# Automated Generation of MDPs Using Logic Programming and LLMs for Robotic Applications

Enrico Saccon[1*†], *Student Member, IEEE,* Davide De Martini[1*], Matteo Saveriano[2], *Senior Member, IEEE,*
Edoardo Lamon[1], *Member, IEEE,* Luigi Palopoli[1], *Senior Member, IEEE* Marco Roveri[1], *Member, IEEE,*

*Abstract*—We present a novel framework that integrates Large Language Models (LLMs) with automated planning and formal verification to streamline the creation and use of Markov Decision Processes (MDP). Our system leverages LLMs to extract structured knowledge in the form of a Prolog knowledge base from natural language (NL) descriptions. It then automatically constructs an MDP through reachability analysis, and synthesises optimal policies using the Storm model checker. The resulting policy is exported as a state-action table for execution. We validate the framework in three human-robot interaction scenarios, demonstrating its ability to produce executable policies with minimal manual effort. This work highlights the potential of combining language models with formal methods to enable more accessible and scalable probabilistic planning in robotics.

*Index Terms*—Planning under Uncertainty, AI-Based Methods, Human-Robot Collaboration

This is the accepted version of the article published in IEEE Robotics and Automation Letters (2025). The final published version will be available via IEEE Xplore.

## I. INTRODUCTION

**T**HE most advanced frontier of robotics is the creation of machines that operate in uncontrolled environments and react to unanticipated conditions within acceptable safety margins. Human-robot interaction falls into this area. Humans can be unpredictable, either because they are not entirely focussed on the task or because they exercise their free will to change their behaviour in ways they consider more convenient. Fortunately, there are situations in which, even though humans can decide freely, their decisions are heavily influenced by

[1] Enrico Saccon, Edoardo Lamon, Luigi Palopoli and Marco Roveri are with Department of Information Engineering and Computer Science, University of Trento, Italy. `name.surname@unitn.it`
[2] Davide De Martini and Matteo Saveriano are with the Department of Industrial Engineering, University of Trento, Italy. `name.surname@unitn.it`
\* These authors contributed equally.
† Corresponding author: `enrico.saccon@unitn.it`.

social conventions. For example, in a collaborative construction task involving two humans, the set of available actions is relatively limited, and in most cases, one action is clearly preferable to the others. Structured behavioural patterns such as this are frequently modelled using probabilistic techniques.

For a robot to operate proficiently in such scenarios, it is essential to adopt a planning framework that defines policies, i.e., closed-loop schemes that associate actions to the perceived state of the system. The definition of the final goal can be complex to specify. For instance, in collaborative robotics we may wish to achieve the objective in a reduced time, leaving *at the same time* the human co-worker with the freedom to choose their favourite course of actions. Our objective in this paper is to derive optimal policies of this kind from an informal specification in Natural Language (NL).

**Related Work.** Markov Decision Processes (MDPs) are commonly used for modeling probabilistic decision making with efficient algorithms for policy generation under full and partial observability [1]. In robotics, MDPs have been applied in various areas such as decision support systems [2], motion planning [3], optimal control [4], planning under uncertainty [5], and coordinated multi-agent behaviour [6], including human adaptation [7] and preference inference [8]. As a result, MDPs are widely used in human-robot interaction [9]. A key challenge in deploying MDPs is creating scalable models for real-time policy generation. This is difficult because human behaviour and dynamic environments are often informally described through language, images, or video, while MDPs require symbolic representations for classical computation.

Reinforcement Learning (RL) [10] offers an alternative, as it learns policies through exploration without requiring prior domain specification. This flexibility has made RL popular in robotics [11]. RL often faces sample inefficiency, reward design challenges, and high computational costs. Incorporating domain knowledge can help, though it's rarely easily obtained from informal sources. Large Language Models (LLMs) offer a promising solution. Initially developed for NL understanding [12], LLMs are now used for planning [13] and robotics applications with uncertainty [14], showing excellent results, especially with few-shot learning, which uses examples to guide task performance [15]. Several authors propose LLMs as standalone planners [16], [17], but they lack predictable, explainable behaviour, especially for complex environments where failure recovery is vital. To enhance reliability, researchers use LLMs to generate explainable representations like PDDL [18], a formal language for planning problems. Extensions such as PPDDL [19] and RDDL [20] enable

probabilistic reasoning with MDP solvers. However, these languages have limited expressiveness and cannot fully capture the complexity of real-world robotic environments [21].

Ontologies and KBs effectively capture human knowledge and common sense. Their automated construction uses tools such as OLAM [22] and interactive learning [23], although these rely on existing plans, lack generalisation, and do not model probabilistic aspects. LLMs offer a scalable solution, synthesising KBs for various scenarios [24].

Logic-based knowledge representations have been used in robotics for planning and automated reasoning [25]. Probabilistic Logic Programming (PLP) [26] adds uncertainty to these representations. For example, distributional clauses [27] facilitate the creation of a probabilistic Monte Carlo planner [28]. However, logic knowledge bases require a complete domain specification. Similarly, the inference mechanism of ProbLog addresses MDPs with value iteration [29], but also requires a fully specified domain. In [30], LLMs are combined with the explainability and reasoning capabilities of logic programming, but it only supports deterministic planning, lacking the probabilistic aspect of human–robot interaction.

**Paper Contribution**. The objective of this paper is to develop an open-source framework[1] for generating probabilistic policies from an informal narrative that describes a robotic scenario. In this work, a scenario refers to a specific scene (e.g., an intersection involving humans and autonomous agents) that we describe in NL and utilise as part of the input. The high-level textual narratives that describe the process or scenario and its desired behaviours are provided by domain experts interacting with the front end of the system, unaware of the underlying system functionality. The system constructs, on the back end, a logical knowledge-base (KB) that encodes entities and actions as predicates. Such KBs are human-readable by system designers and support inference of new facts using logical connectives. The role of system designers is to verify the correctness of the generated KB and tailor the LLM output to the specific problem through methods such as few-shot prompting or fine-tuning. Inspired by PLP, our method links probabilities to predicates to automatically construct an MDP, which a state-of-the-art solver processes to generate a policy. Our framework is designed to augment, not replace, human expertise: domain experts provide high-level conceptual narratives, while system designers ensure formal correctness. The system's interpretable knowledge base enhances productivity, transparency, and trust by allowing experts to inspect and verify its decisions. We demonstrate the efficacy of the framework in three use cases, chosen as representative paradigms of a broader class of applications in which robots operate in uncertain environments. The results show that: i) a KB can be efficiently derived from NL text, and ii) the resulting MDP yields effective policies suitable for industrial applications.

