

# Vision-Language Guided Safe Dynamic Manipulation via Object-Centric Transformers

Saumya Saxena 

Oliver Kroemer

Andrea Bajcsy

Carnegie Mellon University  
`{saumyas,abajcsy,okroemer}@andrew.cmu.edu`

**Abstract:** Safely executing dynamic manipulation tasks—like pulling a box from under a stack—in cluttered, open-world environments is challenging. Robots must identify implicit safety constraints (e.g., avoid collapsing stacks or hitting fragile items) while simultaneously reasoning about the combinatorially large number of interactions and their long-term safety consequences (e.g., aggressively pulling an object from under a stack will cause it to fall and hit nearby objects). In this work, we present *VLTSafe*: *Vision-Language guided Transformers for Safe manipulation*, a framework for safe dynamic manipulation that leverages vision-language models (VLMs) to translate semantic safety concepts into geometric constraints, and object-centric transformers to learn generalizable low-level safe policies. To tractably learn policies that scale well to complex cluttered settings, we: (i) utilize a transformer architecture, representing objects as tokens, enabling a *single* policy to be deployed in variable degrees of clutter; (ii) consider diverse combinations of constraint types during training, enabling generalization to novel test-time constraint compositions; and (iii) optimize a reach-avoid reinforcement learning objective to train a parameterized policy that reasons about long-term safety as well as task completion. At test time, *VLTSafe* uses a VLM to identify relevant geometric constraints from RGB images and a textual task description, enabling open-world constraint specification. In both simulation and hardware experiments with a Franka Panda arm, *VLTSafe* infers nuanced constraints and goals (e.g., soft loofahs can safely be pushed out of the way) that cannot be easily identified with hand-designed heuristics. Furthermore, the learned safe policy shows zero-shot generalization to highly cluttered scenes and novel constraint compositions owing to the transformer’s time-varying attention over *relevant* objects.

**Keywords:** Safe learning, reachability, transformers, reinforcement learning, vision-language models

## 1 Introduction

Consider a robot manipulator that must carefully pull a box from under a stack while surrounded by clutter on a table, then place the box on the far side of the table away from any fragile objects. Accomplishing this task safely requires the robot to reason about many types of interactions and constraints, some of which are directly influenced by the dynamics of the interaction (the top boxes can be flipped over or thrown off the table if the bottom box is pulled aggressively), and others that are semantically informed by the semantic properties of interacting objects – making gentle contact with fragile objects but freely pushing soft objects out of the way to complete the task. All of these interactions must be accounted for, with their corresponding safety constraints respected, while completing the task of pulling a box from under a stack and placing it at a goal location. However, as the environment becomes more complex, explicitly modeling all possible interaction combinations quickly becomes intractable.

In this work, we seek to compute low-level control policies that can safely perform dynamic manipulation tasks in arbitrary cluttered environments and with nuanced robot-object interaction constraints. The core challenge lies in *tractably* computing a policy which preserves *test-time generalization* to novel constraint compositions, object configurations, and degrees of clutter. To tackle this, we introduce **VLTsafe**, framework for safe dynamic manipulation that leverages vision-language models (VLMs) as a test-time constraint specifier and a sim-to-real training recipe for learning generalizable low-level safe policies.

Specifically, **VLTsafe** builds on an object-centric transformer policy architecture, where each object in the scene is represented as a token enriched with constraint attributes derived from a vision-language model’s (VLM) analysis of the scene. We design a custom attention mask based on the key insight that, in cluttered environments, only a small subset of objects are relevant for safe decision-making at any given time. This attention mechanism guides the transformer to focus on critical interactions without being ‘distracted’ by irrelevant context. We train this parameterized policy entirely in simulation using a reach-avoid reinforcement learning (RL) objective [1], which allows the model to reason about satisfying long-horizon safety constraints while pursuing the task goal. During training, we systematically vary both the level of clutter and the composition of constraints, encouraging the policy to learn how different object interactions impact long-term safety across a wide range of scenarios. The result is a single transformer-based policy that generalizes to scenes with varying numbers and types of objects and constraints—without requiring retraining. At test time, given a single image and a task description, an off-the-shelf VLM identifies relevant safety constraints (e.g., “do not touch fragile objects” or “avoid toppling the stack of boxes”) and goal specifications (e.g., “place the box away from fragile objects”). These high-level instructions are translated into geometric constraint parameters, which are then used to condition the transformer policy for deployment in novel environments.

We evaluate **VLTsafe** in controlled simulation experiments and we zero-shot deploy the policy from sim to real on a 7-DoF Franka Panda arm performing the motivating dynamic manipulation task. We find that **VLTsafe** achieves strong zero-shot generalization to novel constraint compositions and significantly higher levels of clutter than those seen during training. It consistently outperforms baseline methods, achieving higher safe success rates and lower failure rates. Our custom attention masking proves crucial for generalization, particularly due to its ability to focus attention on the most relevant objects in cluttered scenes. Moreover, training with diverse constraint types further improves robustness, enabling effective deployment to novel constraint compositions without the need for retraining.

## 2 Related work

**Safe Control for Robotic Manipulation** Impedance and null space control are widely used to ensure compliant interactions in robotic manipulation. Impedance control modulates stiffness and damping properties to enable safe and adaptive interactions with the environment [2, 3], while null space control allows secondary objectives like collision avoidance to be incorporated without interfering with primary tasks [4, 5]. In contrast, contact-aware controllers explicitly reason about contact interactions to keep interaction forces below a safety threshold [6, 7]. Recent work has focused on learning contact-aware complaint controllers using expert demonstrations [8, 9] and RL [10, 11, 12]. Our work builds on these efforts by enabling contact-aware control through reachability-based methods in cluttered environments.

**Reachability-based Safe Control** Hamilton-Jacobi reachability analysis [13, 14, 15] provides theoretical formulations for finding the solution to nonlinear reach-avoid control problems by minimizing the worst case (minimum over time) loss. Compared to approaches that minimize the cumulative loss over time [16, 17], HJR analysis provides rigorous safety assurances, ensuring that the system avoids unsafe states under worst-case scenarios. However, HJR becomes computationally intractable as the dimension of the state space increases [18]. Recent methods [19, 20, 1] take inspiration from RL-based approaches [21, 22] to extend reachability analysis to higher dimensions using a discounted formulation of the reach-avoid bellman equation. Leveraging HJR analysis for

learning safe policies in dynamic, cluttered manipulation remains underexplored due to the high dimensionality of complex multi-body interactions.

**Semantic Safe Planning using Vision-Language Models** Recent works have explored incorporating semantic safety into task planning. [23] integrate safety prompts in the code-as-policies [24] setup, while [25] use LLMs to decompose high-level tasks into subtasks and verify them against LTL specifications. Other approaches infer user preferences through demonstrations and active queries [26]. Semantic constraints derived from VLMs have been used for planning, with simulators verifying feasibility [27, 28]. However, these approaches primarily focus on quasi-static tasks without long-horizon reasoning. Most relevant to our work is [29], which incorporates dynamic constraints within a formal safety framework. However, it requires manually designing barrier functions for each constraint type, and the policy accounts for all constraints simultaneously, making it impractical for cluttered environments. Our approach leverages RARL to learn safety value functions that can be solved for both *liveness* and *safety* based on a subset of relevant task and safety constraints.

### 3 Problem Formulation and the *VLTsafe* Method

In this work, our goal is to compute low-level robot policies that can safely perform dynamic manipulation tasks in cluttered environments with nuanced object interaction constraints. We first formalize the problem mathematically, revealing the underlying challenges. We then present our approach—*VLTsafe*—which breaks down these challenges by training an object-centric transformer policy that is informed about relevant safety constraints at test-time by a vision-language model. Our proposed approach is shown in Fig. 1.

