
SAFEDEC: CONSTRAINED DECODING FOR SAFE AUTOREGRESSIVE GENERALIST ROBOT POLICIES

000
001
002
003
004
005 **Anonymous authors**
006 Paper under double-blind review
007
008
009
010

ABSTRACT

011 Recent advances in end-to-end, multi-task robot policies based on transformer
012 models have demonstrated impressive generalization to real-world embodied AI
013 tasks. Trained on vast datasets of simulated and real-world trajectories, these mod-
014 els map multimodal observations directly to action sequences for physical execu-
015 tion. Despite promising real-world capabilities, these models are still data-driven
016 and, therefore, lack explicit notions of behavioral correctness. We address this
017 gap by introducing **SafeDec**, a constrained decoding framework for autoregres-
018 sive, transformer-based robot policies that enforces invariant safety specifications
019 on candidate action trajectories. Task-specific safety rules are expressed as Signal
020 Temporal Logic (STL) formulas and are enforced at inference time with minimal
021 overhead. Our method ensures that generated actions provably satisfy STL spec-
022 ifications under assumed dynamics at runtime without retraining , while remain-
023 ing agnostic of the underlying policy. We evaluate **SafeDec** on tasks from the
024 CHORES benchmark for state-of-the-art generalist policies (e.g., SPOC, Flare,
025 PoliFormer) across hundreds of procedurally generated environments and show
026 that our decoding-time interventions are useful not only for filtering unsafe actions
027 but also for conditional action generation. Videos are available at **constrained-**
028 **robot-fms.github.io**.
029

1 INTRODUCTION

030 Recent advances in developing large transformer-based models for robotics have enabled general-
031 purpose policies that map multi-modal inputs such as RGB images, natural language instructions,
032 and proprioceptive inputs to action sequences Hu et al. (2023). Shortest Path Oracle Clone (SPOC)
033 Ehsani et al. (2024), PoliFormer Zeng et al. (2025), Flare Hu et al. (2025) and OpenVLA (Kim et al.,
034 2024) exhibit impressive generalization in navigation and manipulation tasks and serve as versatile
035 robot controllers for real-world deployment contexts. However, these models are primarily data-
036 driven and lack any explicit notion of safety. Although these models may implicitly exhibit safety-
037 related behaviors depending on the patterns in their training data, there is no formal guarantee that
038 models will consistently behave safely in all situations. This serves as a limiting factor for deploying
039 these foundation models in the physical world where rule compliance and regulatory safety rule
040 adherence are crucial.
041

042 Formal specifications have long been used to specify safety requirements for robotic deployments
043 (Menghi et al., 2019; Farrell et al., 2018). Temporal logics (TL) (Pnueli, 1977) can capture safety
044 constraints on robot behavior, such as “remain within the permitted region zones and avoid dan-
045 gerous obstacles”. Although TL has seen success in classical robotic planning for safety constraint
046 satisfaction, its use for enforcing safety for large transformer-based robot policies remains limited.
047 Additionally, retraining or fine-tuning these large pre-trained models to directly embed temporal
048 logic specification is challenging (Kapoor et al., 2024). First, retraining models is a costly endeavor
049 in terms of computational resources and data requirements. Moreover, due to the stochastic nature
050 of these models, it is difficult to guarantee strict satisfaction of safety constraints through training
051 alone. Hence, there is a pressing need for methods that can enforce safety specifications efficiently
052 at inference time without disrupting the model’s pre-trained behavior.

053 In the field of natural language processing, syntactic constraints have been successfully enforced
by applying constrained decoding at inference time (Willard & Louf, 2023; Beurer-Kellner et al.,

054 2023; AI, 2023). These approaches typically mask out tokens that violate a syntactic constraint
 055 defined over token sequences. For example, regular expressions (regex) represent a widely used
 056 form of syntactic constraint, requiring that generated token sequences conform to predefined struc-
 057 tural patterns (Willard & Louf, 2023; Beurer-Kellner et al., 2023). Inspired by this line of work,
 058 we extend the paradigm of constrained decoding to enforce safety constraints over action trajec-
 059 tories in dynamical systems and propose *safety specification aligned decoding* (*SafeDec*) for trans-
 060 former based policies that ensures generated action sequences provably satisfy Signal Temporal
 061 Logic (STL) (Maler & Nickovic, 2004) specifications under assumed dynamics. Our key insight is
 062 that decoding-time interventions can be used not just to filter unsafe actions, but to *condition the*
 063 *generation process itself* on specification satisfaction. This conditioning is critical because it steers
 064 the model toward generating safety specification satisfying actions rather than relying on post hoc
 065 rejection. *SafeDec* reduces risk of infeasible outputs while preserving the original action distri-
 066 bution of the model. To enforce such specifications, we leverage the formal semantics of STL to
 067 evaluate candidate actions at runtime and mask those that lead to future violations. Our method
 068 is agnostic to the underlying foundation model, requiring only two properties: (1) access to the
 069 decoding-layer logits during inference, and (2) access to an approximate dynamics model to predict
 070 future states. To efficiently evaluate STL specifications at inference time, we use a high-performance
 071 computational graph based library *STLCG++* (Kapoor et al., 2025b). To the best of our knowledge,
 072 this is the first work to enforce formal safety on transformer-based robotic policies at inference time
 073 using constrained decoding.

074 Our main contributions are as follows: First, we introduce the novel problem of *constrained de-*
 075 *coding for transformer-based policies* under safety STL specifications (Section 3.1). Second, we
 076 propose an inference-time technique that reweights or masks candidate actions using STL satis-
 077 faction scores in (Section 3.2 and 3.3). Finally, we demonstrate the effectiveness of our method on
 078 state-of-the-art object navigation models without modifying model parameters (Section 4).

079 2 PRELIMINARIES

080 2.1 SIGNAL TEMPORAL LOGIC

081 Signal Temporal Logic (STL) is an expressive framework for defining properties and reasoning
 082 over continuous time real-valued signals Maler2004MonitoringTP. Formally, $(s, t) \models \phi$ denotes
 083 that a signal s satisfies the STL formula ϕ at time t . An atomic predicate of an STL formula is
 084 represented by inequalities of the form $\mu(s(t)) > 0$. The truth value of the predicate μ is equivalent
 085 to $\mu(s(t)) > 0$. Note that with slight abuse of notation, μ represents both the predicate and a
 086 function of the trajectory $s(t)$. Any STL formula consists of Boolean and temporal operations on
 087 these predicates, and the syntax of STL formulas is defined recursively as follows:

$$088 \phi := \mu \mid \neg\mu \mid \phi \wedge \psi \mid \phi \vee \psi \mid \mathbf{G}_{[a,b]} \psi \mid \mathbf{F}_{[a,b]} \psi \mid \phi \mathbf{U}_{[a,b]} \psi$$

089 where ψ and ϕ are STL formulae, \mathbf{G} denotes the globally operator, \mathbf{F} the eventually operator, and
 090 \mathbf{U} is the until operator. For example, $s \models \mathbf{G}_{[a,b]} \psi$ specifies that ψ must be in all times in the given
 091 interval, $t \in [a, b]$ of the signal s . Similarly, the operator *until* in $s \models \phi \mathbf{U}_{[a,b]} \psi$ defines that ϕ must
 092 be true until ψ becomes true within a time interval $[a, b]$.