## II. BACKGROUND

**Planning.** We take inspiration from the STRIPS classical planning problem [31], which describes an action $a$ by its preconditions $pre(a)$, which must be satisfied in the current state to be executable, and its effects $eff(a)$, which encode the result of its execution. Taking the lines of PPDDL [19], we consider an action to have more than one effect $eff(a)_i$ associated with a probability $P_i \in [0,1]$ of being applied, such that $\sum_i P_i = 1$. Assuming a finite set of objects, the predicates and the set of actions induce an MDP [32].

**MDPs and Policies.** An MDP is a mathematical framework for modelling sequential decision-making problems in environments characterised by stochasticity. It can be described with a 4-tuple $\langle \mathcal{S}, s_0, \mathcal{A}, \mathcal{P}, \mathcal{R} \rangle$, where: $\mathcal{S} = \{s_1, \ldots, s_N\}$ is a finite set of states representing descriptions of the system; $s_0 \in \mathcal{S}$ denotes the initial state: $\mathcal{A} = \{a_1, \ldots, a_M\}$ is a finite set of actions that cause transitions between states; $\mathcal{P}_a(s, s')$ (transition probability relation) describes the probability of moving from state $s$ to state $s'$ given action $a$; $\mathcal{R}$ is a reward function $\mathcal{R}(s, a)$ that assigns a value for executing action $a$ in state $s$. Given an MDP and a probabilistic temporal property, a policy can be synthesised, which is a mapping of states to actions such that the trajectories (sequences of states satisfying the transition probability relation) satisfy the property [33].

**Prolog.** Prolog is a logic programming language, applied in fields like computational linguistics [34] and AI and robotics [35]. It uses predicates and logic-based rules for symbolic reasoning on simpler facts stored in a KB.

## III. PROPOSED FRAMEWORK

Fig. 1 visually represents the different blocks comprising our framework, all of which are described in detail below.

**Knowledge-base.** In our framework, we use Prolog[2] to formally specify a KB that stores: a) The initial state $s_0$ of the MDP, defined as a list of predicates. b) Grounded predicates that characterise the scenario presented in the input query, i.e., immutable characteristics such as different types of agents, e.g., a robotic arm or a wheeled robot. c) The actions that can be executed to transition from one state to another. d) The reward function to score the transitions. e) Auxiliary functions to write the MDP in the PRISM formalism [33]. The Prolog predicates for reward functions are divided into two categories: necessary and sufficient. The former must be satisfied for the transition to be valid; if even one of them fails, the reward is set to a fixed value to discourage the agent from selecting actions leading to undesired states. If all necessary conditions are satisfied, the final reward is computed by adding the contributions of the sufficient conditions that match the transition. Thus, the reward function is used both to prevent invalid transitions and to promote desired behaviours. The reward functions are pure Prolog code that the LLM generates.

**KB Generation.** This step is performed using an LLM via few-shot prompting (Sec. II). The LLM takes as input both the NL description of the domain elements (Fig. 3), and a set of curated examples, containing some general examples (Fig. 4) and others that are use-case-specific (Fig. 5 to 7). The former are used to give a high-level description of how the LLM has to generate the KB, including both how the desired output should be structured and some focal rules that will be used to correctly
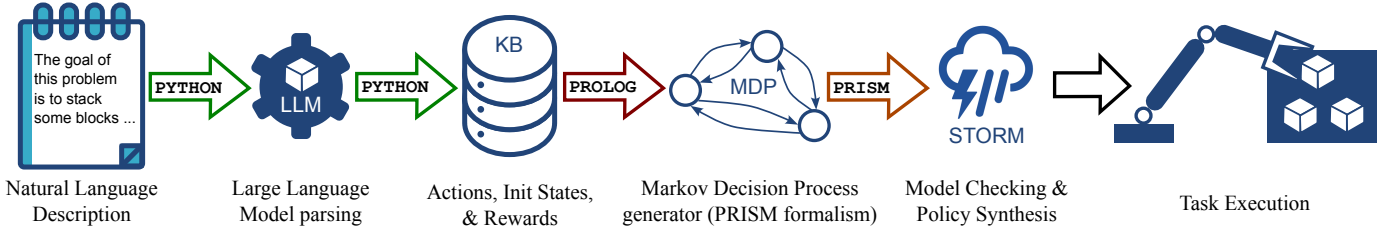
Fig. 1. General diagram of the framework. Given an NL description of the query, it generates a Prolog KB through few-shot prompting with an LLM. The KB is used to extract a PRISM MDP, which Storm uses to synthesize an optima policy. The policy can then be used to execute actions.

```
init_variable_string("p1 : int init 0;\n
    p2 : int init 0;\n p3 : int init 0;\n").
```

Fig. 2. An example of an auxiliary function to generate the PRISM variables: there are three pillars, P1, P2, and P3, which are integers and can have a value between 0 and 4, plus the reward value which is defined as an integer.

```
The use-case consists of an industrial Autonomous
Guided Vehicle (AGV) that has to carry out a task.
The task is completed when the AGV has crossed the
entire factory without emergency stops.

   ... full text in the supplementary material

The goal would be to cross the factory in the least
amount of time possible.
```

Fig. 3. Example of a query passed as input to the LLM to generate the MDP.

```
Q:
 role: 'user'
 content: |
  How should I divide the code?
A:
 role: 'assistant'
 content: |
  When generating an answer, remember to use the
  correct Markdown tags, which are:
  – `kb` for the initial state;
  – `action` for the actions;
  – `reward` for the rewards.
   ... full text in the supplementary material
```

Fig. 4. An example of how the output from the LLM should be formatted.

```
role: 'assistant'
content: |
 Sure. Given the description, we can define the
 initial state a composition of the following
 variables: Section: 1, EStop: false, Delay: 0.
 The initial state can be represented as:
 ```kb
 % Initial state
 initial_state([section(0), estop(false),
 delay(0)]). ```
   ... full text in the supplementary material
```

Fig. 5. An answer used as example for the LLM, containing how to generate the KB from the query shown in Fig. 3.