**Notation** We model the robot’s state as the position and velocity of the end-effector  $s^{\text{EE}} = [x^{\text{EE}}, \dot{x}^{\text{EE}}]$  where  $x^{\text{EE}} \in \mathbb{R}^3$  is the Cartesian position and  $\dot{x}^{\text{EE}} \in \mathbb{R}^3$  is the velocity. Let each object’s state be represented by  $s_i^o = [x_i^o, \dot{x}_i^o]$ ,  $i \in \{1, \dots, N\}$  where  $N$  is the total number of objects in the cluttered scene. We denote the *full state* of the robot and the objects via  $s = [s^{\text{EE}}, s_1^o, \dots, s_N^o] \in \mathcal{S}$  where  $\mathcal{S}$  is the full state space. Finally, let the robot’s actions  $a \in \mathcal{A}$  be represented as the Cartesian displacement of the end effector  $\Delta x^{\text{EE}} \in \mathbb{R}^3$ . In this work, we do not assume access to an analytic dynamics model  $s_{t+1} = f(s_t, a_t)$ , but instead rely on a high-fidelity simulator [30] to evolve the robot and object states as a result of the robot’s actions and physical interaction.

**Safety Constraint & Task Target Representation** We represent safety constraints as a *failure set*,  $\mathcal{F} \subset \mathcal{S}$ , which encodes the forbidden robot-object, object-object or object-environment states. For example, the failure set can prohibit high-velocity contact with fragile objects while permitting arbitrary interactions with soft objects. We seek robot policies that not only comply with this safety specification but are also guaranteed to complete a task. Let the *target set*,  $\mathcal{T} \subset \mathcal{S}$ , be the set of all states that satisfy the robot’s task (i.e., reach a goal). For example, placing a book on an empty shelf.

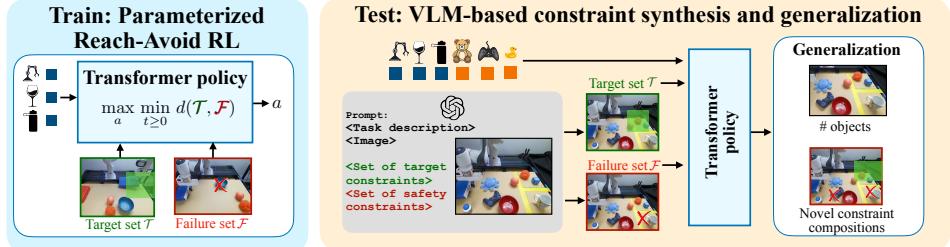
**Goal: A Safe Dynamic Manipulation Policy Adaptable to Test-time Objects and Constraints** Our goal is to find a low-level robot policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  which respects the safety constraints  $\mathcal{F}$  while also achieving the desired task target  $\mathcal{T}$ . To enable test-time adaptation to relevant constraints and targets in the deployment environment, we parameterize the policy via the constraints and targets, denoted by  $\pi_{\mathcal{F}, \mathcal{T}}$ . Formally, we model this problem as a reach-avoid problem [14, 15, 1]. In order to compute the policy which satisfies both properties, we pose a Hamilton-Jacobi reachability problem [13]. Here, the failure set is encoded via the function  $g(\cdot; \mathcal{F})$  and the target set via  $\ell(\cdot; \mathcal{T})$  where

$$g(s; \mathcal{F}) < 0 \iff s \in \mathcal{F}, \quad \ell(s; \mathcal{T}) \geq 0 \iff s \in \mathcal{T}.$$

To jointly account for *multiple* constraints, we define  $g(s; \mathcal{F}) := \min(g_1(s), \dots, g_q(s))$  where  $q$  is the number of safety constraints and to account for *multiple* task goals, we define  $\ell(s; \mathcal{T}) := \min(\ell_1(s), \dots, \ell_p(s))$  where  $p$  is the number of task objectives. Note that both of these functions must be parameterized by the relevant constraint or goal specification.

Following [14, 1], we formulate our policy optimization problem as:

$$V(s; \mathcal{F}, \mathcal{T}) := \max_{\pi_{\mathcal{F}, \mathcal{T}}} \left\{ \max_{t \geq 0} \min \left\{ \ell(\xi_s^\pi(t); \mathcal{T}), \min \{g(\xi_s^\pi(t); \mathcal{F})\} \right\} \right\}. \quad (1)$$



**Figure 1: System overview:** We train a parameterized object-centric transformer policy using reach-avoid reinforcement learning for diverse combination of task and safety constraints. At test time, given an RGB image, task description, and a predefined constraint vocabulary, a vision-language model (VLM) infers the relevant semantic goals and safety constraints for each object. These are used to parameterize the learned policy, enabling zero-shot generalization to novel constraint compositions and highly cluttered scenes.

where  $\xi_s^\pi$  is a state trajectory starting from state  $s$  and executing  $\pi$ , and we index into this trajectory at any discrete time via  $(t)$  notation. The inner optimization (starting with  $\max_{t \geq 0}$ ) “remembers” if the robot was able to reach the target  $\mathcal{T}$  without violating the constraint  $\mathcal{F}$ . The value function will be positive only if both reaching and avoiding happen successfully. The policy’s goal is to maximize this objective, hence reaching the goal without ever violating the specified constraints.

**Challenge: Tractable Computation while Preserving Test-Time Generalization** Solving the above optimization over all possible objects, tasks, and safety constraints is intractable. While restricting to a predefined subset simplifies the problem, it limits safe deployment in novel environments. Therefore, we need an approach that generalizes safe policies across diverse objects and constraint compositions to enable safe deployment in open-world, cluttered environments.

### 3.1 Offline: Safe Policy Learning via Object-Centric Transformers & Reach-Avoid RL

**Policy Representation: Masked Object-centric Transformer Architecture** We propose an object-centric transformer-based model [31] as our policy representation (Fig. 2). This architecture offers two key advantages. First, its input representation is flexible, supporting a variable number of objects and constraint types at test time. Second, by structuring the transformer’s attention mask we can control how the policy attends to objects in extremely cluttered scenes. Our aim is to learn a policy  $\pi_{\mathcal{F}, \mathcal{T}}$  parameterized via constraints and targets, hence we concatenate the input observations with the parameters of the failure  $\phi$  and target  $\tau$  sets, such that state of each object becomes  $[s, \phi, \tau]$ . These observations are tokenized through a two-layer MLP and subsequently processed using two layers of multi-head self-attention, where the attention mechanism is guided by a *custom attention mask*. We design our custom attention mask based on a key insight: in cluttered environments, only a small subset of objects are “relevant” for safe decision-making at any given time step. Consequently, our custom attention mask prioritizes attention to the most relevant neighboring objects. Since we use multiple attention layers, the mask does not fully block information from more distant objects—it simply biases attention toward closer ones. Specifically, all object tokens attend to the end-effector tokens and vice versa, while object-to-object attention is restricted to their  $k$ -nearest neighbors, where  $k$  is task-dependent and can vary across objects (e.g., in our experiments, we choose  $k = 3$ ), while also allowing objects in the stack to always attend to one another.

