

# RT-Sketch: Goal-Conditioned Imitation Learning from Hand-Drawn Sketches

Priya Sundaresan<sup>1,3</sup>, Quan Vuong<sup>2</sup>, Jiayuan Gu<sup>2</sup>, Peng Xu<sup>2</sup>, Ted Xiao<sup>2</sup>, Sean Kirmani<sup>2</sup>, Tianhe Yu<sup>2</sup>, Michael Stark<sup>3</sup>, Ajinkya Jain<sup>3</sup>, Karol Hausman<sup>1,2</sup>, Dorsa Sadigh<sup>\*1,2</sup>, Jeannette Bohg<sup>\*2</sup>, Stefan Schaal<sup>\*3</sup>

\*Equal advising, alphabetical order

<sup>1</sup>Stanford University, <sup>2</sup>Google DeepMind, <sup>3</sup>[Google] Intrinsic

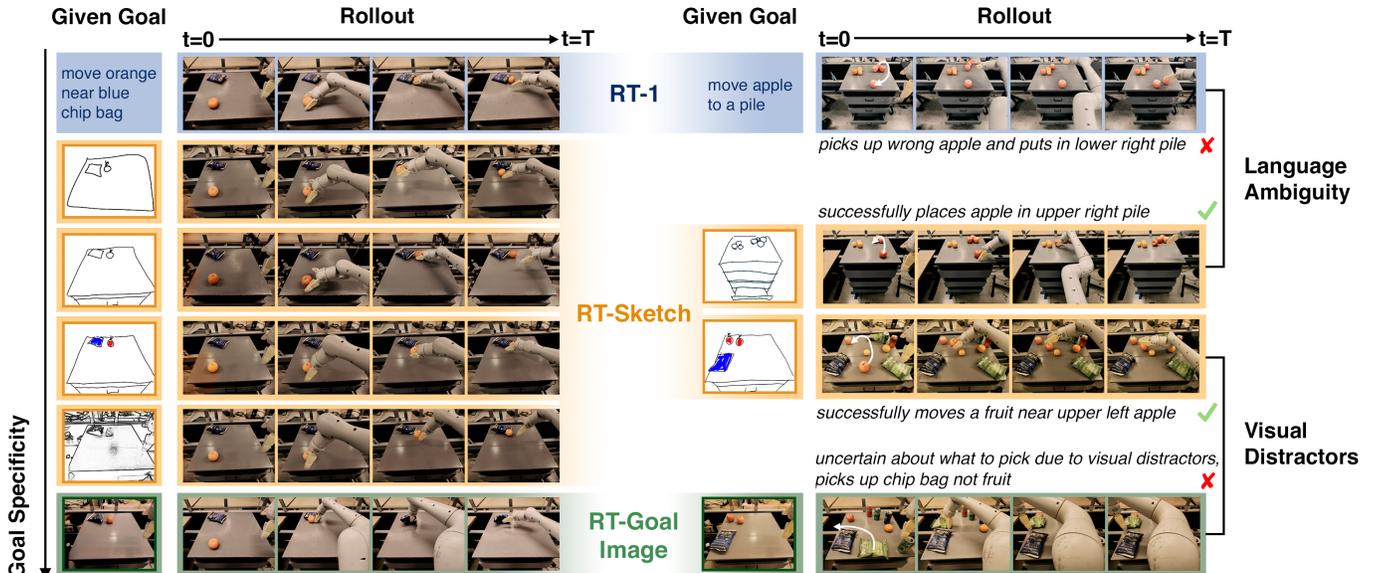


Fig. 1: (Left) Qualitative rollouts comparing RT-Sketch, RT-1, and RT-Goal-Image, (right) highlighting RT-Sketch’s robustness to (top) ambiguous language and (bottom) visual distractors.

**Abstract**—Natural language and images are commonly used as goal representations in goal-conditioned imitation learning (IL). However, natural language can be ambiguous and images can be over-specified. In this work, we study hand-drawn sketches as a modality for goal specification. Sketches are easy for users to provide on the fly like language, but similar to images they can also help a downstream policy to be spatially-aware and even go beyond images to disambiguate task-relevant from task-irrelevant objects. We present RT-Sketch, a goal-conditioned policy for manipulation that takes a hand-drawn sketch of the desired scene as input, and outputs actions. We train RT-Sketch on a dataset of trajectories paired with synthetically generated goal sketches. We evaluate this approach on six manipulation skills involving tabletop object rearrangements on an articulated countertop. Experimentally we find that RT-Sketch is able to perform on a similar level to image or language-conditioned agents in straightforward settings, while achieving greater robustness when language goals are ambiguous or visual distractors are present. Additionally, we show that RT-Sketch has the capacity to interpret and act upon sketches with varied levels of specificity, ranging from minimal line drawings to detailed, colored drawings. For supplementary material and videos, please refer to our [website](#).

## I. INTRODUCTION

Robots operating alongside humans in households, workplaces, or industrial environments have an immense potential for assistance and autonomy, but careful consideration is

needed of what goal representations are easiest *for humans* to convey to robots, and *for robots* to interpret and act upon.

Instruction-following robots attempt to address this problem using the intuitive interface of natural language commands as inputs to language-conditioned imitation learning policies [8, 9, 23, 28, 29]. For instance, imagine asking a household robot to set the dinner table. A language description such as “*put the utensils, the napkin, and the plate on the table*” is under-specified or ambiguous. It is unclear how exactly the utensils should be positioned relative to the plate or the napkin, or whether their distances to each other matter or not. To achieve this higher level of precision, a user may need to give lengthier descriptions such as “*put the fork 2cm to the right of the plate, and 5cm to the leftmost edge of the table.*”, or even online corrections (“*no, you moved too far to the right, move back a bit!*”) [15, 29]. While language is an intuitive way to specify goals, its qualitative nature and ambiguities can make it both inconvenient for humans to provide without lengthy instructions or corrections, and for robot policies to interpret for downstream precise manipulation.

Using goal images to specify objectives and training goal-conditioned imitation learning policies either paired with or without language instructions has shown to be quite successful in recent years [21, 22, 35]. In these settings, an image of the

scene in its desired final state could fully specify the intended goal. However, this has its own shortcomings: access to a goal image is a strong prior assumption, and a pre-recorded goal image can be tied to a particular environment, making it difficult to reuse for generalization.

To summarize: while natural language is highly flexible, it can also be highly ambiguous or require lengthy descriptions. This quickly becomes difficult in long-horizon tasks or those requiring spatial awareness. Meanwhile, goal images overspecify goals in unnecessary detail, leading to the need for internet-scale data for generalization.

To address these challenges, we study *hand-drawn sketches* as a convenient yet expressive modality for goal specification in visual imitation learning. By virtue of being minimal, sketches are still easy for users to provide on the fly like language, but allow for more spatially-aware task specification. Like goal images, sketches readily integrate with off-the-shelf policy architectures that take visual input, but provide an added level of goal abstraction that ignores unnecessary pixel-level details. Finally, the quality and selective inclusion/exclusion of details in a sketch can help a downstream policy distinguish task relevant from irrelevant details.

In this work, we present RT-Sketch, a goal-conditioned policy for manipulation that takes a user-provided hand-drawn sketch of the desired scene as input, and outputs actions. The novel architecture of RT-Sketch modifies the original RT-1 language-to-action Transformer architecture [9] to consume visual goals rather than language, allowing for flexible conditioning on sketches, images, or any other visually representable goals. To enable this, we concatenate a goal sketch and history of observations as input before tokenization, omitting language. We train RT-Sketch on a dataset of 80K trajectories paired with synthetic goal sketches, generated by an image-to-sketch stylization network trained from a few hundred image-sketch pairs.

We evaluate RT-Sketch across six manipulation skills on real robots involving tabletop object rearrangements on a countertop with drawers, subject to a wide range of scene variations. These skills include moving objects near to one another, knocking a can sideways, placing a can upright, closing a drawer, and opening a drawer. Experimentally, we find that RT-Sketch performs on a similar level to image or language-conditioned agents in straightforward settings. When language instructions are ambiguous, or in the presence of visual distractors, we find that RT-Sketch achieves  $\sim 2X$  more spatial precision and alignment scores, as assessed by human labelers, over language or goal image-conditioned policies (see Fig. 1 (right)). Additionally, we show that RT-Sketch can handle different levels of input specificity, ranging from rough sketches to more scene-preserving, colored drawings (see Fig. 1 (left)).

## II. RELATED WORK

In this section, we discuss prior methods for goal-conditioned imitation learning. We also highlight ongoing efforts towards image-sketch conversion, which open new

possibilities for goal-conditioning modalities which are under-explored in robotics.

*a) Goal-Conditioned Imitation Learning:* Despite the similarity in name, our learning of manipulation policies conditioned on hand-drawn sketches of the desired scene is different from the notion of policy sketches [1], symbolic representations of task structure describing its subcomponents. Reinforcement learning (RL) is not easily applicable in our scenario, as it is nontrivial to define a reward objective which accurately quantifies alignment between a provided scene sketch and states visited by an agent during training. We instead focus on imitation learning (IL) techniques, particularly the goal-conditioned setting [17].

Goal-conditioned IL has proven useful in settings where a policy must be able to handle spatial or semantic variations for the same task [3]. These settings include rearrangement of multiple objects [8, 9, 29, 30, 35], kitting [46], folding of deformable objects into different configurations [18], and search for different target objects in clutter [16]. However, these approaches tend to either rely on language [9, 23, 28, 29, 39], or goal images [16] to specify variations. Follow-up work enabled multimodal conditioning on either goal images and language [21], in-prompt images [22], or image embeddings [18, 30, 46]. However, all of these representations are ultimately derived from raw images or language in some way, which overlooks the potential for more abstract goal representations that are easy to specify but preserve spatial awareness, such as sketches.

In addition to their inflexibility in terms of goal representation, goal-conditioned IL tends to overfit to demonstration data and fails to handle even slight distribution shift in new scenarios [37]. For language-conditioning, distribution shift can encompass semantic or spatial ambiguity, novel instructions or phrasing, or unseen objects [9, 21]. Goal-image conditioning is similarly susceptible to out-of-distribution visual shift, such as variations in lighting or object appearances, or unseen background textures [5, 11]. We instead opt for sketches which are minimal enough to combat visual distractors, yet expressive enough to provide unambiguous goals. Prior work, including [4] and [32], have shown the utility of sketches over pure language for navigation and limited manipulation settings. However, the sketches explored in these works are largely intended to guide low-level motion at the joint-level for manipulation, or provide explicit directional cues for navigation. [14] considers sketches amongst other modalities as an input for goal-conditioned manipulation, but does not explicitly train a policy conditioned on sketches. They thus came to the conclusion that the scene image is better than the sketch image at goal specification. Our result is different and complementary, in that policies trained to take sketches as input outperform a scene image conditioned policy, by 1.63x and 1.5x in terms of Likert ratings for perceived spatial and semantic alignment, subject to visual distractors. Most recently, Gu et al. [19] propose an approach to goal-conditioned manipulation via *hindsight-trajectory* sketches. Here, sketches represent 2D paths drawn over an image to indicate the desired

trajectory for the robot to follow. While this work treats sketches as a *motion-centric* representation to specify intended robot trajectories, the sketches in our work are *scene-centric*, representing the desired visual goal state rather than the actions the robot should explicitly take.

*b) Image-Sketch Conversion:* In recent years, sketches have gained increasing popularity within the computer vision community for applications such as object detection [6, 7, 12], visual question answering [25, 33], and scene understanding [13], either in isolation or in addition to text and images. When considering how best to incorporate sketches in IL, an important design choice is whether to take sketches into account (1) at test time (i.e., converting a sketch to another goal modality compatible with a pre-trained policy), or (2) at training time (i.e., explicitly training an IL policy conditioned on sketches). For (1), one could first convert a given sketch to a goal image, and then roll out a vanilla goal-image conditioned policy. This could be based on existing frameworks for sketch-to-image conversion, such as ControlNet [47], GAN-style approaches [24], or text-to-image synthesis, such as Instruct-Pix2Pix [10] or Stable Diffusion [36]. While these models produce photorealistic results under optimal conditions, they do not jointly handle image generation and style transfer, making it unlikely for generated images to match the style of an agent observations. At the same time, these approaches are susceptible to producing hallucinated artifacts, introducing distribution shifts [47].

Based on these challenges, we instead opt for (2), and consider image-to-sketch conversion techniques for hindsight relabeling of terminal images in pre-recorded demonstration trajectories. Recently, Vinker et al. [44, 45] proposes networks for predicting Bezier curve-based sketches of input image objects or scenes. Sketch quality is supervised by a CLIP-based alignment metric. While these approaches generate sketches of high visual fidelity, test-time optimization takes on the order of minutes, which does not scale to the typical size of robot learning datasets (hundreds to thousands of demonstration trajectories). Meanwhile, conditional generative adversarial networks (cGANs) such as Pix2Pix [20] have proven useful for scalable image-to-image translation. Most related to our work is that of Li et al. [26], which trains a Pix2Pix model to produce sketches from given images on a large crowd-sourced dataset of  $5K$  paired images and line drawings. We build on this work to fine-tune an image-to-sketch model on robot trajectory data, and show its utility for enabling downstream manipulation from sketches.

### III. SKETCH-CONDITIONED IMITATION LEARNING

In this section, we will first introduce our problem of learning a sketch-conditioned policy. We will then discuss our approach to train an end-to-end sketch-to-action IL agent. First, in Section III-A, we discuss our instantiation of an auxiliary image-to-sketch translation network which automatically generates sketches from a reference image. In Section III-B, we discuss how we use such a model to automatically hindsight relabel an existing dataset of demonstrations with synthetically

generated goal sketches, and train a sketch-conditioned policy on this dataset.

