

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

Anonymous Author(s)

Affiliation

Address

email

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

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

## 27      1 Introduction

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

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

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

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

## 72 2 Related work

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

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

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

92 **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.

### 103 3 Problem Formulation and the *VLTsafe* Method

104 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  
 105 approach—***VLTsafe***—which breaks down these challenges by training an object-centric trans-  
 106 former policy that is informed about relevant safety constraints at test-time by a vision-language  
 107 model. Our proposed approach is shown in Fig. 1.

108 **Notation** We model the robot’s state as the position and velocity of the end-effector  $s^{\text{EE}} =$   
 109  $[x^{\text{EE}}, \dot{x}^{\text{EE}}]$  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  
 110 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  
 111 cluttered scene. We denote the *full state* of the robot and the objects via  $s = [s^{\text{EE}}, s_1^o, \dots, s_N^o] \in \mathcal{S}$   
 112 where  $\mathcal{S}$  is the full state space. Finally, let the robot’s actions  $a \in \mathcal{A}$  be represented as the Cartesian  
 113 displacement of the end effector  $\Delta x^{\text{EE}} \in \mathbb{R}^3$ . In this work, we do not assume access to an analytic  
 114 dynamics model  $s_{t+1} = f(s_t, a_t)$ , but instead rely on a high-fidelity simulator [30] to evolve the  
 115 robot and object states as a result of the robot’s actions and physical interaction.

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

#### 122 **Goal: A Safe Dynamic Manipulation Policy Adaptable to Test-time Objects and Constraints**

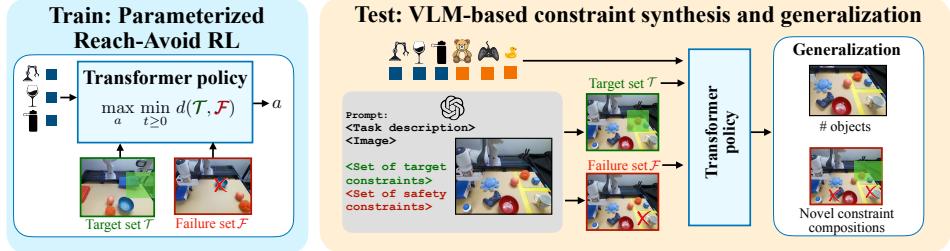
Our goal is to find a low-level robot policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  which respects the safety constraints  $\mathcal{F}$  while also achieving the desired task target  $\mathcal{T}$ . To enable test-time adaptation to relevant constraints and targets in the deployment environment, 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

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

123 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.

124 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.

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

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

### 139 3.1 Offline: Safe Policy Learning via Object-Centric Transformers & Reach-Avoid RL

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

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

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

165

166 where  $\gamma \in [0, 1]$  is the discount factor, which  
 167 can be interpreted as the probability of episode  
 168 continuation. We utilize Deep Deterministic  
 169 Policy Gradient (DDPG) [31] as our RL frame-  
 170 work wherein the bellman update is modified  
 171 to incorporate the time-discounted reach-avoid  
 172 Bellman backup (Equations in Appendix 8.1.)  
 173 We model both the actor and critic using object-  
 174 centric transformer architectures (Fig. 2). The  
 175 actor network encodes observation tokens with  
 176 a transformer, whose output tokens are aggre-  
 177 gated via mean pooling and passed through a  
 178 2-layer MLP policy head to produce the next  
 179 action. Critic network jointly encodes state  
 180 and action tokens using a separate transformer  
 181 with two layers of multi-head self-attention.  
 182 The resulting token embeddings are aggregated  
 183 and concatenated with the action, then passed  
 184 through a two-layer MLP value head to esti-  
 185 mate the Q-value.

