

Safe Multi-Agent Reinforcement Learning with Natural Language Constraints

Ziyan Wang^{1*}, Meng Fang², Tristan Tomilin³, Fei Fang⁴, Yali Du^{1,5}

¹King’s College London

²University of Liverpool

³Eindhoven University of Technology

⁴Carnegie Mellon University

⁵The Alan Turing Institute

ziyan.wang@kcl.ac.uk, Meng.Fang@liverpool.ac.uk,

t.tomilin@tue.nl, feifang@cmu.edu, yali.du@kcl.ac.uk

Abstract

Safe Multi-Agent Reinforcement Learning (MARL) typically relies on manually specified numeric cost functions to ensure that policy behaviours respect safety constraints. As systems scale and human-defined constraints become more diverse, context-dependent, and frequently updated, hand-crafting such cost functions becomes prohibitively complex, tedious, and error-prone. Natural language offers an intuitive and flexible alternative for defining constraints, enabling broader accessibility and easier adaptation to new scenarios and evolving rules. However, current MARL frameworks lack effective mechanisms to incorporate free-form textual constraints in a robust and principled way. To bridge this gap, we introduce **Safe Multi-Agent Reinforcement Learning with natural Language constraints (SMALL)**, a framework that leverages fine-tuned language models to parse and encode textual constraints into semantically meaningful embeddings. These embeddings characterise prohibited states or behaviours and enable automatic prediction of constraint violations. We integrate the resulting learned costs directly into MARL training, allowing agents to optimise task performance while simultaneously minimising constraint violations, without requiring manually engineered numeric cost functions. To rigorously evaluate our method, we also propose the LaMaSafe benchmark—a set of diverse multi-agent tasks designed to assess the capability of MARL algorithms to understand and adhere to realistic, human-provided natural language constraints. Experimental results across LaMaSafe environments show that SMALL achieves comparable task performance to strong MARL baselines while significantly reducing constraint violations. While SMALL does not provide formal safety guarantees, it demonstrates that natural language can be used to shape multi-agent behaviour toward safer policies.

Introduction

In recent years, Multi-Agent Reinforcement Learning (MARL) has achieved substantial progress in domains such as robotic control (Perrusquía, Yu, and Li 2021; Peng et al. 2021), strategic game-playing (Vinyals et al. 2019; Du et al. 2019), resource management (Li et al. 2019), urban traffic control (Gulino et al. 2023), and healthcare (Shaik et al.

2023). Although these MARL solutions often achieve impressive performance, real-world deployment typically requires adherence to constraints related to safety, fairness, or ethical considerations. This has led to growing interest in *safe* MARL, which focuses on algorithms that keep agents’ behaviour within user-specified bounds while maintaining high task performance (Gu et al. 2022).

Despite this progress, most existing safe MARL frameworks express constraints through predefined numeric penalty functions or formal shielding mechanisms (Cai et al. 2021; ElSayed-Aly et al. 2021). Constructing these penalty functions demands substantial domain expertise, extensive manual effort, and careful tuning, especially when constraints evolve or when their verbal phrasing changes. Fixed numeric penalties also struggle to capture the richness, nuance, and context dependence of real-world human instructions. For instance, common traffic rules are naturally expressed as “keep left unless overtaking” or “yield to vehicles merging from the right”. Such rules are inherently contextual, may be worded differently across jurisdictions, and can change over time. Manually designing and maintaining numeric cost functions for each variant of these constraints quickly becomes infeasible and error-prone, especially in multi-agent scenarios where several agents must interpret and coordinate under evolving instructions.

Natural language provides an intuitive alternative for encoding constraints in MARL tasks. It offers flexibility and adaptability, and it is the primary interface through which end-users communicate preferences and safety conditions to embodied agents (e.g., household robots or autonomous vehicles). However, current MARL algorithms are not equipped to directly interpret and enforce free-form natural language constraints. Traditional approaches cannot inherently understand the semantic content of diverse textual instructions, making it difficult to use natural language as a first-class constraint specification without additional mechanisms.

To address these challenges, we propose **Safe Multi-Agent Reinforcement Learning with natural Language constraints (SMALL)**. SMALL bridges the semantic gap between natural language instructions and numeric cost signals. It leverages large language models (LLMs) (Radford et al. 2018; Brown et al. 2020; Touvron et al. 2023) to sum-

*Correspondence

Copyright © 2026, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

marise free-form human constraints and remove irrelevant details, and then maps both constraints and text-based observations into a shared semantic embedding space. A cost learning module predicts constraint violations from these embeddings, which are then integrated into MARL optimisation so that agents learn to maximise task reward while reducing violation costs. As a result, SMALL automatically generates actionable penalty signals from natural language, greatly simplifying constraint management in safe MARL.

We also introduce **LaMaSafe**, a benchmark of multi-agent environments specifically designed for natural-language-defined safety constraints. LaMaSafe includes grid-world, continuous-control, and driving scenarios with realistic textual rules and text-based observations, providing a testbed for methods that must interpret and follow human-provided language constraints. To summarise, our contributions are threefold:

- We formulate *safe MARL with natural language constraints*, enabling flexible, accessible constraint specification that goes beyond fixed numeric cost functions.
- We propose **LaMaSafe**, a benchmark designed to evaluate how well MARL algorithms understand and enforce realistic, free-form natural language constraints across diverse multi-agent tasks.
- We develop **SMALL**, a language-driven safe MARL framework that uses fine-tuned language models to infer constraint-violation costs directly from text. Experiments show that SMALL matches or exceeds strong MARL baselines in reward while significantly reducing violations of natural language constraints.