*a) Problem Statement:* Our goal is to learn a manipulation policy conditioned on a goal *sketch* of the desired scene state and a history of interactions. Formally, we denote such a policy by  $\pi_{\text{sketch}}(a_t|g, \{o_j\}_{j=1}^t)$ , where  $a_t$  denotes an action at timestep  $t$ ,  $g \in \mathbb{R}^{W \times H \times 3}$  is a given goal sketch with width  $W$  and height  $H$ , and  $o_t \in \mathbb{R}^{W \times H \times 3}$  is an observation at time  $t$ . At inference time, the policy takes a given goal sketch along with a history of RGB image observations to infer an action to execute. In practice, we condition  $\pi_{\text{sketch}}$  on a history of  $D$  previous observations rather than all observations from the initial state at  $t = 1$ . To train such a policy, we assume access to a dataset  $\mathcal{D}_{\text{sketch}} = \{g^n, \{(o_t^n, a_t^n)\}_{t=1}^{T^{(n)}}\}_{n=1}^N$  of  $N$  successful demonstrations, where  $T^{(n)}$  refers to the length of the  $n^{\text{th}}$  trajectory in timesteps. Each episode of the dataset consists of a given goal sketch and a corresponding demonstration trajectory, with image observations recorded at each timestep. Our goal is to thus learn the sketch-conditioned imitation policy  $\pi_{\text{sketch}}(a_t|g, \{o_j\}_{j=1}^t)$  trained on this dataset  $\mathcal{D}_{\text{sketch}}$ .

#### A. Image-to-Sketch Translation

Training a sketch-conditioned policy requires a dataset of robot trajectories that are each paired with a sketch of the goal state achieved by the robot. Collecting such a dataset from scratch at scale, including the trajectories themselves and manually drawn sketches, can easily become impractical. Thus, we instead aim to learn an image-to-sketch translation network  $\mathcal{T}(g|o)$  that takes an image observation  $o$  and outputs the corresponding goal sketch  $g$ . This network can be used to post-process an existing dataset of demonstrations  $\mathcal{D} = \{\{(o_t^n, a_t^n)\}_{t=1}^{T^{(n)}}\}_{n=1}^N$  with image observations by appending a synthetically generated goal sketch to each demonstration. This produces a dataset for sketch-based IL:  $\mathcal{D}_{\text{sketch}} = \{g^n, \{(o_t^n, a_t^n)\}_{t=1}^{T^{(n)}}\}_{n=1}^N$ .

*a) RT-1 Dataset:* In this work, we rely on an existing dataset of visual demonstrations collected by prior work [9]. RT-1 is a prior language-to-action imitation learning agent trained on a large-scale dataset ( $80K$  trajectories) of VR-teleoperated demonstrations that include skills such as moving objects near one another, placing cans and bottles upright or sideways, opening and closing cabinets, and performing pick and place on countertops and drawers [9]. Here, we repurpose the RT-1 dataset and further adapt the RT-1 policy architecture to accommodate sketches, detailed in Section III-B.

*b) Assumptions on Sketches:* We acknowledge that there are innumerable ways for a human to provide a sketch corresponding to a given image of a scene. In this work, we make the following assumptions about input sketches for a controlled experimental validation procedure. In particular, we first assume that a given sketch respects the task-relevant contours of an associated image, such that tabletop edges, drawer handles, and task-relevant objects are included and discernible in the sketch. We do not assume contours in the sketch to be edge-aligned or pixel-aligned with those in an

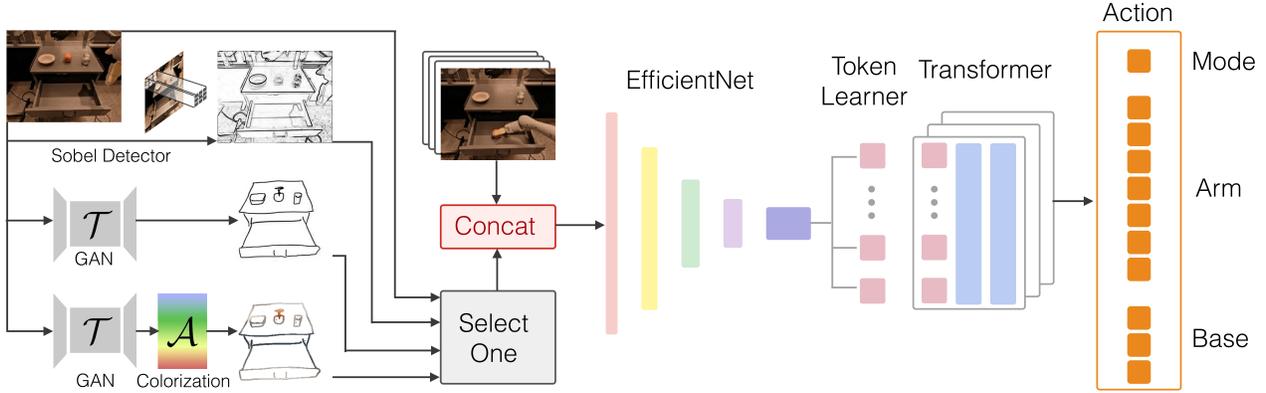


Fig. 2: Architecture of RT-Sketch allowing different kinds of visual input. RT-Sketch adopts the Transformer [43] architecture with EfficientNet [41] tokenization at the input, and outputs bucketized actions.

image. We do assume that the input sketch consists of black outlines at the very least, with shading in color being optional. We further assume that sketches do not contain information not present in the associated image, such as hallucinated objects, scribbles, or textual annotations, but may omit task-irrelevant details that appear in the original image.

c) *Sketch Dataset Generation*: To train an image-to-sketch translation network  $\mathcal{T}$ , we collect a new dataset  $\mathcal{D}_{\mathcal{T}} = \{(o_i, g_i^1, \dots, g_i^{L^{(i)}})\}_{i=1}^M$  consisting of  $M$  image observations  $o_i$  each paired with a set of goal sketches  $g_i^1, \dots, g_i^{L^{(i)}}$ . Those represent  $L^{(i)}$  different representations of the same image  $o_i$ , in order to account for the fact that there are multiple, valid ways of sketching the same scene. To collect  $\mathcal{D}_{\mathcal{T}}$ , we take 500 randomly sampled terminal images from demonstration trajectories in the RT-1 dataset, and manually draw sketches with black lines on a white background capturing the tabletop, drawers, and relevant objects visible on the manipulation surface. While we personally annotate each robot observation with a single sketch only, we add this data to an existing, much larger non-robotic dataset [26]. This dataset captures inter-sketch variation via multiple crowdsourced sketches per image. We do not include the robot arm in our manual sketches, as we find a minimal representation to be most natural. Empirically, we find that our policy can handle such sketches despite actual goal configurations likely having the arm in view. We collect these drawings using a custom digital stylus drawing interface in which a user draws an edge-aligned sketch over the original image (Appendix Fig. 16). The final recorded sketch includes the user’s strokes in black on a white canvas with the original image dimensions.

d) *Image-to-Sketch Training*: We implement the image-to-sketch translation network  $\mathcal{T}$  with the Pix2Pix conditional generative adversarial network (cGAN) architecture, which is composed of a generator  $G_{\mathcal{T}}$  and a discriminator  $D_{\mathcal{T}}$  [20]. The generator  $G_{\mathcal{T}}$  takes an input image  $o$ , a random noise vector  $z$ , and outputs a goal sketch  $g$ . The discriminator  $D_{\mathcal{T}}$  is trained to discriminate amongst artificially generated sketches and ground truth goal sketches. We utilize the standard cGAN supervision loss to train both [20, 26]:

$$\mathcal{L}_{\text{cGAN}} = \min_{G_{\mathcal{T}}} \max_{D_{\mathcal{T}}} \mathbb{E}_{o,g} [\log D_{\mathcal{T}}(o, g)] + \mathbb{E}_{o,g} [\log(1 - D_{\mathcal{T}}(o, G_{\mathcal{T}}(o, g)))] \quad (1)$$

We also add the  $\mathcal{L}_1$  loss to encourage the produced sketches to align with the ground truth sketches as in [26]. To account for the fact that there may be multiple valid sketches for a given image, we only penalize the minimum  $\mathcal{L}_1$  loss incurred across all  $L^{(i)}$  sketches provided for a given image as in Li et al. [26]. This is to prevent wrongly penalizing  $\mathcal{T}$  for producing a valid sketch that aligns well with one example but not another simply due to stylistic differences in the ground truth sketches. The final objective is then a  $\lambda$ -weighted combination of the average cGAN loss and the minimum alignment loss:

$$\mathcal{L}_{\mathcal{T}} = \frac{\lambda}{L^{(i)}} \sum_{k=1}^{L^{(i)}} \mathcal{L}_{\text{cGAN}}(o_i, g_i^{(k)}) + \min_{k \in \{1, \dots, L^{(i)}\}} \mathcal{L}_1(o_i, g_i^{(k)}) \quad (2)$$

In practice, we supplement the 500 manually drawn sketches from  $\mathcal{D}_{\mathcal{T}}$  by leveraging the existing larger-scale Contour Drawing Dataset [26]. We refer to this dataset as  $\mathcal{D}_{\text{CD}}$ , which contains 1000 examples of internet-scraped images containing objects, people, animals from Adobe Stock, paired with  $L^{(i)} = 5$  crowd-sourced black and white outline drawings per image collected on Amazon Mechanical Turk. Visualizations of this dataset are provided in Appendix Fig. 5. We first take a pre-trained image-to-sketch translation network  $\mathcal{T}_{\text{CD}}$  [26] trained on  $\mathcal{D}_{\text{CD}}$ , with  $L^{(i)} = 5$  sketches per image. Then, we fine-tune  $\mathcal{T}_{\text{CD}}$  on  $\mathcal{D}_{\mathcal{T}}$ , with only  $L^{(i)} = 1$  manually drawn sketch per robot observation, to obtain our final image-to-sketch network  $\mathcal{T}$ . Visualizations of the sketches generated by  $\mathcal{T}$  for different robot observations are available in Fig. 6.

## B. RT-Sketch

With a means of translating image observations to black and white sketches via  $\mathcal{T}$  (Section III-A), we can automatically augment the existing RT-1 dataset with goal sketches. This

results in a dataset, which we refer to as  $\mathcal{D}_{\text{sketch}}$ , which can be used for training our algorithm, RT-Sketch.

a) *RT-Sketch Dataset*: The original RT-1 dataset  $\mathcal{D}_{\text{lang}} = \{i^n, \{(o_t^n, a_t^n)\}_{t=1}^{T(n)}\}_{n=1}^N$  consists of  $N$  episodes with a paired natural language instruction  $i$  and demonstration trajectory  $\{(o_t^n, a_t^n)\}_{t=1}^{T(n)}$ . We can automatically hindsight-relabel such a dataset with goal images instead of language goals [2]. Let us denote the last step of a trajectory  $n$  as  $T(n)$ . Then the new dataset with image goals instead of language goals is  $\mathcal{D}_{\text{img}} = \{o_{T(n)}^n, \{(o_t^n, a_t^n)\}_{t=1}^{T(n)}\}_{n=1}^N$ , where we treat the last observation of the trajectory  $o_{T(n)}^n$  as the goal  $g^n$ . To produce a dataset for  $\pi_{\text{sketch}}$ , we can simply replace  $o_{T(n)}^n$  with  $\hat{g}^n = \mathcal{T}(o_{T(n)}^n)$  such that  $\mathcal{D}_{\text{sketch}} = \{\hat{g}^n, \{(o_t^n, a_t^n)\}_{t=1}^{T(n)}\}_{n=1}^N$ .

To encourage the policy to afford different levels of input sketch specificity, we in practice produce goals by  $\hat{g}^n = \mathcal{A}(o_{T(n)}^n)$ , where  $\mathcal{A}$  is a randomized augmentation function.  $\mathcal{A}$  chooses between simply applying  $\mathcal{T}$ ,  $\mathcal{T}$  with colorization during postprocessing (e.g., by superimposing a blurred version of the ground truth RGB image over the binary sketch), a Sobel operator [40] for edge detection, or not applying any operators, which preserves the original ground truth goal image (Fig. 2). By co-training on all representations, we intend for RT-Sketch to handle a spectrum of specificity going from binary sketches; colored sketches; edge detected images; and goal images (Appendix Fig. 6).

b) *RT-Sketch Model Architecture*: In our setting, we consider goals provided as sketches rather than language instructions as was done in RT-1. This change in the input representation necessitates a change in the model architecture. The original RT-1 policy relies on a Transformer architecture backbone [43]. RT-1 first passes a history of  $D = 6$  images through an EfficientNet-B3 model [41] producing image embeddings, which are tokenized, and separately extracts textual embeddings and tokens via FiLM [31] and a Token Learner [38]. The tokens are then fed into a Transformer which outputs bucketized actions. The output action dimensionality is 7 for the end-effector (x, y, z, roll, pitch, yaw, gripper width), 3 for the mobile base, (x, y, yaw), and 1 for a flag that can select amongst base movement, arm movement, and episode termination. To retrain the RT-1 architecture but accommodate the change in input representation, we omit the FiLM language tokenization altogether. Instead, we concatenate a given goal image or sketch with the history of images as input to EfficientNet, and extract tokens from its output, leaving the rest of the policy architecture unchanged. We visualize the RT-Sketch training inputs and policy architecture in Fig. 2. We refer to this architecture when trained only on images (i.e., an image goal-conditioned RT1 policy) as RT-Goal-Image and refer to it as RT-Sketch when it is trained on sketches as discussed in this section.

c) *Training RT-Sketch*: We can now train  $\pi_{\text{sketch}}$  on  $\mathcal{D}_{\pi_{\text{sketch}}}$  utilizing the same procedure as was used to train RT-1 [9], with the above architectural modifications. We fit  $\pi_{\text{sketch}}$  using the behavioral cloning objective function. This aims to minimize the negative log-likelihood of an action provided the

history of observations and a given sketch goal [42]:

$$J(\pi_{\text{sketch}}) = \sum_{n=1}^N \sum_{t=1}^{T(n)} \log \pi_{\text{sketch}}(a_t^n | g^n, \{o_j\}_{j=1}^t)$$

## IV. EXPERIMENTS

We seek to understand the ability of RT-Sketch to perform goal-conditioned manipulation as compared to policies that operate from higher-level goal abstractions like language, or more over-specified modalities, like goal images. To that end, we test the following four hypotheses:

**H1: RT-Sketch is successful at goal-conditioned IL.** While sketches are abstractions of real images, our hypothesis is that they are specific enough to provide manipulation goals to a policy. Therefore, we expect RT-Sketch to perform on a similar level to language goals (RT-1) or goal images (RT-Goal-Image) in straightforward manipulation settings.