### 186 3.2 Online: Adapting via Safety Constraints Inferred from Vision-Language Models

187 The prior section allowed us to pre-compute a low-level robot policy which given the current state  
 188 of the scene  $s$ , a safety specification  $\mathcal{F}$ , and a goal specification  $\mathcal{T}$  safely performs the task. The  
 189 question is, at deployment time, how can the robot know what are the “relevant” safety constraints  
 190  $\mathcal{F}$  and how they may impact which goals  $\mathcal{T}$  are and aren’t allowed? Here, we propose using a  
 191 vision-language model (VLM) as an open-world constraint and target specifier. Given the current  
 192 image of the scene  $\mathcal{I}$  and a text prompt describing the task  $\mathcal{L}_{task}$ , the VLM should use the visual  
 193 and semantic cues to select the relevant safety and targets that are passed into the low-level policy.  
 194 However, a large challenge with using VLMs is the right interface between textual representation  
 195 (that the VLM is trained to output) and the embodied representation (e.g., robot states, objects,  
 196 etc.) that our robot policy needs at test-time to adapt. To bridge this gap, we model the selection  
 197 of constraints/targets as a multiple choice visual question-answering (VQA) problem where the  
 198 multiple choice options are generated via pre-defined geometric constraint and target functions. In  
 199 the future, we envision the possibility of these functions to be written by another LLM based model  
 200 (e.g., as in Code-as-Policies [24]). We first convert the set of all geometric constraints  $\mathcal{T}, \mathcal{F}$  we  
 201 can generate with our functions to corresponding textual descriptions  $\mathcal{L}_{\mathcal{F}}, \mathcal{L}_{\mathcal{T}}$ . For example, a hard  
 202 collision-avoidance constraint between two interacting bodies is represented in text as:

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

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

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

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

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

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

## 209 4 Simulation Experiments

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

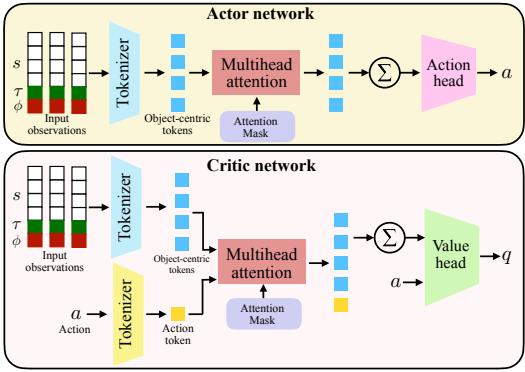


Figure 2: **System architecture:** Actor and critic networks are object-centric transformer-based models. Observations are augmented with parameters of the target  $\tau$  and failure set  $\phi$  and tokenized. These tokens are processed using multi-head self-attention with a custom attention mask.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

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

#### 215 4.1 Experimental Setup

216 **Environment & Task** In all experiments, we use a tabletop setup where a Franka robot interacts  
 217 with household objects (Fig. 5). The task involves quickly but safely pulling a cereal box from  
 218 under another box on a cluttered tabletop and placing it in a goal region. We simulate the robot and  
 219 environment dynamics with MuJoCo [30] and use the Google Scanned Objects Dataset [32, 33, 34]  
 220 to simulate everyday objects in the table arena defined in Robosuite [35].

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

233 **VLM for Test-Time Safety and Target Specification** We use GPT-4o [36] as the VLM. It is  
 234 queried once at the beginning of a test episode to identify the constraint types relevant to all objects  
 235 in the scene and to select the target constraint that is most appropriate for the task. The full prompt  
 236 for the constraint selection is provided in Appendix 8.3

237 **Metrics** We measure (1) *safe success rate* (SafeSucc %) defined as the percentage of trajectories  
 238 that complete the task safely (i.e., satisfy reach and avoid), (2) *stack safety violation rate* (StackFail  
 239 %) defined as the top cereal block moving outside of safety limits, and (3) *object safety violation  
 240 rate* (ObjFail %) percentage of objects not in the stack that are unsafely interacted with.