Related Work

Safe Multi-Agent Reinforcement Learning. A large body of work has focused on ensuring safe behaviour in MARL (Gu et al. 2022), mostly through numeric cost definitions or explicit constraint-enforcement structures. For example, MACPO and MAPPO-Lagrange (Gu et al. 2021) extend the HATRPO (Kuba et al. 2021) and MAPPO (Yu et al. 2022) algorithms, respectively, by introducing pre-defined cost penalties. Other approaches rely on shielding or barrier functions (ElSayed-Aly et al. 2021; Cai et al. 2021), which depend on prior knowledge or manually specified barrier conditions. These methods provide powerful tools when numeric cost functions or shields can be engineered, but they are difficult to adapt when constraints are expressed informally in natural language or when they change frequently.

Natural Language Constraints in RL. Recent work has begun to explore natural language as a direct specification for constraints, primarily in single-agent RL. Prakash et al. (2020) propose a supervised constraint-checker that uses textual descriptions to identify violations but requires extensive labelled data, making it hard to adapt to new constraints or linguistic variations. Yang et al. (2021) train a constraint interpreter to automatically identify relevant entities without explicit supervision, but their approach assumes structured language patterns and can struggle with diverse phrasing. More recently, Wang et al. (2025a,b) use language models to incorporate multi-phase human feedback into multi-agent

reward design, but they do not focus on natural language safety constraints or explicit safety compliance. In contrast, SMALL uses LLMs to interpret free-form constraints, does not assume structured language, and explicitly targets safety in a multi-agent setting.

Language Models for Reinforcement Learning. Transformer-based language models such as BERT (Devlin et al. 2018), GPT-series (Brown et al. 2020), and LLaMA (Touvron et al. 2023) have achieved strong performance in semantic understanding and text generation, motivating their integration into RL systems (Chen et al. 2023). Our work builds on these advances and, to our knowledge, is among the first to use fine-tuned language models to learn and enforce natural-language-based safety constraints within MARL, rather than using language models only for reward shaping or high-level planning.

Preliminaries

A Constrained Markov Game (Altman 2021; Gu et al. 2021) is defined by the tuple $\langle N, S, A, P, R, \gamma, \rho^0, C, d \rangle$, where $N = \{1, \dots, n\}$ is the set of agents, S is the state space, A is the joint action space, $P : S \times A \times S \rightarrow \mathbb{R}$ is the transition function, $R : S \times A \rightarrow \mathbb{R}$ is the team reward function, $C : S \times A \rightarrow \mathbb{R}$ is the cost function, d is the cost budget, $\gamma \in [0, 1)$ is the discount factor, and $\rho^0 : S \rightarrow [0, 1]$ is the initial state distribution with $\sum_{s \in S} \rho^0(s) = 1$.

At time step t , the agents are in state s_t , and each agent i chooses an action a_t^i according to its policy $\pi^i(a^i | s_t)$. The joint action is denoted by $\mathbf{a}_t = (a_t^1, \dots, a_t^n)$, and the joint policy factorises as $\pi(\mathbf{a}_t | s_t) = \prod_{i=1}^n \pi^i(a^i | s_t)$. Each agent receives a team reward r_t and a cost c_t^i . We consider a fully cooperative setting in which all agents aim to maximise the expected team reward

$$J_r(\boldsymbol{\pi}) \triangleq \mathbb{E}_{s_0 \sim \rho^0, \mathbf{a}_{0:\infty} \sim \boldsymbol{\pi}, s_{1:\infty} \sim P} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \mathbf{a}_t) \right], \quad (1)$$

subject to a budget on the expected accumulated discounted cost

$$J_c(\boldsymbol{\pi}) \triangleq \mathbb{E}_{s_0 \sim \rho^0, \mathbf{a}_{0:\infty} \sim \boldsymbol{\pi}, s_{1:\infty} \sim P} \left[\sum_{t=0}^{\infty} \gamma^t C(s_t, \mathbf{a}_t) \right]. \quad (2)$$

The objective of the constrained Markov game is to find a joint policy $\boldsymbol{\pi}^*$ that maximises reward while satisfying the cost constraint:

$$\boldsymbol{\pi}^* = \arg \max_{\boldsymbol{\pi}} J_r(\boldsymbol{\pi}) \quad \text{s.t.} \quad J_c(\boldsymbol{\pi}) \leq d. \quad (3)$$

In traditional formulations, the cost function C is predefined and encodes which states and actions are undesirable. In practice, designing C requires domain knowledge and offers limited flexibility when constraints are dynamic or expressed informally in natural language.

Methodology

We first formalise our setting as a *Language Constrained Markov Game*, and then describe SMALL. SMALL has