**Policy Optimization: Reach-Avoid Reinforcement Learning** With our policy architecture setup, we now turn to policy optimization. Recall the reach-avoid optimization problem we formulated in Equation 1. Exact grid-based dynamic programming solvers are intractable for the high-dimensional state space representation we use in this work. Thus, to tractably compute a reach-avoid policy, we use a reinforcement-learning based relaxation of this optimization. Specifically, we adopt the time-discounted reach-avoid Bellman backup from [1] which discounts the safety value function to induce a contraction mapping, yeilding the reach-avoid problem compatible with off-the-shelf reinforcement learning solvers.

$$Q(s, a; \mathcal{F}, \mathcal{T}) = (1-\gamma) \min \{ \ell(s; \mathcal{T}), g(s; \mathcal{F}) \} + \gamma \min \left\{ g(s; \mathcal{F}), \max \{ \ell(s; \mathcal{T}), Q(s', a'; \mathcal{F}, \mathcal{T}) \} \right\}$$

where  $\gamma \in [0, 1]$  is the discount factor, which can be interpreted as the probability of episode continuation. We utilize Deep Deterministic Policy Gradient (DDPG) [32] as our RL framework wherein the bellman update is modified to incorporate the time-discounted reach-avoid Bellman backup (Equations in Appendix 8.1.). We model both the actor and critic using object-centric transformer architectures (Fig. 2). The actor network encodes observation tokens with a transformer, whose output tokens are aggregated via mean pooling and passed through a 2-layer MLP policy head to produce the next action. Critic network jointly encodes state and action tokens using a separate transformer with two layers of multi-head self-attention. The resulting token embeddings are aggregated and concatenated with the action, then passed through a two-layer MLP value head to estimate the Q-value.

### 3.2 Online: Adapting via Safety Constraints Inferred from Vision-Language Models

The prior section allowed us to pre-compute a low-level robot policy which given the current state of the scene  $s$ , a safety specification  $\mathcal{F}$ , and a goal specification  $\mathcal{T}$  safely performs the task. The question is, at deployment time, how can the robot know what are the “relevant” safety constraints  $\mathcal{F}$  and how they may impact which goals  $\mathcal{T}$  are and aren’t allowed? Here, we propose using a vision-language model (VLM) as an open-world constraint and target specifier. Given the current image of the scene  $\mathcal{I}$  and a text prompt describing the task  $\mathcal{L}_{task}$ , the VLM should use the visual and semantic cues to select the relevant safety and targets that are passed into the low-level policy.

However, a large challenge with using VLMs is the right interface between textual representation (that the VLM is trained to output) and the embodied representation (e.g., robot states, objects, etc.) that our robot policy needs at test-time to adapt. To bridge this gap, we model the selection of constraints/targets as a multiple choice visual question-answering (VQA) problem where the multiple choice options are generated via pre-defined geometric constraint and target functions. In the future, we envision the possibility of these functions to be written by another LLM based model (e.g., as in Code-as-Policies [24]). We first convert the set of all geometric constraints  $\mathcal{T}, \mathcal{F}$  we can generate with our functions to corresponding textual descriptions  $\mathcal{L}_{\mathcal{F}}, \mathcal{L}_{\mathcal{T}}$ . For example, a hard collision-avoidance constraint between two interacting bodies is represented in text as:

$$Math Representation of \mathcal{F}: \quad g(s_{EE}, s_{cup}) = ||x_{EE} - x_{cup}|| - \epsilon \quad (2)$$

$$Text Representation of \mathcal{F}: \quad \mathcal{L}_{\mathcal{F}}(<EE>, <cup>) = <\text{no contact}> \quad (3)$$

For the target set, the geometric objective placing a book on an empty book shelf can be written as:

$$Math Representation of \mathcal{T}: \quad \ell(x_{book}) = (||x_{book} - x_{shelf}|| - \epsilon) \quad (4)$$

$$Text Representation of \mathcal{T}: \quad \mathcal{L}_{\mathcal{T}}(<\text{book}>) = <\text{second shelf}>, \quad (5)$$

if the VLM detects the second shelf to be empty. Note that without a VLM, such open-world specifications would need to be manually specified by an expert designer. Once the set of semantic constraints are selected by the VLM, they can be converted back to their geometric formulations and used to parameterize the learned policy for execution. Constraints are identified at the beginning of each episode ( $t = 0$ ) and kept fixed for the duration of the rollout.

## 4 Simulation Experiments

We first devised a series of simulation experiments to carefully test the impact of each of our design decisions on overall task performance. Specifically, we ask the questions (1) When the training and

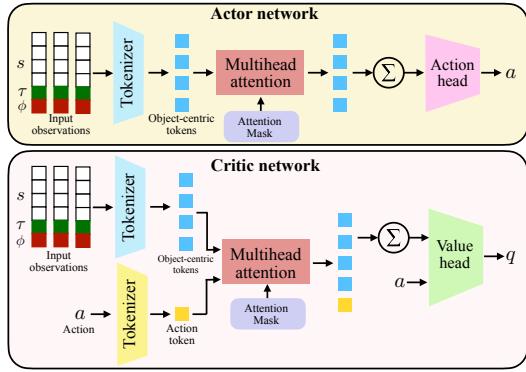


Figure 2: **System architecture:** Actor and critic networks are object-centric transformer-based models. Observations are augmented with parameters of the target  $\tau$  and failure set  $\phi$  and tokenized. These tokens are processed using multi-head self-attention with a custom attention mask. Output tokens are aggregated and used to output action (actor) and value estimate (critic).

test environments are the same, how much does the transformer architecture influence task success?, (2) How well can our approach generalize to varying degrees of novel, test-time object clutter?, and (3) How well does our approach generalize to novel, test-time constraint compositions?

#### 4.1 Experimental Setup

**Environment & Task** In all experiments, we use a tabletop setup where a Franka robot interacts with household objects (Fig. 5). The task involves quickly but safely pulling a cereal box from under another box on a cluttered tabletop and placing it in a goal region. We simulate the robot and environment dynamics with MuJoCo [30] and use the Google Scanned Objects Dataset [33, 34, 35] to simulate everyday objects in the table arena defined in Robosuite [36].

**Safety and Target Representation** Intuitively, the safety constraints capture how the stack of boxes should remain stable as the robot pulls the cereal box from underneath, modeling critical robot-object interactions. Specifically,  $\mathcal{F}$  enforces that the top cereal box does not displace beyond a threshold, the end-effector avoids contact with *fragile* objects, may make limited-velocity contact with *soft* objects, can freely interact with *durable* objects, and must not move over *sensitive* items like a laptop. Since these constraints depend on the semantic properties of objects, the VLM assigns constraint types via the parameterization  $\phi \in \{\text{no-contact}, \text{soft-contact}, \text{any-contact}, \text{do not move over}\}$ . The target set  $\mathcal{T}$  defines successful task completion as the bottom cereal box being fully separated from the top box and placed in a goal region at the far end of the table. This yields two possible targets—*top-goal* and *bottom-goal*—which the VLM selects at deployment time via the parameterization  $\tau \in \{\text{top-goal}, \text{bottom-goal}\}$ . Additional implementation details are provided in Appendix 8.2.

**VLM for Test-Time Safety and Target Specification** We use GPT-4o [37] as the VLM. It is queried once at the beginning of a test episode to identify the constraint types relevant to all objects in the scene and to select the target constraint that is most appropriate for the task. The full prompt for the constraint selection is provided in Appendix 8.3

**Metrics** We measure (1) *safe success rate* (SafeSucc %) defined as the percentage of trajectories that complete the task safely (i.e., satisfy reach and avoid), (2) *stack safety violation rate* (StackFail %) defined as the top cereal block moving outside of safety limits, and (3) *object safety violation rate* (ObjFail %) percentage of objects not in the stack that are unsafely interacted with.

