

## 583 A Appendix

### 584 A.1 APIs for VoxPoser

585 Central to VoxPoser is an LLM generating Python code that is executed by a Python interpreter. Besides exposing  
586 NumPy package to the LLM, we provide the following environment APIs that LLMs can choose to invoke:

587 **detect(\*obj\_names)**: Takes in a (list of) object name(s) and returns a list of dictionaries, where each  
588 dictionary corresponds to one instance of the matching object containing center position, occupancy grid,  
589 and mean normal vector.

590 **execute(movable, affordance\_map, avoidance\_map, rotation\_map, velocity\_map, gripper\_map)**: Takes in “movable” (a dictionary returned by `detect`) and (optionally) a list of value maps and  
591 invokes the motion planner to execute the trajectory. Note that in MPC settings, “movable” and the input  
592 value maps are functionals that can be invoked to reflect the latest environment observation.

594 **index2cm(index, direction)**: Takes in an integer and a direction vector and returns the distance in  
595 centimeters in world coordinates displaced by the integer in voxel coordinates.

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

598 **set\_voxel\_by\_radius(voxel\_map, voxel\_xyz, radius\_cm, value)**: Assigns “value” to voxels within  
599 “radius\_cm” from “voxel\_xyz”.

600 **get\_empty\_affordance\_map()**: Returns a default affordance map initialized with 0.

601 **get\_empty\_avoidance\_map()**: Returns a default affordance map initialized with 0.

602 **get\_empty\_rotation\_map()**: Returns a default affordance map initialized with current end-effector  
603 quaternion.

604 **get\_empty\_gripper\_map()**: Returns a default affordance map initialized with current gripper action.

605 **get\_empty\_velocity\_map()**: Returns a default affordance map initialized with 1.

606 **reset\_to\_default\_pose()**: Reset to rest pose.

### 607 A.2 Real-World Environment Setup

608 We use a Franka Emika Panda robot with a tabletop setup. We use Operational Space Controller with  
609 impedance from Deoxys [114]. We mount two RGB-D cameras (Azure Kinect) at two opposite ends of  
610 the table: bottom right and top left from the top down view. At the start of each rollout, both cameras start  
611 recording and return the real-time RGB-D observations at 20 Hz.

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

617 We compare to a variant of Code as Policies [71] as a baseline that an LLM with action primitives. The  
618 primitives include: `move_to_pos`, `rotate_by_quat`, `set_vel`, `open_gripper`, `close_gripper`. We do  
619 not provide primitives such as pick-and-place as they would be tailored for a particular suite of tasks that  
620 we do not constrain to in our study (similar to the control APIs for VoxPoser).

#### 621 A.2.1 Tasks

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

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

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

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

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

### 628 **A.3 Simulated Environment Setup**

629 We implement a tabletop manipulation environment with a Franka Emika Panda robot in SAPIEN [115]. The  
630 controller takes as input a desired end-effector 6-DoF pose, calculates a sequence of interpolated waypoints  
631 using inverse kinematics, and finally follows the waypoints using a PD controller. We use a set of 10 colored  
632 blocks and 10 colored lines in addition to an articulated cabinet with 3 drawers. They are initialized differently  
633 depending on the specific task. The lines are used as visual landmarks and are not interactable. For perception,  
634 a total of 4 RGB-D cameras are mounted at each end of the table pointing at the center of the workspace.

#### 635 **A.3.1 Tasks**

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

##### 641 **Seen Attributes:**

642 **[pos]:** [“back left corner of the table”, “front right corner of the table”, “right side of the table”, “back side  
643 of the table”]

644 **[obj]:** [“blue block”, “green block”, “yellow block”, “pink block”, “brown block”]

645 **[preposition]:** [“left of”, “front side of”, “top of”]

646 **[deixis]:** [“topmost”, “second to the bottom”]

647 **[dist]:** [3, 5, 7, 9, 11]

648 **[region]:** [“right side of the table”, “back side of the table”]

649 **[velocity]:** [“faster speed”, “a quarter of the speed”]

650 **[line]:** [“blue line”, “green line”, “yellow line”, “pink line”, “brown line”]

##### 651 **Unseen Attributes:**

652 **[pos]:** [“back right corner of the table”, “front left corner of the table”, “left side of the table”, “front side  
653 of the table”]

654 **[obj]:** [“red block”, “orange block”, “purple block”, “cyan block”, “gray block”]

655 **[preposition]:** [“right of”, “back side of”]

656 **[deixis]:** [“bottommost”, “second to the top”]

657 **[dist]:** [4, 6, 8, 10]

658 **[region]:** [“left side of the table”, “front side of the table”]

659 **[velocity]:** [“slower speed”, “3x speed”]

660 **[line]:** [“red line”, “orange line”, “purple line”, “cyan line”, “gray line”]

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. “SA” indicates seen attributes and “UA” indicates unseen attributes.