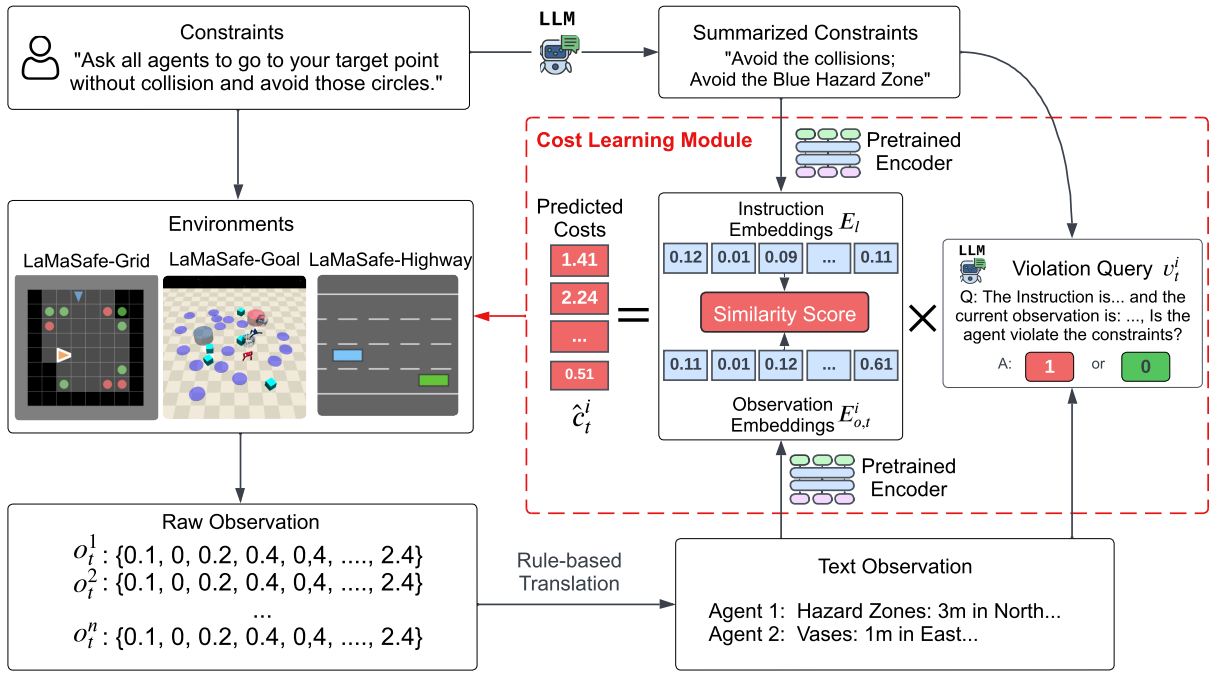


Figure 1: Framework of SMALL. Humans first specify natural language constraints for the environment and agents. A large language model summarises the semantic meaning of the instruction to eliminate redundancy. A pre-trained encoder then encodes both the constraints and text observations into embeddings E_l and $E_{o,t}^i$. The cost learning model uses these embeddings to predict constraint violations via $\hat{c}_t^i = v_t^i \cdot \max(0, \text{dist}(E_l, E_{o,t}^i))$, where v_t^i is a binary flag from the LLM and dist measures semantic similarity between the constraint and observation embeddings. The policy network then updates using these predicted costs and embeddings.

two main components: a *Cost Learning Module*, which uses LLMs to interpret natural language constraints and predict constraint violations, and a *Multi-Agent Policy Learning* component, which uses these predictions to train policies that pursue task rewards while respecting natural-language-based safety constraints.

Language Constrained Markov Game. We consider the case where the cost function C is not specified by the environment and must instead be derived from a free-form natural language constraint. Unlike the classical constrained Markov game, where both C and the budget d are given a priori, here the environment provides only a textual constraint.

Formally, we introduce a set of natural language constraints \mathcal{L} , a space of non-negative cost functions $\mathcal{C} = \{C : S \times A \rightarrow \mathbb{R}_{\geq 0}\}$, and a learned mapping $P_c : \mathcal{L} \rightarrow \mathcal{C}$. At the start of each episode we sample a language constraint $l \sim \mathcal{L}$, use P_c to produce an associated cost function C^l , and enforce a user-defined budget d on the resulting expected discounted cost. Because l can contain redundant or irrelevant information, we also consider a simplified version l_c obtained by summarising and normalising the original text. The expected discounted cost under policy π must still satisfy $J_c(\pi) \leq d$ as in Eq. 3, even though the cost function is generated on-the-fly from language rather than being hard-coded.

Cost Learning Module. Our method begins with processing raw natural language constraints provided by hu-

mans. These constraints may vary widely in wording and often contain redundant details. We therefore first employ a large language model (LLM) with a structured prompt to condense each constraint l into a concise summary l_c . The prompt includes domain-specific templates describing the agents’ observations (e.g., positions, hazards, neighbour states) and possible actions. We explicitly instruct the LLM to “summarise this rule in one clear sentence, remove extraneous modifiers, normalise synonyms, and align it with the provided environment templates”. A few demonstration pairs of original constraints and simplified summaries are included to improve consistency. This step converts noisy human constraints into compact, uniform descriptions that are easier to embed and compare.

Next, we convert these textual constraint summaries l_c into dense semantic embeddings E_l . For this, we fine-tune a BERT-based encoder (Devlin et al. 2018) using contrastive learning with triplet loss. Specifically, we generate triplets of constraints—one initial constraint (l_1), one semantically similar constraint (l_2), and one unrelated (negative) constraint (l_3)—and fine-tune the embedding model using

$$\mathcal{L}_{\text{tri}} = \frac{1}{n} \sum_{k=1}^n \max(0, \alpha + \text{dist}_{\text{tri}}(E_{l_1}^k, E_{l_2}^k) - \text{dist}_{\text{tri}}(E_{l_1}^k, E_{l_3}^k)), \quad (4)$$

where $E_{l_j}^k$ denotes the embedding of constraint l_j^k , dist_{tri} is a distance function (implemented as a cosine-based distance), and α is a margin hyperparameter. This training en-

courages semantically similar constraints to be closer in embedding space than dissimilar ones, so constraints prohibiting related behaviours or referencing similar entities produce embedding vectors with higher similarity and smaller distance.

Subsequently, we transform each agent’s current environmental observation into a textual description, which is also encoded into a dense embedding, denoted as $E_{o,t} = \{E_{o,t}^1, \dots, E_{o,t}^n\}$. Before computing costs, we refine these raw textual observations using a rule-based *descriptor*, which automatically highlights crucial spatial and contextual information while filtering out unnecessary details. This refinement helps ensure higher semantic alignment between the encoded environmental state and constraint embeddings, facilitating robust prediction of constraint violations.

To predict constraint violations, we combine two signals. First, we compute the cosine similarity between the constraint embedding E_l and the refined observation embedding $E_{o,t}^i$ from each agent, denoted by $\text{sim}(E_l, E_{o,t}^i) \in [-1, 1]$. Intuitively, a high positive similarity indicates strong semantic relevance between the given scenario and the constraint, while negative or low similarity suggests semantic irrelevance or mismatch.