093 Given a signal s_t representing a signal starting at time t , the Boolean semantics of satisfaction of
 094 $s_t \models \phi$ are defined inductively as follows:

$$\begin{aligned} 095 s_t \models \mu &\iff \mu(s(t)) > 0 \\ 096 s_t \models \neg\varphi &\iff \neg(s_t \models \varphi) \\ 097 s_t \models \varphi_1 \wedge \varphi_2 &\iff (s_t \models \varphi_1) \wedge (s_t \models \varphi_2) \\ 098 s_t \models \mathbf{F}_{[a,b]}(\varphi) &\iff \exists t' \in [t + a, t + b] \text{ s.t. } s_{t'} \models \varphi \\ 099 s_t \models \mathbf{G}_{[a,b]}(\varphi) &\iff \forall t' \in [t + a, t + b] \text{ s.t. } s_{t'} \models \varphi \end{aligned}$$

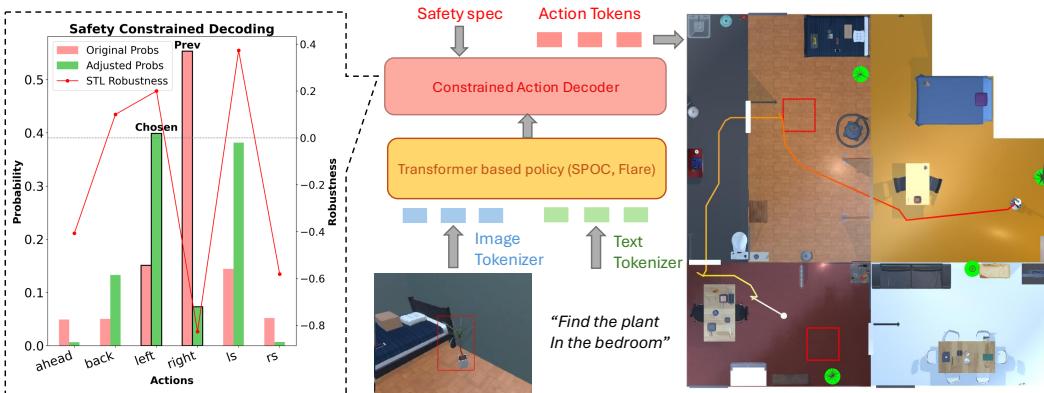
100 Apart from the Boolean semantics, quantitative semantics are defined for a signal to compute a real-
 101 valued metric indicating *robustness*, i.e., the strength of satisfaction or violation. For the sake of
 102 brevity, the definition of robustness is provided in Appendix A.1.

108 **2.2 CONSTRAINED DECODING IN TRANSFORMERS**
 109

110 A large variety of autoregressive transformer-based models generate final outputs by producing a
 111 probability distribution over the model vocabulary at each timestep. This distribution is generated
 112 by performing a softmax operation over the model’s last hidden layer. Then, through the process of
 113 *decoding*, tokens are selected to maximize the overall likelihood of an output sequence. In standard
 114 decoding, this maximization can be performed by either greedily selecting the most probable token
 115 at each step or by using a beam search to maintain multiple high-likelihood candidates. However,
 116 this often leads to degenerate output sequences that are repetitive (Holtzman et al., 2019). A common
 117 approach is to use sampling strategies like top k (Fan et al., 2018), and nucleus sampling (Holtzman
 118 et al., 2019) that introduce stochasticity to encourage more diverse outputs. Constrained decoding
 119 (Hokamp & Liu, 2017) modifies this probabilistic selection by pruning invalid tokens to ensure
 120 that the generated sequences satisfy predefined constraints. These constraints are often syntactic,
 121 such as regular expressions, JSON formatting, or programming language grammars (Welleck et al.,
 122 2024). There is also recent work on enforcing *semantic* constraints that ensure coherence of the
 123 output or alignment with specific knowledge bases (Peyrard et al., 2024). Formally, constrained
 124 decoding can be seen as maximizing the probability of the output sequence subject to a constraint
 \mathcal{C} : $\arg \max_{y \in \mathcal{Y}_C} P(y | x)$ where \mathcal{Y}_C is the set of sequences satisfying \mathcal{C} .

125

126 **3 SPECIFICATION-GUIDED CONSTRAINED DECODING**
 127



128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

Figure 1: Overview of our specification aligned decoding framework. Given multimodal inputs (e.g., RGB images and natural language instructions), a pretrained transformer-based robot policy (e.g., SPOC) generates candidate actions. These actions are filtered or reweighted by the constrained decoding technique based on STL safety constraints. In the figure, green markers denote target object locations, while red zones represent regions to avoid.

In this section, we introduce a novel problem formulation for SafeDec in autoregressive transformer-based robot policies. First, we highlight the challenge in specification checking for these policies in contrast to traditional syntactical constraint checking adopted by LLMs, and our solution to remedy it. Then, we propose two novel inference-time techniques for specification aligned decoding: *Hard Constrained Decoding* (HCD) and *Robustness Constrained Decoding* (RCD).

154 **3.1 PROBLEM STATEMENT**

155

156 As highlighted in the background section, existing techniques in constrained decoding for language
 157 models enforce *syntactic constraints* defined over tokens, such as conforming to a context-free grammar
 158 or matching a regular expression. In these setups, constraint checking can be performed in the
 159 model’s token space.

160 In contrast, transformer-based policies operate in a physical environment and constraints (captured
 161 via temporal logic) are defined over state variables. Since a large class of end-to-end policies solely
 propose action sequences, specification checking can only be performed as actions are executed

162 and the environment is simulated forward. In this case, constraint checking cannot be done solely
 163 in the token space and requires environmental feedback or a dynamics stepping function. Hence,
 164 we leverage an approximate first order dynamics function to compute specification satisfaction of
 165 different action sequences proposed by these policies.

