom left**), and estimating physical properties of objects (**bottom right**).

647 Emergent capabilities refer to unpredictable phenomena that are only present in large mod-
648 els [125]. As VoxPoser uses pre-trained LLMs as backbone, we observe similar embodied emergent
649 capabilities driven by the rich world knowledge of LLMs. In particular, we focus our study on the
650 *behavioral capabilities* that are unique to VoxPoser. We observe the following capabilities:

- 651 • **Behavioral Commonsense Reasoning:** During a task where robot is setting the table, the
652 user can specify behavioral preferences such as “I am left-handed”, which requires the
653 robot to comprehend its meaning in the context of the task. VoxPoser decides that it should
654 move the fork from the right side of the bowl to the left side.
- 655 • **Fine-grained Language Correction:** For tasks that require high precision such as “cov-
656 ering the teapot with the lid”, the user can give precise instructions to the robot such as
657 “you’re off by 1cm”. VoxPoser similarly adjusts its action based on the feedback.
- 658 • **Multi-step Visual Program** [126, 127]: Given a task “open the drawer precisely by half”
659 where there is insufficient information because object models are not available, VoxPoser
660 can come up with multi-step manipulation strategies based on visual feedback that first
661 opens the drawer fully while recording handle displacement, then close it back to the mid-
662 point to satisfy the requirement.
- 663 • **Estimating Physical Properties:** Given two blocks of unknown mass, the robot is tasked
664 to conduct physics experiments using an existing ramp to determine which block is heav-
665 ier. VoxPoser decides to push both blocks off the ramp and choose the block traveling the
666 farthest as the heavier block. Interestingly, this mirrors a common human oversight: in
667 an ideal, frictionless world, both blocks would traverse the same distance under the influ-
668 ence of gravity. This serves as a lighthearted example that language models can exhibit
669 limitations similar to human reasoning.

670 A.2 APIs for VoxPoser

671 Central to VoxPoser is an LLM generating Python code that is executed by a Python interpreter.
672 Besides exposing NumPy [16] and the [Transforms3d](#) library to the LLM, we provide the following
673 environment APIs that LLMs can choose to invoke:

674 **detect(obj_name)**: Takes in an object name and returns a list of dictionaries, where each dictionary
675 corresponds to one instance of the matching object, containing center position, occupancy grid, and
676 mean normal vector.

677 **execute(movable, affordance_map, avoidance_map, rotation_map, velocity_map, gripper_map)**:
678 Takes in an “entity of interest” as “movable” (a dictionary returned by detect) and (optionally)
679 a list of value maps and invokes the motion planner to execute the trajectory. Note that in MPC
680 settings, “movable” and the input value maps are functions that can be re-evaluated to reflect the
681 latest environment observation.

682 **cm2index(cm, direction)**: Takes in a desired offset distance in centimeters along direction and
683 returns 3-dim vector reflecting displacement in voxel coordinates.

684 **index2cm(index, direction)**: Inverse of cm2index. Takes in an integer “index” and a “direction”
685 vector and returns the distance in centimeters in world coordinates displaced by the “integer” in
686 voxel coordinates.

687 **pointat2quat(vector)**: Takes in a desired pointing direction for the end-effector and returns a
688 satisfying target quaternion.

689 **set_voxel_by_radius(voxel_map, voxel_xyz, radius_cm, value)**: Assigns “value” to voxels
690 within “radius_cm” from “voxel_xyz” in “voxel_map”.

691 **get_empty_affordance_map()**: Returns a default affordance map initialized with 0, where a high
692 value attracts the entity.

693 **get_empty_avoidance_map()**: Returns a default avoidance map initialized with 0, where a high
694 value repulses the entity.

695 **get_empty_rotation_map()**: Returns a default rotation map initialized with current end-effector
696 quaternion.

697 **get_empty_gripper_map()**: Returns a default gripper map initialized with current gripper action,
698 where 1 indicates “closed” and 0 indicates “open”.

699 **get_empty_velocity_map()**: Returns a default affordance map initialized with 1, where the number
700 represents scale factor (e.g., 0.5 for half of the default velocity).

701 **reset_to_default_pose()**: Reset to robot rest pose.

702 A.3 Real-World Environment Setup

703 We use a Franka Emika Panda robot with a tabletop setup. We use Operational Space Controller
704 with impedance from Deoxys [128]. We mount two RGB-D cameras (Azure Kinect) at two opposite
705 ends of the table: bottom right and top left from the top down view. At the start of each rollout, both
706 cameras start recording and return the real-time RGB-D observations at 20 Hz.

707 For each task, we evaluate each method on two settings: without and with disturbances. For tasks
708 with disturbances, we apply three kinds of disturbances to the environment, which we pre-select a
709 sequence of them at the start of the evaluation: 1) random forces applied to the robot, 2) random
710 displacement of task-relevant and distractor objects, and 3) reverting task progress (e.g., pull drawer
711 open while it’s being closed by the robot). We only apply the third disturbances to tasks where
712 “entity of interest” is an object or object part.