```
role: 'assistant'
content: |
 Sure. The actions that we need to perform are two:
 one for waiting and one for proceeding. ...
   ...
 The resulting actions will be:
 ```actions
 % Action: wait
 % Preconditions: AGV in section (1 to 4), no estop
 % Effects: delay plus 20, change section, no estop
 action(wait(Section, Delay), ```
   ... full text in the supplementary material
```

Fig. 6. An answer used as example for the LLM, containing how to generate the actions. The query contains both the KB (Fig. 5) and the NL text (Fig. 3).

```
role: 'assistant'
content: |
 Sure. The rewards for the AGV use-case can be
 defined based on the provided description. From the
 description, there are two type of rewards:
 – the necessary one is that the AGV will not crash
 with the workers, and
 – the sufficient ones are that the delay is
 minimized and that the task is completed.
 % Sufficient condition: complete task
 suff_condition(complete_task, _State, [Section, _,
 _], _Action, Reward) :–
    Section = 5, Reward is 100. ```
     ... full text in the supplementary material
```

Fig. 7. An answer used as example for the LLM, containing how to generate the rewards functions from the NL text. The query for this answer contains the previously generated KB, actions and the NL description shown in Fig. 3.

kb (Fig. 5), actions (Fig. 6), and rewards (Fig. 7).

The general examples describe technical aspects, for example, how an action should be structured, as in using the predicate action with arguments the name with the action arguments, and the lists of preconditions and effects, how the initial state should be formatted, i.e., as a list contained in the predicate initial_state, and the labels for Storm, since we want to either maximize the probability of success of a label or minimize the reward value.

The use-case-specific examples define how the NL should be converted into a KB for that specific task. To provide the LLM with this information, we break down the query and examples into different parts, each corresponding to the components of an MDP: the actions, the reward, and the states. The examples explain both how the KB shall be formatted (see Fig. 5 to 7), and also include counterexamples that illustrate wrong modelling choices. In Fig. 5, we see that the LLM shall generate the initial state in the PRISM formalism, as well as the other predicates, e.g., get_state and get_printable_state that are problem dependent and hence must change depending on the input query.

The LLM generates all the components of the KB: initial state, grounded predicates, actions, reward functions, and auxiliary functions to generate the PRISM file. These parts of the KB are produced in multiple steps and, to ensure consistency in the terminology and structure of the predicates, the framework feeds back the previously generated information into the prompt at each subsequent step. We adopted this incremental approach rather than generating the entire KB in a single pass, as preliminary experiments indicated that it would provide more accurate and coherent outputs from the LLM. For example, we begin with the query presented in Fig. 3, which is then passed to the LLM along with a series of examples. The LLM returns a well-formed output containing

parse the NL description in order to generate a correct and coherent KB. The Python parser requires the LLM output to contain Markdown-like tags that allow for easily parsing the output and capturing the required information. The tags are

the KB within the tags `kb`, which is then parsed and fed-back to the LLM with the same original query, and we ask the LLM to generate the actions. As in the previous step, we extract the relevant part for the actions from the output and then ask the LLM to generate the final part for the rewards, feeding both the KB and the actions, plus the original query.

Importantly, the reward functions are *automatically generated by the LLM*. Drawing on its internal knowledge and the few-shot examples, the model not only infers the appropriate reward structure for the given scenario but also distinguishes between *necessary* and *sufficient* conditions. As mentioned before, the formers correspond to constraints not to be violated, whereas the latter correspond to desirable but non-mandatory objectives. This classification is entirely handled by the model during the reward generation step, based on the patterns and instructions provided in the examples and its internal logic.

Finally, a human operator is in charge of checking the quality of the produced KB and eventually fixing the mistakes.
**MDP Generator.** This module, implemented in Prolog, generates the MDP from the KB and the specific query with the following four steps. *Step 1: Graph generation.* The MDP generator module computes a graph $G = \langle S, E \rangle$ starting from an initial state corresponding to the initial condition. The set of nodes $S$ comprises all possible states of the MDP, and the edges $E$ correspond to the transitions: each transition (edge) is associated with an action $a$, with the probabilistic effects $E_i \in eff(a)$ applied during the transition, and with the corresponding probability $P_i$. The algorithm recursively checks if an action $a$ from the action space $\mathcal{A}$ can be executed. If, in the current state $s$, the action's preconditions $pre(a)$ are satisfied, then the probabilistic effects are applied, enabling the algorithm to generate new states. This is repeated by recursively considering each set $E_i$ in the probabilistic effects $eff(a)$ and by checking that the delete effects within $E_i$ have grounded arguments. If they have, the effects are grounded and applied, generating one or more new states. It is important to note that a single set of effects $E_i$ can transition to multiple states as the values of the lifted predicates are recursively and exhaustively assigned. For example, let $s_c$ be the current state in which it is possible to apply the action $a$=`move_block` with the set of effects $E_i = add(block(B))$, and let the KB contain the predicates `block(b1)`, `block(b2)`. The action $a$ from state $s_c$ leads to two different states $s_{n1}, s_{n2}$, one containing `block(b1)` and the other containing `block(b2)`. This is done by exhaustively grounding all the lifted predicates in the lists of effects of the action. If any of the delete effects of $E_i$ cannot be grounded in the current state, this means that $E_i$ cannot be applied in the current state (we expect them to hold in the current state; otherwise, there is no reason to delete them). However, unlike preconditions that prevent an action from being applied, here we add a self-transition $edge_{\langle a, E_i \rangle}(s, s, P)$ to $G$. This interprets the case as a no-op (do nothing). Upon the generation of a new state $s'$, the algorithm evaluates whether it is already included within $S$ of $G$. If included, the algorithm merely adds a new transition from $s$ to $s'$. Conversely, if the state is absent, the new state $s'$ is incorporated into $S$ prior to the addition of the transition. The function is then invoked