166 Consider a discrete dynamical system with states $x_t \in \mathbb{R}^n$ and actions $a_t \in \mathbb{A}$ at time step t . The
 167 system's dynamics are defined by $x_{t+1} = f(x_t, a_t)$ where $f : \mathbb{R}^n \times \mathbb{A} \rightarrow \mathbb{R}^n$ maps the current state
 168 ($x_t \in \mathbb{R}^n$) and a discrete action ($a_t \in A$) to the next state $x_{t+1} \in \mathbb{R}^n$. This system is controlled by
 169 a policy that selects an action a_t at each time step based on observations and task context such as
 170 user-provided natural-language instructions or goal waypoints.

171 In this work, we focus on robot policies that generate actions based on multi-modal inputs, including
 172 sensor observations (e.g., RGB, depth, LiDAR) and natural language instructions. Let \mathcal{I}_t represent
 173 the aggregated input at timestep t . Given the history of encoded inputs up to time t , a transformer-
 174 based policy parameterized by θ predicts embeddings for the next $T-t$ actions:
 175

$$\{\hat{e}_{a_{t+k}}\}_{k=1}^{T-t} = \text{Transformer}_\theta(\{e_{\mathcal{I}_\tau}\}_{\tau=0}^t).$$

176 Each predicted action embedding is decoded into an action $\hat{a}_{t+k} \in \mathbb{A}$, resulting in a predicted action
 177 sequence $\{\hat{a}_{t+1}, \dots, \hat{a}_T\}$, where T denotes the planning horizon.

178 Now, consider that the system is required to satisfy requirements encoded using an STL formula
 179 φ defined over the state variables of the system. Formally, the goal is to ensure that the resulting
 180 trajectory satisfies the specification φ :

$$\{(x_0, \hat{a}_0), \dots, (x_T, \hat{a}_T)\} \models \varphi$$

181 Most of the existing techniques for specification enforcement perform posthoc manipulation of pro-
 182 posed actions through filtering or rejecting action sequences that violate the specification φ . Al-
 183 though manipulation after sampling can ensure specification satisfaction, it can lead to distorting
 184 the model's learned distribution, producing low-likelihood output. This undermines the inductive
 185 biases learned during pretraining and leads to degenerate, brittle behaviors. A similar problem was
 186 highlighted when ensuring compliance with logical constraints for large language models (LLM) in
 187 Park et al. (2024). Additionally, these models decode actions sequentially, where each action a_t is
 188 conditioned on previously generated tokens $a_{<t}$. Posthoc manipulation can disrupt this causal chain
 189 and lead to a mismatch between the model's internal hidden state and the executed sequence. Hence,
 190 we propose the following problem statement:
 191

192 *How can we enforce temporal logic constraints during action generation in robot
 193 foundation models such that the output sequence (1) satisfies an STL specification
 194 φ , and (2) remains faithful to the model's autoregressive distribution $\pi(a_{1:T} |$
 195 $\mathcal{I}_{1:T})$?*

196 Let $\pi(a_{0:T})$ be the *unconstrained* action-sequence distribution produced by the transformer-based
 197 policy's decoder (e.g. the softmax over logits generated by the Transformer). We define the ideal
 198 constrained distribution over action sequences as:

$$Q_{\pi, \varphi}(a_{0:T}) = \frac{\pi(a_{0:T}) \cdot \mathbf{1}[(x_{1:T}, a_{0:T}) \models \varphi]}{\sum_{a'_{0:T}} \pi(a'_{0:T}) \cdot \mathbf{1}[(x'_{1:T}, a'_{0:T}) \models \varphi]} \quad (1)$$

199 where $x_{1:T}$ denotes the state trajectory induced by the system dynamics under actions $a_{0:T}$ and
 200 $\mathbf{1}[\cdot]$ is the indicator function that returns 1 iff the trajectory-action pair satisfies the specification.
 201 Equation 1 is the exact Bayesian conditioning of π on the event that the generated rollout satisfies φ .
 202 Hence, sampling from $Q_{\pi, \varphi}$ would give sequences that (i) inherit the original model's preferences
 203 encoded in π and (ii) *guarantee* specification satisfaction.

204 In this work, we propose a technique to overcome the drawbacks of post-hoc safety enforcement
 205 methods (such as filtering) by leveraging constrained decoding techniques. Specifically, we propose
 206 SafeDec : A constrained decoding strategy that integrates STL specifications into the foundation
 207 model action selection process itself, ensuring satisfaction without distorting the model's distribu-
 208 tion.

216 3.2 HARD CONSTRAINED DECODING
 217

218 As highlighted in the background section, in the final layer, predictions are detokenized and a pro-
 219 jection layer converts the embeddings into logits over the vocabulary space. These logits are further
 220 converted into a probability distribution using a softmax operation. In prior work, for structured out-
 221 put generation in LLMs, some invalid tokens are masked based on syntactical constraints or other
 222 criteria (Welleck et al., 2024; Park et al., 2024). This is done by setting their logit value as $-\infty$
 223 before the softmax operation is applied. For HCD, we use a similar approach as constrained decod-
 224 ing literature (Welleck et al., 2024) and mask out predicted action tokens that violate our given STL
 225 specification φ during sequential generation. Formally, to enforce the STL specification φ during
 226 sequential generation, we adjust the logits at each timestep $t + k$ as follows:
 227

228 Let \mathbf{z}_{t+k} denote the logits at timestep $t + k$. For each action choice i at timestep $t + k$, we define:

$$229 \quad z_{t+k}^{(i)} = \begin{cases} -\infty, & \text{if } \hat{x}_{t+k}^{(i)} = f(x_{t+k-1}, \hat{a}_{t+k}^{(i)}) \text{ violates } \varphi \\ z_{t+k}^{(i)}, & \text{otherwise} \end{cases}$$

230

232 Here t is the current decision step, k is an index for the look-ahead step $t + k$ within a planning
 233 horizon of length T ($k \in [1..T]$), $\hat{a}_{t+k}^{(i)}$ is the action mapping to the token i and $\hat{x}_{t+k}^{(i)}$ is the next
 234 state value upon taking this action. This next state is elicited using a simple dynamics model (f) as
 235 highlighted in the previous section. Adjusting logits in this fashion ensures that any invalid token
 236 with respect to the safety specification will have zero probability of being selected after applying the
 237 softmax function.