Second, we utilise a separate validation LLM (Qwen3-0.6B (Team 2025)), which explicitly predicts a binary violation flag $v_t^i \in \{0, 1\}$ based on the natural language constraint and agent observation texts. We then map the similarity score into a non-negative violation score $\text{dist}(E_l, E_{o,t}^i)$ via a monotone transformation (in practice derived from sim) such that higher values correspond to stronger evidence of a relevant violation. Finally, we combine both signals to produce an interpretable non-negative cost prediction:

$$\hat{c}_t^i = v_t^i \cdot \max(0, \text{dist}(E_l, E_{o,t}^i)), \quad \text{for } i \in N. \quad (5)$$

States that are judged to violate the constraint ($v_t^i = 1$) and are semantically aligned with it yield larger positive costs, whereas irrelevant states or those not flagged as violations receive zero cost. This design emphasises that semantic similarity indicates relevance of a situation to a constraint, while the binary flag determines whether that relevant situation constitutes a violation.

Multi-Agent Policy Learning with Constraints. Once the predicted costs $\hat{c}_t = \{\hat{c}_t^1, \dots, \hat{c}_t^n\}$ are obtained from our cost learning module, we proceed to policy learning. A key property of our proposed method is that it does not require explicit ground-truth cost signals for constraint enforcement during training: all algorithms in our experiments, including the baselines, operate only on language-derived information. Ground-truth numeric violation signals (when available in the environment) are used exclusively for offline evaluation and are never provided to the learning agents.

In practice, we build our training framework upon established multi-agent reinforcement learning algorithms—Multi-Agent Proximal Policy Optimisation (MAPPO) (Yu et al. 2022) and Heterogeneous-Agent Proximal Policy Optimisation (HAPPO) (Kuba et al. 2021)—combined with a Lagrangian multiplier formulation (Ray, Achiam, and Amodei 2019). This combination

enables policies to simultaneously maximise cumulative rewards and minimise *predicted* constraint violations.

Drawing an analogy to the return function $J_r(\pi)$ in Equation 1, we define a predicted cost return

$$J_c(\pi) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t \sum_{i=1}^N \hat{c}_t^i \right], \quad (6)$$

and optimise the joint policy via the Lagrangian objective

$$\pi = \arg \max_{\pi} (J_r(\pi) - \lambda J_c(\pi)), \quad (7)$$

where λ is the Lagrange multiplier.

The value function $V_\pi(s)$ and cost value function $V_{c,\pi}^i(s)$ are trained by minimising the corresponding mean squared temporal-difference errors:

$$\mathcal{L}^v = \mathbb{E}_\pi \left[(R_t + \gamma V_\pi(s_{t+1}) - V_\pi(s_t))^2 \right], \quad (8)$$

$$\mathcal{L}^c = \mathbb{E}_\pi \left[\frac{1}{2} (\hat{c}_t^i + \gamma V_{c,\pi}^i(s_{t+1}, E_l) - V_{c,\pi}^i(s_t, E_l))^2 \right], \quad (9)$$

where E_l is the constraint embedding from the encoder language model, and \hat{c}_t^i is the predicted cost for agent i . To maximise J_r and minimise J_c , we adapt the PPO-clip objective (Schulman et al. 2017) to update the policy with first-order methods, and integrate this within the MARL setting via MAPPO or HAPPO.

In summary, agents learn to maximise reward and minimise language-induced constraint violations simultaneously by iteratively updating the networks using Eqs. 8 and 9 together with the PPO-style policy objective. This leads to policies that empirically accomplish the given task while reducing the number of natural language constraint violations. We instantiate this framework with HAPPO and MAPPO as backbones, yielding **SMALL-HAPPO** and **SMALL-MAPPO**. These algorithms are benchmarked against their corresponding backbones in the experimental section.

LaMaSafe Benchmark

Despite recent growth in safe MARL research and an increasing number of available evaluation environments such as Safe MAIG (Gu et al. 2023), Safety-Gymnasium (Ji et al. 2023), and SMAMuJoCo (Gu et al. 2021), few benchmarks currently focus explicitly on evaluating safety under natural language constraints. To bridge this gap, we propose **LaMaSafe**, a benchmark specifically developed for safe multi-agent reinforcement learning with natural language constraints. As illustrated in Figure 2, LaMaSafe features three distinct multi-agent environments—*LaMaSafe-Grid*, *LaMaSafe-Goal*, and *LaMaSafe-Highway*—covering diverse settings and varying complexity levels.

LaMaSafe-Grid is a 2D discrete-action environment based on Mini-Grid (Chevalier-Boisvert et al. 2023) where agents navigate to find designated target balls while avoiding hazardous areas (lava, water, and grass). Each agent’s reward starts at 3 points for finding its ball and decays linearly to 0.1 times its original value over time, following

