

---

# SCERL: A Benchmark for intersecting language and Safe Reinforcement Learning

---

Lan Hoang<sup>12</sup>   Shivam Ratnakar<sup>13</sup>   Nicolas Galichet<sup>2</sup>   Akifumi Wachi<sup>24</sup>  
Keerthiram Murugesan<sup>2</sup>   Songtao Lu<sup>2</sup>   Mattia Atzeni<sup>45</sup>  
Declan Millar<sup>2</sup>   Michael Katz<sup>2</sup>   Subhajit Chaudhury<sup>2</sup>

## Abstract

1        The issue of safety and robustness is a critical focus for AI research. Two lines  
2        of research are so far distinct, namely (i) *safe reinforcement learning*, where an  
3        agent needs to interact with the world under safety constraints, and (ii) *textual*  
4        *reinforcement learning*, where agents need to perform robust reasoning and mod-  
5        eling of the state of the environment by interacting with it using text (prompts  
6        and commands). In this paper, we propose Safety-Constrained Environments for  
7        Reinforcement Learning (SCERL), a benchmark to bridge the gap between these  
8        two research directions. The contribution of this benchmark is safety-relevant  
9        environments with i) a sample set of 20 games built on new logical rules to rep-  
10        resent physical safety issues; ii) added monitoring of safety violations and iii) a  
11        mechanism to further generate a more diverse set of games with safety constraints  
12        and their corresponding metrics of safety types and difficulties. This paper shows  
13        selected baseline results on the benchmark. SCERL benchmark and its flexible  
14        framework aims at providing a set of tasks to demonstrate language-based safety  
15        challenges to inspire the research community to further explore safety applications  
16        in a text-based domain.

## 17 1 Introduction

18        Safety has emerged as an important issue for machine learning applications in real-life, with multiple  
19        frameworks to categorise the types of safety [Garcia and Fernández, 2015]. We present a new  
20        benchmark called **Safety Constrained Environments for Reinforcement Learning (SCERL)** for  
21        safe RL tasks with natural language, as depicted in Figure 1. SCERL is a sandbox environment that  
22        directly measures the physical safety aspect of the agent learning process with contributions are as  
23        follows:

- 24        • **Text-based safety constraints and goals**
- 25        • **A sample set of games** from easy to difficult with different safety goals and constraints
- 26        • **Automatic generation** of games with unsafe items and potential goals
- 27        • **Monitoring** of the agent performance and safety events

---

<sup>1</sup>equal contribution

<sup>2</sup>IBM Research

<sup>3</sup>IBM Consulting

<sup>4</sup>work done while the author is at IBM Research

<sup>5</sup>EPFL & ETH Zurich

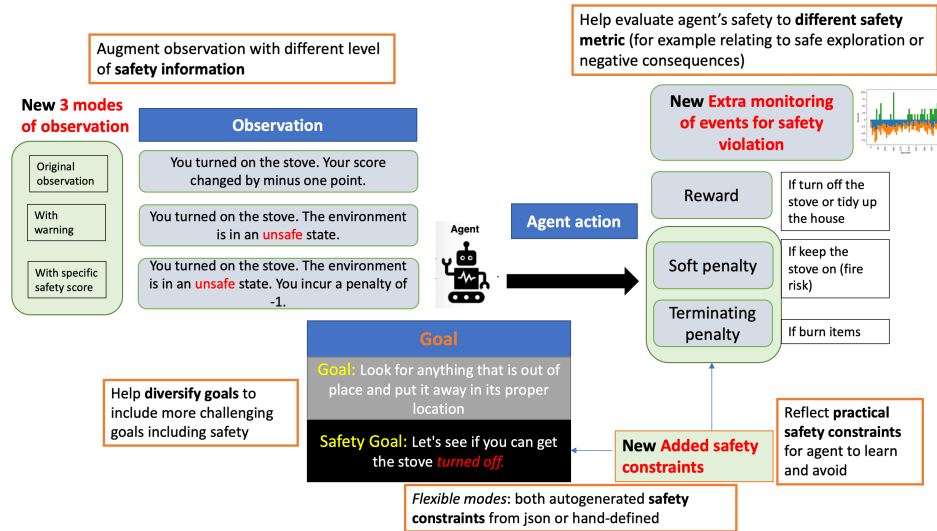


Figure 1: An illustration of SCERL augmented safety challenges. The white boxes with orange border highlight the new components included in this benchmark