238
 239 3.3 ROBUSTNESS CONSTRAINED DECODING
 240

241 HCD ensures compliance but can lead to compromising task success, which can be undesirable. A
 242 similar tradeoff was observed by Liu et al. (2021) when probability space-steering preserved model
 243 fluency while reducing toxic continuations compared with hard-filtering strategies that inflated per-
 244 perplexity and eroded diversity. Hence, we propose an alternative approach, called RCD, where we
 245 leverage the quantitative semantics of STL specifications (robustness). Unlike HCD, which applies
 246 hard masking to completely remove unsafe actions, RCD softly guides the model toward safer ac-
 247 tions by incorporating robustness scores that reflect the degree of satisfaction of φ . This is similar
 248 to the approach proposed in Liu et al. (2021) where the next-token distributions were re-weighted
 249 based on the utility scores provided by another language model. Our utility scores are quantified by
 250 the robustness function($\rho(x_t, \varphi)$) that returns a real-valued score indicating how well a predicted
 251 state satisfies the specification. Positive robustness values denote specification satisfaction, while
 252 negative values capture the degree of violation. A formal definition of robustness in line with the
 253 STL quantitative semantics is provided in Appendix A.1.

254 First, we compute a robustness score for each candidate action: $r_{t+k}^{(i)} = \rho(\hat{x}_{t+k}^{(i)}, \varphi)$ where $\rho(\cdot, \varphi)$ is
 255 the STL robustness metric, and $\hat{x}_{t+k}^{(i)}$ is the predicted next state under action $\hat{a}_{t+k}^{(i)}$. This robustness
 256 score $r_{t+k}^{(i)}$ quantifies how well each candidate action satisfies the specification φ . These scores
 257 are then converted into weights using exponential scaling: $w_{t+k,i} = \exp(\alpha \cdot r_{t+k,i})$ where α is
 258 a temperature parameter that adjusts the sharpness of the bias. We use these weights to shift the
 259 original logits: $\tilde{z}_{t+k,i} = z_{t+k,i} + \beta \cdot w_{t+k,i}$ where β is a hyperparameter that modulates the trade-
 260 off between specification adherence and the original task objective. Finally, we obtain the action
 261 distribution by applying softmax over the adjusted logits: $p_{t+k} = \text{softmax}(\tilde{\mathbf{z}}_{t+k})$

262 This approach allows for graded preferences that improve flexibility and robustness to dynamics ap-
 263 proximation errors. By ensuring that all actions preserve a non-zero probability, the policy remains
 264 capable of recovering from errors arising from an imperfect dynamics model. Concretely, if the
 265 predicted successor \hat{x}_{t+1} is off by ϵ , an action that appeared marginally unsafe can be safe in the
 266 true system, and vice versa. Retaining a weighted down probability for this action gives the sampler
 267 a fall-back option whereas HCD would completely rule this action out due to 0 probability. Since
 268 we are shifting the probability mass for unsatisfying actions, it is possible that they are still chosen
 269 and lead to a violation. However, this is a tradeoff that we allow to achieve a given task objective.
 We note that this still ensures higher STL satisfaction than unconstrained actions.

270 4 EVALUATION
271

272 4.1 IMPLEMENTATIONAL DETAILS
273

274 We comprehensively evaluate our constrained decoding framework on procedurally generated AI2-
275 THOR (Kolve et al., 2022) indoor scenes with diverse objects and layouts using three state-of-
276 the-art (SOTA) generalist robot policies: Shortest Path Oracle Clone (SPOC) (Ehsani et al., 2024),
277 PoliFormer (Zeng et al., 2025) and Flare (Hu et al., 2025). All three are large transformer-based
278 embodied agents trained on extensive language-conditioned robot trajectory datasets. These models
279 achieve strong zero-shot generalization for a vast variety of navigation tasks that span open vocab-
280 ular object-goal navigation (“find a mug”), room-to-room traversal (“visit all rooms”), waypoint-
281 based navigation (“move three meters forward and stop near the red rug”), and attribute-conditioned
282 variants (“locate the chair closest to the refrigerator in the kitchen”). These models also demonstrate
283 reliable zero-shot transfer to real-world environments, achieving robust task satisfaction.

284 In addition, these models capture three different training paradigms for generalist robot policies.
285 SPOC is trained purely with imitation learning from shortest-path rollouts. Poliformer employs a
286 hybrid approach that combines reinforcement learning and imitation learning, enabling it to learn
287 long-horizon structure while retaining expert priors. Flare adopts a large-scale pretraining plus fine-
288 tuning on embodied navigation data in line with recent foundation model training paradigms. This
289 diversity in training paradigms allows us to evaluate the applicability of SafeDec across different
290 learning regimes.

291 In this work, we address safety specifications for robotics and, therefore, select those most rel-
292 evant to real-world deployment. In particular, we focus on an important class of safety spec-
293 ifications called *invariants*, which are specifications that must be enforced at every reachable
294 state of the system (e.g., “always avoid an unsafe region”). We enforce *geofencing* and *ob-
295 stacle avoidance* by encoding them as invariant specifications in STL. Specifically, we generate
296 random regions in the configuration space that the robot must either avoid (obstacle zones)
297 or remain within (safe zones), and apply these constraints in real time during execution. The
298 specifications used are: $\varphi_{\text{geofence}} = \mathbf{G}\left(\bigvee_{i=1}^N \left(x_i^{\text{L}} \leq x \leq x_i^{\text{U}} \wedge z_i^{\text{L}} \leq z \leq z_i^{\text{U}}\right)\right)$, $\varphi_{\text{avoid}} =$
299 $\mathbf{G}\left(\bigwedge_{i=1}^N \neg\left(x_i^{\text{L}} \leq x \leq x_i^{\text{U}} \wedge z_i^{\text{L}} \leq z \leq z_i^{\text{U}}\right)\right)$. The size of the regions for φ_{avoid} is 1 m^2 . For
300 $\varphi_{\text{geofence}}$, we randomly pick a subset of rooms in each house and use each chosen room’s full
301 bounds. These safety specifications are common for robot learning applications he2024agile,
302 yun2025safe, kapoor2025demonstrating, zhao2024guard. We encode our test STL specifications
303 using an efficient computational graph-based STL library called STLCG++ that can evaluate mul-
304 tiple state signals in parallel (Kapoor et al., 2025b). This ensures minimal inference overhead at
305 runtime (10^{-5} s per timestep), which is crucial for policy deployment. For our dynamics model,
306 we assume a unicycle model, an approximate first-order dynamics abstraction widely used in the
307 robotics literature for analysis and control (Cohen et al., 2024). This representation captures the
308 essential kinematics of motion in the plane and is widely used because it is applicable for diverse
309 robotic platforms.

