

485 **Appendix**

486 **Imitation Learning and RT-1**

487 MOO builds upon a language-conditioned imitation learning setup. The goal of language-  
 488 conditioned imitation learning is to learn a policy  $\pi(a | \ell, o)$ , where  $a$  is a robot action that should be  
 489 applied given the current observation  $o$  and task instruction  $\ell$ . To learn a language-conditioned policy  
 490  $\pi$ , we build on top of RT-1 [24], a recent robotics transformer-based model that achieves high lev-  
 491 els of performance across a wide variety of manipulation tasks. RT-1 uses behavioral cloning [53],  
 492 which optimizes  $\pi$  by minimizing the negative log-likelihood of an action  $a$  given the image ob-  
 493 servations seen so far in the trajectory and the language instruction, using a demonstration dataset  
 494 containing  $N$  demonstrations:

$$J(\pi) := \sum_{n=1}^N \sum_{t=1}^{T^{(n)}} \log \pi(a_t^{(n)} | \ell^{(n)}, \{o_j^{(n)}\}_{j=1}^t). \quad (1)$$

495 **Vision-Language Models**

496 In recent years, there has been a growing interest in developing models that can detect objects  
 497 in images based on natural language queries. These models, known as vision-language models  
 498 (VLMs), are enabling detectors to identify a wide range of objects based on natural language queries.  
 499 Typically the text queries are tokenized and embedded in a high-dimensional space by a pre-trained  
 500 language encoder, and the image is processed by a separate network to extract image features into  
 501 the same embedding space as the text features. The language and image representations are then  
 502 combined to make predictions of the bounding boxes and segmentation masks. Given a natural  
 503 language query,  $q$ , and an image observation on which to run detection,  $o$ , these models aim to  
 504 produce a set of embeddings for the image  $f_i(o)$  with shape (height, width, feature dim) and an  
 505 embedding of the language query  $f_l(q)$  with shape feature dim such that logits =  $f_i(o) \cdot f_l(q)$  gives  
 506 a logit score map and is maximized at regions in  $o$  which correspond to the queries in  $q$ . Each  
 507 image embedding location within  $f_i(o)$  is also associated with a predicted bounding box or mask  
 508 indicating the spatial extent of that object corresponding to  $f_i(o)$ . In this work, we use the Owl-ViT  
 509 detector [54], which we discuss further in Sec. 3.4.

510 **Datasets**

511 We collect a focused collection of teleoperated demonstration data that focuses on increasing object  
 512 diversity for the most efficient skill to collect data for, the picking task. Detailed dataset statistics  
 513 across objects are shown in Appendix Figure 9.

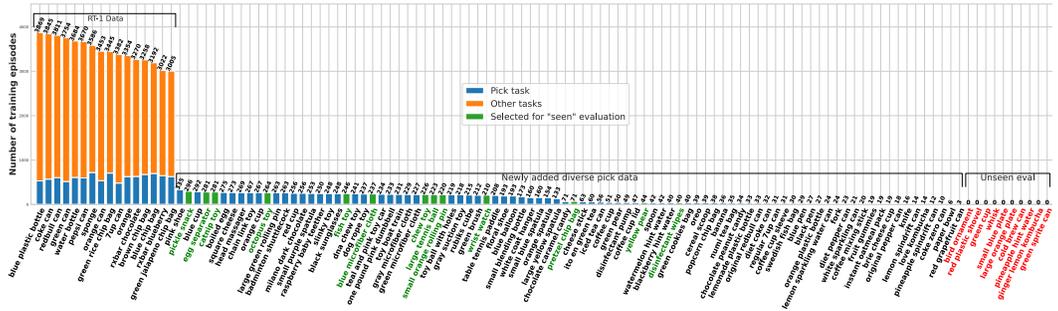


Figure 9: Distribution of training objects for “pick” episodes and other skills. The data on the left was what was used by [24]. We augmented RT-1 data with a large number of diverse pick episodes in order to demonstrate strong generalization to unseen objects. Blue and green bars represent “pick” episodes and orange bars represent other tasks like “move near” or “knock.” “Green” bars were the objects we randomly selected for “seen” evaluations. All randomly selected “unseen” objects are shown to the right.