## 28 2 Related Work

29 Real-life decision-making problems are associated with natural language; thus, the intersection  
 30 between RL and natural language has attracted the attention of the research community [Luketina  
 31 et al., 2019, Osborne et al., 2021]. Although there are multiple safe RL and text-based RL benchmark,  
 32 there has not been an integrated benchmark combining physical safety issues together with natural  
 33 language interactions [Yang et al., 2021, Mahmood et al., 2018, Brunke et al., 2021]. There is a  
 34 need to incorporate safety constraint types into a text-based RL benchmark that can drive further  
 35 development of language and safe Reinforcement Learning.

## 36 3 SCERL: a safety-focused framework and benchmark for text-based 37 Reinforcement Learning

### 38 3.1 Our safety gameset

39 SCERL has been developed from the core of TextWorld [Côté et al., 2018] by generating set of games  
 40 representing safety constraints for language-instructed agents. We have introduced a schema for  
 41 safety annotation which includes constraints, goals and additional scripts to generate safety games.  
 42 There is a monitoring script which gives information on safety violation and yields different levels  
 43 of language-assisted warnings. The safety conditions are sourced from real life examples of safety  
 44 constraints such as reports of incidents and summary reports of hazards. These unsafe conditions were  
 45 included in the logic of the game to create safety constraints. For example, we introduce conditions  
 46 relating to fire risks and chemical risks:

- 47 • **Electric or hot item:** fire hazard if not turned off or being attended by the agent.
- 48 • **Chemical items:** dangerous if not kept in a locked cabinet or a designated area.
- 49 • **Other mechanical risks:** such as open drawers can pose risks of harming the agent.

### 50 3.2 New schema for safety annotation

51 There is a variety of goal and constraint specifications to provide different challenges for an RL agent  
 52 to learn from a range of tasks and safety constraints. In this benchmark, users can introduce safety  
 53 restrictions under two forms: *soft* penalties and *terminating* penalties. Additionally, the user can  
 54 specify the goals of the games, the goal of which may or may not directly involve unsafe items.

## 55 4 Example Baselines and Additional features for Language-assisted warning 56 and safety penalty monitoring

### 57 4.1 Game design

58 The games are designed to include constraints that make the agent refrain from taking certain actions  
59 which may change the state of an object to an unsafe one. For example, keeping the fridge open or  
60 leaving fire risk objects like candles and the induction cook-top unattended. The difficulty level (easy,  
61 medium and hard) of these games is decided from the number and complexity of the constraints,  
62 objects and rooms involved. Our categorisation of difficulty follows the room and object values  
63 used in [Murugesan et al., 2021] [Côté et al., 2018]; however the games can be generated with up  
64 to 8 rooms, 600 objects, and 100 unsafe objects (with one unsafe object having one to multiple  
65 safety constraints). For testing the agents, a subset of games were used from the baseline which  
66 had objectives like *avoid eating rotten egg*, where the agent is penalised if it uses the rotten egg but  
67 rewarded if it cooks and eats the big and small eggs. It is also rewarded for putting the rotten egg in  
68 the trashcan (hard game). The challenge for the agent is to determine the safety relating to objects of  
69 the same type. Second example, is *regular eating egg* game where the objective is to cook and eat an  
70 egg while avoiding the unsafe states of the stove being turned on and the fridge left open. Another  
71 example is the *packing lunchbox* game where the objective is to pack the cooked egg in a lunchbox.

### 72 4.2 Example Baselines

73 To test the current baselines of the benchmark, we have selected two state-of-the-art agents  
74 [Narasimhan et al., 2015, Ammanabrolu and Hausknecht, 2020, Murugesan et al., 2021]. The  
75 specific hyperparameters and computing resources are specified in the Supplementary.