**Training details** Model architectures and hyperparameters are provided in the Appendix 8.2).

Method	Always Same Constraint ( <i>AllFragileConst</i> )			Random Constraints per Episode ( <i>RandomConst</i> )		
	SafeSucc (%) ↑	StackFail (%) ↓	ObjFail (%) ↓	SafeSucc (%) ↑	StackFail (%) ↓	ObjFail (%) ↓
MLP	82	9	9	40.5	43.5	15.5
<b>VLTSafe-NoMask</b>	85.5	3.5	8.5	60.3	29.6	10.1
<b>VLTSafe</b>	<b>95.5</b>	<b>0.5</b>	<b>4</b>	<b>78.5</b>	<b>11.5</b>	<b>4</b>

Table 1: **Policy Architecture vs. Performance.** Comparison of our method against baselines in two training domains, *AllFragileConst* and *RandomConst*

#### 4.2 How much does the transformer policy architecture influence task success?

First, we study how much the transformer design influences the overall task success, assuming the training and test environments are the same.

**Methods** We compare three policy architectures: **VLTSafe** which is our object-centric transformer architecture *with* masking, **VLTSafe-NoMask** which has no masking, **MLP** which is a MLP policy architecture. To ensure a fair comparison, we kept the representation capacity (number of trainable parameters) roughly the same across all models (see Appendix 8.2 for details).

**Training & Test Setup** We train each policy in two training domains. First, we train policies in the *AllFragileConst* domain where all objects are assigned the same constraint type during all episodes:  $\phi = \text{no-contact}$  (i.e. all objects are assumed to be fragile). Next, we train policies in the *RandomConst* domain. Here, constraints for each object are randomly sampled (from the set of

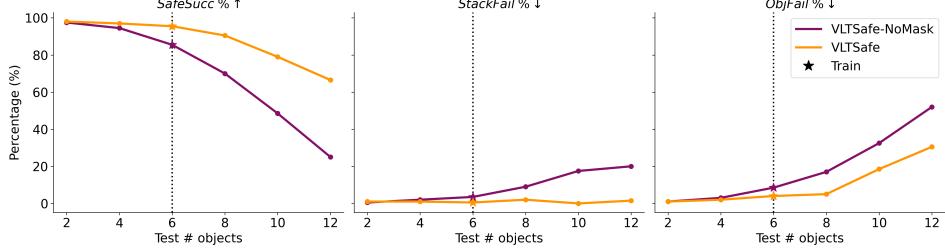


Figure 3: **Generalization: Varying Degrees of Clutter.** Performance variation as the number of objects is varied in the test scenario. Both **VLTsafe** vs **VLTsafe-NoMask** are trained in the *AllFragileConst* domain with six objects. Black dotted line represents the number of objects seen during training.

constraints) at the start of each episode during training. We train with only six objects in the scene and then deploy with six objects in the scene, and the training and the test domain are identical.

**Results** We report all metrics described Sec. 4.1 in Table 1. We find that **VLTsafe** outperforms both baselines with the highest safe success rate and lowest failures rates. Specifically, we find that the use of masking significantly helps the policy. We hypothesize that, without masking, the transformer attends equally to all objects in the scene, making the learning problem unnecessarily difficult. By restricting attention to only the most “relevant” objects through masking, the model can focus on critical near-term interactions—such as between the two cereal boxes — enabling more effective learning of safer policies. This effect is further highlighted in *RandomConst*, where **VLTsafe-NoMask** performs noticeably worse than **VLTsafe**. Furthermore, even with the same representation capacity, the MLP model is unable to learn safe policies that satisfy complex interactive constraints. This further highlights the advantage of using a transformer-based architecture.

#### 4.3 How well can our approach generalize to varying degrees of test-time clutter?

**Methods** Next, we study the ability of the policy to generalize to novel object/clutter amounts at test-time. Since the MLP cannot handle arbitrary number of object at test-time, we only compare **VLTsafe** and **VLTsafe-NoMask**.

**Training & Test Setup** During training, we only consider the *AllFragileConst* domain trained with six objects, as in Section 4.2. At test-time, we deploy each policy in varying degrees of clutter—with 2,4,6,8,10, and 12 objects—to test generalization to number of scene objects.

**Results** In Fig. 3, our quantitative results show that masking plays a large role in enabling the policy to maintain performance and generalize to environments with varying degrees of clutter. Specifically, with **VLTsafe** we see a 29% decrease in the safe success rate when we double the number of objects from 6 to 12, while with **VLTsafe-NoMask** we see a 60.5% decrease. Looking deeper into the failure modes, we see that the main gains come from how **VLTsafe** maintains a near zero-percent StackFail rate as we increase the number of objects. In contrast, for **VLTsafe-NoMask** the StackFail rate increases to 20%. We hypothesize that, similar to 4.2, this is again due to the capability of **VLTsafe** to focus on relevant objects and critical near-term interactions without being ‘distracted’ by the apparent complexity of the task.

#### 4.4 How well does **VLTsafe** generalize to novel, test-time constraint compositions?

**Methods** Next, we want to study how well our approach generalizes to novel constraint compositions at test-time. Here we only test **VLTsafe**.

**Training & Test Setup** We train our policy in two settings: *RandomConst*, where constraints vary randomly between episodes, and *FixedConst*, where constraints are fixed for each object. Training is limited to 6-object scenes, and at test time, we evaluate generalization in the *RandomConst* domain with 6, 8, and 10 objects and novel constraint combinations.

**Results** Recall that at training-time (with six objects), **VLTsafe** achieves a safe success rate of 79.5% in the *FixedConst* domain and a safe success rate of 78.5% in *RandomConst* (star icon in Fig. 4). Note that this test scenario is novel for the policy trained in *FixedConst* even with



**Figure 5: Real-world deployment:** Zero-shot sim-to-real transfer of a policy trained in simulation (*RandomConst*, 6 objects) to a real-world setting with 8 objects. Constraints are first inferred using a VLM. As the robot pulls the blue box from under the red one, it makes allowable contact with the blue plush toy, then lifts to avoid the bowl and places the box safely in the goal region, away from the fragile porcelain mug.

$N=6$ , since during training the policy never saw randomized object constraints. This explains the drop in performance (12.5 %) when this policy is tested in *RandomConst*. For the policy trained with *RandomConst*, only the 8 and 10-object environments are novel and, since the number of objects is larger, the constraint composition is novel. In Fig. 4, we see that **VLTsafe** trained with random combinations of constraints exhibits better generalization capabilities, compared to training with *FixedConst*, since it sees a larger variation in constraint compositions during training. However, we note that this performance gap quickly diminishes as the environment complexity increases and the task becomes significantly more challenging than the train setting.

## 5 Real World Experiments

**Hardware Setup and State Estimation** We use a 7-DoF Franka Emika Panda arm and an Azure Kinect RGB-D camera for both state estimation and image capture for VLM-based constraint inference. State estimation relies on TapNet [38, 39], a point tracking algorithm applied to RGB-D inputs. At the start of each trajectory, we manually annotate object keypoints, a step that could be automated in future work using VLMs. TapNet tracks object states at 17 Hz.

**Sim-to-Real** We deploy policies trained entirely in simulation directly in the real world, without any real-world fine-tuning. The experimental setup—including the robot’s state space and the object configurations on the table—closely matches that of the simulation environment, enabling zero-shot transfer of policies trained in simulation using object-centric observations, directly in the real world. We deploy the policy trained in the *RandomConst* simulation domain. This policy was trained considering six objects in the scene, however we test in scenes containing up to eight objects.