514 **Experiments**

515 We show a visualization of our 7-DoF manipulation robot in Figure 10.

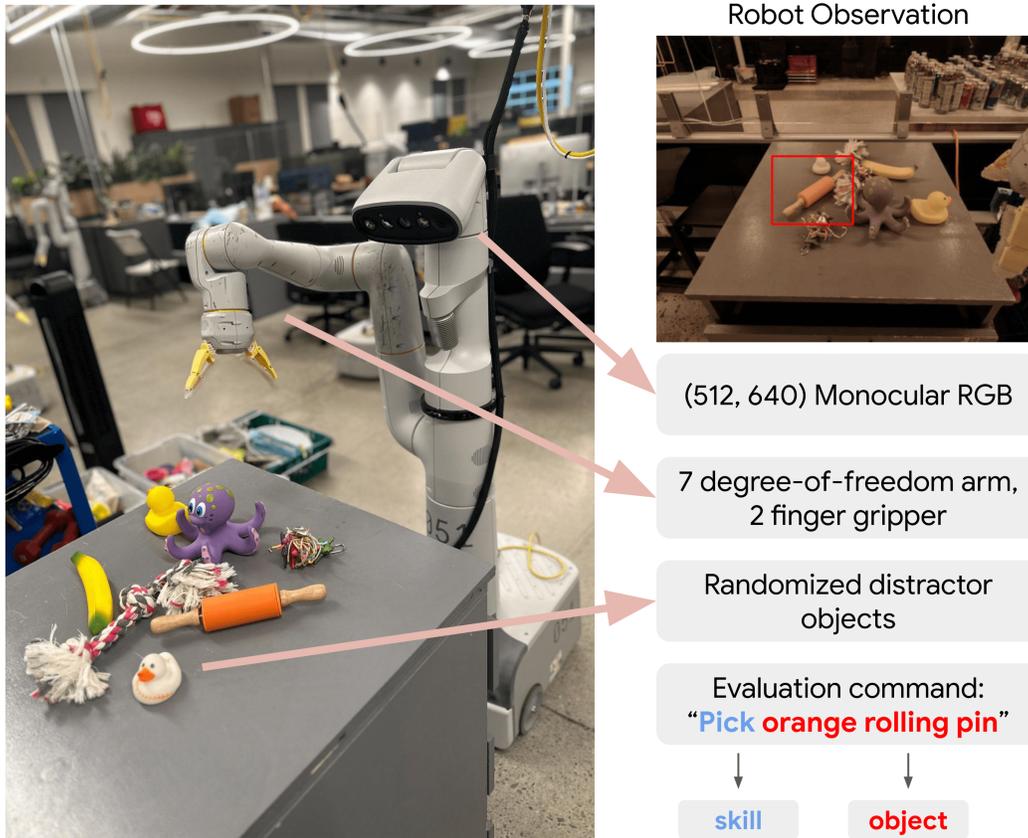


Figure 10: Image of our robot hardware and evaluation setting.

516 **Skills.** Our experiments evaluate the percent of successfully completed manipulation commands  
 517 which include five skills: “pick”, “move near”, “knock,” “place upright,” and “place into” across  
 518 a set of evaluation episodes. The definition of the tasks follows RT-1 [24]: For “pick”, success is  
 519 defined as (1) grasping the specified object and (2) lifting the object at least 6 inches from the table  
 520 top. For “move near”, success is defined as (1) grasping the specified object and (2) placing it within  
 521 6 inches of the specified target object. For “knock”, success is defined as placing the specified object  
 522 from an “upright” position onto its side. “Place upright” tasks are the inverse of “knock” and involve  
 523 placing an object from its side into an upright position. Finally, “place into” tasks involve placing  
 524 one object into another, such as an apple into a bowl.

525 **Robustness evaluation details.** We evaluate the robustness of MOO on a variety of visually chal-  
 526 lenging scenarios with drastically different furniture and backgrounds, as shown in Figure 6; the  
 527 results are reported in Figure 5. The first set of these difficult evaluation scenes introduces six evalu-  
 528 ations across five additional open-world objects that correspond to various household items that have  
 529 not been seen at any point during training. The second set of difficult scenes introduces 14 evalu-  
 530 ations across two patterned tablecloths; these tablecloth textures are significantly more challenging  
 531 than the plain gray counter-tops seen in the training demonstration dataset. Finally, the last set of  
 532 difficult scenes include 14 evaluations across three new environments in natural kitchen and office  
 533 spaces that were never present training. These new scenes simultaneously change the counter-top  
 534 materials, backgrounds, lighting conditions, and distractor items.