---

**Algorithm 1:** The MDP generation functions

```
Function generate_mdp(s):
    while a ← select_action(a) do
        check_preconditions(a);
        for (Prob, E_i) ∈ eff(a) do
            apply_effects(s, s', a, E_i, Prob);

Function apply_effects(s, s', a, E_i, Prob):
    for del(Predicate) in E_i do
        if ground(Predicate) then
            if s' is ¬Visited then
                add_state(s');
                generate_mdp(s');
            add_edge(edge(s, s', a, E_i, Prob));
        else
            add_edge(edge(s, s, a, E_i, Prob));
```

---

**Algorithm 2:** The transition probabilities refinement function

```
Function refine_probabilities(G = (S, E)):
    Initialize changed_edges ← [];
    foreach e in E do
        if e ∉ changed_edges then
            ℰ ← findall(edge(e.s, e.s', e.a, e.eff, _));
            new_prob ← e.prob/|ℰ|;
            change_prob(new_prob, ℰ);
            changed_edges.append(ℰ);
```

---

recursively on it ensuring a complete reachability analysis, which, therefore, ends when every possible combination of actions, states, and grounding possibilities has been examined. *Step 2: Generation of the Transition Probabilities.* After $G$ has been generated, the next step is to distribute the probabilities for its transitions. The probabilities $P_1, \ldots, P_n$ associated in KB with an action's set of effects, $eff(a) = \{E_1, \ldots, E_n\}$, do not depend on the number of states generated by applying that action's effects. This means that each probability $P_i$ must be distributed over the possible transitions created by the action. To perform this task, the framework does a complete traversal of $G$ and, for each edge, if 1) it has the same starting state, 2) it is associated with the same action $a$, and 3) it originates from the same set of effects $E_i$, but 4) leads to different new states, then its probability $P_i$ is distributed across the created transitions. For an illustration, see Fig. 8 and Algorithm 2. This step generates the graph $G'$, which is associated with $G$ and encodes the transition probabilities. For instance, consider the following toy example: the current state has predicate `position(1, 0)`, `position(2,0)`, meaning that there are two pillars that have no blocks. If we consider the action `bi`, meaning that we propose a block of type base and one of type intermediate on the tray, with probabilities 0.75 and 0.25 of being taken by the user, respectively, then the effect `0.75:[del(position(Pillar, 0)), add(position(Pillar, 1))]` corresponding to the user choosing the block of type base, has two possible outcomes: `Pillar={1, 2}`, each with probability $0.75/2$, and each leading to a different state.

*Step 3: Generation of the rewards.* To assign reward

values to the transitions of the MDP graph $G'$, we distinguish between necessary and sufficient rewards. Necessary rewards correspond to strict constraints that must always hold. If even one necessary condition is violated, the transition is immediately assigned a fixed large negative penalty, and no further conditions are evaluated. This ensures that behaviours violating mandatory requirements are strongly discouraged.

If all necessary conditions are satisfied, the system then evaluates the sufficient rewards. These represent desirable but non-mandatory properties and contribute non-negative values (i.e., positive or zero) that accumulate when conditions are met. The more sufficient conditions are satisfied, the higher the reward of the transition. This formulation allows us to separate hard constraints from soft preferences, enabling policies that are both correct and optimised for user-specified objectives.

For example, consider the *structure building* scenario with the policy "Build the pillars by layers, while maximizing the user's freedom". The requirement that no two pillars differ in height by more than one block constitutes a *necessary reward*: any violation results in a severe negative penalty. Conversely, providing the user with blocks that differ from one another represents a *sufficient reward*: its fulfillment yields a positive bonus, but its absence does not incur any penalty.

*Step 4: PRISM Script Compilation.* The compilation of the graph $G'$ into a PRISM script is performed in three steps. First, the PRISM header is created using the LLM-generated strings as auxiliary functions for variable initialisation (see Paragraph *Knowledge-base* in Section III). Second, the transition model is built by traversing the graph (as an MDP) and generating the transitions, formatted using the LLM-generated predicates. Finally, it reports the labels generated by the LLM from the the NL input, and the complete MDP is saved in a file for processing by Storm. As mentioned earlier, the LLM system generates two labels for the MDP: one to maximize the probability of success (doneP) and another to minimize the reward (doneR). This distinction exists because the optimal final state may differ depending on the objective. For example, in the blocks-world use-case, both labels coincide since the task requires all three pillars to be built. In contrast, the AGV use-case has different labels: doneP requires reaching the final section without emergency stops, while doneR focuses on reaching it as quickly as possible, regardless of stops. When compiling Storm, it is up to the user to choose which label is more appropriate, depending on whether they want to prioritize the successful execution of the task (doneP), or whether they want to optimize a reward value doneR. The two labels are Storm functionalities, which will then be used in the model checker to find the optimal solution for the provided label.

**Policy synthesis and execution.** A policy for a given property is synthesized using the Storm solver [36]. Storm provides a Python wrapper for easy integration. To obtain a policy from Storm, specify an optimization goal as a property. The common property is Pmax=? [ F "doneP" ], which asks Storm to find a policy maximizing the expected probability of reaching states where "doneP" holds. For reward optimization, minimizing the reward is achieved by setting the property to: Rmin=? [ F "doneR" ]. The synthesized policy is stored as a *state-action table*, which can be implemented in
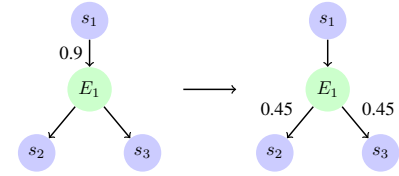


Fig. 8. Example of how to distribute the probability: assume that $s_i$ are the states, $E_1$ is an effect of the action with an associated probability of 0.9. Assuming a uniform distribution, the states $s_2$ and $s_3$, generated by applying the same lifted effects $E_1$, have the same probability of being reached; hence, the probability 0.9 is split between the two, as shown on the right.
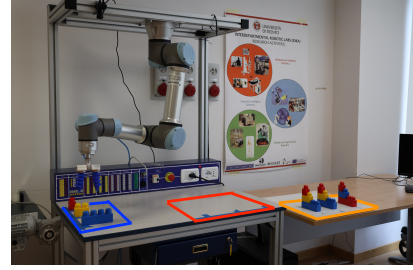


Fig. 9. Experimental setup for the blocks-world scenario.

any robotic runtime environment, such as the adopted ROS2.