**Results** Fig. 5 shows a real world rollout of our learned policy in a scene with eight objects in clutter. First, a VLM identifies the goal and object constraints: Soft objects are typically assigned the `<any contact>` constraint, while mugs and electronic items—considered safety-critical—are labeled as fragile `<no contact>`. When prompted to choose between two target sets—one containing soft toys, and the other containing a porcelain mug—the VLM selects the former, as interactions with soft objects are safer. During execution, the robot briefly contacts the blue plush toy while extracting the blue box—an acceptable interaction—then lifts the end-effector to avoid colliding with a nearby bowl. It ultimately places the box safely atop the orange loofah and blue ball, both of which are deemed safe for contact. Additional real world experiments are included in the Appendix 8.4.

## 6 Conclusion

We introduce **VLTsafe**, an approach for learning low-level policies to safely perform dynamic tasks in cluttered environments. By leveraging vision-language models (VLMs) for open-world constraint synthesis, our object-centric transformer-based policy demonstrates zero-shot generalization to high degrees of clutter and novel constraint compositions. Training with the reach-avoid reinforcement learning (RARL) objective enables long-horizon safety reasoning and task accomplishment. **VLTsafe** outperforms baseline approaches, achieving high safe success rates, low failure rates, and strong generalization across novel scenarios.



**Figure 4: Generalization: Constraint Compositions.** Evaluating policies trained in *FixedConst* and *RandomConst* domains with six objects, in *RandomConst* domain with increasing number of objects.

## 7 Limitations

Our approach has several limitations. It relies on off-the-shelf state estimation methods to enable learning of object-centric representations. While this enables strong generalization and seamless sim-to-real transfer, it may limit performance in certain scenarios. Future work will explore visual object proposals, which can implicitly capture both geometric and semantic information, reducing reliance on traditional state estimation methods. Additionally, our experiments were constrained to the Cartesian state space, but we aim to extend our approach to 6-DoF manipulation for more dexterous and dynamic tasks.

## References

- [1] K.-C. Hsu, V. Rubies-Royo, C. J. Tomlin, and J. F. Fisac. Safety and liveness guarantees through reach-avoid reinforcement learning. *arXiv preprint arXiv:2112.12288*, 2021. [2](#), [3](#), [4](#), [1](#)
- [2] N. Hogan. Impedance control: An approach to manipulation. In *1984 American control conference*, pages 304–313. IEEE, 1984. [2](#)
- [3] F. J. Abu-Dakka and M. Saveriano. Variable impedance control and learning—a review. *Frontiers in Robotics and AI*, 7:590681, 2020. [2](#)
- [4] H. Sadeghian, L. Villani, M. Keshmiri, and B. Siciliano. Task-space control of robot manipulators with null-space compliance. *IEEE Transactions on Robotics*, 30(2):493–506, 2013. [2](#)
- [5] A. De Luca and L. Ferrajoli. Exploiting robot redundancy in collision detection and reaction. In *2008 IEEE/RSJ international conference on intelligent robots and systems*, pages 3299–3305. IEEE, 2008. [2](#)
- [6] Z. Li, M. Zamora, H. Zheng, and S. Coros. Embracing safe contacts with contact-aware planning and control. *arXiv preprint arXiv:2308.04323*, 2023. [2](#)
- [7] X. Zhu, W. Lian, B. Yuan, C. D. Freeman, and M. Tomizuka. Allowing safe contact in robotic goal-reaching: Planning and tracking in operational and null spaces. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 8120–8126. IEEE, 2023. [2](#)
- [8] L. Peternel, T. Petrič, and J. Babič. Human-in-the-loop approach for teaching robot assembly tasks using impedance control interface. In *2015 IEEE international conference on robotics and automation (ICRA)*, pages 1497–1502. IEEE, 2015. [2](#)
- [9] F. J. Abu-Dakka, L. Rozo, and D. G. Caldwell. Force-based learning of variable impedance skills for robotic manipulation. In *2018 IEEE-RAS 18th International Conference on Humanoid Robots (Humanoids)*, pages 1–9. IEEE, 2018. [2](#)
- [10] X. Zhang, C. Wang, L. Sun, Z. Wu, X. Zhu, and M. Tomizuka. Efficient sim-to-real transfer of contact-rich manipulation skills with online admittance residual learning. In *Conference on Robot Learning*, pages 1621–1639. PMLR, 2023. [2](#)
- [11] X. Zhu, S. Kang, and J. Chen. A contact-safe reinforcement learning framework for contact-rich robot manipulation. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2476–2482. IEEE, 2022. [2](#)
- [12] R. Martín-Martín, M. A. Lee, R. Gardner, S. Savarese, J. Bohg, and A. Garg. Variable impedance control in end-effector space: An action space for reinforcement learning in contact-rich tasks. In *2019 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 1010–1017. IEEE, 2019. [2](#)

- [13] S. Bansal, M. Chen, S. Herbert, and C. J. Tomlin. Hamilton-jacobi reachability: A brief overview and recent advances. In *2017 IEEE 56th Annual Conference on Decision and Control (CDC)*, pages 2242–2253. IEEE, 2017. 2, 3
- [14] K. Margellos and J. Lygeros. Hamilton–jacobi formulation for reach–avoid differential games. *IEEE Transactions on automatic control*, 56(8):1849–1861, 2011. 2, 3
- [15] I. M. Mitchell, A. M. Bayen, and C. J. Tomlin. A time-dependent hamilton–jacobi formulation of reachable sets for continuous dynamic games. *IEEE Transactions on automatic control*, 50(7):947–957, 2005. 2, 3
- [16] D. Bertsekas. *Dynamic programming and optimal control: Volume I*, volume 4. Athena scientific, 2012. 2
- [17] B. D. Anderson and J. B. Moore. *Optimal control: linear quadratic methods*. Courier Corporation, 2007. 2
- [18] J. Darbon and S. Osher. Algorithms for overcoming the curse of dimensionality for certain hamilton–jacobi equations arising in control theory and elsewhere. *Research in the Mathematical Sciences*, 3(1):19, 2016. 2
- [19] J. F. Fisac, N. F. Lugovoy, V. Rubies-Royo, S. Ghosh, and C. J. Tomlin. Bridging hamilton–jacobi safety analysis and reinforcement learning. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 8550–8556. IEEE, 2019. 2, 1
- [20] A. K. Akametalu, S. Ghosh, J. F. Fisac, V. Rubies-Royo, and C. J. Tomlin. A minimum discounted reward hamilton–jacobi formulation for computing reachable sets. *IEEE Transactions on Automatic Control*, 2023. 2
- [21] R. S. Sutton. Learning to predict by the methods of temporal differences. *Machine learning*, 3:9–44, 1988. 2, 1
- [22] C. J. Watkins and P. Dayan. Q-learning. *Machine learning*, 8:279–292, 1992. 2, 1
- [23] J. Varley, S. Singh, D. Jain, K. Choromanski, A. Zeng, S. B. R. Chowdhury, A. Dubey, and V. Sindhwani. Embodied ai with two arms: Zero-shot learning, safety and modularity. *arXiv preprint arXiv:2404.03570*, 2024. 3
- [24] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng. Code as policies: Language model programs for embodied control. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9493–9500. IEEE, 2023. 3, 5
- [25] Z. Yang, S. S. Raman, A. Shah, and S. Tellec. Plug in the safety chip: Enforcing constraints for llm-driven robot agents. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 14435–14442. IEEE, 2024. 3
- [26] H. Wang, N. Chin, G. Gonzalez-Pumariega, X. Sun, N. Sunkara, M. A. Pace, J. Bohg, and S. Choudhury. Apricot: Active preference learning and constraint-aware task planning with llms. *arXiv preprint arXiv:2410.19656*, 2024. 3
- [27] N. Kumar, F. Ramos, D. Fox, and C. R. Garrett. Open-world task and motion planning via vision-language model inferred constraints. *arXiv preprint arXiv:2411.08253*, 2024. 3
- [28] A. Curtis, N. Kumar, J. Cao, T. Lozano-Pérez, and L. P. Kaelbling. Trust the proc3s: Solving long-horizon robotics problems with llms and constraint satisfaction. *arXiv preprint arXiv:2406.05572*, 2024. 3
- [29] L. Brunke, Y. Zhang, R. Römer, J. Naimer, N. Staykov, S. Zhou, and A. P. Schoellig. Semantically safe robot manipulation: From semantic scene understanding to motion safeguards. *arXiv preprint arXiv:2410.15185*, 2024. 3