**H2: RT-Sketch is able to handle varying levels of specificity.** There are as many ways to sketch a scene as there are people. Because we have trained RT-Sketch on sketches of varying levels of specificity, we expect it to be robust against variations of the input sketch for the same scene.

**H3: Sketches enable better robustness to distractors than goal images.** Sketches focus on task-relevant details of a scene. Therefore, we expect RT-Sketch to provide robustness against distractors in the environment that are not included in the sketch compared to RT-Goal-Image that operates on detailed image goals.

**H4: Sketches are favorable when language is ambiguous.** We expect RT-Sketch to provide a higher success rate compared to ambiguous language inputs when using RT-1.

### A. Experimental Setup

a) *Policies*: We compare RT-Sketch to the original language-conditioned agent RT-1 [9], and RT-Goal-Image, a policy identical in architecture to RT-Sketch, but taking a goal image as input rather than a sketch. All policies are trained on a multi-task dataset of  $\sim 80\text{K}$  real-world trajectories manually collected via VR teleoperation using the setup from Brohan et al. [9]. These trajectories span a suite of common office and kitchen tasks such as picking and placing objects, re-orienting cups and bottles upright or sideways, opening and closing drawers, and rearranging objects between drawers or a countertop.

b) *Evaluation protocol*: To ensure fair comparison, we control for the same initial and goal state of the environment across different policy rollouts via a catalog of well-defined evaluation scenarios that serve as references for human robot operators. For each scenario, we record an initial image (RGB observation) of the scene, the goal image (with objects manually rearranged as desired), a natural language task string describing the desired agent behavior to achieve the goal, and a set of hand-drawn sketches corresponding to the recorded goal image. At test time, a human operator retrieves a particular evaluation scenario from the catalog, aligns the physical robot and scene according to a reference image using a custom

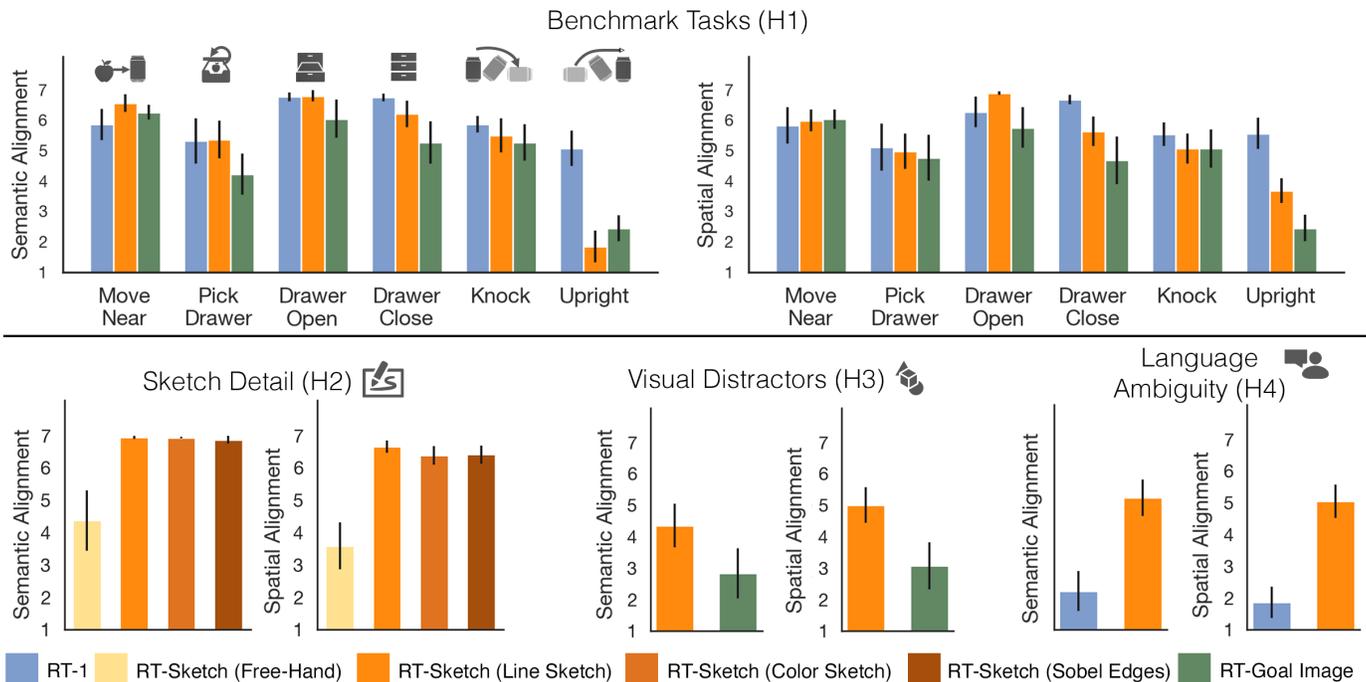


Fig. 3: **Goal Alignment Results:** Average Likert scores for different policies rating perceived semantic alignment (**Q1**) and spatial alignment (**Q2**) to a provided goal. For straightforward benchmark manipulation tasks, RT-Sketch performs comparably and in some cases better than RT-1 and RT-Goal-Image in terms of both metrics, for 5 out of 6 skills (**H1**). RT-Sketch further exhibits the ability to handle sketches of different levels of detail (**H2**), while achieving better goal alignment than baselines when the visual scene is distracting (**H3**) or language would be ambiguous (**H4**). Error bars indicate standard error across labeler ratings.

visualization utility, and places the relevant objects in their respective locations. Finally, the robot selects one of the goal representations (language, image, sketch, etc.) for the scenario as input to a policy. We record a video of the policy rollout for downstream evaluation (see Section IV-B). We perform all experiments using the Everyday Robot, which contains a mobile base, an overhead camera, and a 7-DoF manipulator arm with a parallel jaw gripper. All sketches for evaluation are collected with a custom manual drawing interface by a single human annotator on a tablet with a digital stylus.

c) *Performance Metrics:* Defining a standardized, automated evaluation protocol for goal alignment is non-trivial. Since binary task success is too coarse-grained and image-similarity metrics like frame-differencing or CLIP [34] tend to be brittle, we measure performance with two more targeted metrics. First, we quantify policy precision as the distance (in pixels) between object centroids in achieved and ground truth goal states, using manual keypoint annotations. Although leveraging out-of-the box object detectors to detect object centroids is a possibility, we want to avoid conflating errors in object detection (imprecise bounding box, wrong object, etc.) from manipulation error of the policy itself. Second, we gather human-provided assessments of perceived goal alignment, following the commonly-used Likert [27] rating scheme from 1 (Strongly Disagree) to 7 (Strongly Agree), for:

- (**Q1**) *The robot achieves semantic alignment with the given goal during the rollout.*

- (**Q2**) *The robot achieves spatial alignment with the given goal during the rollout.*

For **Q1**, we present labelers with the policy rollout video along with the given ground-truth language task description. We expect reasonably high ratings across all methods for straightforward manipulation scenarios (**H1**). Sketch-conditioned policies should yield higher scores than a language-conditioned policy when a task string is ambiguous (**H4**). **Q2** is instead geared at measuring to what degree a policy can spatially arrange objects as desired. For instance, a policy can achieve semantic alignment for the instruction *place can upright* as long as the can ends up in the right orientation. For **Q2**, we visualize a policy rollout side-by-side with a given visual goal (ground truth image, sketch, etc.) to assess perceived spatial alignment. We posit that all policies should receive high ratings for straightforward scenarios (**H1**), with a slight edge for visual-conditioned policies which implicitly have stronger spatial priors encoded in goals. We further expect that as the visual complexity of a scene increases, sketches may be able to better attend to pertinent aspects of a goal and achieve better spatial alignment than image-conditioned agents (**H3**), even for different levels of sketch specificity (**H4**). We provide a visualization of the assessment interface for **Q1** and **Q2** in Appendix Fig. 17. We note that we perform these human assessment surveys across 62 individuals (non-expert, unfamiliar with our system), where we assign between 8 and 12 people to evaluate each of the 6 different

Skill	Spatial Precision (RMSE in px.)			Failure Occurrence (Excessive Retrying)		
	RT-1	RT-Sketch	RT-Goal-Image	RT-1	RT-Sketch	RT-Goal-Image
Move Near	5.43 ± 2.15	<b>3.49 ± 1.38</b>	3.89 ± 1.16	<b>0.00</b>	0.06	0.33
Pick Drawer	5.69 ± 2.90	4.77 ± 2.78	<b>4.74 ± 2.01</b>	<b>0.00</b>	0.13	0.20
Drawer Open	4.51 ± 1.55	<b>3.34 ± 1.08</b>	4.98 ± 1.16	<b>0.00</b>	<b>0.00</b>	0.07
Drawer Close	<b>2.69 ± 0.93</b>	3.02 ± 1.35	3.71 ± 1.67	<b>0.00</b>	<b>0.00</b>	0.07
Knock	7.39 ± 1.77	<b>5.36 ± 2.74</b>	5.63 ± 2.60	<b>0.00</b>	0.13	0.40
Upright	7.84 ± 2.37	5.08 ± 2.08	<b>4.18 ± 1.54</b>	0.06	<b>0.00</b>	0.27
Visual Distractors	-	<b>4.78 ± 2.17</b>	7.95 ± 2.86	-	<b>0.13</b>	0.67
Language Ambiguity	8.03 ± 2.52	<b>4.45 ± 1.54</b>	-	0.40	<b>0.13</b>	-

TABLE I: **Spatial Precision and Failure Occurrence** : Left: We report the level of spatial precision achieved across policies, measured in terms of RMSE of the centroids of manipulated objects in achieved vs. given reference goal images. Darker shading indicates higher precision (lower centroid distance). Fig. 8 contains visualizations illustrating the degree of visual alignment that different RMSE values correspond to. Right: We report the proportion of rollouts in which different policies exhibit *excessive retrying* behavior. Bolded numbers indicate the most precise and least failure-prone policy for each skill.

skills considered below.

### B. Experimental Results

In this section, we present our findings related to the hypotheses of Section IV. Tables I and II measure the spatial precision achieved by policies in terms of pixelwise distance, while Fig. 3 shows the results of human-perceived semantic and spatial alignment, based on a 7-point Likert scale rating.

**H1:** We evaluate 6 skills from the RT-1 benchmark [9]: *move X near Y*, *place X upright*, *knock X over*, *open the X drawer*, *close the X drawer*, and *pick X from Y*. For each skill, we record 15 different catalog scenarios, varying both objects (16 unique in total) and their placements.

In general, we find that RT-Sketch performs on a comparable level to RT-1 and RT-Goal-Image for both semantic (Q1) and spatial alignment (Q2), achieving ratings in the ‘Agree’ to ‘Strongly Agree’ range on average for nearly all skills (Fig. 3 (top)). A notable exception is *upright*, where RT-Sketch essentially fails to accomplish the goal semantically (Q1), albeit with some degree of spatial alignment (Q2). Both RT-Sketch and RT-Goal-Image tend to position cans or bottles appropriately and then terminate, without realizing the need for reorientation (Appendix Fig. 9). This behavior results in low centroid-distance to the goal (darker gray in Table I (left)). RT-1, on the other hand, reorients cans and bottles successfully, but at the expense of higher error (Appendix Fig. 9, light color in Table I (left)). In our experiments, we also observe the occurrence of *excessive retrying behavior*, in which a policy attempts to align the current scene with a given goal with retrying actions such as grasping and placing. However, performing these low-level actions with a high degree of precision is challenging, and thus excessive retrying can actually disturb the scene leading to knocking objects off the table or undoing task progress. In Table I, we report the proportion of rollouts in which we observe this behavior across all policies. We note that RT-Goal-Image is most susceptible to this failure mode, as a result of over-attending to pixel-level details and trying in excess to match a given goal exactly. Meanwhile, RT-Sketch and RT-1 are far less vulnerable, since both sketches and language provide a higher level of goal abstraction.

**H2:** We next assess RT-Sketch’s ability to handle input sketches of varied levels of detail (free-hand, edge-aligned line sketch, colored line sketch, and a Sobel edge-detected image as an upper bound). Free-hand sketches are drawn with a reference image next to a blank canvas, while line sketches are drawn on a semi-transparent canvas overlaid on the image (see Appendix Fig. 16). We find such a UI to be convenient and practical, as an agent’s current observations are typically available and provide helpful guides for sketching lines and edges. Across 5 trials each of the *move near* and *open drawer* skills, we see in Table II that all types of sketches produce reasonable levels of spatial precision. As expected, Sobel edges incur the least error, but even free-hand sketches, which do not necessarily preserve perspective projection, and line sketches, which are far sparser in detail, are not far behind. This is also reflected in the corresponding Likert ratings (Fig. 3 (left, bottom)). Free-hand sketches already garner moderate ratings (around 4) of perceived spatial and semantic alignment, but line sketches result in a marked performance improvement to nearly 7, on par with the upper bound of providing an edge-detected goal image. Adding color does not improve performance further, but leads to interesting qualitative differences in behavior (see Appendix Fig. 10).