713 We compare to a variant of Code as Policies [75] as a baseline that uses an LLM with ac-
714 tion primitives. The primitives include: `move_to_pos`, `rotate_by_quat`, `set_vel`, `open_gripper`,
715 `close_gripper`. We do not provide primitives such as pick-and-place as they would be tailored
716 for a particular suite of tasks that we do not constrain to in our study (similar to the control APIs
717 for VoxPoser specified in Sec. A.2).

718 A.3.1 Tasks

719 **Move & Avoid:** “Move to the top of [obj1] while staying away from [obj2]”, where [obj1] and [obj2]
720 are randomized everyday objects selected from the list: apple, banana, yellow bowl, headphones,
721 mug, wood block.

722 **Set Up Table:** “Please set up the table by placing utensils for my pasta”.

723 **Close Drawer:** “Close the [deixis] drawer”, where [deixis] can be “top” or “bottom”.

724 **Open Bottle:** “Turn open the vitamin bottle”.

725 **Sweep Trash:** “Please sweep the paper trash into the blue dustpan”.

726 A.4 Simulated Environment Setup

727 We implement a tabletop manipulation environment with a Franka Emika Panda robot in
728 SAPIEN [112]. The controller takes as input a desired end-effector 6-DoF pose, calculates a se-
729 quence of interpolated waypoints using inverse kinematics, and finally follows the waypoints using
730 a PD controller. We use a set of 10 colored blocks and 10 colored lines in addition to an articulated
731 cabinet with 3 drawers. They are initialized differently depending on the specific task. The lines are
732 used as visual landmarks and are not interactable. For perception, a total of 4 RGB-D cameras are
733 mounted at each end of the table pointing at the center of the workspace.

734 A.4.1 Tasks

735 We create a custom suite of 13 tasks shown in Table 4. Each task comes with a templated instruction
736 (shown in Table 4) where there may be one or multiple attributes randomized from the pre-defined
737 list below. At reset time, a number of objects are selected (depending on the specific task) and are
738 randomized across the workspace while making sure that task is not completed at reset and that task
739 completion is feasible. A complete list of attributes can be found below, divided into “seen” and
740 “unseen” categories:

741 Seen Attributes:

- 742 • **[pos]**: [“back left corner of the table”, “front right corner of the table”, “right side of the
743 table”, “back side of the table”]
- 744 • **[obj]**: [“blue block”, “green block”, “yellow block”, “pink block”, “brown block”]
- 745 • **[preposition]**: [“left of”, “front side of”, “top of”]
- 746 • **[deixis]**: [“topmost”, “second to the bottom”]
- 747 • **[dist]**: [3, 5, 7, 9, 11]
- 748 • **[region]**: [“right side of the table”, “back side of the table”]
- 749 • **[velocity]**: [“faster speed”, “a quarter of the speed”]
- 750 • **[line]**: [“blue line”, “green line”, “yellow line”, “pink line”, “brown line”]

751 Unseen Attributes:

- 752 • **[pos]**: [“back right corner of the table”, “front left corner of the table”, “left side of the
753 table”, “front side of the table”]
- 754 • **[obj]**: [“red block”, “orange block”, “purple block”, “cyan block”, “gray block”]
- 755 • **[preposition]**: [“right of”, “back side of”]
- 756 • **[deixis]**: [“bottommost”, “second to the top”]
- 757 • **[dist]**: [4, 6, 8, 10]
- 758 • **[region]**: [“left side of the table”, “front side of the table”]
- 759 • **[velocity]**: [“slower speed”, “3x speed”]
- 760 • **[line]**: [“red line”, “orange line”, “purple line”, “cyan line”, “gray line”]

761 **A.4.2 Full Results on Simulated Environments**

Tasks	U-Net + MP		LLM + Prim.		VoxPoser	
	SA	UA	SA	UA	SA	UA
move to the [preposition] the [obj]	95.0%	0.0%	85.0%	60.0%	90.0%	55.0%
move to the [pos] while staying on the [preposition] the [obj]	100.0%	10.0%	80.0%	30.0%	95.0%	50.0%
move to the [pos] while moving at [velocity] when within [dist]cm from the obj	80.0%	0.0%	10.0%	0.0%	100.0%	95.0%
close the [deixis] drawer by pushing	0.0%	0.0%	60.0%	60.0%	80.0%	80.0%
push the [obj] along the [line]	0.0%	0.0%	0.0%	0.0%	65.0%	30.0%
grasp the [obj] from the table at [velocity]	35.0%	0.0%	75.0%	70.0%	65.0%	40.0%
drop the [obj] to the [pos]	70.0%	10.0%	60.0%	100.0%	60.0%	100.0%
push the [obj] while letting it stay on [region]	0.0%	5.0%	10.0%	0.0%	50.0%	50.0%
move to the [region]	5.0%	0.0%	100.0%	95.0%	100.0%	100.0%
move to the [pos] while staying at least [dist]cm from the [obj]	0.0%	0.0%	15.0%	20.0%	85.0%	90.0%
move to the [pos] while moving at [velocity] in the [region]	0.0%	0.0%	90.0%	45.0%	85.0%	85.0%
push the [obj] to the [pos] while staying away from [obstacle]	0.0%	0.0%	0.0%	10.0%	45.0%	55.0%
push the [obj] to the [pos]	0.0%	0.0%	20.0%	25.0%	80.0%	75.0%