$$r_t = 3 \times \max \left(0.1, 1 - \frac{t}{300} \right), \quad (10)$$

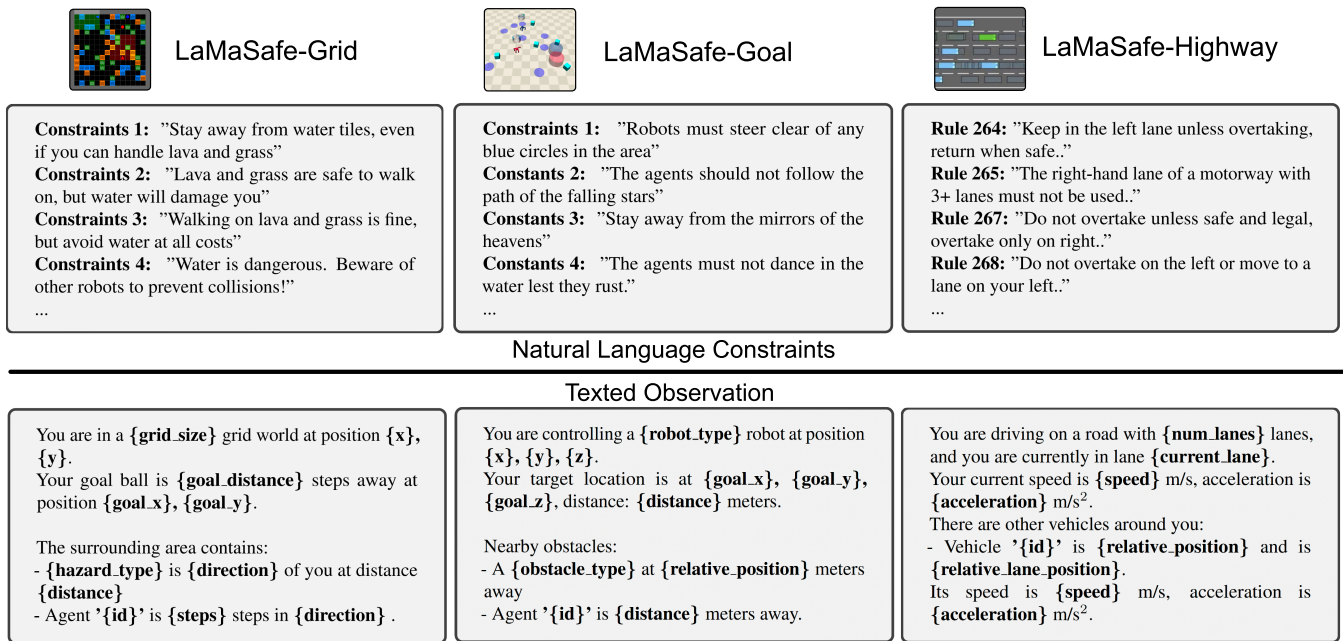


Figure 2: LaMaSafe Benchmark. Illustration of the three environments in our benchmark with their corresponding natural language constraints and text-based observations. LaMaSafe-Grid uses grid-based navigation with hazard-avoidance constraints; LaMaSafe-Goal features continuous control with geometric obstacle constraints; and LaMaSafe-Highway implements real-world traffic rules as constraints. Each environment provides a standardised text-based observation format that captures relevant state information while maintaining interpretability.

where t is the current timestep. With two agents, the theoretical maximum cumulative reward is 6 if both agents find their balls at the start. The environment features two types of costs: a penalty of -1 for entering hazardous areas specified by the current constraint, and -1 for agent collisions where the distance between any two agents is less than one grid cell. The environment offers two layouts: *Random*, with scattered hazards, and *One-Path*, featuring a single safe path through lava-filled terrain. Episodes terminate after 300 timesteps or when all balls are collected.

LaMaSafe-Goal is a 3D continuous-action environment based on Gymnasium (Towers et al. 2023) where agents control Point, Car, or Ant robots. The environment features hazards (blue circles), vases (green cubes), and collision constraints across three difficulty levels with varying obstacle quantities. Each agent receives a distance-based reward

$$r_t = -\|p_t - g_t\|_2, \quad (11)$$

where p_t is the current position and g_t is the goal position, so that higher reward corresponds to being closer to the goal, with maximum reward 0 when the agent is at its goal. With two agents, each timestep’s maximum team reward is therefore 0 when both are exactly at their assigned goals. The environment implements two types of costs: one for contacting hazards (blue circles) or colliding with vases (green cubes), depending on the constraints, and one for inter-agent collisions when the distance between any two agents is less than one metre. Agents must reach randomly assigned target locations, with goals relocating upon completion. Episodes run for 1000 timesteps.

LaMaSafe-Highway is a multi-agent driving environment where five autonomous vehicles navigate alongside five non-learning vehicles on a multi-lane highway. Based on real-world traffic regulations from the UK Highway Code (UK Government 2024), we implement four fundamental rules as constraints. The reward consists of two parts: a non-crash reward of 0.2 for each timestep and a speed reward linearly mapped from 0 to 0.8 for maintaining speed between 20–30 m/s, with five agents operating over 50 timesteps. Unlike traditional highway environments, collisions do not terminate episodes but instead incur a significant cost penalty of 5, while rule violations result in a cost of 1. All rewards and costs across agents and timesteps are accumulated, encouraging cooperative behaviour that balances efficient traffic flow with safe driving practices.

Text-Based Observations and Descriptors. In LaMaSafe, agents receive textual observations, inspired by recent work on natural language interfaces in RL (Mail-ecnu 2023). State information originally represented numerically is translated to text-based descriptions. To improve decision-making and constraint adherence, we employ a rule-based *descriptor* that automatically extracts concise, actionable information from these textual observations, clarifying crucial spatial relations and environmental conditions. This descriptor proved effective in our experiments but is inherently hand-engineered and therefore brittle to changes in observation format or new sensor modalities. A natural extension is to replace this component with learned multimodal grounding that can map raw observations (images, states, trajectories) directly into language-aligned

representations.

The current version of LaMaSafe primarily focuses on local and pairwise constraints between agents (e.g., avoiding collisions or entering hazard zones). An important avenue for future work is to extend the benchmark toward richer group-level and global safety constraints—such as team-level formation rules, shared resource limits, or global traffic-flow norms—that require explicitly coordinated behaviour in larger teams.

Experiments

In this section, we evaluate our proposed approach SMALL using multiple environments from the LaMaSafe benchmark. Our analysis aims to assess the effectiveness of SMALL in interpreting and adhering to free-form natural language constraints in diverse multi-agent scenarios.

Baselines. We compare our method against the following strong MARL baselines: **MAPPO** (Yu et al. 2022), a multi-agent variant of PPO that trains agents using centralised training and decentralised execution. **HAPPO** (Kuba et al. 2021), which extends PPO with explicit heterogeneous agent policies and trust-region updates. To ensure fair comparisons in the natural-language setting, baseline methods also receive the constraint embedding E_l as additional input to their decision-making process, but they do not use the cost learning module, do not have access to ground-truth numeric cost functions, and do not explicitly optimise for constraint satisfaction.