We also evaluate whether RT-Sketch can generalize to sketches drawn by different individuals and handle stylistic variations. We first collect 30 sketches drawn by 6 different annotators using line sketching (tracing) on 5 goal images from the *move near* evaluation scenarios. We obtain the resulting rollouts produced by RT-Sketch with these sketches as input. Across 22 human evaluators who report their perceived spatial alignment via Likert ratings, we find that RT-Sketch achieves high spatial alignment on sketches drawn by other annotators. Notably, there is no significant dropoff in performance between sketches drawn by different annotators, or in the policy performance as compared to when our original sketches are the input (Fig. 4).

**H3:** Next, we compare the robustness of RT-Sketch and RT-Goal-Image to the presence of visual distractors. We re-use 15 *move X near Y* trials from the catalog, but introducing 5 – 9 distractor objects into the initial visual scene after alignment.

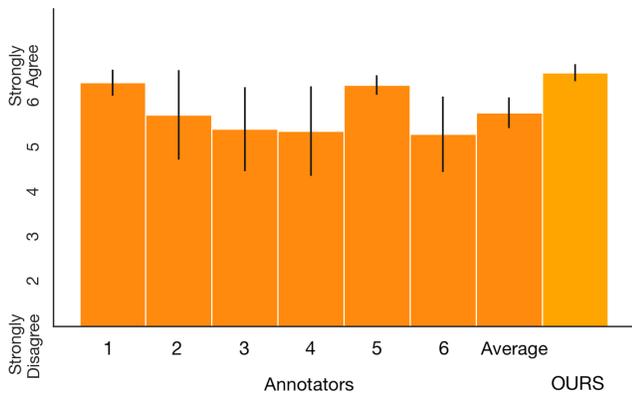


Fig. 4: **Perceived Spatial Alignment for Sketches Drawn by Other Annotators (H2)**: Across line sketches drawn by 6 annotators who are not represented in the training dataset for RT-Sketch, we record policy rollouts with these sketches as input for the *move near* skill. We evaluate the resulting rollouts across 22 human evaluators who provide Likert ratings measuring spatial alignment between the achieved goal state and given sketch. RT-Sketch’s performance on these new input sketches is on par with policy performance on our original sketches (OURS), and with no significant dropoff between sketches drawn by different annotators.

Skill	Free-Hand	Line Sketch	Color Sketch	Sobel Edges
Move Near	7.21 ± 2.76	3.49 ± 1.38	3.45 ± 1.03	<b>3.36 ± 0.66</b>
Drawer Open	3.75 ± 1.63	3.34 ± 1.08	2.48 ± 0.50	<b>2.13 ± 0.25</b>

TABLE II: **RT-Sketch Spatial Precision across Sketch Types (RMSE (centroid-distance) in px.** We report the spatial precision achieved by RT-Sketch subject to different input modalities. As expected, for less detailed and more rough sketches, RT-Sketch achieves lower precision (lighter shading), and for richer representations RT-Sketch is more precise (bolded, darker shading). Still, there is a relatively small difference in performance between line, color, and edge-detected representations, indicating RT-Sketch’s ability to afford different levels of input specificity.

This testing procedure is adapted from RT-1 generalization experiments referred to as *medium-high* difficulty [9]. In Table I (left, bottom), we see that RT-Sketch exhibits far lower spatial errors on average, while producing higher semantic and spatial alignment scores over RT-Goal-Image( Fig. 3 (middle, bottom)). RT-Goal-Image is easily confused by the distribution shift introduced by distractor objects, and often cycles between picking up and putting down the wrong object. RT-Sketch, on the other hand, ignores task-irrelevant objects not captured in a sketch and completes the task in most cases (see Appendix Fig. 11).

**H4**: Finally, we evaluate whether sketches as a representation are favorable when language goals alone are ambiguous. We collect 15 scenarios encompassing 3 types of ambiguity in language instructions: instance ambiguity (**T1**) (e.g., *move apple near orange* when multiple orange instances are present), somewhat out-of-distribution (OOD) language (**T2**) (e.g., *move left apple near orange*), and highly OOD language (**T3**) (e.g., *complete the rainbow*) (see Appendix Fig. 12). While the latter two qualifications should intuitively help resolve ambiguities, they were not explicitly made part of the original RT-1 training [9], and hence only provide limited utility. In Table I (left, bottom), RT-Sketch achieves nearly half the error of RT-1, and a 2.39-fold and 2.79-fold score increase for semantic and spatial alignment, respectively (Fig. 3 (right, bottom)). For **T1** and **T2** scenarios, RT-1 often tries to pick

up an instance of any object mentioned in the task string, but fails to make progress beyond that (Appendix Fig. 13). This further suggests the utility of sketches to express new, unseen goals with minimal overhead, when language could otherwise be opaque or difficult to express with only in-distribution vocabulary (Appendix Fig. 14).

a) *Limitations and Failure Modes*: Firstly, the image-to-sketch generation network used in this work is fine-tuned on a dataset of sketches provided by a single human annotator. Although we empirically show that despite this, RT-Sketch can handle sketches drawn by other annotators, we have yet to investigate the effects of training RT-Sketch at scale with sketches produced by different people. Secondly, we note that RT-Sketch shows some inherent biases towards performing certain skills it was trained on, and occasionally performs the wrong skill. For a detailed breakdown of RT-Sketch’s limitations and failure modes, please see Appendix A).

## V. CONCLUSION

We propose RT-Sketch, a goal-conditioned policy for manipulation that takes a hand-drawn sketch of the desired scene as input, and outputs actions. To enable such a policy, we first develop a scalable way to generate paired sketch-trajectory training data via an image-to-sketch translation network, and modify the existing RT-1 architecture to take visual information as an input. Empirically, we show that RT-Sketch not only performs on a comparable level to existing language or goal-image conditioning policies for a number of manipulation skills, but is amenable to different degrees of sketch fidelity, and more robust to visual distractors or ambiguities. Future work will focus on extending hand-drawn sketches to more structured representations, like schematics or diagrams for assembly tasks. While powerful, sketches are not without their own limitations – namely ambiguity due to omitted details or poor quality sketches. In the future, we are excited by avenues for multimodal goal specification that can leverage the benefits of language, sketches, and other modalities to jointly resolve ambiguity from any single modality alone.

## REFERENCES

- [1] Jacob Andreas, Dan Klein, and Sergey Levine. Modular multitask reinforcement learning with policy sketches. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 166–175. PMLR, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/andreas17a.html>.
- [2] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, Pieter Abbeel, and Wojciech Zaremba. Hind-sight experience replay. In *31st Conference on Neural Information Processing Systems (NIPS 2017)*, 2017.
- [3] Brenna D Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. A survey of robot learning from demonstration. *Robotics and autonomous systems*, 57(5):469–483, 2009.

- [4] Christine M Barber, Robin J Shucksmith, Bruce MacDonald, and Burkhard C Wünsche. Sketch-based robot programming. In *2010 25th International Conference of Image and Vision Computing New Zealand*, pages 1–8. IEEE, 2010.
- [5] Suneel Belkhale, Yuchen Cui, and Dorsa Sadigh. Data quality in imitation learning. *arXiv preprint arXiv:2306.02437*, 2023.
- [6] Ayan Kumar Bhunia, Viswanatha Reddy Gajjala, Subhadeep Koley, Rohit Kundu, Aneeshan Sain, Tao Xiang, and Yi-Zhe Song. Doodle it yourself: Class incremental learning by drawing a few sketches. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2293–2302, 2022.
- [7] Ayan Kumar Bhunia, Subhadeep Koley, Amandeep Kumar, Aneeshan Sain, Pinaki Nath Chowdhury, Tao Xiang, and Yi-Zhe Song. Sketch2saliency: Learning to detect salient objects from human drawings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2733–2743, 2023.
- [8] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence, Chuyuan Fu, Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Kehang Han, Karol Hausman, Alex Herzog, Jasmine Hsu, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Lisa Lee, Tsang-Wei Edward Lee, Sergey Levine, Yao Lu, Henryk Michalewski, Igor Mordatch, Karl Pertsch, Kanishka Rao, Krista Reymann, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Pierre Sermanet, Jaspier Singh, Anikait Singh, Radu Soricut, Huong Tran, Vincent Vanhoucke, Quan Vuong, Ayzaan Wahid, Stefan Welker, Paul Wohlhart, Jialin Wu, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *arXiv preprint arXiv:2307.15818*, 2023.
- [9] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alexander Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jor-nell Quiambao, Kanishka Rao, Michael S Ryoo, Grecia Salazar, Pannag R Sanketi, Kevin Sayed, Jaspier Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan H Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. RT-1: Robotics Transformer for Real-World Control at Scale. In *Proceedings of Robotics: Science and Systems*, Daegu, Republic of Korea, July 2023. doi: 10.15607/RSS.2023.XIX.025.
- [10] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023.
- [11] Kaylee Burns, Tianhe Yu, Chelsea Finn, and Karol Hausman. Robust manipulation with spatial features. In *CoRL 2022 Workshop on Pre-training Robot Learning*, 2022.
- [12] Pinaki Nath Chowdhury, Ayan Kumar Bhunia, Aneeshan Sain, Subhadeep Koley, Tao Xiang, and Yi-Zhe Song. What can human sketches do for object detection? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15083–15094, 2023.
- [13] Pinaki Nath Chowdhury, Ayan Kumar Bhunia, Aneeshan Sain, Subhadeep Koley, Tao Xiang, and Yi-Zhe Song. Scenetrilogy: On human scene-sketch and its complementarity with photo and text. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10972–10983, 2023.
- [14] Yuchen Cui, Scott Niekum, Abhinav Gupta, Vikash Kumar, and Aravind Rajeswaran. Can foundation models perform zero-shot task specification for robot manipulation? In *Learning for Dynamics and Control Conference*, pages 893–905. PMLR, 2022.
- [15] Yuchen Cui, Siddharth Karamcheti, Raj Palleti, Nidhya Shivakumar, Percy Liang, and Dorsa Sadigh. No, to the right: Online language corrections for robotic manipulation via shared autonomy. In *Proceedings of the 2023 ACM/IEEE International Conference on Human-Robot Interaction*, pages 93–101, 2023.
- [16] Michael Danielczuk, Andrey Kurenkov, Ashwin Balakrishna, Matthew Matl, David Wang, Roberto Martín-Martín, Animesh Garg, Silvio Savarese, and Ken Goldberg. Mechanical search: Multi-step retrieval of a target object occluded by clutter. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 1614–1621. IEEE, 2019.
- [17] Yiming Ding, Carlos Florensa, Pieter Abbeel, and Mariano Phielipp. Goal-conditioned imitation learning. *Advances in neural information processing systems*, 32, 2019.
- [18] Aditya Ganapathi, Priya Sundaresan, Brijen Thananjeyan, Ashwin Balakrishna, Daniel Seita, Jennifer Grannen, Minho Hwang, Ryan Hoque, Joseph E Gonzalez, Nawid Jamali, et al. Learning dense visual correspondences in simulation to smooth and fold real fabrics. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11515–11522. IEEE, 2021.
- [19] Jiayuan Gu, Sean Kirmani, Paul Wohlhart, Yao Lu, Montserrat Gonzalez Arenas, Kanishka Rao, Wenhao Yu, Chuyuan Fu, Keerthana Gopalakrishnan, Zhuo Xu, et al. Rt-trajectory: Robotic task generalization via hindsight trajectory sketches. *arXiv preprint arXiv:2311.01977*, 2023.