## IV. USE-CASES AND EXPERIMENTS

In this section, we present three use-cases on which we tested the system. For each, we first outline the scenario, then provide a detailed explanation of the implementation and a description of the devised experiments. We carried out five tests for each use-case to validate the framework's ability to generate a correct policy and to highlight the generalisation capabilities of the domain generation step.

### A. Structure Building

Inspired by the blocks-world scenario [37] widely used in task planning, we devised a proof-of-concept scenario in which a human agent and a robotic manipulator cooperate to construct a structure. The use-case was tested both in simulation and in a real execution environment.

**Problem description.** A human operator must build $N$ towers using different blocks (orange rectangle in Fig. 9). A robotic manipulator places the blocks on a tray (red rectangle) at predefined positions, all reachable but at varying distances from the human's hands, hence with increasing probability of choosing a block based on proximity. The operator selects one block from the tray, discarding the others. The robot can draw any block from an infinite storage area (blue rectangle).

**KB.** Fig. 10 shows an excerpt of the KB for the scenario. Fig. 10 illustrates the initial state with no blocks on pillars. The right side displays an action placing two top blocks (t) and one intermediate block (i) on the tray, with the top blocks nearer to the human. An action includes a name with arguments, preconditions, and probabilistic effects. The tti action needs at least a pillar of height 2 or 1.

**Reward value.** We want to encourage freedom of choice, hence the reward values are chosen so that the system encourages diversity of the blocks in the tray (as far each block corresponds to an admissible choice).

```
init_state([          action(tti(P1, P2),                              % action_name(args)
    pil(1,0),            [],                                            % AND Preconditions
    pil(2,0),            [pil(P1, 2), pil(P2, 1)],                      % OR preconditions
    pil(3,0)]).          [],                                            % Predicates to verify
                         [0.9:[del(pil(P1,2)),add(pil(P1,3))], 0.1:[del(pil(P2,1)),add(pil(P2,2))]]).  % Effects
```

Fig. 10. Initial state (left): no block has been placed. An example of actions (right) with the list of preconditions and effects.

```
init_state([          action(proceed(S),                               % action_name(args)
  section(1),           [section(S), estop(0)],                        % AND preconditions
  estop(0),             [],                                            % OR preconditions
  delay(0)]).           [Pno is 1-S/10, Pyes is S/10, NS is S+1, S<5]  % Predicates to verify
                        [Pyes:[del(estop(0)),add(estop(1))], Pno:[del(section(S)),add(section(NS))]]).  % Effects
```

Fig. 11. Initial state (left): AGV is at starting state without any emergency stop and delay. An example of actions (right) in which the predicates to verify plays a crucial role to also compute the probabilities of the different effects.

**Experiments.** We performed five tests for this use-case: 1) There are a total of 3 types of blocks and the robotic arm places 3 blocks on the tray. The operator has to create 3 towers picking blocks from the provided ones. This case is the one we also used in the examples for the LLM. 2) The probabilities of picking the blocks on the tray are changed with values not present in the examples. 3) The number of pillars increases to five. 4) The number of block types increases to four, and the number of pillars reduces to two. 5) We considered two pillars, three types of blocks, and a new action that can be executed only at the end to place an architrave on top of the pillars. Tests are conducted to highlight the ability of the framework to generalise in the scenario and extract consistent policies.

**Real-world experiment.** Once the state-action table is extracted, it can be used in any ad-hoc system. We developed a ROS2 workspace to conduct real-world experiments.

### B. Traffic management of AGVs in an industrial scenario.

**Problem description.** An Autonomous Guided Vehicle (AGV) must navigate a factory floor, each route is divided into sections and in each section there are workers who might disrupt its path. More workers increase the chance of emergency stops to prevent collisions. The objective is to traverse all sections without stops. The AGV has two actions: `wait`, which slows down to avoid crashes but delays the task, or `proceed`, which moves ahead without delay, risking an emergency stop. If the AGV stops, it ends negatively, as indefinite waiting causes delays that reduce factory efficiency.

**KB.** The initial state includes three predicates: tracking the section, reporting an emergency stop, and expressing delay (left of Fig. 11). Two actions model this use-case by evaluating predicates to determine effect probabilities, incrementing the section, and increasing the delay.

**Reward value.** The goal is to reduce the likelihood of entering emergency stop mode while ensuring factory efficiency. Mandatory properties in the NL text include human safety, while optional ones like minimizing time spent in a section help compute reward values. The reward function penalizes AGV emergency stops, and rewards for transitions without stops are inversely related to accumulated delay.

**Experiments.** For this use-case, we rely on simulation to validate the proposed approach since we have no access to the necessary hardware and resources. We considered the following scenarios. 1) There are a total of 5 sections that the robot has travel across. The two only actions that it can do are: `wait` and `proceed`, as aforementioned. 2) It extends the first test by increasing the section number to 10 sections in total. 3) It decreases the delay caused by the wait action by 15. 4) It combines the second and third tests together, and also changes the `wait` action. The robot now stops and it can only advance to the next section if `proceed` is called. 5) A third previously unseen action, namely `speed-up`, has to be generate by the LLM. This action increases the velocity of the AGV setting the delay to 0, but it also increases the probability of an emergency stop. `proceed` now allows for moving to the next section with a reduced probability of collision.

### C. Gripper domain

In the gripper domain [38], a robot with two grippers moves balls between rooms. This domain tests our system's generalisation, serving as a balance between the first two use-cases by considering different settings with probabilistic actions, multi-stage transitions, resource constraints, and value prioritization.

1) A robot with two grippers moves 3 balls from Room A to Room B. The move is certain, but picking succeeds at 0.9, and dropping at 0.95. 2) Two labelled balls are involved: one is in Room A, the other is in Room B, requiring return to Room A before delivery. Movement succeeds at 0.98, picking at 0.88, and dropping at 0.94. 3) In a three-room setup, the robot moves balls from Room A through a Corridor to Room B. Initial setup includes 2 balls in Room A and 2 in the Corridor. Move success is 0.97 from Room A to Corridor, and 0.95 from Corridor to Room B. Picking succeeds at 0.90, dropping at 0.93. 4) Energy is limited: 4 fragile balls need delivery with 10 energy units. Each move, pick, or drop uses 1 unit of energy. Pick and drop success rates are 0.92 and 0.97. Task fails if energy depletes before completion. 5) Delivery prioritization has 4 balls in Room A, with 2 as high-value. Movement is certain. High-value balls pick at 0.85, drop at 0.90; standard balls pick at 0.90, drop at 0.96. Carrying both high-value balls reduces drop success by 0.03.