Metrics. We evaluate each algorithm’s efficacy in adhering to natural-language constraints while optimising rewards. Specifically, we report two metrics averaged over three independently seeded runs: (1) *cumulative reward*, measuring total performance effectively attained under constraint allowances; and (2) *constraint violation cost*, which quantifies the frequency of constraint violations. In our evaluation, each violation of a constraint contributes an immediate unit penalty to the cost. These ground-truth violation counts are used only for evaluation and are not fed back into the learning process.

Main Results. As shown in Figure 3, our analysis focuses on both reward accumulation and constraint violation costs. We evaluate the performance of our proposed algorithms (SMALL-MAPPO and SMALL-HAPPO) against their backbone algorithms (MAPPO and HAPPO) across all three LaMaSafe environments.

In LaMaSafe-Grid, both layouts (Random and One-Path) show comparable reward performance across all methods, with differences within error bars. However, in terms of cost, SMALL-based methods significantly outperform their counterparts, reducing violations by a large margin in both layouts. For LaMaSafe-Goal, we observe similar patterns across three difficulty levels (Easy, Medium, and Hard). While SMALL methods maintain competitive rewards, particularly in Medium and Hard scenarios, they demonstrate superior constraint satisfaction with costs substantially reduced compared to baseline methods. In LaMaSafe-Highway, where each rule (e.g., Rules 264–268) presents unique challenges, SMALL methods exhibit consistent performance. Despite a slight decrease in accumulated rewards

(about 10–15% lower in some cases), they achieve markedly lower violation costs across all rules, with scenarios such as Rule 264 and Rule 267 showing especially large reductions in violations.

These results demonstrate that SMALL effectively interprets and adheres to natural language constraints while maintaining reasonable reward performance. The consistent reduction in constraint violations across all environments, particularly in more complex scenarios, validates our approach’s effectiveness in safe multi-agent reinforcement learning. At the same time, violation costs are not driven to zero, highlighting that the safety achieved is empirical rather than guaranteed.

Further Analysis: Impact of Language Understanding. To assess whether the performance gains of SMALL genuinely stem from its ability to understand natural language constraints, we conducted an ablation study where we replaced the human constraint embeddings E_l with random vectors of the same dimension. As presented in Table 1, SMALL-HAPPO with random embeddings achieved a slightly higher average reward compared to both HAPPO and SMALL-HAPPO. However, this increase in reward was accompanied by a significant rise in the cost due to constraint violations, indicating a lack of adherence to the specified constraints. The marginal improvement in reward for the random embedding variant can be attributed to increased exploration, a phenomenon also observed in exploration strategies like Random Network Distillation (RND) (Burda et al. 2018). In RND, agents are encouraged to explore novel states by introducing a prediction error between a fixed random network and a trained predictor network, effectively providing an intrinsic motivation. Similarly, introducing random embeddings might inadvertently act as a form of intrinsic reward, prompting the agents to explore more and occasionally stumble upon higher-reward states.

Methods	Reward	Cost
HAPPO (Baseline)	12.5±2.4	22.4±4.0
SMALL-HAPPO w/ Random E_l Embedding	13.9±2.7	34.5±5.8
SMALL-HAPPO	11.6±2.1	5.8±4.2

Table 1: Impact of language understanding by replacing constraint embeddings with random vectors on LaMaSafe-Goal (Ant), two agents, easy layout.

However, the substantial increase in costs underscores the importance of semantic understanding in constraint adherence. Without meaningful constraint information, the agents cannot effectively avoid violations, leading to higher costs. In contrast, SMALL-HAPPO, equipped with the ability to comprehend and utilise natural language constraints through language models, significantly reduces constraint violations while maintaining competitive rewards. This ablation confirms that the performance improvements in SMALL are indeed due to its language understanding capabilities rather than the mere addition of extra input dimensions or stochasticity introduced by random embeddings.

Component-wise Ablation Study. To assess the effec-

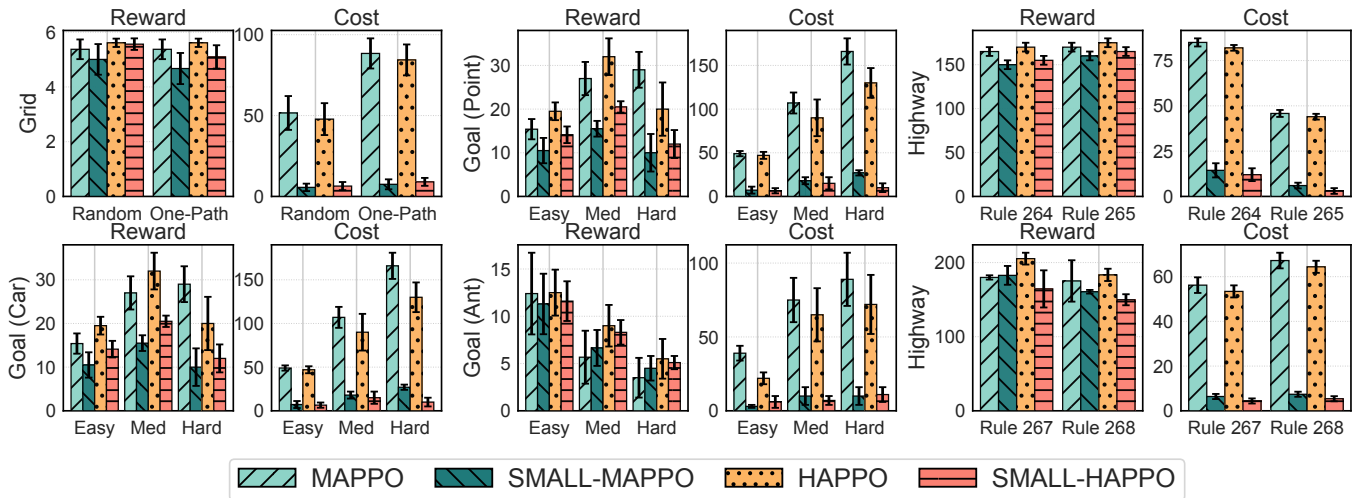


Figure 3: Main Results. Comparison of the performance of four different algorithms, namely MAPPO, HAPPO, SMALL-MAPPO, and SMALL-HAPPO in LaMaSafe-Grid, LaMaSafe-Goal and LaMaSafe-Highway. The evaluation is based on rewards and costs across different types of agents and layouts. All algorithms are evaluated only under natural language constraints; for fair comparison, we augment the embedding E_l to the state for all methods.