- [20] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017.
- [21] Eric Jang, Alex Irpan, Mohi Khansari, Daniel Kappler, Frederik Ebert, Corey Lynch, Sergey Levine, and Chelsea Finn. Bc-z: Zero-shot task generalization with robotic imitation learning. In *Conference on Robot Learning*, pages 991–1002. PMLR, 2022.
- [22] Yunfan Jiang, Agrim Gupta, Zichen Zhang, Guanzhi Wang, Yongqiang Dou, Yanjun Chen, Li Fei-Fei, Anima Anandkumar, Yuke Zhu, and Linxi Fan. Vima: General robot manipulation with multimodal prompts. *arXiv preprint arXiv:2210.03094*, 2022.
- [23] Siddharth Karamcheti, Suraj Nair, Annie S Chen, Thomas Kollar, Chelsea Finn, Dorsa Sadigh, and Percy Liang. Language-driven representation learning for robotics. *arXiv preprint arXiv:2302.12766*, 2023.
- [24] Subhadeep Koley, Ayan Kumar Bhunia, Aneeshan Sain, Pinaki Nath Chowdhury, Tao Xiang, and Yi-Zhe Song. Picture that sketch: Photorealistic image generation from abstract sketches. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6850–6861, 2023.
- [25] Zixing Lei, Yiming Zhang, Yuxin Xiong, and Siheng Chen. Emergent communication in interactive sketch question answering. *arXiv preprint arXiv:2310.15597*, 2023.
- [26] Mengtian Li, Zhe Lin, Radomir Mech, Ersin Yumer, and Deva Ramanan. Photo-sketching: Inferring contour drawings from images. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1403–1412. IEEE, 2019.
- [27] Rensis Likert. A technique for the measurement of attitudes. *Archives of Psychology*, 1932.
- [28] Corey Lynch and Pierre Sermanet. Language conditioned imitation learning over unstructured data. *arXiv preprint arXiv:2005.07648*, 2020.
- [29] Corey Lynch, Ayzaan Wahid, Jonathan Tompson, Tianli Ding, James Betker, Robert Baruch, Travis Armstrong, and Pete Florence. Interactive language: Talking to robots in real time. *IEEE Robotics and Automation Letters*, 2023.
- [30] Lucas Manuelli, Wei Gao, Peter Florence, and Russ Tedrake. kpm: Keypoint affordances for category-level robotic manipulation. In *The International Symposium of Robotics Research*, pages 132–157. Springer, 2019.
- [31] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- [32] David Porfirio, Laura Stegner, Maya Cakmak, Allison Sauppé, Aws Albarghouthi, and Bilge Mutlu. Sketching robot programs on the fly. In *Proceedings of the 2023 ACM/IEEE International Conference on Human-Robot Interaction, HRI '23*, page 584–593, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9781450399647. doi: 10.1145/3568162.3576991. URL <https://doi.org/10.1145/3568162.3576991>.
- [33] Shuwen Qiu, Sirui Xie, Lifeng Fan, Tao Gao, Jungseock Joo, Song-Chun Zhu, and Yixin Zhu. Emergent graphical conventions in a visual communication game. *Advances in Neural Information Processing Systems*, 35:13119–13131, 2022.
- [34] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [35] Moritz Reuss, Maximilian Li, Xiaogang Jia, and Rudolf Lioutikov. Goal-conditioned imitation learning using score-based diffusion policies. *Robotics: Science and Systems (RSS)*, 2023.
- [36] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, June 2022.
- [37] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 627–635. JMLR Workshop and Conference Proceedings, 2011.
- [38] Michael Ryoo, AJ Piergiovanni, Anurag Arnab, Mostafa Dehghani, and Anelia Angelova. Tokenlearner: Adaptive space-time tokenization for videos. *Advances in Neural Information Processing Systems*, 34:12786–12797, 2021.
- [39] Lin Shao, Toki Migimatsu, Qiang Zhang, Karen Yang, and Jeannette Bohg. Concept2robot: Learning manipulation concepts from instructions and human demonstrations. In *Proceedings of Robotics: Science and Systems (RSS)*, 2020.
- [40] Irwin Sobel. An isotropic 3x3 image gradient operator. *Presentation at Stanford A.I. Project 1968*, 1968.
- [41] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- [42] Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral cloning from observation. *arXiv preprint arXiv:1805.01954*, 2018.
- [43] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [44] Yael Vinker, Yuval Alaluf, Daniel Cohen-Or, and Ariel Shamir. Clipscene: Scene sketching with differ-

ent types and levels of abstraction. *arXiv preprint arXiv:2211.17256*, 2022.

- [45] Yael Vinker, Ehsan Pajouheshgar, Jessica Y Bo, Roman Christian Bachmann, Amit Haim Bermano, Daniel Cohen-Or, Amir Zamir, and Ariel Shamir. Clipasso: Semantically-aware object sketching. *ACM Transactions on Graphics (TOG)*, 41(4):1–11, 2022.
- [46] Kevin Zakka, Andy Zeng, Johnny Lee, and Shuran Song. Form2fit: Learning shape priors for generalizable assembly from disassembly. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9404–9410. IEEE, 2020.
- [47] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023.

## APPENDIX

In this section, we provide further details on the visual goal representations RT-Sketch sees at train and test time (Appendix A), qualitative visualizations of experimental roll-outs (Appendix A), limitations (Appendix A) of RT-Sketch, as well as the interfaces used for data annotation, evaluation, and human assessment (Appendix A).

Since the main bottleneck to training a sketch-to-action policy like RT-Sketch is collecting a dataset of paired trajectories and goal sketches, we first train an image-to-sketch translation network  $\mathcal{T}$  mapping image observations  $o_i$  to sketch representations  $g_i$ , discussed in Section III. To train  $\mathcal{T}$ , we first take a pre-trained network for sketch-to-image translation [26] trained on the ContourDrawing dataset of paired images and edge-aligned sketches (Fig. 5). This dataset contains  $L^{(i)} = 5$  crowdsourced sketches per image for 1000 images. By pre-training on this dataset, we hope to embed a strong prior in  $\mathcal{T}$  and accelerate learning on our much smaller dataset. Next, we finetune  $\mathcal{T}$  on a dataset of 500 manually drawn line sketches for RT-1 robot images. We visualize a few examples of our manually sketched goals in Fig. 6 under ‘Line Drawings’.

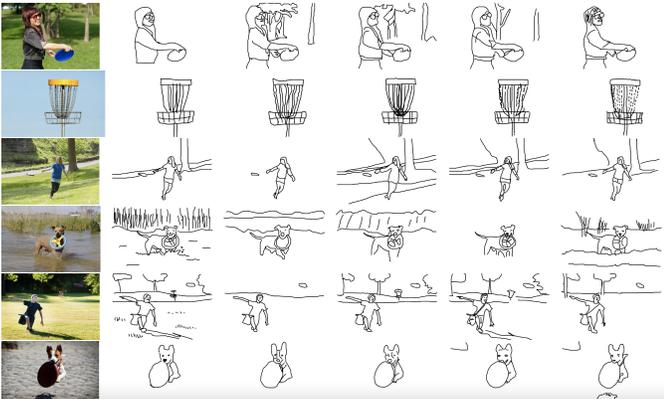


Fig. 5: **ContourDrawing Dataset:** We visualize 6 samples from the ContourDrawing Dataset from [26]. For each image, 5 separate annotators provide an edge-aligned sketch of the scene by outlining on top of the original image. As depicted, annotators are encouraged to preserve main contours of the scene, but background details or fine-grained geometric details are often omitted. Li et al. [26] then train an image-to-sketch translation network  $\mathcal{T}$  with a loss that encourages aligning with at least one of the given reference sketches (Eq. (2)).

Notably, while we only train  $\mathcal{T}$  to map an image to a black-and-white line sketch  $\hat{g}_i$ , we consider various augmentations  $\mathcal{A}$  on top of generated goals to simulate sketches with varied colors, affine and perspective distortions, and levels of detail. Fig. 6 visualizes a few of these augmentations, such as automatically colorizing black-and-white sketches by superimposing a blurred version of the original RGB image, and treating an edge-detected version of the original image as a generated sketch to simulate sketches with a lot of details. We generate a dataset for training RT-Sketch by ‘sketchifying’ hind-sight relabeled goal images via  $\mathcal{T}$  and  $\mathcal{A}$ .

Although RT-Sketch is only trained on generated line sketches, colorized line sketches, edge-detected images, and

goal images, we find that it is able to handle sketches of even greater diversity. This includes non-edge aligned free-hand sketches and sketches with color infills, like those shown in Fig. 6.

### A. Alternate Image-to-Sketch Techniques

The choice of image-to-sketch technique we use is critical to the overall success of the RT-Sketch pipeline. We experiment with various other techniques before converging on the above approach.

Recently, two recent works, CLIPAsso [45] and CLIPAScene [44] explore methods for automatically generating a sketch from an image. These works pose sketch generation as inferring the parameters of Bezier curves representing ‘strokes’ in order to produce a generated sketch with maximal CLIP-similarity to a given input image. These methods perform a per-image optimization to generate a plausible sketch, rather than a global batched operation across many images, limiting their scalability. Additionally, they are fundamentally more concerned with producing high-quality, aesthetically pleasing sketches which capture a lot of extraneous details.

We, on the other hand, care about producing a minimal but reasonable-quality sketch. The second technique we explore is trying the pre-trained Photosketching GAN [26] on internet data of paired images and sketches. However, this model output does not capture object details well, likely due to not having been trained on robot observations, and contains irrelevant sketch details. Finally, by finetuning this PhotoSketching GAN on our own data, the outputs are much closer to real, hand-drawn human sketches that capture salient object details as minimally as possible. We visualize these differences in Fig. 7.

To further interpret RT-Sketch’s performance, we provide visualizations of the precision metrics and experimental roll-outs. In Fig. 8, we visualize the degree of alignment RT-Sketch achieves, as quantified by the pixelwise distance of object centroids in achieved vs. given goal images. In Fig. 9, Fig. 10, Fig. 11, and Fig. 13, we visualize each policy’s behavior for H1, H2, H3 and H4, respectively. Fig. 12 visualizes the four tiers of difficulty in language ambiguity that we analyze for H4.

While RT-Sketch is performant at several manipulation benchmark skills, capable of handling different levels of sketch detail, robust to visual distractors, and unaffected by ambiguous language, it is not without failures and limitations.

In Fig. 15, we visualize the failure modes of RT-Sketch. One failure mode we see with RT-Sketch is occasionally re-trying excessively, as a result of trying to align the scene as closely as possible. For instance, in the top row, Rollout Image 3, the scene is already well-aligned, but RT-Sketch keeps shifting the chip bag which causes some misalignment in terms of the chip bag orientation. Still, this kind of failure is most common with RT-Goal-Image (Table I), and is not nearly as frequent for RT-Sketch. We posit that this could be due to the fact that sketches enable high-level spatial reasoning without over-attending to pixel-level details.

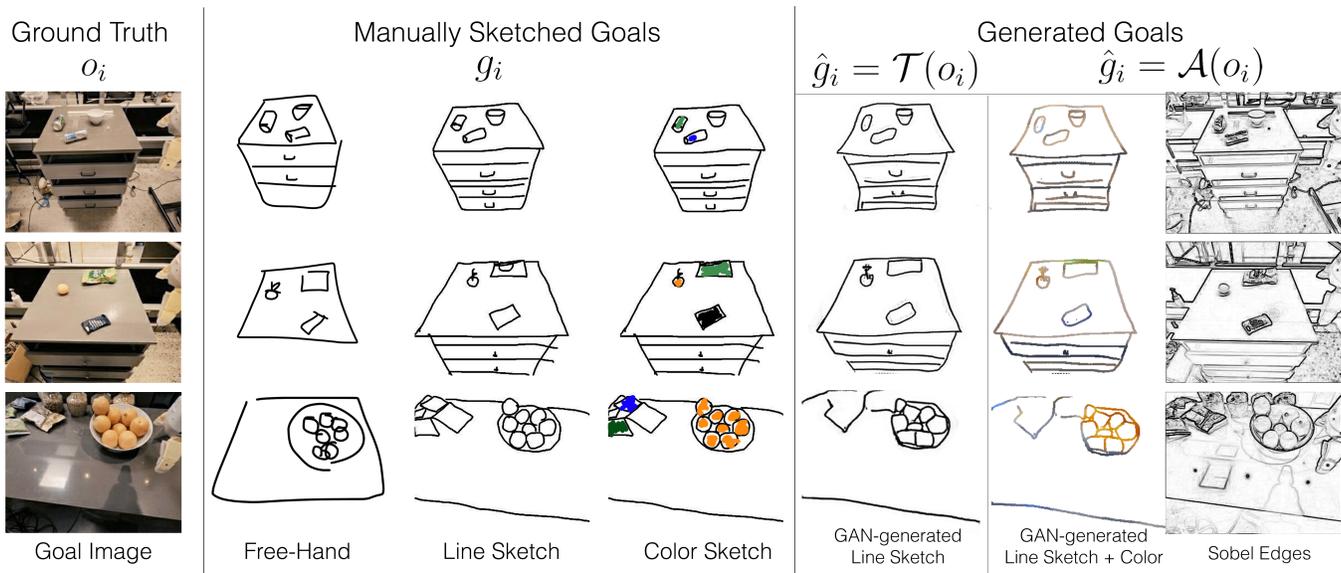


Fig. 6: **Visual Goal Diversity**: RT-Sketch is capable of handling a variety of visual goals at both train and test time. RT-Sketch is trained on generated and augmented images like those shown on the right below ‘Generated Goals’. But it can also interpret free-hand, line sketches, and colored sketches at test time such as those on the left below ‘Manually Sketched Goals’.

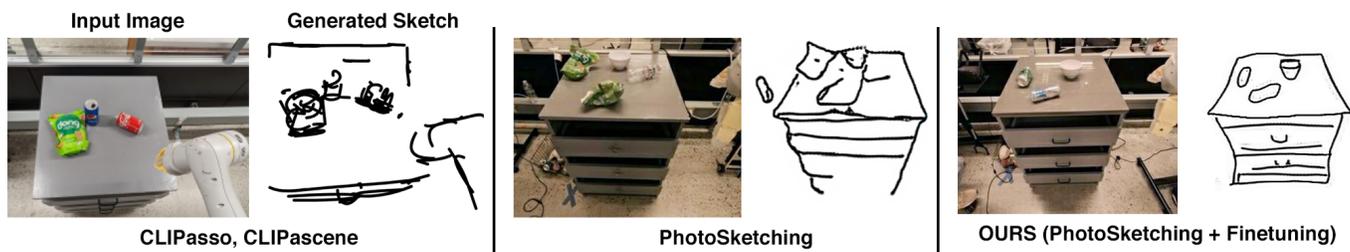


Fig. 7: **Alternate Image-to-Sketch Techniques**

One consequence of spatial reasoning at such a high level, though, is an occasional lack of precision. This is noticeable when RT-Sketch orients items incorrectly (second row) or positions them slightly off, possibly disturbing other items in the scene (third row). This may be due to the fact that sketches are inherently imperfect, which makes it difficult to reason with such high precision.