- [30] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033. IEEE, 2012. doi:10.1109/IROS.2012.6386109. 3, 6
- [31] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 4, 3
- [32] T. Lillicrap. Continuous control with deep reinforcement learning. *International Conference on Learning Representations*, 2016. 5, 1
- [33] L. Downs, A. Francis, N. Koenig, B. Kinman, R. Hickman, K. Reymann, T. B. McHugh, and V. Vanhoucke. Google scanned objects: A high-quality dataset of 3d scanned household items. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 2553–2560. IEEE, 2022. 6
- [34] L. Downs, A. Francis, N. Koenig, B. Kinman, R. Hickman, K. Reymann, T. B. McHugh, and V. Vanhoucke. Google scanned objects: A high-quality dataset of 3d scanned household items, 2022. URL <https://arxiv.org/abs/2204.11918>. 6
- [35] K. Zakka. Scanned Objects MuJoCo Models, 7 2022. URL [https://github.com/kevinzakka/mujoco\\_scanned\\_objects](https://github.com/kevinzakka/mujoco_scanned_objects). 6
- [36] Y. Zhu, J. Wong, A. Mandlekar, R. Martín-Martín, A. Joshi, S. Nasiriany, Y. Zhu, and K. Lin. robosuite: A modular simulation framework and benchmark for robot learning. In *arXiv preprint arXiv:2009.12293*, 2020. 6
- [37] e. a. OpenAI Team. Gpt-4o system card, 2024. URL <https://arxiv.org/abs/2410.21276>. 6
- [38] C. Doersch, Y. Yang, M. Vecerik, D. Gokay, A. Gupta, Y. Aytar, J. Carreira, and A. Zisserman. TAPIR: Tracking any point with per-frame initialization and temporal refinement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10061–10072, 2023. 8
- [39] C. Doersch, A. Gupta, L. Markeeva, A. Recasens, L. Smaira, Y. Aytar, J. Carreira, A. Zisserman, and Y. Yang. TAP-vid: A benchmark for tracking any point in a video. *Advances in Neural Information Processing Systems*, 35:13610–13626, 2022. 8
- [40] R. Bellman, R. E. Kalaba, et al. *Dynamic programming and modern control theory*, volume 81. Citeseer, 1965. 1

## 8 Appendix

### 8.1 Background

#### 8.1.1 Hamilton-Jacobi Reachability Analysis

Hamilton-Jacobi Reachability (HJR) analysis provides a formal approach for computing safe policies that guarantee constraint satisfaction by finding an optimal solution to a reach-avoid problem. Consider a discrete-time dynamical system  $s_{t+1} = f(s_t, a_t)$ , where  $t$  is the current time step,  $s_t \in \mathcal{S} \subset \mathbb{R}^n$ ,  $a \in \mathcal{U} \subset \mathbb{R}^m$ ,  $\mathcal{U}$  is compact and dynamics  $f$  are bounded and Lipschitz continuous. Target set  $\mathcal{T} \subset \mathcal{S}$  describes *reach* states that satisfy  $\{s : l(s) > 0\}$  and failure set  $\mathcal{F} \subset \mathcal{S}$  describes *avoid* states that satisfy  $\{s : g(s) < 0\}$ .  $\mathcal{T}$  and  $\mathcal{F}$  are closed sets and  $l(s), g(s) : s \rightarrow \mathbb{R}$  are Lipschitz continuous functions. The safety problem is to find the *reach-avoid* set,  $\mathcal{RA}(\mathcal{T}; \mathcal{F})$ , which is the set of states from which a controller can drive the system to  $\mathcal{T}$  while avoiding  $\mathcal{F}$  at all times  $t$ . It was demonstrated in [19] that states belonging to the set  $\mathcal{RA}(\mathcal{T}; \mathcal{F})$  satisfy  $V(s) > 0$ , where  $V(s)$  is the solution to the following fixed-point Reach-Avoid Bellman equation (RABE):

$$\begin{aligned} V(s) &= \min [g(s), \max_{a \in \mathcal{U}} (l(s), \sup V(f(s, a)))] \\ V(s) &> 0 \Leftrightarrow s \in \mathcal{RA}(\mathcal{T}; \mathcal{F}) \end{aligned} \quad (6)$$

It is important to note that, unlike the Bellman updates in dynamic programming based approaches [40], (6) does not have a time-discounting term, thus it does not induce a contraction mapping and cannot converge to a fixed point using value iteration.

#### 8.1.2 Reach-Avoid Reinforcement Learning (RARL)

Based on time-discounting in temporal difference learning [21] and Q-learning [22] based methods, [19] introduced a time-discounting term in RABE (6), thus inducing contraction mapping in value function learning, which was extended to the reach-avoid setting and deep Q-learning by [1]. The Discounted Reach-Avoid Bellman Equation (DRABE) [1] can be written as:

$$V(s) = \gamma \min [g(s), \max_{a \in \mathcal{U}} (l(s), \sup V(f(s, a)))] + (1 - \gamma) \min (l(s), g(s))$$

where  $\gamma$  is the discount factor, which can be interpreted as the probability of episode continuation.

#### 8.1.3 Offline Parameterized Safe Policy Learning

We utilize Deep Deterministic Policy Gradient (DDPG) [32] as our reinforcement learning framework to learn a parameterized safe policy in continuous action spaces. The Q-function update is modified to the discounted reach-avoid formulation (7). Suppose  $Q_\phi$  is the Q-network,  $\mu_\theta$  the policy network, and  $Q_{\phi_{\text{targ}}}$  and  $\mu_{\theta_{\text{targ}}}$  are the target Q-network and target policy networks respectively. Given a sample  $(s, a, r, s', d)$  from the replay buffer  $\mathcal{D}$ , the Mean-squared Bellman Equation (MSBE) in DDPG is as follows:

$$L(\phi, \mathcal{D}) = \underset{(s, a, r, s', d) \sim \mathcal{D}}{\mathbb{E}} \left[ \left( Q_\phi(s, a) - (r + \gamma(1 - d)Q_{\phi_{\text{targ}}}(s', \mu_{\theta_{\text{targ}}}(s')) \right)^2 \right]$$

where  $r = \min(l(s), g(s))$  is the reward and  $d$  indicates whether state  $s'$  is terminal. The Discounted Reach-Avoid Bellman Error for the DDPG algorithm can be written as:

$$\begin{aligned} L(\phi, \mathcal{D})_{\text{DRABE}} &= \underset{(s, a, r, s', d) \sim \mathcal{D}}{\mathbb{E}} \left[ Q_\phi(s, a) - \left[ (1 - \gamma) \min(l(s), g(s)) \right. \right. \\ &\quad \left. \left. + \gamma \min \left\{ g(s), \max \{l(s), Q_{\phi_{\text{targ}}}(s', \mu_{\theta_{\text{targ}}}(s'))\} \right\} \right] \right] \end{aligned}$$

## 8.2 Experimental details

We denote the *full state* of the robot and the objects via  $s = [s^{\text{EE}}, s^{\text{cbot}}, s^{\text{ctop}}, s_1^o, \dots, s_N^o] \in \mathcal{S}$  where  $\mathcal{S}$  is the full state space,  $s^{\text{cbot}}$  is the state of the bottom cereal box,  $s^{\text{ctop}}$  is the state of the top cereal box and  $s_i^o$  is the state of the  $i^{\text{th}}$  object in clutter on the table. We model the robot’s state as the position and velocity of the end-effector  $s^{\text{EE}} = [x^{\text{EE}}, \dot{x}^{\text{EE}}]$  where  $x^{\text{EE}} \in \mathbb{R}^3$  is the Cartesian position and  $\dot{x}^{\text{EE}} \in \mathbb{R}^3$  is the velocity. Let each object’s state be represented by  $s_i^o = [x_i^o, \dot{x}_i^o], i \in \{1, \dots, N\}$  where  $N$  is the total number of cluttered objects on the table, that is, in addition to the two cereal boxes.