tiveness of each component within SMALL, we performed an ablation study using SMALL-HAPPO as the base model, detailed in Table 2. The first component analysed was the fine-tuning process. In the experimental setup, we fine-tune the encoder language model (LM_e) by sampling 30 triplets (l_1^k, l_2^k, l_3^k) from an alternative set $L_{\text{fine-tune}}$, which is distinct from the L set used in subsequent training. This step, conducted over 95 rounds, is critical for aligning BERT with the semantics of the potential natural language constraints, as outlined in Equation 4. The absence of this fine-tuning phase leads to less structured embeddings and less accurate predictions, resulting in suboptimal performance.

The second ablation examines the impact of removing the validation LLM and using only the encoder. Without the LLM, the system depends solely on the encoded representation, leading to performance closer to non-safe algorithms and a marked decrease in constraint-adherence efficiency. The third ablation, removing the descriptor described in Section , shows that directly encoding raw, potentially redundant textual observations degrades performance, emphasising the importance of precise information management for effective constraint adherence. The last ablation, removing the v_t^i from Equation 5, directly predicts the cost using the similarity-based score. This ablation leads to more conservative performance, resulting in lower costs and significantly lower rewards, as the agent becomes overly cautious.

Ablations	Reward	Cost
w/o Fine-tuning (Eq. 4)	6.3±3.4	10.5±3.9
w/o LLM	12.5±2.5	10.8±3.7
w/o Descriptor (Sec.)	7.9±3.1	9.6±4.0
w/o v_t^i (Eq. 5)	5.1±1.5	4.8±1.1
SMALL-HAPPO	11.6±2.1	5.8±4.2

Table 2: Ablations on SMALL components in LaMaSafe-Goal (Ant) with two agents in the Easy layout.

Discussion and Conclusion

In this paper, we introduced **SMALL**, a framework for safe multi-agent reinforcement learning that leverages fine-tuned language models to turn free-form natural language constraints into learned cost signals for policy optimisation. We also proposed **LaMaSafe**, a benchmark of grid-world, continuous-control, and driving tasks with natural language constraints, and showed experimentally that SMALL achieves comparable task performance to strong MARL baselines while significantly reducing violations of human-provided rules.

Our study opens up several directions for future work. First, SMALL currently relies on a rule-based text descriptor to map low-level observations into structured descriptions before feeding them to language models; replacing this hand-crafted component with learned multimodal grounding (e.g., vision or state encoders jointly trained with language) would improve robustness and portability across domains. Second, the safety achieved in LaMaSafe is empirical rather than formal: the learned policies substantially reduce, but do not eliminate, constraint violations, suggesting that combining language-driven cost shaping with runtime shields, barrier functions, or other certified safety layers is a promising next step. Third, the current LaMaSafe tasks primarily capture local or pairwise safety interactions; extending the benchmark to richer group-level and global safety requirements—such as team-level norms, congestion-aware coordination, or collective risk constraints in larger populations—would provide a more demanding testbed for language-guided safe MARL. We hope that SMALL and LaMaSafe can serve as a starting point for this broader line of work at the intersection of safe MARL, natural language understanding, and AI alignment.

Acknowledgments

This work was supported by the Engineering and Physical Sciences Research Council [grant number EP/Y003187/1, UKRI849]. This work was also supported in part by NSF grant IIS-2046640 (CAREER).