535 **Input modality demonstration details.** We explore the ability of MOO to incorporate object-  
 536 centric mask representations that are generated via different processes than the one used during  
 537 training. During training, an OWL-ViT generates mask visual representations from textual prompts,  
 538 as described in Section 3.2. We study whether MOO can successfully accomplish manipulation  
 539 tasks given (1) a mask generated from a text caption from a generative VLM, (2) a mask generated  
 540 from an image query instead of a text query, or (3) a mask directly provided by a human via a

Dataset Filtering		Pick	
Objects	Episodes per Object	Seen objects	Unseen objects
100%	100%	<b>98</b>	<b>79</b>
50%	100%	92	75
18%	100%	88	19
100%	50%	46	38
100%	10%	23	0

Table 1: Performance of MOO in percentage of success relative to the amount of data used for training. Both data scale and data diversity are important.

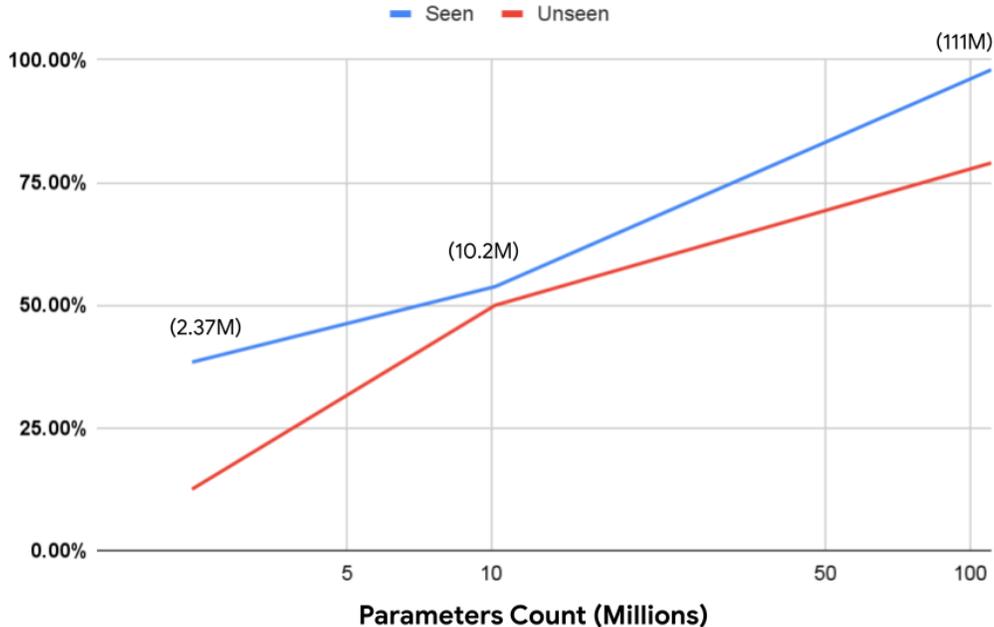


Figure 11: Pick success vs. model size. We see continuous improvements on both seen and unseen objects as we increase the number of parameters of our model architecture while keeping the data set size fixed. In comparison to our main model, we scaled down layer widths and depth by the same constant multiplier. We expect more performance gains at larger model capacity, yet are currently unable to scale further due to real time inference constraints on our robot.

541 GUI. For each of these cases, we implement different procedures for generating the object mask  
542 representation, which are then fed to the frozen MOO policy.

543 **Training data ablation.** We ablate the amount of data used to train MOO, and find that both data  
544 diversity and data scale are important, as shown in Table 1.

#### 545 Prompts used

546 We use the following prompts to OWL-ViT detect our objects. All prompts were prefixed with the  
547 phrase “An image of a”.