241 **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*

#### 242 4.2 How much does the transformer policy architecture influence task success?

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

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

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

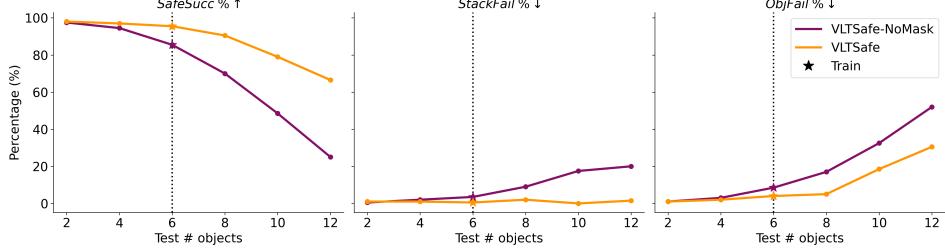


Figure 3: **Generalization: Varying Degrees of Clutter.** Performance variation as the number of objects is varied in the test scenario. Both **VLTsafe** vs **VLTsafe-NoMask** are trained in the *AllFragileConst* domain with six objects. Black dotted line represents the number of objects seen during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

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

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

#### 265 4.3 How well can our approach generalize to varying degrees of test-time clutter?

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

269 **Training & Test Setup** During training, we only consider the *AllFragileConst* domain trained with  
270 six objects, as in Section 4.2. At test-time, we deploy each policy in varying degrees of clutter—with  
271 2,4,6,8,10, and 12 objects—to test generalization to number of scene objects.

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

#### 281 4.4 How well does **VLTsafe** generalize to novel, test-time constraint compositions?

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

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

288 **Results** Recall that at training-time (with six objects), **VLTsafe** achieves a safe success rate of  
289 79.5% in the *FixedConst* domain and a safe success rate of 78.5% in *RandomConst* (star icon  
290 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.

291 N=6, since during training the policy never saw randomized object constraints. This explains  
 292 the drop in performance (12.5 %) when this policy is tested in *RandomConst*. For the policy  
 293 trained with *RandomConst*, only the 8 and 10-object environments are novel and, since the num-  
 294 ber of objects is larger, the constraint composition is novel. In Fig. 4, we see that **VLTsafe**  
 295 trained with random combinations of constraints exhibits better generalization capabilities, com-  
 296 pared to training with *FixedConst*, since it sees a larger variation in constraint compositions dur-  
 297 ing training. However, we note that this performance gap quickly diminishes as the environment  
 298 complexity increases and the task becomes significantly more challenging than the train setting.  
 299

## 300 5 Real World Experiments

301 **Hardware Setup and State Estimation** We use a 7-  
 302 DoF Franka Emika Panda arm and an Azure Kinect RGB-  
 303 D camera for both state estimation and image capture for  
 304 VLM-based constraint inference. State estimation relies  
 305 on TapNet [37, 38], a point tracking algorithm applied to  
 306 RGB-D inputs. At the start of each trajectory, we man-  
 307 nually annotate object keypoints, a step that could be auto-  
 308 mated in future work using VLMs. TapNet tracks object  
 309 states at 17 Hz.

310 **Sim-to-Real** We deploy policies trained entirely in sim-  
 311 ulation directly in the real world, without any real-world  
 312 fine-tuning. The experimental setup—including the robot’s state space and the object configura-  
 313 tions on the table—closely matches that of the simulation environment, enabling zero-shot transfer of  
 314 policies trained in simulation using object-centric observations, directly in the real world. We de-  
 315 ploy the policy trained in the *RandomConst* simulation domain. This policy was trained consider-  
 316 ing six objects in the scene, however we test in scenes containing up to eight objects.