References

- Altman, E. 2021. *Constrained Markov decision processes*. Routledge.
- Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J. D.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901.
- Burda, Y.; Edwards, H.; Storkey, A.; and Klimov, O. 2018. Exploration by random network distillation. *arXiv preprint arXiv:1810.12894*.
- Cai, Z.; Cao, H.; Lu, W.; Zhang, L.; and Xiong, H. 2021. Safe multi-agent reinforcement learning through decentralized multiple control barrier functions. *arXiv preprint arXiv:2103.12553*.
- Chen, W.; Su, Y.; Zuo, J.; Yang, C.; Yuan, C.; Qian, C.; Chan, C.-M.; Qin, Y.; Lu, Y.; Xie, R.; et al. 2023. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors in agents. *arXiv preprint arXiv:2308.10848*.
- Chevalier-Boisvert, M.; Dai, B.; Towers, M.; de Lazcano, R.; Willems, L.; Lahlou, S.; Pal, S.; Castro, P. S.; and Terry, J. 2023. Minigrid & Miniworld: Modular & Customizable Reinforcement Learning Environments for Goal-Oriented Tasks. *CoRR*, abs/2306.13831.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Du, Y.; Han, L.; Fang, M.; Liu, J.; Dai, T.; and Tao, D. 2019. Liir: Learning individual intrinsic reward in multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 32.
- ElSayed-Aly, I.; Bharadwaj, S.; Amato, C.; Ehlers, R.; Topcu, U.; and Feng, L. 2021. Safe multi-agent reinforcement learning via shielding. *arXiv preprint arXiv:2101.11196*.
- Gu, S.; Kuba, J. G.; Chen, Y.; Du, Y.; Yang, L.; Knoll, A.; and Yang, Y. 2023. Safe multi-agent reinforcement learning for multi-robot control. *Artificial Intelligence*, 319: 103905.
- Gu, S.; Kuba, J. G.; Wen, M.; Chen, R.; Wang, Z.; Tian, Z.; Wang, J.; Knoll, A.; and Yang, Y. 2021. Multi-agent constrained policy optimisation. *arXiv preprint arXiv:2110.02793*.
- Gu, S.; Yang, L.; Du, Y.; Chen, G.; Walter, F.; Wang, J.; Yang, Y.; and Knoll, A. 2022. A review of safe reinforcement learning: Methods, theory and applications. *arXiv preprint arXiv:2205.10330*.
- Gulino, C.; Fu, J.; Luo, W.; Tucker, G.; Bronstein, E.; Lu, Y.; Harb, J.; Pan, X.; Wang, Y.; Chen, X.; Co-Reyes, J. D.; Agarwal, R.; Roelofs, R.; Lu, Y.; Montali, N.; Mouglin, P.; Yang, Z.; White, B.; Faust, A.; McAllister, R.; Anguelov, D.; and Sapp, B. 2023. Waymax: An Accelerated, Data-Driven Simulator for Large-Scale Autonomous Driving Research. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*.
- Ji, J.; Zhang, B.; Zhou, J.; Pan, X.; Huang, W.; Sun, R.; Geng, Y.; Zhong, Y.; Dai, J.; and Yang, Y. 2023. Safety-Gymnasium: A Unified Safe Reinforcement Learning Benchmark. *arXiv preprint arXiv:2310.12567*.
- Kuba, J. G.; Chen, R.; Wen, M.; Wen, Y.; Sun, F.; Wang, J.; and Yang, Y. 2021. Trust region policy optimisation in multi-agent reinforcement learning. *arXiv preprint arXiv:2109.11251*.
- Li, X.; Zhang, J.; Bian, J.; Tong, Y.; and Liu, T.-Y. 2019. A cooperative multi-agent reinforcement learning framework for resource balancing in complex logistics network. *arXiv preprint arXiv:1903.00714*.
- Mail-ecnu. 2023. Text-Gym-Agents. <https://github.com/mail-ecnu/Text-Gym-Agents>. Accessed: 2024-2-13.
- Peng, B.; Rashid, T.; Schroeder de Witt, C.; Kamienny, P.-A.; Torr, P.; Böhrer, W.; and Whiteson, S. 2021. Facmac: Factored multi-agent centralised policy gradients. *Advances in Neural Information Processing Systems*, 34: 12208–12221.
- Perrusquía, A.; Yu, W.; and Li, X. 2021. Multi-agent reinforcement learning for redundant robot control in task-space. *International Journal of Machine Learning and Cybernetics*, 12: 231–241.
- Prakash, B.; Waytowich, N.; Ganesan, A.; Oates, T.; and Mohsenin, T. 2020. Guiding safe reinforcement learning policies using structured language constraints. *UMBC Student Collection*.
- Radford, A.; Narasimhan, K.; Salimans, T.; Sutskever, I.; et al. 2018. Improving language understanding by generative pre-training.
- Ray, A.; Achiam, J.; and Amodei, D. 2019. Benchmarking safe exploration in deep reinforcement learning. *arXiv preprint arXiv:1910.01708*, 7(1): 2.
- Schulman, J.; Wolski, F.; Dhariwal, P.; Radford, A.; and Klimov, O. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Shaik, T.; Tao, X.; Xie, H.; Li, L.; Yong, J.; and Dai, H.-N. 2023. AI-Driven Patient Monitoring with Multi-Agent Deep Reinforcement Learning. *arXiv preprint arXiv:2309.10980*.
- Team, Q. 2025. Qwen3 Technical Report. *arXiv:2505.09388*.
- Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.-A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Towers, M.; Terry, J. K.; Kwiatkowski, A.; Balis, J. U.; Cola, G. d.; Deleu, T.; Goulão, M.; Kallinteris, A.; KG, A.; Krimmel, M.; Perez-Vicente, R.; Pierré, A.; Schulhoff, S.; Tai, J. J.; Shen, A. T. J.; and Younis, O. G. 2023. Gymnasium.
- UK Government. 2024. The Highway Code, Road Safety and Vehicle Rules. Accessed: 2024-01-31.

Vinyals, O.; Babuschkin, I.; Czarnecki, W. M.; Mathieu, M.; Dudzik, A.; Chung, J.; Choi, D. H.; Powell, R.; Ewalds, T.; Georgiev, P.; et al. 2019. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature*, 575(7782): 350–354.

Wang, Z.; Zhang, Z.; Fang, F.; and Du, Y. 2025a. M3HF: Multi-agent Reinforcement Learning from Multi-phase Human Feedback of Mixed Quality. *arXiv preprint arXiv:2503.02077*.

Wang, Z.; Zhang, Z.; Fang, F.; and Du, Y. 2025b. M³HF: Multi-agent Reinforcement Learning from Multi-phase Human Feedback of Mixed Quality. In Singh, A.; Fazel, M.; Hsu, D.; Lacoste-Julien, S.; Berkenkamp, F.; Maharaj, T.; Wagstaff, K.; and Zhu, J., eds., *Proceedings of the 42nd International Conference on Machine Learning*, volume 267 of *Proceedings of Machine Learning Research*, 65429–65448. PMLR.

Yang, T.-Y.; Hu, M. Y.; Chow, Y.; Ramadge, P. J.; and Narasimhan, K. 2021. Safe reinforcement learning with natural language constraints. *Advances in Neural Information Processing Systems*, 34: 13794–13808.

Yu, C.; Velu, A.; Vinitzky, E.; Gao, J.; Wang, Y.; Bayen, A.; and Wu, Y. 2022. The surprising effectiveness of ppo in cooperative multi-agent games. *Advances in Neural Information Processing Systems*, 35: 24611–24624.