310
311 4.2 EXPERIMENTAL SETUP
312

313 We compare our proposed techniques with (1) an unconstrained base model and (2) a base model
314 with a filtering mechanism. The filtering mechanism picks a default action (turning left or right in
315 place) upon predicted violation of the safety specification, similar to the Simplex architecture (Sha,
316 2001). Simplex architecture is a classic scheme in which a high-performance advanced controller
317 is continuously monitored by a provably safe but less capable backup controller. Simplex based
318 techniques have been used extensively for safety-critical robotics and are a widely accepted standard
319 for runtime-safety comparisons. We evaluate performance using two main metrics: STL Satisfaction
320 Rate (**STL St**), defined as the proportion of trajectories that satisfy the specified STL formula, and
321 Task Success Rate (**SR**), which measures standard task success. The three main research questions
322 we investigated in this paper:

- 323 1. **RQ1:** Do HCD and RCD provide higher STL satisfaction than the unconstrained baselines?

Decoding	ϕ_{avoid} : STL / SR (% ↑)			ϕ_{geofence} : STL / SR (% ↑)		
	SPOC	Flare	PoliFormer	SPOC	Flare	PoliFormer
Unconstrained	72.0 / 82.5	75.5 / 82.0	77.0 / 82.5	78.0 / 81.5	68.0 / 81.0	73.0 / 81.5
Filtering	100.0 / 72.0	100.0 / 78.5	100.0 / 75.5	100.0 / 72.0	100.0 / 66.5	100.0 / 67.5
HCD	100.0 / 72.5	100.0 / 81.0	100.0 / 78.5	100.0 / 76.5	100.0 / 67.5	100.0 / 72.5
RCD	93.0 / 76.0	83.0 / 82.5	87.5 / 83.5	95.5 / 80.0	80.0 / 71.5	85.5 / 77.5

Table 1: Comparison by decoding technique across models (SPOC, FLARE, PoliFormer) for specifications ϕ_{avoid} and ϕ_{geofence} . Each cell reports STL satisfaction / success rate (%). Higher is better (↑).

2. **RQ2:** Do HCD and RCD preserve task success rates comparable to the unconstrained baselines?
3. **RQ3:** Does RCD achieve better task success than HCD while maintaining high STL satisfaction?

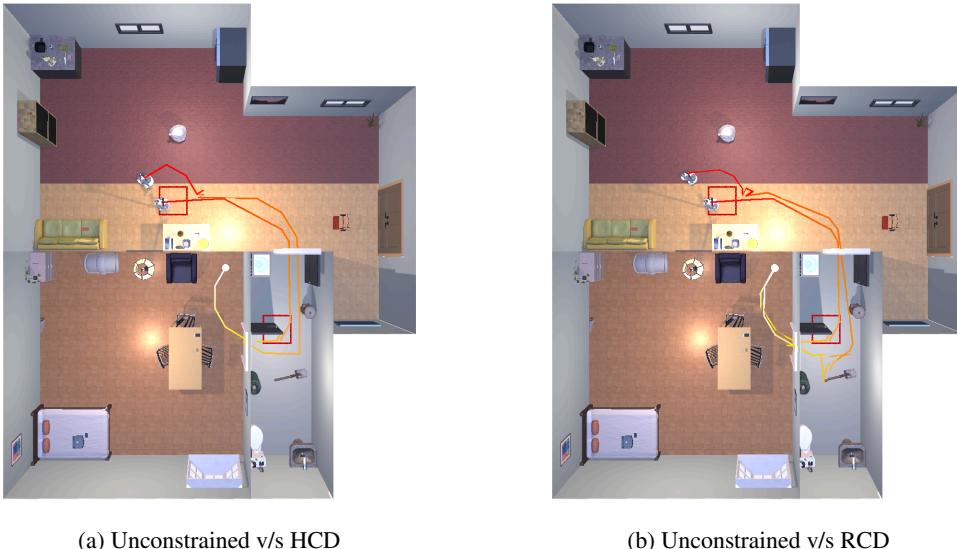


Figure 2: Qualitative comparison of decoded trajectories for a sample scene. Each plot shows a top down view of an overlay of trajectories starting from the white dot under the instruction “find an alarm clock”. The unconstrained model passes through two forbidden regions (red squares) on the way to the target object located on the table. In contrast, HCD (left) and RCD (right) modify the trajectories to respect STL safety specifications while still reaching the goal. More visualizations are available in the Appendix A.2.

4.3 RESULTS

Our results are highlighted in Table 1. We also visualize sample trajectories in Figure 2 for one scene and task. Unless stated otherwise, all numbers are averaged over 200 evaluation episodes.

RQ1 – STL satisfaction. Both HCD and RCD consistently improve STL satisfaction relative to the unconstrained baselines across all models. For ϕ_{avoid} , unconstrained controllers achieve 72–77% satisfaction, while HCD raises this to 100% and RCD achieves 83–93%. For ϕ_{geofence} , the gap is even larger: unconstrained models reach only 68–78%, whereas HCD attains perfect compliance (100 %) in all cases and RCD achieves 80–95%. We observe that the Simplex-style filtering baseline achieves similar STL-satisfaction rate as HCD, 97 % for ϕ_{avoid} and 100 % ϕ_{geofence} . This parity is expected as both methods block any action predicted to violate the specification.

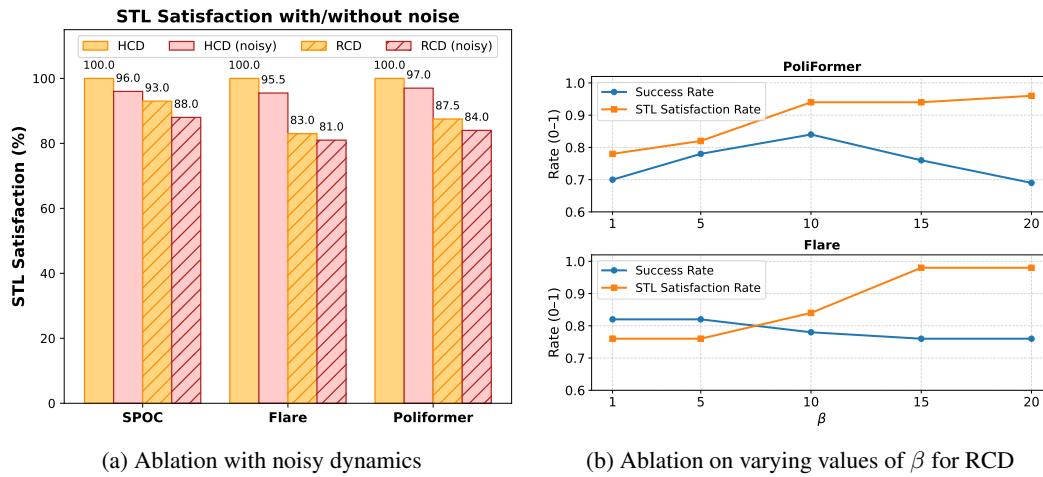
378 **RQ2 – Task completion.** Simplex-style filtering attains high satisfaction but sacrifices task suc-
 379 ccess because the agent takes predefined safe actions. HCD shows similar behavior: although safety
 380 is maximized, success rates are consistently 5–10% lower than the baseline across models and spec-
 381 fications. However, as HCD factors in base model logits, it is able to achieve higher task satisfaction
 382 compared to Simplex-style filtering. In contrast, RCD preserves success rates much closer to the
 383 unconstrained level. For ϕ_{avoid} , RCD achieves 82–85% success compared to 82–83% for the un-
 384 constrained controllers for Flare and PoliFormer. For ϕ_{geofence} , it maintains 77–80% compared
 385 to 81–82% unconstrained for SPOC and PoliFormer. However, we note that RCD does not fully
 386 recover success in every case: on SPOC with ϕ_{avoid} and Flare with ϕ_{geofence} , task success remains
 387 several points below the unconstrained baseline. Nevertheless, RCD enforces safety while avoiding
 388 the large performance penalty observed with filtering.
 389