317 **Results:** Fig. 5 shows a real world rollout of our learned policy in a scene with eight objects in  
 318 clutter. First, a VLM identifies the goal and object constraints: Soft objects are typically assigned  
 319 the `<any contact>` constraint, while mugs and electronic items—considered safety-critical—are  
 320 labeled as fragile `<no contact>`. When prompted to choose between two target sets—one contain-  
 321 ing soft toys, and the other containing a porcelain mug—the VLM selects the former, as interactions  
 322 with soft objects are safer. During execution, the robot briefly contacts the blue plush toy while ex-  
 323 tracting the blue box—an acceptable interaction—then lifts the end-effector to avoid colliding with  
 324 a nearby bowl. It ultimately places the box safely atop the orange loofah and blue ball, both of which  
 325 are deemed safe for contact. Additional real world experiments are included in the Appendix 8.5.

## 326 6 Conclusion

327 We introduce **VLTsafe**, an approach for learning low-level policies to safely perform dynamic tasks  
 328 in cluttered environments. By leveraging vision-language models (VLMs) for open-world constraint  
 329 synthesis, our object-centric transformer-based policy demonstrates zero-shot generalization to high  
 330 degrees of clutter and novel constraint compositions. Training with the reach-avoid reinforcement  
 331 learning (RARL) objective enables long-horizon safety reasoning and task accomplishment. **VLT-**  
 332 **Safe** outperforms baseline approaches, achieving high safe success rates, low failure rates, and  
 333 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.

334 **7 Limitations**

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

342 **References**

- 343 [1] K.-C. Hsu, V. Rubies-Royo, C. J. Tomlin, and J. F. Fisac. Safety and liveness guarantees  
344 through reach-avoid reinforcement learning. *arXiv preprint arXiv:2112.12288*, 2021. 2, 3, 4
- 345 [2] N. Hogan. Impedance control: An approach to manipulation. In *1984 American control*  
346 *conference*, pages 304–313. IEEE, 1984. 2
- 347 [3] F. J. Abu-Dakka and M. Saveriano. Variable impedance control and learning—a review. *Frontiers*  
348 *in Robotics and AI*, 7:590681, 2020. 2
- 349 [4] H. Sadeghian, L. Villani, M. Keshmiri, and B. Siciliano. Task-space control of robot manip-  
350 ulators with null-space compliance. *IEEE Transactions on Robotics*, 30(2):493–506, 2013.  
351 2
- 352 [5] A. De Luca and L. Ferrajoli. Exploiting robot redundancy in collision detection and reaction. In  
353 *2008 IEEE/RSJ international conference on intelligent robots and systems*, pages 3299–3305.  
354 IEEE, 2008. 2
- 355 [6] Z. Li, M. Zamora, H. Zheng, and S. Coros. Embracing safe contacts with contact-aware plan-  
356 ning and control. *arXiv preprint arXiv:2308.04323*, 2023. 2
- 357 [7] X. Zhu, W. Lian, B. Yuan, C. D. Freeman, and M. Tomizuka. Allowing safe contact in robotic  
358 goal-reaching: Planning and tracking in operational and null spaces. In *2023 IEEE Interna-*  
359 *tional Conference on Robotics and Automation (ICRA)*, pages 8120–8126. IEEE, 2023. 2
- 360 [8] L. Peternel, T. Petrič, and J. Babič. Human-in-the-loop approach for teaching robot assembly  
361 tasks using impedance control interface. In *2015 IEEE international conference on robotics*  
362 *and automation (ICRA)*, pages 1497–1502. IEEE, 2015. 2
- 363 [9] F. J. Abu-Dakka, L. Rozo, and D. G. Caldwell. Force-based learning of variable impedance  
364 skills for robotic manipulation. In *2018 IEEE-RAS 18th International Conference on Hu-*  
365 *manoid Robots (Humanoids)*, pages 1–9. IEEE, 2018. 2
- 366 [10] X. Zhang, C. Wang, L. Sun, Z. Wu, X. Zhu, and M. Tomizuka. Efficient sim-to-real transfer  
367 of contact-rich manipulation skills with online admittance residual learning. In *Conference on*  
368 *Robot Learning*, pages 1621–1639. PMLR, 2023. 2
- 369 [11] X. Zhu, S. Kang, and J. Chen. A contact-safe reinforcement learning framework for contact-  
370 rich robot manipulation. In *2022 IEEE/RSJ International Conference on Intelligent Robots*  
371 *and Systems (IROS)*, pages 2476–2482. IEEE, 2022. 2
- 372 [12] R. Martín-Martín, M. A. Lee, R. Gardner, S. Savarese, J. Bohg, and A. Garg. Variable  
373 impedance control in end-effector space: An action space for reinforcement learning in  
374 contact-rich tasks. In *2019 IEEE/RSJ international conference on intelligent robots and sys-*  
375 *tems (IROS)*, pages 1010–1017. IEEE, 2019. 2