Finally, we see that RT-Sketch occasionally manipulates the wrong object (rows 4 and 5). Interestingly, we see that a fairly frequent pattern of behavior is to manipulate the wrong object (orange in row 4) to the right target location (near green can in row 4). This may be due to the fact that the sketch-generating GAN has occasionally hallucinated artifacts or geometric details missing from the actual objects. Having been trained on some examples like these, RT-Sketch can mistakenly perceive the wrong object to be aligned with an object drawn in the sketch. However, the sketch still indicates the relative desired spatial positioning of objects in the scene, so in this case RT-Sketch still attempts to align the incorrect object with the proper place.

Finally, the least frequent failure mode is manipulating the wrong object to the wrong target location (i.e. opening the

wrong drawer handle). This is most frequent when the input is a free-hand sketch, and could be mitigated by increasing sketch detail (Table II).

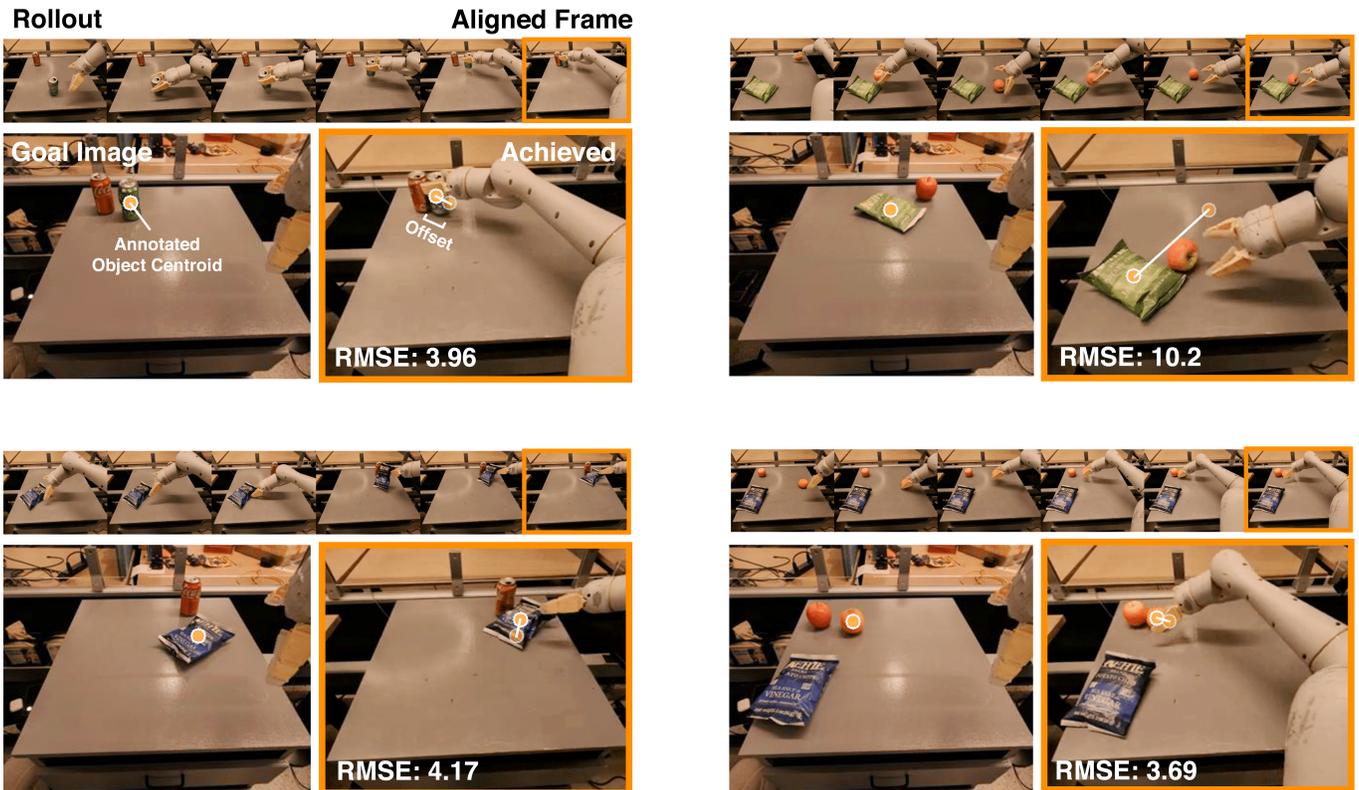


Fig. 8: **Spatial Precision Visualization:** We visualize four trials of RT-Sketch on the Move Near skill, along with the measured spatial precision in terms of RMSE. To evaluate spatial precision, we have a human annotator annotate the frame that is visually most aligned, and then keypoints for the object that was moved in this frame and in the provided reference goal image. For each of the four trials, we visualize the rollout frames until alignment is achieved, along with the labeled object centroids and the offset in achieved vs. desired positions. The upper right example shows a failure of RT-Sketch in which the apple is moved instead of the chip bag, incurring a high RMSE. These visualizations are intended to better contextualize the numbers from [Table I](#).

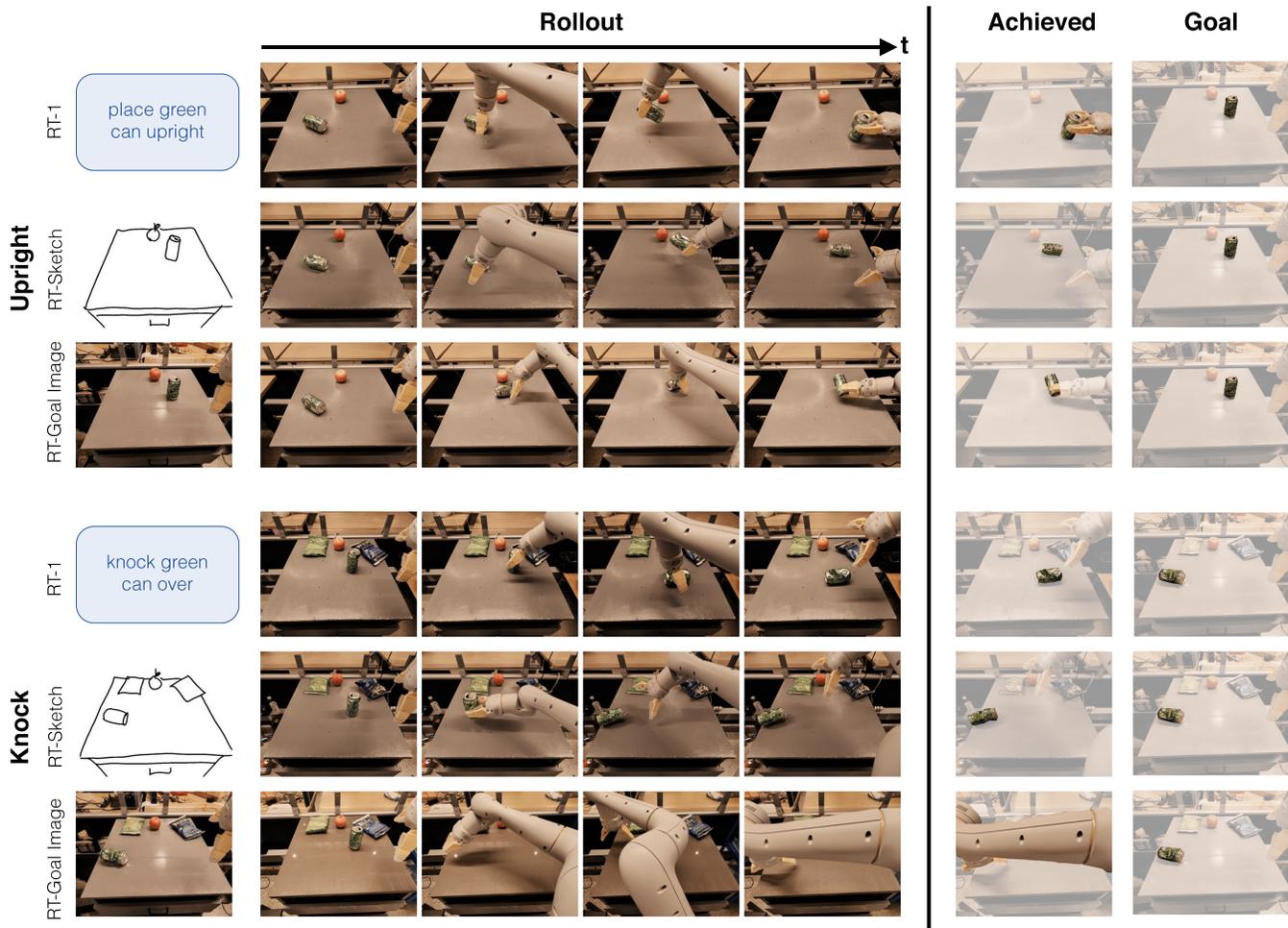


Fig. 9: **H1 Rollout Visualization:** We visualize the performance of RT-1, RT-Sketch, and RT-Goal-Image on two skills from the RT-1 benchmark (*upright* and *knock*). For each skill, we visualize the goal provided as input to each policy, along with the policy rollout. We see that for both skills, RT-1 obeys the semantic task at hand by successfully placing the can upright or sideways, as intended. Meanwhile, RT-Sketch and RT-Goal-Image struggle with orienting the can upright, but successfully knock it sideways. Interestingly, both RT-Sketch and RT-Goal-Image are able to place the can in the desired location (disregarding can orientation) whereas RT-1 does not pay attention to where in the scene the can should be placed. This is indicated by the discrepancy in position of the can in the achieved versus goal images on the right. This trend best explains the anomalous performance of RT-Sketch and RT-Goal-Image in perceived Likert ratings for the upright task (Fig. 3), but validates their comparably higher spatial precision compared to RT-1 across all benchmark skills (Table I).

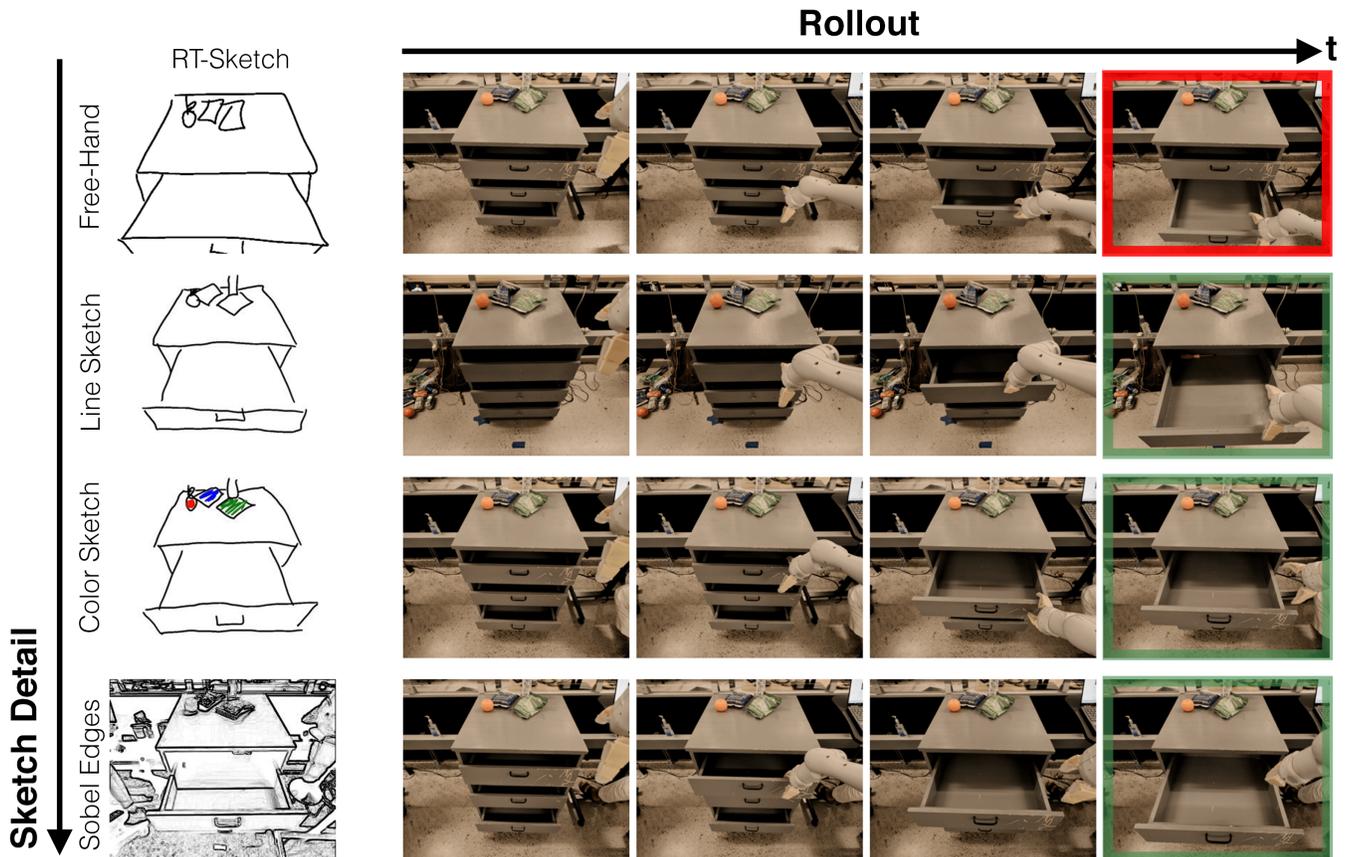


Fig. 10: **H2 Rollout Visualization**: For the *open drawer* skill, we visualize four separate rollouts of RT-Sketch operating from different input types. Free-hand sketches are drawn without outlining over the original image, such that they can contain marked perspective differences, partially obscured objects (drawer handle), and roughly drawn object outlines. Line sketches are drawn on top of the original image using the sketching interface we present in Appendix Fig. 16. Color sketches merely add color infills to the previous modality, and Sobel Edges represent an upper bound in terms of unrealistic sketch detail. We see that RT-Sketch is able to successfully open the correct drawer for any sketch input except the free-hand sketch, without a noticeable performance gain or drop. For the free-hand sketch, RT-Sketch still recognizes the need for opening a drawer, but the differences in sketch perspective and scale can occasionally cause the policy to attend to the wrong drawer, as depicted.