548 7up can → “white can of soda”

549 banana → “banana”

550 black pen → “black pen”

551 blue chip bag → “blue bag of chips”

552 blue pen → “blue pen”

553 brown chip bag → “brown bag of chips”

554 cereal scoop → “cereal scoop”

555 chocolate peanut candy → “bag of candy snack”

556 coffee cup → “coffee cup”

557 coke can → “red can of soda”

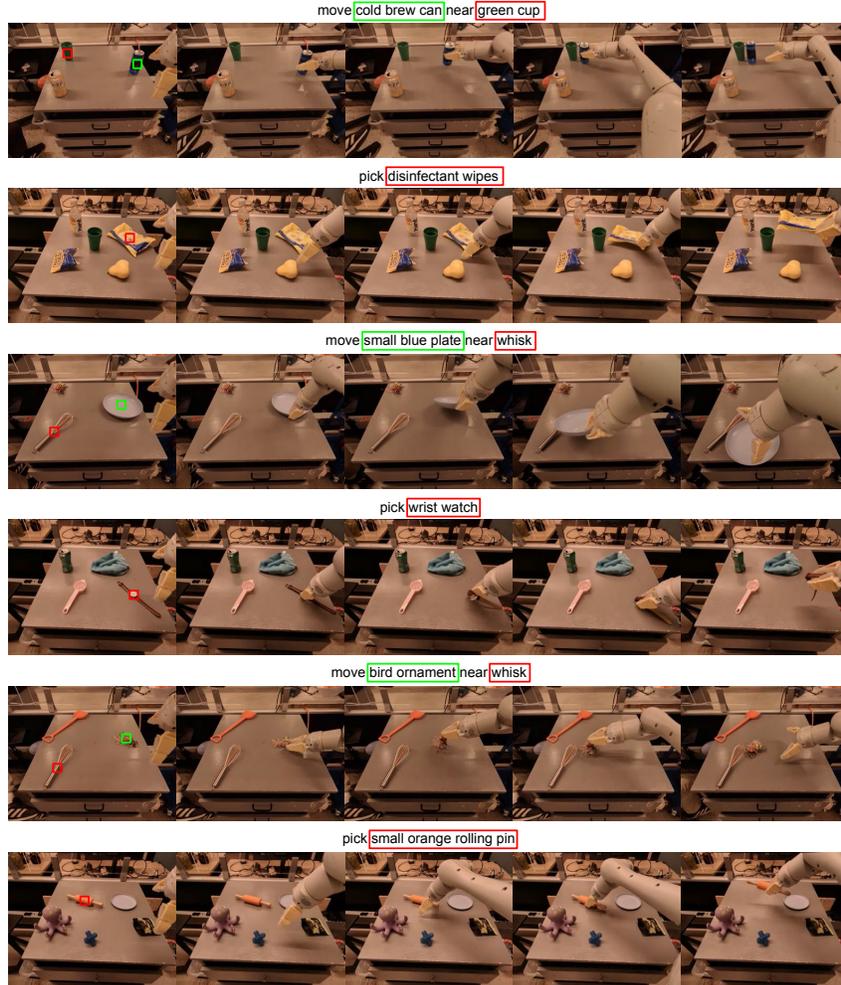


Figure 12: Example images of our policy detecting and grasping objects not seen during training time. The object detections are colored in correspondence to the text above the image, and the images are ordered left to right across time.

- 558 coke zero can → “can of soda”
- 559 disinfectant pump → “bottle”
- 560 fork → “fork”
- 561 green can → “green aluminum can”
- 562 green cookies bag → “green snack food bag”
- 563 green jalapeno chip bag → “green bag of chips”
- 564 green sprite can → “green soda can”
- 565 knife → “knife”
- 566 orange can → “orange aluminum can”
- 567 orange plastic bottle → “orange bottle”
- 568 oreo → “cookie snack food bag”
- 569 pepsi can → “blue soda can”
- 570 popcorn chip bag → “bag of chips”
- 571 pretzel chip bag → “bag of chips”
- 572 red grapefruit can → “red aluminum can”
- 573 redbull can → “skinny silver can of soda”
- 574 rxbar blueberry → “small blue rectangular snack food bar”
- 575 spoon → “spoon”
- 576 swedish fish bag → “bag of candy snack food”
- 577 water bottle → “clear plastic waterbottle with white cap”
- 578 white sparkling can → “aluminum can”