## 662 A.4 Prompts

663 **planner**: Takes in a user instruction  $\mathcal{L}$  and generates a sequence of sub-tasks  $\ell_i$  which is fed into “composer”.

664 simulation: [voxposer-anon.github.io/sim\\_planner\\_prompt.txt](https://voxposer-anon.github.io/sim_planner_prompt.txt).

665 real-world: [voxposer-anon.github.io/real\\_planner\\_prompt.txt](https://voxposer-anon.github.io/real_planner_prompt.txt).

666 **composer**: Takes in sub-task instruction  $\ell_i$  and invokes necessary value map LMPs to compose affordance  
667 maps and constraint maps.

668 simulation: [voxposer-anon.github.io/sim\\_composer\\_prompt.txt](https://voxposer-anon.github.io/sim_composer_prompt.txt).

669 real-world: [voxposer-anon.github.io/real\\_composer\\_prompt.txt](https://voxposer-anon.github.io/real_composer_prompt.txt).

670 **parse\_query\_obj**: Takes in a text query of object/part name and returns a list of dictionaries, where each  
671 dictionary corresponds to one instance of the matching object containing center position, occupancy grid,  
672 and mean normal vector.

673 simulation: [voxposer-anon.github.io/sim\\_parse\\_query\\_obj\\_prompt.txt](https://voxposer-anon.github.io/sim_parse_query_obj_prompt.txt).

674 real-world: [voxposer-anon.github.io/real\\_parse\\_query\\_obj\\_prompt.txt](https://voxposer-anon.github.io/real_parse_query_obj_prompt.txt).

675 **get\_target\_map**: Takes in natural language parametrization from composer and returns a NumPy array  
676 for task affordances.

677 simulation: [voxposer-anon.github.io/sim\\_get\\_target\\_map\\_prompt.txt](https://voxposer-anon.github.io/sim_get_target_map_prompt.txt).

678 real-world: [voxposer-anon.github.io/real\\_get\\_target\\_map\\_prompt.txt](https://voxposer-anon.github.io/real_get_target_map_prompt.txt).

679 **get\_avoidance\_map**: Takes in natural language parametrization from composer and returns a NumPy array  
680 for end-effector rotation.

681 simulation: [voxposer-anon.github.io/sim\\_get\\_avoidance\\_map\\_prompt.txt](https://voxposer-anon.github.io/sim_get_avoidance_map_prompt.txt).

682 real-world: [voxposer-anon.github.io/real\\_get\\_avoidance\\_map\\_prompt.txt](https://voxposer-anon.github.io/real_get_avoidance_map_prompt.txt).

683 **get\_rotation\_map**: Takes in natural language parametrization from composer and returns a NumPy array  
684 for end-effector rotation.

685 simulation: [voxposer-anon.github.io/sim\\_get\\_rotation\\_map\\_prompt.txt](https://voxposer-anon.github.io/sim_get_rotation_map_prompt.txt).

686 real-world: [voxposer-anon.github.io/real\\_get\\_rotation\\_map\\_prompt.txt](https://voxposer-anon.github.io/real_get_rotation_map_prompt.txt).

687 **get\_gripper\_map**: Takes in natural language parametrization from composer and returns a NumPy array  
688 for gripper action.

689 simulation: [voxposer-anon.github.io/sim\\_get\\_gripper\\_map\\_prompt.txt](https://voxposer-anon.github.io/sim_get_gripper_map_prompt.txt).

690 real-world: [voxposer-anon.github.io/real\\_get\\_gripper\\_map\\_prompt.txt](https://voxposer-anon.github.io/real_get_gripper_map_prompt.txt).

691 **get\_velocity\_map**: Takes in natural language parametrization from composer and returns a NumPy array  
692 for end-effector velocity.

693 simulation: [voxposer-anon.github.io/sim\\_get\\_velocity\\_map\\_prompt.txt](https://voxposer-anon.github.io/sim_get_velocity_map_prompt.txt).

694 real-world: [voxposer-anon.github.io/real\\_get\\_velocity\\_map\\_prompt.txt](https://voxposer-anon.github.io/real_get_velocity_map_prompt.txt).