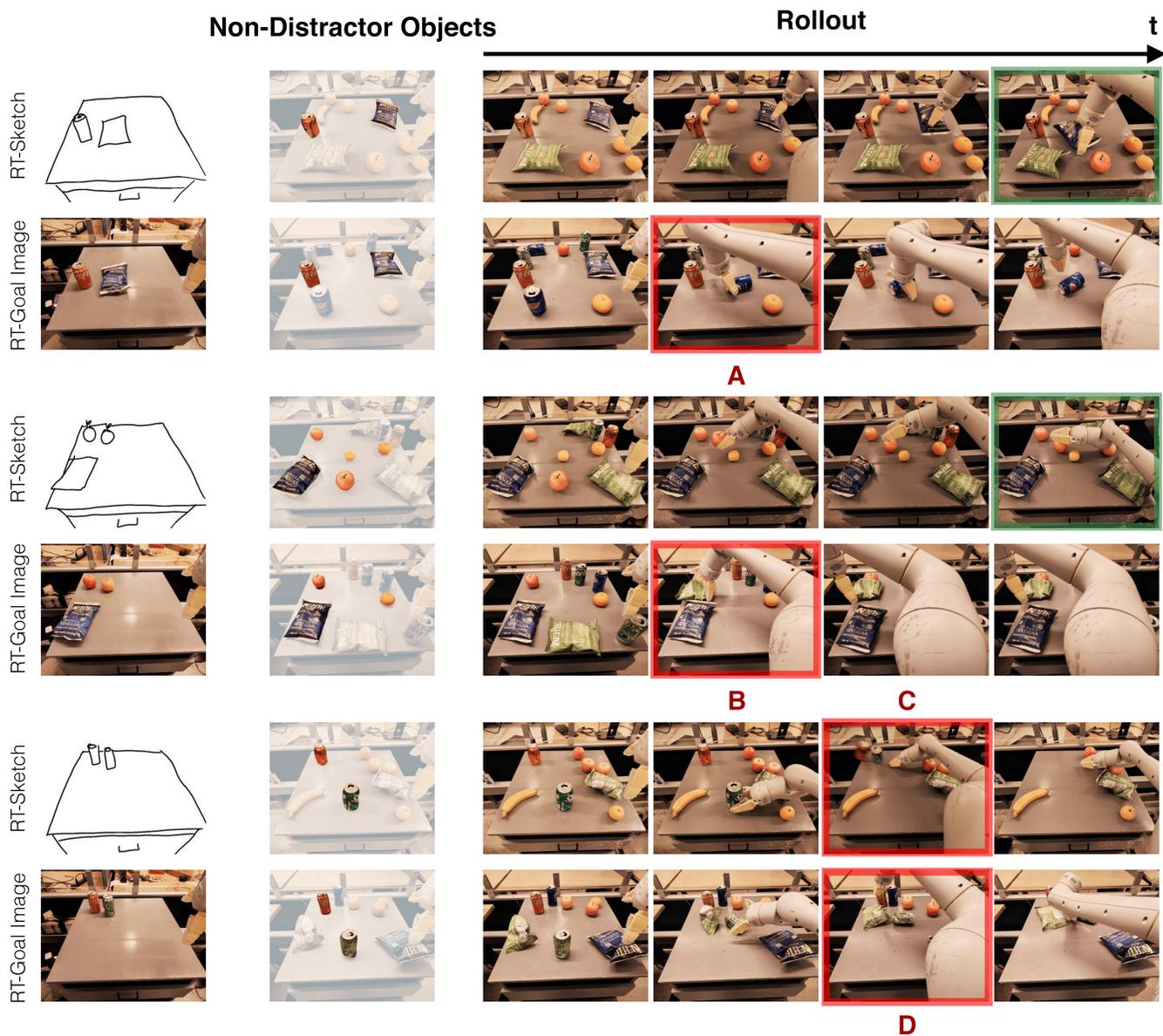


Fig. 11: **H3 Rollout Visualization:** We visualize qualitative rollouts for RT-Sketch and RT-Goal-Image for 3 separate trials of the *move near* skill subject to distractor objects. In Column 2, we highlight the relevant non-distractor objects that the policy must manipulate in order to achieve the given goal. In Trial 1, we see that RT-Sketch successfully attends to the relevant objects and moves the blue chip bag near the coke can. Meanwhile, RT-Goal-Image is confused about which blue object to manipulate, and picks up the blue pepsi can instead of the blue chip bag (A). In Trial 2, RT-Sketch successfully moves an apple near the fruit on the left. A benefit of sketches is their ability to capture instance multimodality, as any of the fruits highlighted in Column 2 are valid options to move, whereas this does not hold for an overspecified goal image. RT-Goal-Image erroneously picks up the green chip bag (B) instead of a fruit. Finally, Trial 3 shows a failure for both policies. While RT-Sketch successfully infers that the green can must be moved near the red one, it accidentally knocks over the red can (C) in the process. Meanwhile, RT-Goal-Image prematurely drops the green can and instead tries to pick the green chip bag (D).

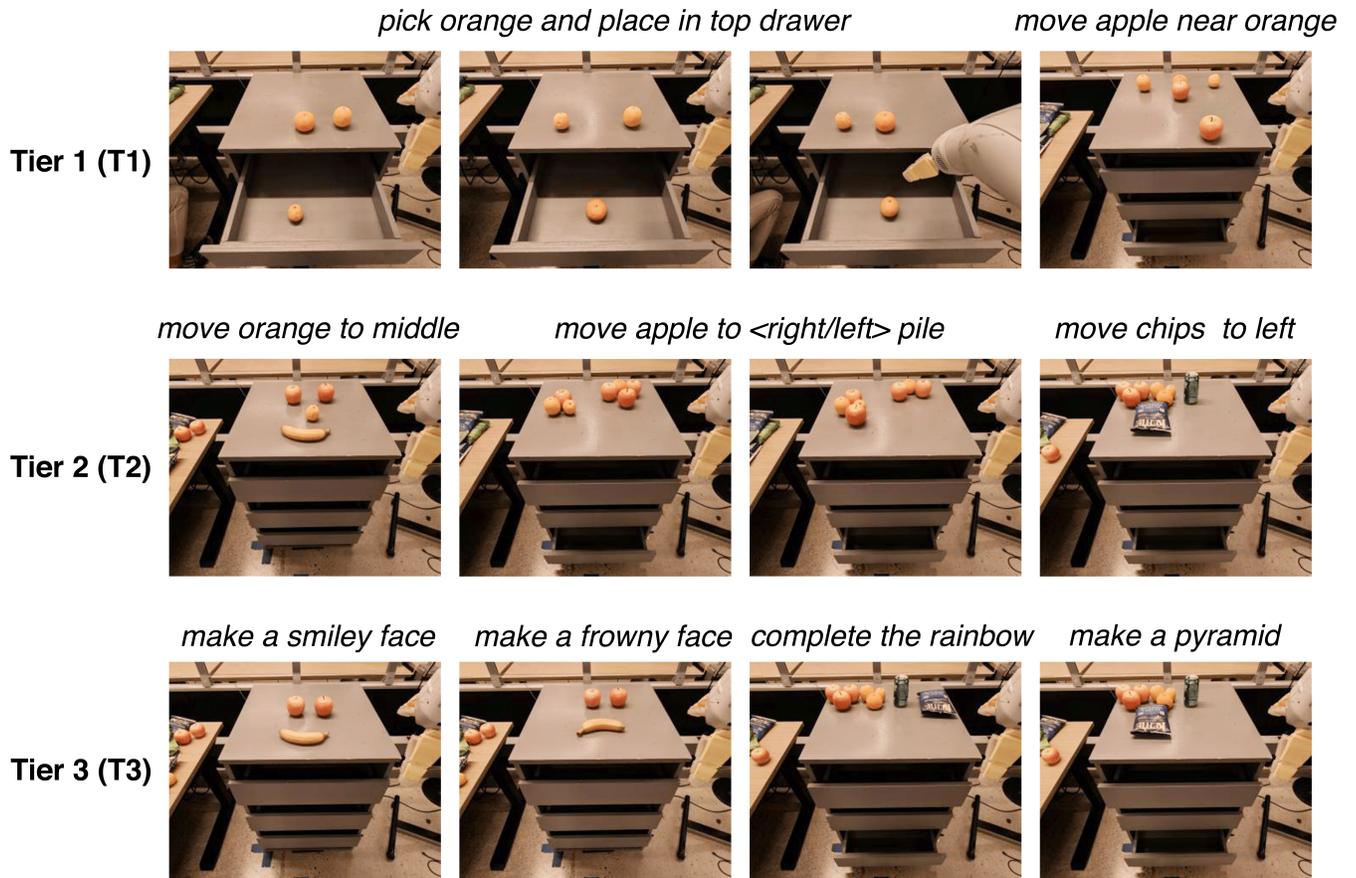


Fig. 12: **H4 Tiers of Difficulty**: To test **H4**, we consider language instructions that are either ambiguous due to the presence of multiple similar object instances (**T1**), are somewhat out-of-distribution for RT-1 (**T2**), or are far out-of-distribution and difficult to specify concretely without lengthier descriptions (**T3**). Each image represents the ground truth goal image paired with the task description.

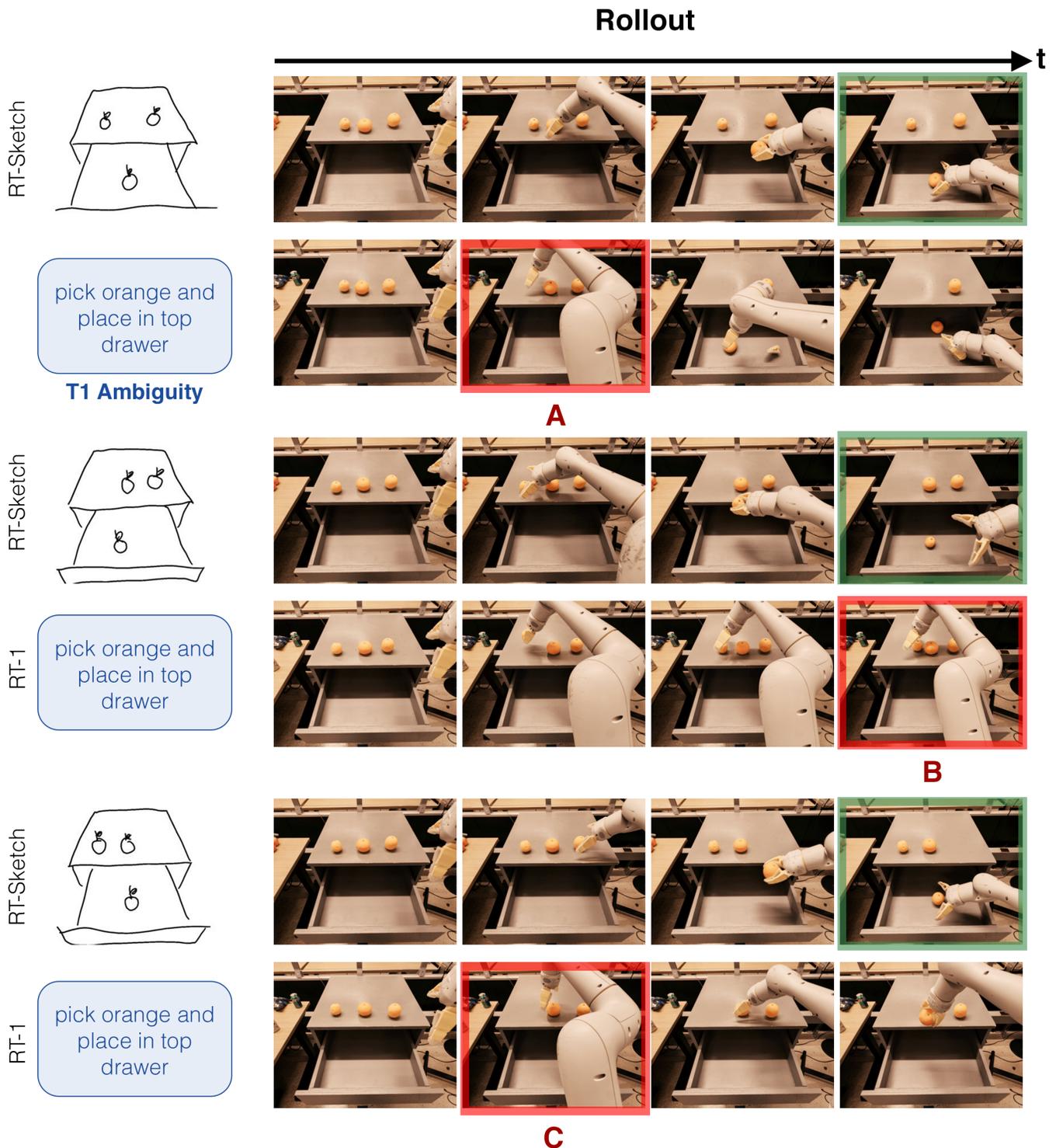


Fig. 13: **H4 Rollout Visualization (T1 as visualized in Fig. 12)**: One source of ambiguity in language descriptions is mentioning an object for which there are multiple instances present. For example, we can easily illustrate three different desired placements of an orange in the drawer via a sketch, but an ambiguous instruction cannot easily specify which orange is relevant to pick and place. In all rollouts, RT-Sketch successfully places the correct orange in the drawer, while RT-1 either picks up the wrong object (A), fails to move to the place location (B), or knocks off one of the oranges (C). Although in this case, the correct orange to manipulate could easily be specified with a spatial relation like *pick up the  $\langle \text{left/middle/right} \rangle$  orange*, we show below in Appendix Fig. 14 that this type of language is still out of the realm of RT-1’s semantic familiarity.