- 376 [13] S. Bansal, M. Chen, S. Herbert, and C. J. Tomlin. Hamilton-jacobi reachability: A brief  
 377 overview and recent advances. In *2017 IEEE 56th Annual Conference on Decision and Control*  
 378 (*CDC*), pages 2242–2253. IEEE, 2017. 2, 3
- 379 [14] K. Margellos and J. Lygeros. Hamilton–jacobi formulation for reach–avoid differential games.  
 380 *IEEE Transactions on automatic control*, 56(8):1849–1861, 2011. 2, 3
- 381 [15] I. M. Mitchell, A. M. Bayen, and C. J. Tomlin. A time-dependent hamilton–jacobi formulation  
 382 of reachable sets for continuous dynamic games. *IEEE Transactions on automatic control*, 50  
 383 (7):947–957, 2005. 2, 3
- 384 [16] D. Bertsekas. *Dynamic programming and optimal control: Volume I*, volume 4. Athena scien-  
 385 tific, 2012. 2
- 386 [17] B. D. Anderson and J. B. Moore. *Optimal control: linear quadratic methods*. Courier Corpo-  
 387 ration, 2007. 2
- 388 [18] J. Darbon and S. Osher. Algorithms for overcoming the curse of dimensionality for certain  
 389 hamilton–jacobi equations arising in control theory and elsewhere. *Research in the Mathemat-  
 390 ical Sciences*, 3(1):19, 2016. 2
- 391 [19] J. F. Fisac, N. F. Lugovoy, V. Rubies-Royo, S. Ghosh, and C. J. Tomlin. Bridging hamilton-  
 392 jacobi safety analysis and reinforcement learning. In *2019 International Conference on*  
 393 *Robotics and Automation (ICRA)*, pages 8550–8556. IEEE, 2019. 2
- 394 [20] A. K. Akametalu, S. Ghosh, J. F. Fisac, V. Rubies-Royo, and C. J. Tomlin. A minimum dis-  
 395 counted reward hamilton–jacobi formulation for computing reachable sets. *IEEE Transactions*  
 396 *on Automatic Control*, 2023. 2
- 397 [21] R. S. Sutton. Learning to predict by the methods of temporal differences. *Machine learning*,  
 398 3:9–44, 1988. 2
- 399 [22] C. J. Watkins and P. Dayan. Q-learning. *Machine learning*, 8:279–292, 1992. 2
- 400 [23] J. Varley, S. Singh, D. Jain, K. Choromanski, A. Zeng, S. B. R. Chowdhury, A. Dubey, and  
 401 V. Sindhwani. Embodied ai with two arms: Zero-shot learning, safety and modularity. *arXiv*  
 402 *preprint arXiv:2404.03570*, 2024. 3
- 403 [24] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng. Code  
 404 as policies: Language model programs for embodied control. In *2023 IEEE Interna-  
 405 tional Conference on Robotics and Automation (ICRA)*, pages 9493–9500. IEEE, 2023. 3, 5
- 406 [25] Z. Yang, S. S. Raman, A. Shah, and S. Tellec. Plug in the safety chip: Enforcing constraints for  
 407 llm-driven robot agents. In *2024 IEEE International Conference on Robotics and Automation*  
 408 (*ICRA*), pages 14435–14442. IEEE, 2024. 3
- 409 [26] H. Wang, N. Chin, G. Gonzalez-Pumariega, X. Sun, N. Sunkara, M. A. Pace, J. Bohg, and  
 410 S. Choudhury. Apricot: Active preference learning and constraint-aware task planning with  
 411 llms. *arXiv preprint arXiv:2410.19656*, 2024. 3
- 412 [27] N. Kumar, F. Ramos, D. Fox, and C. R. Garrett. Open-world task and motion planning via  
 413 vision-language model inferred constraints. *arXiv preprint arXiv:2411.08253*, 2024. 3
- 414 [28] A. Curtis, N. Kumar, J. Cao, T. Lozano-Pérez, and L. P. Kaelbling. Trust the proc3s: Solv-  
 415 ing long-horizon robotics problems with llms and constraint satisfaction. *arXiv preprint*  
 416 *arXiv:2406.05572*, 2024. 3
- 417 [29] L. Brunke, Y. Zhang, R. Römer, J. Naimer, N. Staykov, S. Zhou, and A. P. Schoellig. Se-  
 418 mantically safe robot manipulation: From semantic scene understanding to motion safeguards.  
 419 *arXiv preprint arXiv:2410.15185*, 2024. 3