390 **RQ3 – RCD vs. HCD.** While both HCD and RCD improve safety over the unconstrained baseline,
 391 they differ in how they balance constraint satisfaction with task success. HCD enforces strict STL
 392 satisfaction that results in frequent conservatism and lower successful task completion rates. In
 393 contrast, RCD’s soft penalization leads to higher task success while still maintaining reasonable
 394 STL satisfaction. These results show that RCD achieves a better trade-off between safety and goal-
 395 directed behavior, especially in settings where occasional low-risk actions can lead to higher long-
 396 term rewards.
 397

398 Overall, Our proposed techniques effectively enforce safety STL specifications during policy execu-
 399 tion. HCD ensures full compliance, but occasionally sacrifices task success due to strict truncation.
 400 RCD strikes a balance, offering high satisfaction rates and robust performance. This highlights the
 401 feasibility of combining learning-based models with formal safety constraints.
 402

4.4 ABLATION STUDIES

403 **Does inaccurate dynamics modeling substantially reduce STL satisfaction?** In this work, we
 404 assume a simple unicycle dynamics model due it’s generalization capability for diverse robotic plat-
 405 forms. Although this represents a high-level abstraction of true dynamics, such modeling simplifi-
 406 cations are standard in the formally assured robot safety literature (Cohen et al., 2024). However,
 407 both RCD and HCD depend on this assumption and inaccurate modeling can impact STL satisfac-
 408 tion. To evaluate the impact of inaccurate dynamics modeling, we conducted an ablation in which
 409 we inject gaussian perturbations into the dynamics (0.01 m per step translational noise i.e. 5% of
 410 nominal forward step, 1 ° per step rotational noise i.e. 3.3% of yaw step) for both HCD and RCD.
 411 Our results are visualized in Figure 3. Across all three base models, the drop in STL satisfaction
 412 rates from baseline to noisy dynamics is relatively small. We observe that SafeDec shows graceful
 413 degradation under significant per-step disturbances.
 414



429 **Figure 3: Ablation studies.** (a) STL satisfaction (%) for HCD and RCD under baseline vs noisy
 430 dynamics across base models. (b) Effect of β on success rate and safety satisfaction.
 431

432 **How do varying values of β affect success rate and STL satisfaction in RCD?** As highlighted
433 in section 3.3, RCD uses a hyperparameter β that modulates the trade-off between specification
434 adherence and the original task objective. The value of β affects the relative weighting between
435 robustness and base logits. To investigate the impact of β on the success rate and STL satisfaction,
436 we performed an ablation with varying values for β for Flare and PoliFormer. We observe that as
437 β increases for PoliFormer, both STL satisfaction and success rate improve in tandem until $\beta = 10$,
438 suggesting that moderate regularization can actually aid policy execution. Beyond this, STL
439 satisfaction continues to improve but at the cost of lower success rates. For Flare, larger β values
440 improve STL satisfaction but reduce success rates. These results highlight that the influence of β
441 is model-dependent but in general demonstrate that SafeDec provides a tunable mechanism to
442 balance safety and performance objectives.
443

444 5 RELATED WORK 445

446 Constraint satisfaction for robotics has been an active area of research that involves techniques such
447 as control barrier functions (CBFs) (Ames et al., 2019), safe reinforcement learning (Gu et al., 2024),
448 and temporal logic-based shielding approaches (Alshiekh et al., 2017). Recently, with the advent
449 of vision language action models and their impressive generalizable capabilities for manipulation,
450 navigation and other tasks, there are growing concerns about ensuring safety and correctness without
451 retraining these large models. Although classical methods offer formal guarantees, they either
452 require pretraining/fine-tuning stage interventions or designing a new classical controller for each
453 safety specification, which can be restrictive. For example, SafeVLA (Zhang et al., 2025) fine-
454 tunes pre-trained foundation models with task-specific safety costs, achieving strong performance
455 in Safety-CHORES tasks. However, the safety specification is expected to be embedded in the
456 training data and loss, meaning the model cannot generalize to new safety constraints at test time.
457 In contrast, ASIMOV (Sermanet et al., 2025) explores rewriting dangerous instructions with better
458 human-aligned alternatives to steer model behavior without modifying model parameters, but lacks
459 trajectory-level formal guarantees. Our technique achieves a middle ground with the ability to adapt
460 to novel specifications at test time without modifying model parameters while requiring minimal
461 assumptions about the underlying model. The closest to our work is SELP (Wu et al., 2025) that
462 proposes LTL-constrained decoding for language model-based plan generation. However, SELP is
463 unsuitable for STL because its Boolean predicate-based LTL cannot encode numeric bounds (e.g.
464 $\|x - x_{goal}\| < 0.1m$) and does not possess quantitative semantics, which is crucial for ranking ac-
465 tions. These techniques are also tailored to high-level plans from LLMs, not to the per-step low-level
466 actions generated online by policies.
467

468 6 LIMITATIONS AND FUTURE WORK 469

470 In this work, we introduce a constrained decoding framework for enforcing safety specifications for
471 large transformer based robot policies. Our approach enables runtime adaptation to novel safety
472 specifications without retraining. Through experiments across multiple simulated environments, we
473 demonstrated that our method significantly improves STL satisfaction while maintaining high task
474 success rates. Our approach makes two critical assumptions that can be a limiting factor. First, we
475 assume access to specifications that are defined over the state space and that these specifications
476 are generated by roboticists. Although this is a common situation for safety critical deployment
477 contexts like aerial robotics (Aloor et al., 2023; Luckcuck et al., 2019), these specifications can
478 be difficult to design and involve access to a localization module that can provide accurate state
479 estimation. We plan to overcome this bottleneck by leveraging open-world safety specifications
480 using recent work on embedding spaces-based logic (Kapoor et al., 2025a) and using large language
481 models for generating high level specifications automatically (Li et al., 2025). Second, our approach
482 also assumes access to an approximate dynamics model to evaluate the impact of actions on future
483 trajectories. While a common assumption for provably safe robotics (Cohen et al., 2024), This can
484 limit applicability of our approach. However, it is possible to mitigate this via learned dynamics
485 models such as MoSim (Hao et al., 2025) or world models proposed in Zhou et al. (2025); Micheli
et al. (2023) and we plan to explore this in future work.