### 8.2.1 Safety and Target Representation

The safety constraints capture how the stack of boxes should not topple too far when the robot is pulling the cereal box from underneath, and models robot-object interactions. Specifically,  $\mathcal{F}$  reasons about the following constraints: 1) at any time step  $t$ , the top cereal block should not be excessively displaced from its initial position  $g_{\text{stack}}(s_t) = d_{\text{thresh}} - ||x_{t=t}^{\text{ctop}} - x_{t=0}^{\text{ctop}}||$ . For instance, if the bottom box is pulled very aggressively, the top cereal box will be flipped over or if the bottom block is pulled too slowly, the top block will move with it due to friction. Both these interactions should be avoided; 2) the bottom cereal box should not make contact with any of the *fragile* objects  $g_{\text{fragile}}(s) = ||x^{\text{cbot}} - x_{\text{fragile}}^o|| - d_{\text{thresh}}$ ; 3) the bottom cereal box should make soft contact with the *soft* objects  $g_{\text{soft}}(s) = d_{\text{thresh}} - ||\dot{x}^{\text{cbot}} - \dot{x}_{\text{soft}}^o||$ . This is a relative velocity constraint, which requires long-horizon reasoning, so that the cereal box can slow down before making contact with the object; 4) the bottom cereal box should not move over *sensitive* objects  $g_{\text{sensitive}}(s) = ||x^{\text{cbot}}[:, 2] - x_{\text{sensitive}}^o[:, 2]|| - d_{\text{thresh}}$  i.e. their x-y positions should not overlap; 5) any contact is allowed with *durable* objects  $g_{\text{durable}}(s) = +\inf$ . Note that since the safety constraint depends on the semantic properties of the object in the scene, during deployment the VLM specifies the constraints types via the parameterization  $\phi \in \{\text{no-contact}, \text{soft-contact}, \text{any-contact}, \text{do not move over}\}$ . During training, we consider all combinations of object properties and respective constraints.

The target set represents the desired cereal box reaching the correct goal location. Specifically,  $\mathcal{T}$  always requires 1) the bottom cereal box to be completely separated from the top cereal box  $l_{\text{slide}}(s) = ||x^{\text{cbot}} - x^{\text{ctop}}|| - d_{\text{thresh}}$ , and 2) placed in one of two regions near the far end of the table  $l_{\text{goal}}(s) = d_{\text{thresh}} - ||x^{\text{cbot}} - x^{\text{goal } 1, 2}||$ . This induces two possible target sets that the VLM can choose from at deployment time via the target parameterization  $\tau \in \{\text{top-goal}, \text{bottom-goal}\}$ . During training, we randomly choose between the two goal regions.

### 8.2.2 Custom attention mask

The custom attention mask used in **VLTSafe** is illustrated in Fig. 8.2.2. Each element  $m_j^i$  of the mask indicates whether object  $i$  attends to object  $j$ ; that is,  $m_j^i = 1$  if object  $i$  attends to object  $j$ , and  $m_j^i = 0$  otherwise. As noted in Fig. 8.2.2, the end-effector ( $s^{\text{EE}}$ ), bottom cereal box ( $s^{\text{cbot}}$ ) and top cereal box ( $s^{\text{ctop}}$ ) attend to each other. Specifically,  $m_j^{\text{EE}} = 1$  if object  $j$  is among the top-k nearest objects to the end-effector. The same logic applies for the other pairs.

	EE	cbot	ctop	$o_1$	...	$o_N$
EE	1	1	1	$m_1^{\text{EE}}$	...	$m_N^{\text{EE}}$
cbot	1	1	1	$m_1^{\text{cbot}}$	...	$m_N^{\text{cbot}}$
ctop	1	1	1	$m_1^{\text{ctop}}$	...	$m_N^{\text{ctop}}$
$o_1$	$m_1^{\text{EE}}$	$m_1^{\text{cbot}}$	$m_1^{\text{ctop}}$	1	...	$m_N^1$
:	:	:	:	:	...	:
$o_N$	$m_N^{\text{EE}}$	$m_N^{\text{cbot}}$	$m_N^{\text{ctop}}$	$m_1^N$	...	$m_N^N$

Figure 6: Custom attention mask

Hyperparameter	Value
Batch size	200
# Training steps	3M
Q learning rate	1e-4
Policy learning rate	1e-4
Max. episode length	300
Optimizer	AdamW
Eps	1e-8
Weight decay	0.01
Clip grad norm	5
Replay buffer size	100k
Discount factor	0.99
Polyak	0.995
Action noise	0.1
Collision threshold	1 cm

Table 2: Training hyperparameters

### 8.2.3 Architecture and training details

All policies are trained using the same target and failure sets as defined in Sec. 8.2.1 and optimized via the update rule discussed in Sec. 8.1.3. The training hyperparameters are kept consistent across all methods and are shown in Table 2. The policies are trained using one NVIDIA RTX A6000 GPU and an Intel i9-10900K CPU.

In **VLTsafe**, both the actor and critic networks are composed of their respective tokenizers, multihead self-attention layers and output heads. Each tokenizer comprises two hidden layers with 64 units each and SiLU activation, followed by an output layer that projects inputs into a 128-dimensional token embedding space. The core transformer backbone follows the standard architecture of Vaswani et al. [31], with two layers of multi-head self-attention, each employing four attention heads, as shown in Fig. 7. The feedforward block is composed of a single hidden layer of size 256 and an output layer that projects back to the token embedding dimension. The action and value prediction heads are implemented as two-layer MLPs with 64 hidden units per layer. In total, **VLTsafe** contains approximately 620k trainable parameters. The actor and critic networks in the **MLP** baseline are implemented simply as five-layer MLPs with 256 hidden units per layer and ReLU activation resulting in 540k total trainable parameters.