- 420 [30] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In  
421 *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–  
422 5033. IEEE, 2012. doi:10.1109/IROS.2012.6386109. 3, 6
- 423 [31] T. Lillicrap. Continuous control with deep reinforcement learning. *International Conference  
424 on Learning Representations*, 2016. 5
- 425 [32] L. Downs, A. Francis, N. Koenig, B. Kinman, R. Hickman, K. Reymann, T. B. McHugh,  
426 and V. Vanhoucke. Google scanned objects: A high-quality dataset of 3d scanned household  
427 items. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 2553–  
428 2560. IEEE, 2022. 6
- 429 [33] L. Downs, A. Francis, N. Koenig, B. Kinman, R. Hickman, K. Reymann, T. B. McHugh, and  
430 V. Vanhoucke. Google scanned objects: A high-quality dataset of 3d scanned household items,  
431 2022. URL <https://arxiv.org/abs/2204.11918>. 6
- 432 [34] K. Zakka. Scanned Objects MuJoCo Models, 7 2022. URL [https://github.com/kevinzakka/mujoco\\_scanned\\_objects](https://github.com/kevinzakka/mujoco_scanned_objects). 6
- 433 [35] Y. Zhu, J. Wong, A. Mandlekar, R. Martín-Martín, A. Joshi, S. Nasiriany, Y. Zhu, and K. Lin.  
434 robosuite: A modular simulation framework and benchmark for robot learning. In *arXiv  
435 preprint arXiv:2009.12293*, 2020. 6
- 436 [36] e. a. OpenAI Team. Gpt-4o system card, 2024. URL <https://arxiv.org/abs/2410.21276>. 6
- 437 [37] C. Doersch, Y. Yang, M. Vecerik, D. Gokay, A. Gupta, Y. Aytar, J. Carreira, and A. Zisserman.  
438 TAPIR: Tracking any point with per-frame initialization and temporal refinement. In *Proceedings  
439 of the IEEE/CVF International Conference on Computer Vision*, pages 10061–10072,  
440 2023. 8
- 441 [38] C. Doersch, A. Gupta, L. Markeeva, A. Recasens, L. Smaira, Y. Aytar, J. Carreira, A. Zisser-  
442 man, and Y. Yang. TAP-vid: A benchmark for tracking any point in a video. *Advances in  
443 Neural Information Processing Systems*, 35:13610–13626, 2022. 8