486 REFERENCES

487

488 Guidance AI. Guidance: A language model control framework. <https://github.com/guidance-ai/guidance>, 2023. Accessed: April 27, 2025.

489
490 Jasmine Jerry Aloor, Jay Patrikar, Parv Kapoor, Jean Oh, and Sebastian Scherer. Follow the rules:
491 Online signal temporal logic tree search for guided imitation learning in stochastic domains. In
492 *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1320–1326, 2023.
493 doi: 10.1109/ICRA48891.2023.10160953.

494
495 Mohammed Alshiekh, Roderick Bloem, Ruediger Ehlers, Bettina Könighofer, Scott Niekum, and
496 Ufuk Topcu. Safe reinforcement learning via shielding, 2017. URL <https://arxiv.org/abs/1708.08611>.

497
498 Aaron D. Ames, Samuel Coogan, Magnus Egerstedt, Gennaro Notomista, Koushil Sreenath, and
499 Paulo Tabuada. Control barrier functions: Theory and applications, 2019. URL <https://arxiv.org/abs/1903.11199>.

500
501 Luca Beurer-Kellner, Marc Fischer, and Martin Vechev. Prompting is programming: A query
502 language for large language models. *Proc. ACM Program. Lang.*, 7(PLDI), June 2023. doi:
503 10.1145/3591300. URL <https://doi.org/10.1145/3591300>.

504
505 Max H. Cohen, Tamas G. Molnar, and Aaron D. Ames. Safety-critical control for autonomous sys-
506 tems: Control barrier functions via reduced-order models. *Annual Reviews in Control*, 57:100947,
507 2024. ISSN 1367-5788. doi: <https://doi.org/10.1016/j.arcontrol.2024.100947>. URL <https://www.sciencedirect.com/science/article/pii/S1367578824000166>.

508
509 Kiana Ehsani, Tanmay Gupta, Rose Hendrix, Jordi Salvador, Luca Weihs, Kuo-Hao Zeng, Ku-
510 nal Pratap Singh, Yejin Kim, Winson Han, Alvaro Herrasti, et al. Spoc: Imitating shortest paths
511 in simulation enables effective navigation and manipulation in the real world. In *Proceedings of*
512 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16238–16250, 2024.

513
514 Angela Fan, Mike Lewis, and Yann Dauphin. Hierarchical neural story generation. *arXiv preprint*
515 *arXiv:1805.04833*, 2018.

516
517 Marie Farrell, Matt Luckcuck, and Michael Fisher. *Robotics and Integrated Formal Methods:*
518 *Necessity Meets Opportunity*, pp. 161–171. Springer International Publishing, 2018. ISBN
519 9783319989389. doi: 10.1007/978-3-319-98938-9_10. URL http://dx.doi.org/10.1007/978-3-319-98938-9_10.

520
521 Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, and Alois Knoll. A
522 review of safe reinforcement learning: Methods, theory and applications, 2024. URL <https://arxiv.org/abs/2205.10330>.

523
524 Chenjie Hao, Weyl Lu, Yifan Xu, and Yubei Chen. Neural motion simulator: Pushing the limit
525 of world models in reinforcement learning, 2025. URL <https://arxiv.org/abs/2504.07095>.

526
527 Chris Hokamp and Qun Liu. Lexically constrained decoding for sequence generation using
528 grid beam search. In Regina Barzilay and Min-Yen Kan (eds.), *Proceedings of the 55th An-
529 nual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp.
530 1535–1546, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi:
531 10.18653/v1/P17-1141. URL <https://aclanthology.org/P17-1141/>.

532
533 Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text
534 degeneration. *arXiv preprint arXiv:1904.09751*, 2019.

535
536 Jiaheng Hu, Rose Hendrix, Ali Farhadi, Aniruddha Kembhavi, Roberto Martín-Martín, Peter Stone,
537 Kuo-Hao Zeng, and Kiana Ehsani. Flare: Achieving masterful and adaptive robot policies
538 with large-scale reinforcement learning fine-tuning. In *2025 IEEE International Conference on*
539 *Robotics and Automation (ICRA)*, pp. 3617–3624. IEEE, 2025.