### D. Implementation details and experimental setup.

All experiments were executed on a desktop PC running Ubuntu 22.04 with an AMD Ryzen 7 7700X CPU and 64GB of DDR5 memory. We used Swipl version 9.2.9, and stormpy version 1.9.0. For the generation of the KB, we used GPT-4o and GPT-5-mini as the LLM. To reducestochasticity, we set the temperature to 0 and fixed the seed (to 42) beforehand. Although this choice does not produce repeatable behaviours, it significantly alleviates the unpredictability of the LLM.

TABLE I

RESULTS OF THE EXPERIMENTS ON THE PROPOSED USE-CASES AVERAGED OVER 100 TRIALS. IN KB GENERATION, WE REPORT ✓ IF THE LLM OUTPUT WAS CORRECT, OR WE USE $X(N, M)$, WHERE $N$ IS THE NUMBER OF LOGICAL ERRORS AND $M$ THE NUMBER OF CORRECTIONS. MDP EXTRACTION SHOWS THE TIME TO GENERATE THE MDP, REFINE PROBABILITIES (SEC. III), AND WRITE TO FILE. POLICY EXTRACTION REPORTS THE AVERAGE TIMES (OVER 10000 RUNS) TO EXTRACT THE OPTIMAL POLICY FOR PROBABILITY MAXIMISATION (TOP) AND REWARD MAXIMISATION (BOTTOM).

| | STRUCTURE BUILDING | | | | | AGV | | | | | GRIPPERS | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 1 | 2 | 3 | 4 | 5 | 1 | 2 | 3 | 4 | 5 |
| # STATES | 64 | 64 | 1024 | 125 | 17 | 35 | 120 | 48 | 78 | 194 | 14 | 17 | 104 | 2027 | 263 |
| # ACTIONS | 27 | 27 | 27 | 16 | 28 | 2 | 2 | 2 | 2 | 3 | 59 | 5 | 12 | 18 | 19 |
| KB GPT-4O | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | X(2,2) | X(1,1) |
| KB GPT-5-MINI | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| MDP EXTRACTION | 0.073 | 0.074 | 28.928 | 0.055 | 0.006 | 0.001 | 0.005 | 0.001 | 0.003 | 0.014 | 0.008 | 0.002 | 0.021 | 3.184 | 0.080 |
| | 0.788 | 0.792 | 780.792 | 0.270 | 0.053 | 0.000 | 0.001 | 0.000 | 0.001 | 0.004 | 0.009 | 0.001 | 0.004 | 0.999 | 0.070 |
| | 0.404 | 0.407 | 374.034 | 0.140 | 0.016 | 0.000 | 0.001 | 0.000 | 0.001 | 0.002 | 0.005 | 0.000 | 0.004 | 0.298 | 0.028 |
| POLICY EXTRACTION | 0.175 | 0.180 | 11.701 | 0.148 | 0.036 | 0.016 | 0.027 | 0.017 | 0.021 | 0.046 | 0.042 | 0.016 | 0.044 | 0.015 | 0.171 |
| | 0.193 | 0.184 | 18.702 | 0.172 | 0.038 | 0.049 | 0.027 | 0.017 | 0.021 | 0.048 | 0.044 | 0.017 | 0.046 | 2.584 | 0.185 |

## V. EXPERIMENTAL RESULTS

We analyse the results obtained in the three domains described in Sec. IV. Table I presents results for the various steps of our framework, while Table II shows policy execution.

### A. KB Generation

The results in Table I show promising generalisation abilities from the LLM, demonstrated across 3 different domains. **Structure Building** The LLM generated the correct KB in all the five cases, across varying numbers of pillars, blocks and sections. The examples used for the few-shot prompt were similar to the tested use-case. Notably, it correctly formulated the action to place the architrave on top of the pillars, an action that can be executed only at the end in Case 5.
**AGV** The LLM was able to generate consistent KBs from the examples. Only one error occurred in scenario 4, in which the LLM added an additional non instantiated predicate to the *verify predicate* list. This is considered a minor error, since the Swipl interpreter immediately identified the error, and the correction consisted in removing such predicate from the list.

**Grippers** The LLM generated the correct KB in three out of five cases, without any new example. The errors committed by GPT-4o were not conceptual errors, but rather simple syntax mistakes that could be easily detected and corrected by an expert. For instance, in both Case 4 and 5, the error consisted in using `ball3` instead of `ball3_position`. Such errors could also be automatically identified and corrected by integrating syntax and and consistency-checking tools [39].

Finally, given the nature of the observed errors, we also conducted additional tests using GPT-5-mini. This model exhibited improved generalisation capabilities and achieved a perfect score in the generation of the KB, confirming the robustness and scalability of our approach.

### B. MDP and Policy Extraction

The results of the generation of the MDP and for the extraction of the policy are shown in Table I. For the considered examples, the generation of KB and MDP took a negligible time, which allowed us to make multiple queries in a small time. This does not mean that the framework could be used for online generation in industrial scale applications.

Overall, performance varies significantly with the complexity of the MDP, i.e., the number of states and actions.