579 blue plastic bottle → “clear plastic waterbottle with white cap”  
580 diet pepper can → “can of soda”  
581 disinfectant wipes → “yellow and blue pack”  
582 green rice chip bag → “green bag of chips”  
583 orange → “round orange fruit”  
584 paper bowl → “round bowl”  
585 rxbar chocolate → “small black rectangular snack food bar”  
586 sponge → “scrub sponge”  
587 blackberry hint water → “clear plastic bottle with white cap”  
588 pineapple hint water → “clear plastic bottle with white cap”  
589 watermelon hint water → “clear plastic bottle with white cap”  
590 regular 7up can → “can of soda”  
591 lemonade plastic bottle → “clear plastic bottle with white cap”  
592 diet coke can → “silver can of soda”  
593 yellow pear → “yellow pear”  
594 green pear → “green pear”  
595 instant oatmeal pack → “flat brown pack of instant oatmeal”  
596 coffee mixing stick → “small thin flat wooden popsicle stick”  
597 coffee cup lid → “round disposable coffee cup lid”  
598 coffee cup sleeve → “brown disposable coffee cup sleeve”  
599 numi tea bag → “small flat packet of tea”  
600 fruit gummies → “small blue bag of snacks”  
601 chocolate caramel candy → “small navy bag of candy”  
602 original redbull can → “can of energy drink with dark blue label”  
603 cold brew can → “blue and black can”  
604 ginger lemon kombucha → “yellow and tan aluminum can with brown writing”  
605 large orange plate → “circular orange plate”  
606 small blue plate → “circular blue plate”  
607 love kombucha → “white and orange can of soda”  
608 original pepper can → “dark red can of soda”  
609 ito en green tea → “light green can of soda”  
610 iced tea can → “black can of soda”  
611 cheese stick → “yellow cheese stick in wrapper”  
612 brie cheese cup → “small white cheese cup with wrapper”  
613 pineapple spindrift can → “white and cyan can of soda”  
614 lemon spindrift can → “white and brown can of soda”  
615 lemon sparkling water can → “yellow can of soda”  
616 milano dark chocolate → “white pack of snacks”  
617 square cheese → “small orange rectangle packet”  
618 boiled egg → “small white egg in a plastic wrapper”  
619 pickle snack → “small black and green snack bag”  
620 red cup → “plastic red cup”  
621 blue cup → “plastic blue cup”  
622 orange cup → “plastic orange cup”  
623 green cup → “plastic green cup”  
624 head massager → “metal head massager with many wires”  
625 chew toy → “blue and yellow toy with orange polka dots”  
626 wrist watch → “wrist watch”  
627 small orange rolling pin → “small orange rolling pin with wooden handles”  
628 large green rolling pin → “large green rolling pin with wooden handles”  
629 rubiks cube → “rubiks cube”  
630 blue microfiber cloth → “blue cloth”  
631 gray microfiber cloth → “gray cloth”  
632 green microfiber cloth → “green cloth”  
633 small blending bottle → “small turquoise and brown bottle”  
634 large tennis ball → “large tennis ball”  
635 table tennis paddle → “table tennis paddle”  
636 octopus toy → “purple toy octopus”  
637 pink shoe → “pink shoe”  
638 floral shoe → “red and blue shoe”  
639 whisk → “whisk”  
640 orange spatula → “orange spatula”  
641 small blue spatula → “small blue spatula”  
642 large yellow spatula → “large yellow spatula”  
643 egg separator → “large pink cooking spoon”

644 green brush → “green brush”  
645 small purple spatula → “small purple spatula”  
646 badminton shuttlecock → “shuttlecock”  
647 black sunglasses → “black sunglasses”  
648 toy ball with holes → “toy ball with holes”  
649 red plastic shovel → “red plastic shovel”  
650 bird ornament → “colorful ornament with blue and yellow confetti”  
651 blue balloon → “blue balloon animal”  
652 catnip toy → “small dark blue plastic cross toy”  
653 raspberry baby teether → “red and green baby pacifier”  
654 slinky toy → “gray metallic cylinder slinky”  
655 dna chew toy → “big orange spring”  
656 gray suction toy → “gray suction toy”  
657 teal and pink toy car → “teal and pink toy car”  
658 two pound purple dumbbell → “purple dumbbell”  
659 one pound pink dumbbell → “pink dumbbell”  
660 three pound brown dumbbell → “brown dumbbell”  
661 dog rope toy → “white pink and gray rope with knot”  
662 fish toy → “fish”  
663 chain link toy → “skinny green rectangular toy”  
664 toy boat train → “plastic toy boat”  
665 white coat hanger → “white coat hanger”  
666