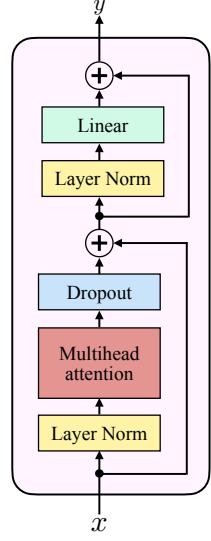


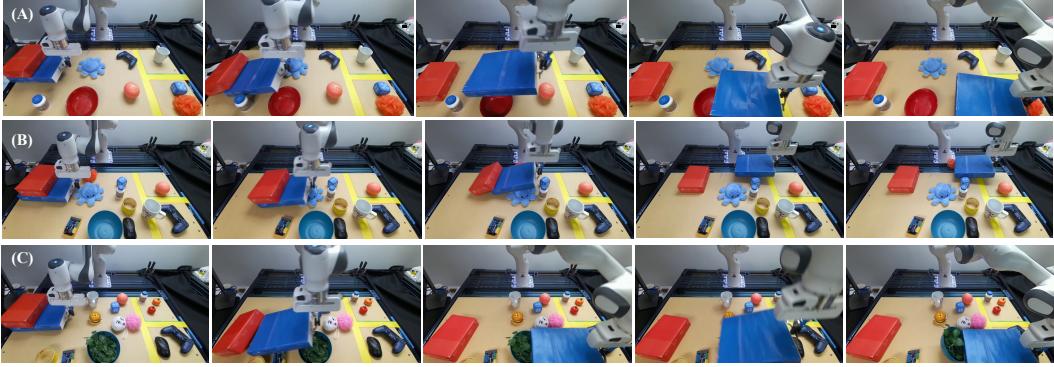
Figure 7: **Attention block**

## 8.3 VLM for Test-Time Safety and Target Specification

**Agent role prompt:** You are an excellent safe planning agent for dynamic tasks. You are given a task description and an image showing the robot and objects on a table. The robot is trying to slide the blue box from under the red box and to the goal regions to the right of the table, without damaging other objects along the way.

**Safety constraint selection prompt:** Each object on the table can potentially come in contact with the end-effector. You need to decide the safe interaction type for each object on the table from the list of constraint types. Here the description of the constraint types: `no_contact` implies that there should absolutely be no contact with a certain object. `soft_contact` implies that you can softly interact with that object, push it softly, etc. `any_contact` implies that any kind of interaction including aggressive impact is allowed. `no_over` implies that the robot is not allowed to move over (on top of) the object. Some hints on how to decide on the constraint type for an object: If an object is soft or made of durable material, and softly pushing it or moving it without toppling it is okay, `soft_contact` can be allowed with that object. If an object is very durable, and pushing it aggressively will not damage it, `any_contact` can be allowed with that object. If an object is fragile, and contacting it might damage it, `no_contact` should be allowed with that object. If an object is very sensitive like an open laptop or a bowl of food, and moving over it might be risky, `no_over` should be constrained for that object. Usually objects such as cups, wine glasses, bowls, electronics, etc are considered fragile and should be `no_contact`. Plastic objects such as bottles, plastic cans, tubes can be allowed `soft_contact`. Soft and non-critical objects such as toys, clothing, etc. are soft and can be ignored and allowed `any_contact`. Provide brief explanation, for choosing a specific constraint type for an object. In `image_description` briefly describe the scene and features relevant to the task.

**Target constraint selection prompt:** The robot is trying to slide the blue box from under the red box and place it in one of the two goal regions to the right of the table indicated by the two yellow squares on the table. The square closer to the robot is `bottom_goal` and the square further away is `top_goal`. Choose the target region for the blue box from the two goal regions such that it is safe to slide the blue box to that target region. Also provide a brief explanation for choosing that target region.



**Figure 8: Real-world deployment:** Zero-shot sim-to-real transfer of a policy trained in simulation (*Random-Const*, 6 objects) to a real-world settings with (A) 8, (B) 10 and (C) 14 objects. Constraints are first inferred using a VLM. In (A) and (B) the VLM is asked to specify constraints based on common sense reasoning while in (C) the VLM is directed to strictly avoid food shaped objects (see Sec. 8.4.1). As the robot pulls the blue box from under the red one, it makes allowable contact with the soft/durable objects, moves around or lifts up to avoid fragile objects and places the box safely in the goal regions deemed safe by the VLM.

**Structured Output** We employ the structured output capabilities of OpenAI’s Python API to force a desired structure on what is output by GPT4o. Below is the `create_vlm_constraint_response` function used in the implementation of **VLTsafe**

```

1 def create_vlm_constraint_response(object_list, safety_constraint_list,
2 , target_constraint_list):
3
4     fields = {}
5     for i in range(len(object_list)):
6         fields[f"explanation_obj_{object_list[i]}"] = (str, ...)
7         fields[f"{object_list[i]}"] = (safety_constraint_list, ...)
8
9     SafetyConstraintsPerObject = create_model(
10     "SafetyConstraintsPerObject", **fields)
11
12     class ConstraintResponse(BaseModel):
13         safety_constraints_per_object: SafetyConstraintsPerObject
14         target_region: target_constraint_list
15         target_description: str
16         image_description: str
17
18     return ConstraintResponse

```

Code Listing 1: The `create_vlm_constraint_response` function used to structure output from GPT4o.

The `create_vlm_constraint_response` function takes as input the list of objects in the scene, list of semantic targets (`bottom_goal`, `top_goal`) and the list of valid safety constraints for each object (`no_contact`, `soft_contact`, `any_contact`, `no_over`). These lists are used to populate the member variables of `ConstraintResponse` ensuring that both target and safety constraint outputs are restricted to valid options when querying the OpenAI API, thereby helping to prevent hallucinations.

## 8.4 Real-world Experiments

We control the Franka robot using Cartesian impedance control, which takes the desired end-effector pose as input at each time step. Our learned policy outputs the desired displacement of the end-effector, which is converted into an absolute pose and passed to the controller. Real-world rollouts for increasing number of objects than seen during training are shown in Fig. 8

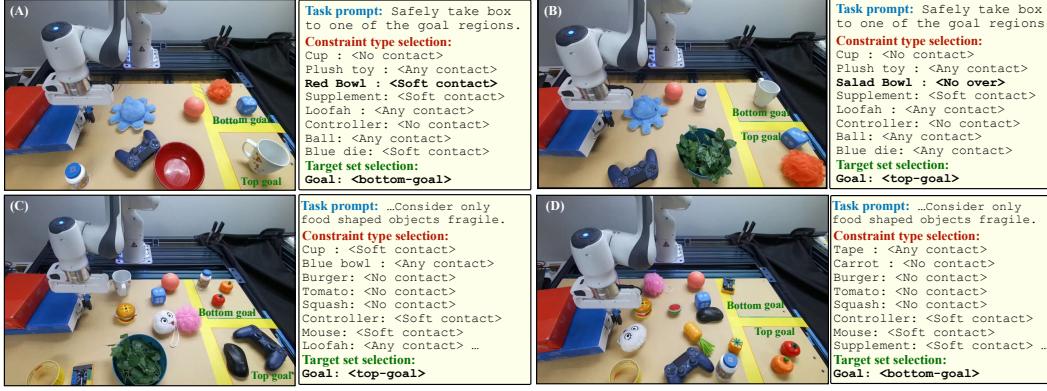


Figure 9: Qualitative results for the VLM selecting safety constraint types and target regions based on the semantics of the objects as observed in the input image as well as the task description.

#### 8.4.1 Qualitative examples of VLM constraint selection

In Fig. 9(A) and (B), the VLM is prompted to specify target regions and safety constraints based on commonsense reasoning about the material properties of objects. As a result, electronics and porcelain cups are treated as fragile. Accordingly, the goal region in both (A) and (B) is selected to avoid including the porcelain mug. When the bowl contains salad, as in (B), it is assigned the constraint no\_over; otherwise, as in (A), it is assigned soft\_contact.

In Fig. 9(C) and (D), the VLM is given explicit instructions to treat certain objects in specific ways. In this case, it is directed to consider only food-shaped objects as fragile. Consequently, these objects are assigned the constraint no\_contact, and the goal regions are selected to avoid them. Notably, electronics—previously considered no\_contact in (A) and (B)—are now assigned soft\_contact in (C) and (D), in accordance with the user’s updated instructions. This highlights the VLM’s capacity to interpret user intent and adapt task constraints accordingly, which is particularly valuable for nuanced and complex tasks requiring human-in-the-loop guidance.