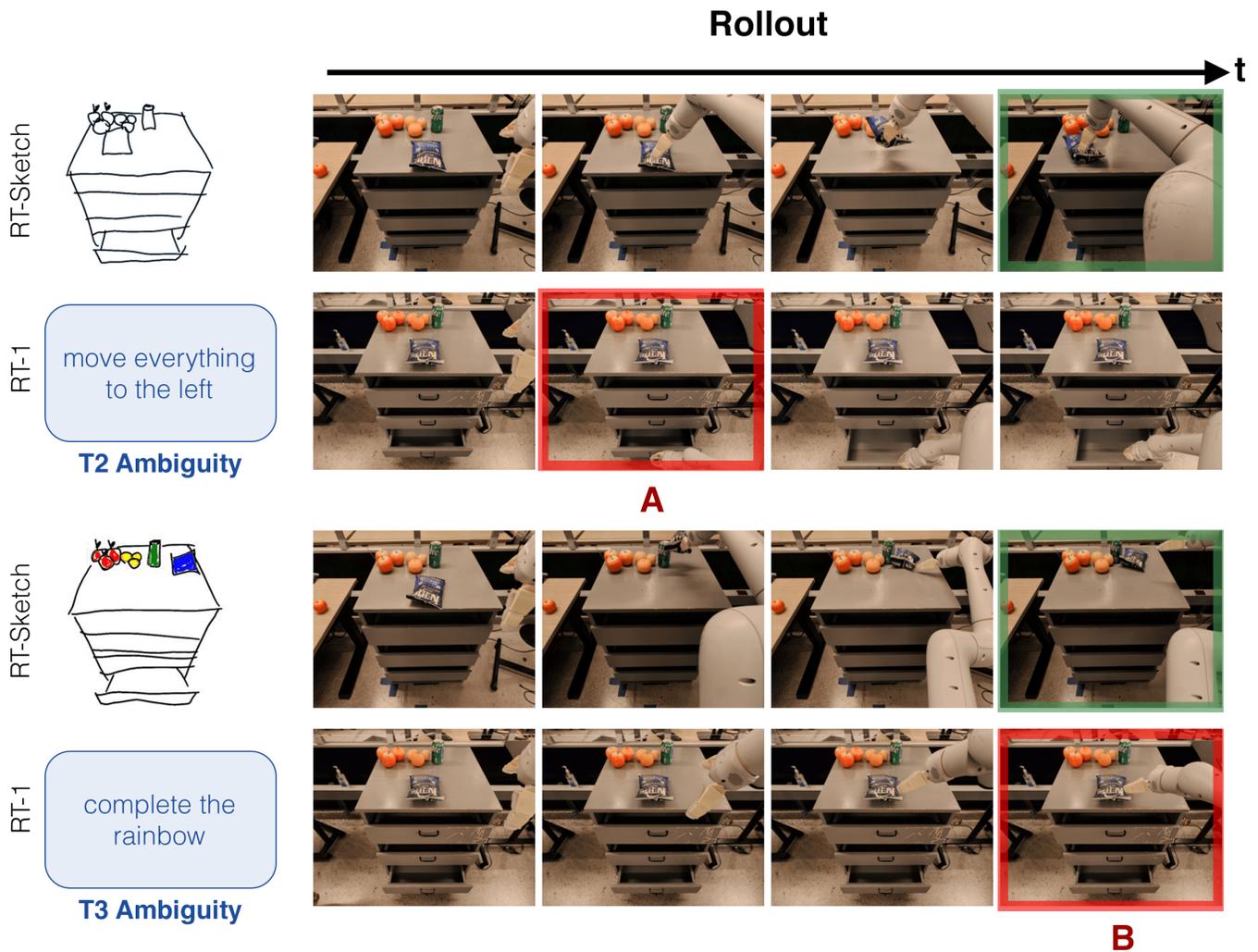


Fig. 14: **H4 Rollout Visualization (T2-3 as visualized in Fig. 12)**: For T2, we consider language with spatial cues that intuitively should help the policy disambiguate in scenarios like the oranges in Fig. 13. However, we find that RT-1 is not trained to handle such spatial references, and this kind of language causes a large distribution shift leading to unwanted behavior. Thus, for the top rollout of trying to move the chip bag to the left where there is an existing pile, RT-Sketch completes the skill without issues, but RT-1 attempts to open the drawer instead of even attempting to rearrange anything on the countertop (A). For T3, we consider language goals that are even more abstract in interpretation, without explicit objects mentioned or spatial cues. Here, sketches are advantageous in their ability to succinctly communicate goals (i.e. visual representation of a rainbow), whereas the corresponding language task string is far too underspecified and OOD for the policy to handle (B).

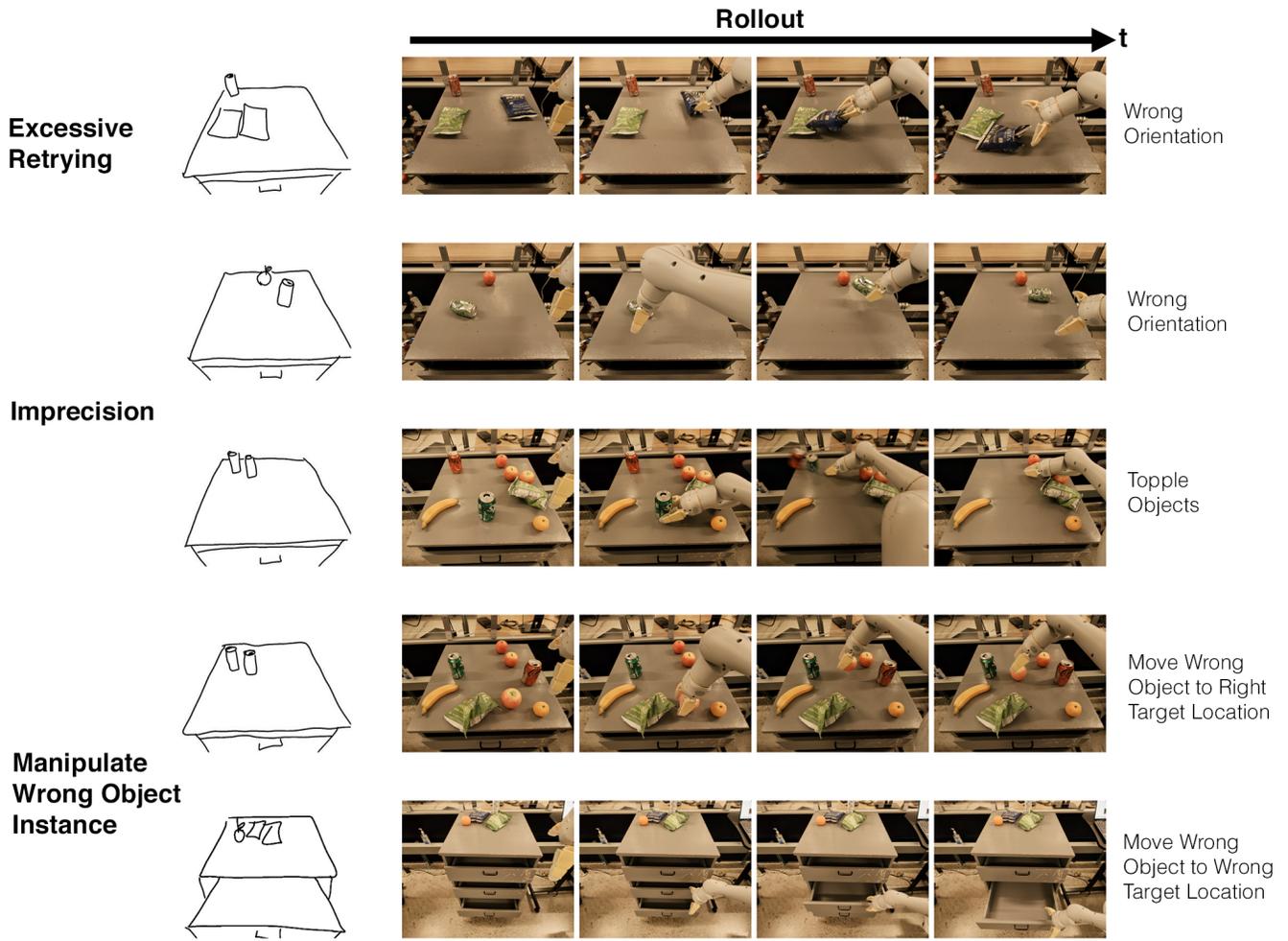


Fig. 15: RT-Sketch Failure Modes

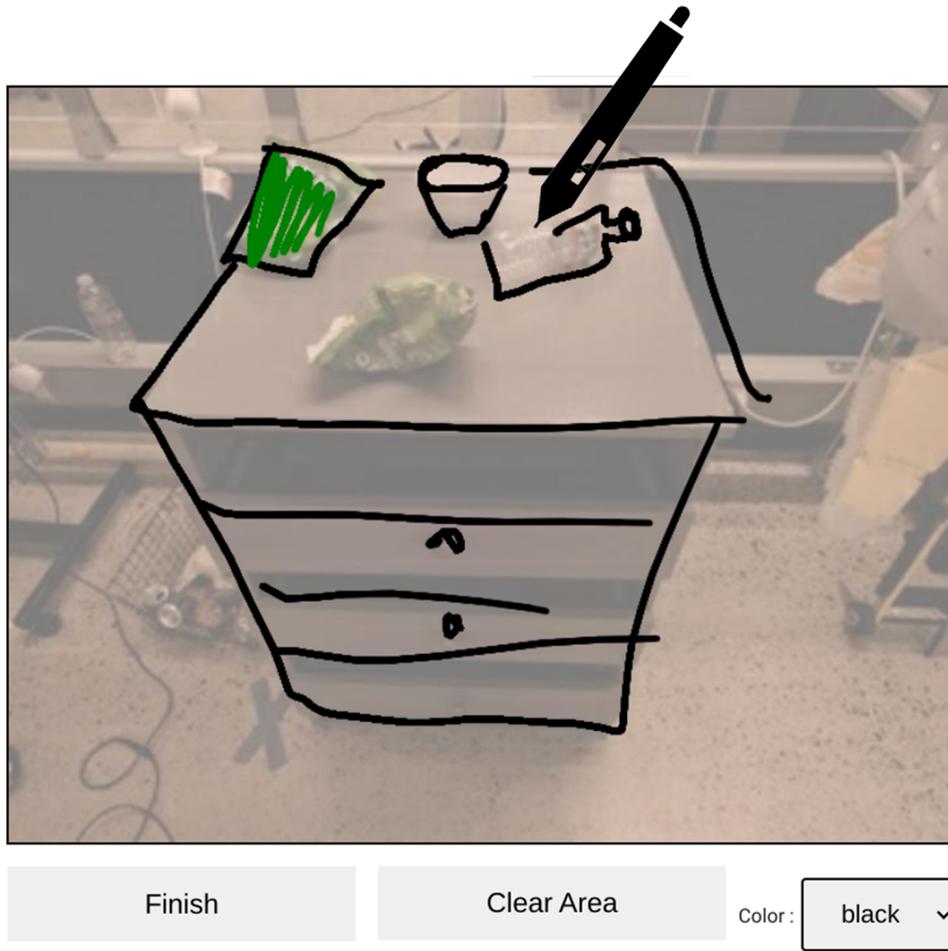


Fig. 16: **Sketching UI**: We design a custom sketching interface for manually collecting paired robot images and sketches with which to train  $\mathcal{T}$ , and for sketching goals for evaluation. The interface visualizes the current robot observation, and provides the ability to draw on a digital screen with a stylus. The interface supports different colors and erasure. We note that intuitively, drawing on top of the image is not an unreasonable assumption to make, since current agent observations are far more readily available than a goal image, for instance. Additionally, the overlay is intended to make the sketching interface easy for the user to provide, without having to eyeball edges for the drawers or handles blindly.

## Q1

### Reference Instruction + Actual Rollout



The robot achieves *semantic alignment* with the goal during the rollout. \*

1    2    3    4    5    6    7

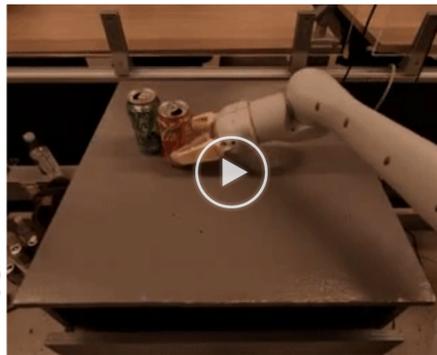
Strongly Disagree                                Strongly Agree

## Q2

### Reference Goal



### Actual Rollout



The robot achieves *spatial alignment* with the goal during the rollout. \*

1    2    3    4    5    6    7

Strongly Disagree                                Strongly Agree

Fig. 17: **Assessment UI**: For all skills and methods, we ask labelers to assess semantic and spatial alignment of the recorded rollout relative to the ground truth semantic instruction and visual goal. We show the interface above, where labelers are randomly assigned to skills and methods (anonymized). The results of these surveys are reported in Fig. 3.