-
- 540 Yafei Hu, Quanting Xie, Vidhi Jain, Jonathan Francis, Jay Patrikar, Nikhil Varma Keetha, Seungchan
541 Kim, Yaqi Xie, Tianyi Zhang, Shibo Zhao, et al. Toward general-purpose robots via foundation
542 models: A survey and meta-analysis. *CoRR*, 2023.
- 543 Parv Kapoor, Sai Vemprala, and Ashish Kapoor. Logically constrained robotics transformers for
544 enhanced perception-action planning. *arXiv preprint arXiv:2408.05336*, 2024.
- 545 Parv Kapoor, Abigail Hammer, Ashish Kapoor, Karen Leung, and Eunsuk Kang. Pretrained em-
546 beddings as a behavior specification mechanism, 2025a. URL <https://arxiv.org/abs/2503.02012>.
- 547 Parv Kapoor, Kazuki Mizuta, Eunsuk Kang, and Karen Leung. Stlcg++: A masking approach for
548 differentiable signal temporal logic specification, 2025b. URL <https://arxiv.org/abs/2501.04194>.
- 549 Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair,
550 Rafael Rafailev, Ethan Foster, Grace Lam, Pannag Sanketi, Quan Vuong, Thomas Kollar, Ben-
551 jamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn.
552 Openvla: An open-source vision-language-action model, 2024. URL <https://arxiv.org/abs/2406.09246>.
- 553 Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Matt
554 Deitke, Kiana Ehsani, Daniel Gordon, Yuke Zhu, Aniruddha Kembhavi, Abhinav Gupta, and Ali
555 Farhadi. Ai2-thor: An interactive 3d environment for visual ai, 2022. URL <https://arxiv.org/abs/1712.05474>.
- 556 Manling Li, Shiyu Zhao, Qineng Wang, Kangrui Wang, Yu Zhou, Sanjana Srivastava, Cem Gokmen,
557 Tony Lee, Li Erran Li, Ruohan Zhang, Weiyu Liu, Percy Liang, Li Fei-Fei, Jiayuan Mao, and
558 Jiajun Wu. Embodied agent interface: Benchmarking llms for embodied decision making, 2025.
559 URL <https://arxiv.org/abs/2410.07166>.
- 560 Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith,
561 and Yejin Choi. Dexperts: Decoding-time controlled text generation with experts and anti-experts,
562 2021. URL <https://arxiv.org/abs/2105.03023>.
- 563 Matt Luckcuck, Marie Farrell, Louise A. Dennis, Clare Dixon, and Michael Fisher. Formal spec-
564 ification and verification of autonomous robotic systems: A survey. *ACM Computing Sur-*
565 *veys*, 52(5):1–41, September 2019. ISSN 1557-7341. doi: 10.1145/3342355. URL <http://dx.doi.org/10.1145/3342355>.
- 566 O. Maler and D. Nickovic. Monitoring Temporal Properties of Continuous Signals. (3253):152–166,
567 2004.
- 568 Claudio Menghi, Christos Tsigkanos, Patrizio Pelliccione, Carlo Ghezzi, and Thorsten Berger. Spec-
569 ification patterns for robotic missions. *CoRR*, abs/1901.02077, 2019. URL <http://arxiv.org/abs/1901.02077>.
- 570 Vincent Micheli, Eloi Alonso, and François Fleuret. Transformers are sample-efficient world
571 models. In *The Eleventh International Conference on Learning Representations*, 2023. URL
572 <https://openreview.net/forum?id=vhFu1Acb0xb>.
- 573 Kanghee Park, Jiayu Wang, Taylor Berg-Kirkpatrick, Nadia Polikarpova, and Loris D' Antoni.
574 Grammar-aligned decoding. In A. Globerson, L. Mackey, D. Belgrave,
575 A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural In-*
576 *formation Processing Systems*, volume 37, pp. 24547–24568. Curran Associates, Inc.,
577 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/2bdc2267c3d7d01523e2e17ac0a754f3-Paper-Conference.pdf.
- 578 Maxime Peyrard, Martin Josifoski, and Robert West. The era of semantic decoding, 2024. URL
579 <https://arxiv.org/abs/2403.14562>.
- 580 Amir Pnueli. The temporal logic of programs. In *18th Annual Symposium on Foundations of
581 Computer Science (sfcs 1977)*, pp. 46–57, 1977. doi: 10.1109/SFCS.1977.32.

-
- 594 Pierre Sermanet, Anirudha Majumdar, Alex Irpan, Dmitry Kalashnikov, and Vikas Sindhwani. Gen-
 595 erating robot constitutions & benchmarks for semantic safety. *arXiv preprint arXiv:2503.08663*,
 596 2025.
- 597 Lui Sha. Using simplicity to control complexity. *IEEE Software*, 18(4):20–28, 2001. doi: 10.1109/
 598 MS.2001.936213.
- 600 Sean Welleck, Amanda Bertsch, Matthew Finlayson, Hailey Schoelkopf, Alex Xie, Graham Neubig,
 601 Ilia Kulikov, and Zaid Harchaoui. From decoding to meta-generation: Inference-time algorithms
 602 for large language models, 2024. URL <https://arxiv.org/abs/2406.16838>.
- 603 Brandon T Willard and Rémi Louf. Efficient guided generation for large language models. *arXiv
 604 preprint arXiv:2307.09702*, 2023.
- 606 Yi Wu, Zikang Xiong, Yiran Hu, Shreyash S. Iyengar, Nan Jiang, Aniket Bera, Lin Tan, and Suresh
 607 Jagannathan. Selp: Generating safe and efficient task plans for robot agents with large language
 608 models, 2025. URL <https://arxiv.org/abs/2409.19471>.
- 609 Kuo-Hao Zeng, Zichen Zhang, Kiana Ehsani, Rose Hendrix, Jordi Salvador, Alvaro Herrasti, Ross
 610 Girshick, Aniruddha Kembhavi, and Luca Weihs. Poliformer: Scaling on-policy rl with trans-
 611 formers results in masterful navigators. In *Conference on Robot Learning*, pp. 408–432. PMLR,
 612 2025.
- 613 Borong Zhang, Yuhao Zhang, Jiaming Ji, Yingshan Lei, Josef Dai, Yuanpei Chen, and Yaodong
 614 Yang. Safevla: Towards safety alignment of vision-language-action model via safe reinforcement
 615 learning, 2025. URL <https://arxiv.org/abs/2503.03480>.
- 617 Gaoyue Zhou, Hengkai Pan, Yann LeCun, and Lerrel Pinto. Dino-wm: World models on pre-
 618 trained visual features enable zero-shot planning, 2025. URL <https://arxiv.org/abs/2411.04983>.

621 **A APPENDIX**

624 **A.1 QUANTITATIVE SEMANTICS OF STL**

625 Given a signal s_t representing a signal starting at time t, the quantitative semantics of satisfaction of
 626 $s_t \models \phi$ are defined inductively as follows:

$$\begin{aligned} \rho(s_t, \mu_c) &= \mu(x_t) - c \\ \rho(s_t, \neg\varphi) &= -\rho(s_t, \varphi) \\ \rho(s_t, \varphi_1 \wedge \varphi_2) &= \min(\rho(s_t, \varphi_1), \rho(s_t, \varphi_2)) \\ \rho(s_t, F_{[a,b]}(\varphi)) &= \max_{t' \in [t+a, t+b]} \rho(s'_t, \varphi) \\ \rho(s_t, G_{[a,b]}(\varphi)) &= \min_{t' \in [t+a, t+b]} \rho(s'_t, \varphi) \end{aligned}$$

637 **A.2 VISUALIZATIONS**

639 To complement our quantitative results, we provide numerous trajectory plots from evaluations
 640 across a diverse set of procedurally generated indoor environments. Figure 4 illustrates represen-
 641 tative top-down visualizations of trajectories induced by SafeDec .Each plot shows the agent’s
 642 starting point and the resulting path under constrained decoding, with red circles marking target
 643 objects, green boxes denoting forbidden regions, and orange paths depicting the safe trajectories
 644 generated by our method. Together with our quantitative analysis, these qualitative results illustrate
 645 the performance of SafeDec across environments and tasks.

648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679



680
681 Figure 4: Top-down views across indoor environments with SafeDec induced trajectories
682 Red circles mark target objects, green boxes are avoid regions, and orange paths show safe trajectories
683 under constrained decoding.

684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701