TABLE II

SIMULATION RESULTS. FOR THE STRUCTURE BUILDING USE-CASE, WE SHOW THE TOTAL NUMBER OF ACTIONS AND THE NUMBER OF ACTIONS THAT PROPOSE 2 OR 3 EQUAL BLOCKS. FOR THE AGV USE-CASE, WE COMPARE THE OPTIMAL POLICY WITH A FAULTY ONE IN WHICH THERE IS A PROBABILITY $p_f = 0.4$ OF CHOOSING THE WRONG ACTION.

| | | STRUCTURE BUILDING | | | AGV | |
|---|---|---|---|---|---|---|
| | | | # ACTIONS | | SUCCESS RATIO | |
| | TEST | TOT | 2 | 3 | OPT | FAULTY |
| OPTIMAL POLICY | 1 | 90000 | 67754 | 22246 | 1.0000 | 0.8024 |
| | 2 | 90000 | 67273 | 22727 | 1.0000 | 0.8591 |
| | 3 | 150000 | 116596 | 33404 | 1.0000 | 0.8090 |
| | 4 | 120000 | 20268 | - | 0.0148 | 0.0163 |
| | 5 | 70000 | 43454 | 16546 | 0.0697 | 0.0516 |
| ε-GREEDY POLICY | 1 | 90000 | 20000 | 3257 | 0.0380 | 0.0602 |
| | 2 | 90000 | 0 | 24053 | 0.0003 | 0.0007 |
| | 3 | 150000 | 113349 | 0 | 0.0126 | 0.0232 |
| | 4 | 120000 | 0 | - | 0.0126 | 0.0175 |
| | 5 | 70000 | 0 | 1133 | 0.0166 | 0.0185 |

The framework solves the tests for the AGV use-case–which includes between 35 and 189 states and 2 or 3 actions– consistently faster than the tests for the structure building use-case, which ranges from 17 to 1024 states and 16 or 27 actions. Notably, the time required to refine transition probabilities increases with the number of available actions per state. This factor has a substantial impact, often dominating the total computation time (Table I). Extracted policies respect the criteria written in the input queries both when maximising the probabilities of success or minimising the reward values.

### C. Policy execution

To provide a brief evaluation of policy behaviour, we tested the framework in three representative scenarios. In the structure building task, the probability-maximising policy `doneP` often used repeat block combinations (over 75% duplicates in test 1), while the reward-oriented policy `doneR` created more diverse sets, demonstrating our approach's ability to bias policies toward reliability or diversity. In the AGV use case, the probability-maximising policy consistently achieved goals under ideal conditions (100% success in tests 1–3) and performed well under faults (over 80% success), though its performance decreased in tests 4–5 due to limited action effectiveness. Conversely, the reward-minimising policy struggled even without faults, with faults further reducing success, highlighting its fragility. These experiments show that our framework supports different optimisation objectives and

exposes their trade-offs in robustness and diversity.

**Real-world experiment.** In the real-world experiment, we used a UR5e robotic arm, equipped with a 2-finger soft gripper (Figure 9), and involved a total of four subjects, including two authors and two additional participants, each of whom completed the test twice. In the first round, the user was instructed to build three pillars using blocks from the robotic arm. In the second round, they were informed of the building policy, i.e., building by layer. All participants successfully constructed the structure each time, confirming the correctness and robustness of the extracted policies. A video of the experiment can be found in the multimedia material.

## VI. DISCUSSION AND CONCLUSIONS

In this paper, we presented a complete framework that automates the generation of MDP policies from natural language queries. The proposed approach leverages state-of-art prompt-engineering techniques to construct a Prolog KB, which is subsequently used to derive an MDP and its corresponding policy through Storm. The resulting policy can be conveniently stored in a state–action table, enabling straightforward deployment in real-world scenarios. Overall, this study provides a foundational step toward integrating natural language interfaces with formal decision-making frameworks, paving the way for more accessible and adaptive autonomous systems.

Despite its promising results, the framework has several limitations. The most critical limitation arises from LLMs inherent tendency to generate incorrect or inconsistent outputs. At present, the system does not include an automated mechanism for error detection or correction, and therefore primarily functions as a support tool for domain experts. Future work will focus on introducing a feedback loop for self-correction when the parser fails or after MDP analysis. Additionally, employing LLM ensembles or fine-tuned domain-specific models may help mitigate such inaccuracies. A second limitation concerns the LLM's inability to autonomously infer action probabilities, which currently requires expert intervention. To address this, we are investigating methods to estimate these probabilities from prior task demonstrations, such as video data or raw sensory inputs. Furthermore, while the framework is designed for expert users, its initial learning curve can be steep. We are developing comprehensive usage guidelines, prompt templates, and an intuitive user interface to improve accessibility and usability. Another challenge involves the generation of few-shot examples to guide the LLM in creating effective knowledge bases. Although the model demonstrates strong generalisation capabilities across domains, curating high-quality examples remains essential for robust performance. To this end, we are preparing standardised templates and documentation to facilitate this process. As future work, we also aim to extend the framework to support partially observable MDPs and hidden Markov models, thereby enhancing its ability to model uncertainty and improve interaction with human operators.

## REFERENCES

[1] L. P. Kaelbling, M. L. Littman, and A. R. Cassandra, "Planning and acting in partially observable stochastic domains," *Artificial Intelligence*, vol. 101, no. 1-2, pp. 99–134, 1998.

[2] S. Doltsinis, P. Ferreira, M. M. Mabkhot, and N. Lohse, "A Decision Support System for rapid ramp-up of industry 4.0 enabled production systems," *Computers in Industry*, vol. 116, p. 103190, 2020.

[3] H. Kurniawati, Y. Du, D. Hsu, and W. S. Lee, "Motion planning under uncertainty for robotic tasks with long time horizons," *The International Journal of Robotics Research*, vol. 30, no. 3, pp. 308–323, 2011.

[4] X. Ding, S. L. Smith, C. Belta, and D. Rus, "Optimal Control of Markov Decision Processes With Linear Temporal Logic Constraints," *IEEE Transactions on Automatic Control*, vol. 59, no. 5, pp. 1244–1257, 2014.

[5] N. t. Gopalan, "Planning with abstract markov decision processes," in *ICAPS*, vol. 27, 2017, pp. 480–488.

[6] K. Bogert and P. Doshi, "Interacting partially observable markov decision processes for modeling coordinated multi-agent behavior," *Autonomous Agents and Multi-Agent Systems*, vol. 28, pp. 684–705, 2014.

[7] S. Nikolaidis, S. Nath, A. D. Procaccia, and S. S. Srinivasa, "Game-theoretic modeling of human adaptation in human–robot collaboration," in *ACM/IEEE Int. Conf. on Human-Robot Inter.*, 2017, pp. 323–331.