Table 4: Full experimental results in simulation on [seen tasks](#) and [unseen tasks](#). “SA” indicates seen attributes and “UA” indicates unseen attributes. Each entry represents success rate averaged across 20 episodes.

762 A.5 Prompts

763 Prompts used in Sec. 4.1 and Sec. 4.2 can be found below.

764 **planner**: Takes in a user instruction \mathcal{L} and generates a sequence of sub-tasks ℓ_i which is fed into
765 “composer” (Note that planner is not used in simulation as the evaluated tasks consist of a single
766 manipulation phase).

767 real-world: [voxboser-anon.github.io/prompts/real_planner_prompt.txt](https://github.com/voxboser-anon/prompts/real_planner_prompt.txt).

768 **composer**: Takes in sub-task instruction ℓ_i and invokes necessary value map LMPs to compose
769 affordance maps and constraint maps.

770 simulation: [voxboser-anon.github.io/prompts/sim_composer_prompt.txt](https://github.com/voxboser-anon/prompts/sim_composer_prompt.txt).

771 real-world: [voxboser-anon.github.io/prompts/real_composer_prompt.txt](https://github.com/voxboser-anon/prompts/real_composer_prompt.txt).

772 **parse_query_obj**: Takes in a text query of object/part name and returns a list of dictionaries, where
773 each dictionary corresponds to one instance of the matching object containing center position, oc-
774 cupancy grid, and mean normal vector.

775 simulation: [voxboser-anon.github.io/prompts/sim_parse_query_obj_prompt.txt](https://github.com/voxboser-anon/prompts/sim_parse_query_obj_prompt.txt).

776 real-world: [voxboser-anon.github.io/prompts/real_parse_query_obj_prompt.txt](https://github.com/voxboser-anon/prompts/real_parse_query_obj_prompt.txt).

777 **get_affordance_map**: Takes in natural language parametrization from composer and returns a
778 NumPy array for task affordance map.

779 simulation: [voxboser-anon.github.io/prompts/sim_get_affordance_map_prompt.txt](https://github.com/voxboser-anon/prompts/sim_get_affordance_map_prompt.txt).

780 real-world: [voxboser-anon.github.io/prompts/real_get_affordance_map_prompt.txt](https://github.com/voxboser-anon/prompts/real_get_affordance_map_prompt.txt).

781 **get_avoidance_map**: Takes in natural language parametrization from composer and returns a
782 NumPy array for task avoidance map.

783 simulation: [voxboser-anon.github.io/prompts/sim_get_avoidance_map_prompt.txt](https://github.com/voxboser-anon/prompts/sim_get_avoidance_map_prompt.txt).

784 real-world: [voxboser-anon.github.io/prompts/real_get_avoidance_map_prompt.txt](https://github.com/voxboser-anon/prompts/real_get_avoidance_map_prompt.txt).

785 **get_rotation_map**: Takes in natural language parametrization from composer and returns a NumPy
786 array for end-effector rotation map.

787 simulation: [voxboser-anon.github.io/prompts/sim_get_rotation_map_prompt.txt](https://github.com/voxboser-anon/prompts/sim_get_rotation_map_prompt.txt).

788 real-world: [voxboser-anon.github.io/prompts/real_get_rotation_map_prompt.txt](https://github.com/voxboser-anon/prompts/real_get_rotation_map_prompt.txt).

789 **get_gripper_map**: Takes in natural language parametrization from composer and returns a NumPy
790 array for gripper action map.

791 simulation: [voxboser-anon.github.io/prompts/sim_get_gripper_map_prompt.txt](https://github.com/voxboser-anon/prompts/sim_get_gripper_map_prompt.txt).

792 real-world: [voxboser-anon.github.io/prompts/real_get_gripper_map_prompt.txt](https://github.com/voxboser-anon/prompts/real_get_gripper_map_prompt.txt).

793 **get_velocity_map**: Takes in natural language parametrization from composer and returns a NumPy
794 array for end-effector velocity map.

795 simulation: [voxboser-anon.github.io/prompts/sim_get_velocity_map_prompt.txt](https://github.com/voxboser-anon/prompts/sim_get_velocity_map_prompt.txt).

796 real-world: [voxboser-anon.github.io/prompts/real_get_velocity_map_prompt.txt](https://github.com/voxboser-anon/prompts/real_get_velocity_map_prompt.txt).