[8] A. Fern, S. Natarajan, K. Judah, and S. Kambhampati, "Modeling and learning human user strategies in planning systems," in *AAAI*, 2010.

[9] A. Karami and A. Mouaddib, "A Decision Model of Adaptive Interaction Selection for a Robot Companion," in *ECMR*, 2011, pp. 83–88.

[10] L. P. Kaelbling, M. L. Littman, and A. W. Moore, "Reinforcement learning: A survey," *JAIR*, vol. 4, pp. 237–285, 1996.

[11] B. Singh, R. Kumar, and V. P. Singh, "Reinforcement learning in robotic applications: a comprehensive survey," *Artificial Intelligence Review*, vol. 55, no. 2, pp. 945–990, 2022.

[12] L. Yuan, Y. Chen, G. Cui, H. Gao, F. Zou, X. Cheng, H. Ji, Z. Liu, and M. Sun, "Revisiting out-of-distribution robustness in nlp: Benchmarks, analysis, and llms evaluations," in *Advances in Neural Information Processing Systems*, vol. 36, 2023, pp. 58 478—-58 507.

[13] M. Aghzal, E. Plaku, G. J. Stein, and Z. Yao, "A survey on large language models for automated planning," 2025.

[14] R. Mon-Williams and G. t. Li, "Embodied large language models enable robots to complete complex tasks in unpredictable environments," *Nature Machine Intelligence*, vol. 7, no. 4, pp. 592–601, 2025.

[15] T. t. Brown, "Language models are few-shot learners," in *NeurIPS*, vol. 33, 2020, pp. 1877—-1901.

[16] M. Ahn, N. Brohan, N. Brown, Y. C. Hong, C. Cremer, C. Finn, S. Levine, J. Lu, R. Martens, and S. t. Saxena, "Code as policies: Language model-driven robot task planning," in *Conference on Robot Learning*, 2022, pp. 1–21.

[17] N. Brohan, K. Jang, R. Joshi, T. Xiao, E. Irpan, C. Cremer, W. Jang, N. Brown, J. Burgner-Kahrs, and K. t. Daniilidis, "RT-2: Vision-Language-Action Models Learn To Follow Instructions, Reason, and Plan," in *ICRA*, 2023, pp. 3337–3349.

[18] J. Oswald, K. Srinivas, H. Kokel, J. Lee, M. Katz, and S. Sohrabi, "Large language models as planning domain generators," in *ICAPS*, vol. 34, no. 1, 2024, pp. 423–431.

[19] H. L. Younes and M. L. Littman, "Ppddl1. 0: An extension to pddl for expressing planning domains with probabilistic effects," *Techn. Rep. CMU-CS-04-162*, vol. 2, p. 99, 2004.

[20] S. t. Sanner, "Relational Dynamic Influence Diagram Language (RDDL): Language Description," *Unpublished ms. Australian National University*, vol. 32, p. 27, 2010.

[21] D. Rao, Z. Jiang, Y. Wen, and J. Li, "Performance of the rddl planners," in *IEEE Conf. of Online Analysis & Comp. Sci.*, 2016, pp. 115–118.

[22] L. Lamanna, A. Saetti, L. Serafini, A. Gerevini, and P. t. Traverso, "Online learning of action models for pddl planning." in *IJCAI*, 2021, pp. 4112–4118.

[23] D. t. Meli, "Inductive learning of answer set programs for autonomous surgical task planning: Application to a training task for surgeons," *Machine Learning*, vol. 110, no. 7, pp. 1739–1763, 2021.

[24] D. t. Xu, "Large language models for generative information extraction: A survey," *Frontiers of Computer Science*, vol. 18, no. 6, p. 186357, 2024.

[25] D. t. Meli, "Logic programming for deliberative robotic task planning," *Artificial Intelligence Review*, vol. 56, no. 9, pp. 9011–9049, 2023.

[26] F. Riguzzi and T. Swift, "A survey of probabilistic logic programming," *Declarative Logic Programming: Theory, Systems, and Applications*, pp. 185–228, 2018.

[27] D. t. Fierens, "Inference and learning in probabilistic logic programs using weighted boolean formulas," *Theory and Practice of Logic Programming*, vol. 15, no. 3, pp. 358–401, 2015.

[28] D. Nitti, V. Belle, T. De Laet, and L. De Raedt, "Planning in hybrid relational MDPs," *Machine Learning*, vol. 106, pp. 1905–1932, 2017.

[29] T. P. Bueno, D. D. Mauá, L. N. De Barros, and F. G. Cozman, "Markov decision processes specified by probabilistic logic programming: representation and solution," in *Braz. Conf. on Intel. Sys.*, 2016, pp. 337–342.

[30] E. t. Saccon, "When Prolog Meets Generative Models: a New Approach for Managing Knowledge and Planning in Robotic Applications," in *IEEE ICRA*, 2024, pp. 17 065–17 071.

[31] M. Ghallab, D. S. Nau, and P. Traverso, *Automated planning - theory and practice*. Elsevier, 2004.

[32] M. L. Puterman, "Markov decision processes," in *Handbooks in OR and MS*. Elsevier, 1990, ch. 8, pp. 331–434.

[33] V. Forejt, M. Z. Kwiatkowska *et al.*, "Automated verification techniques for probabilistic systems," in *SFM*, vol. 6659, 2011, pp. 53–113.

[34] C. Matthews, *An introduction to natural language processing through Prolog*. Routledge, 2016.

[35] B. G. Batchelor and R. Hack, "Robot vision system programmed in Prolog," in *Machine Vision Applications, Architectures, and Systems Integration*, vol. 2597, 1995, pp. 239–252.

[36] C. Hensel, S. Junges, J.-P. Katoen, T. Quatmann, and M. Volk, "The Probabilistic Model Checker Storm," 2020.

[37] N. Gupta and D. S. Nau, "On the complexity of blocks-world planning," *Artificial Intelligence*, vol. 56, no. 2, pp. 223–254, 1992.

[38] D. V. McDermott, "The 1998 AI planning systems competition," *AI Mag.*, vol. 21, no. 2, pp. 35–55, 2000.

[39] D. D. *et al.*, "Llm-driven knowledge extraction in temporal and description logics," in *EKAW*, vol. 15370, 2024, pp. 190–208.