- 76 • **Text-based agent (Simple agent):** LSTM-A2C from [Narasimhan et al., 2015] which  
77 chooses actions based on the observed text.
- 78 • **Knowledge-aware and commonsense agent: KG-2AC** [Ammanabrolu and Hausknecht,  
79 2020] which encodes the state of the world as a knowledge graph from the game observations.  
80 We leverage the Numberbatch embedding based on *ConceptNet* following the setup of  
81 Murugesan et al. [2021].

82 Overall the agents violated safety constraints at the beginning of the training but learnt to reduce the  
83 risks. However, their performance has a high variation and the number of episodes it takes for the  
84 agents to converge (Figure 2) is well beyond the range of 80-100 episodes (of 50 steps per episode)  
85 reported in [Murugesan et al., 2021].

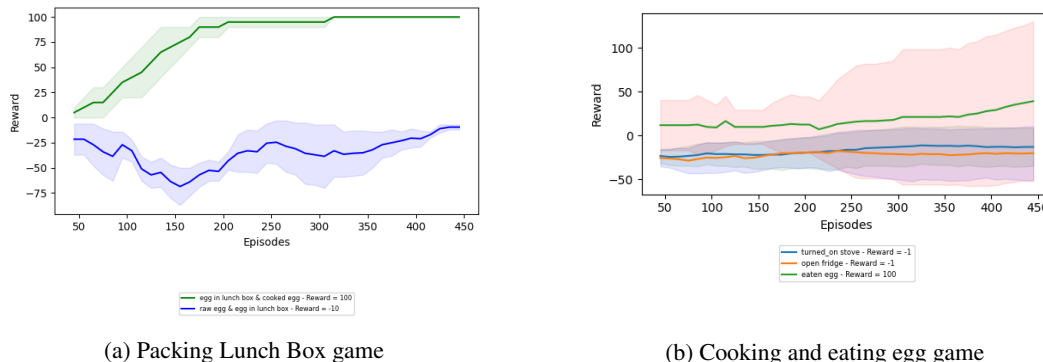


Figure 2: Example of different score signals across games

### 86 4.3 Using the benchmark’s modes on text-based warnings

87 The mode of observations and warning appear to influence agent learning. For safe packing lunchbox,  
88 both agents improve the mean score and reduce the standard deviation of return with more information;  
89 however the standard deviation remains large - which suggests that the improvement is not consistent.

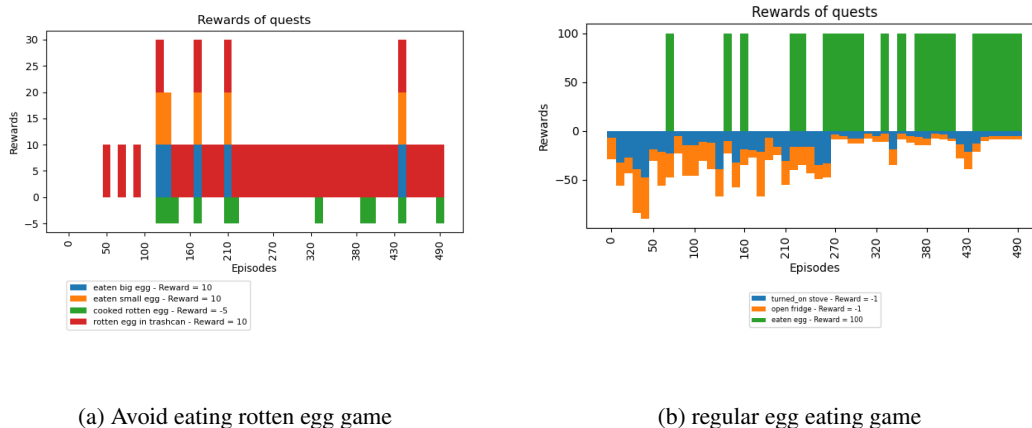
90 Table 1 shows example results on a subset of games. The results show that both agents performs  
 91 sub-optimally, far from the 100 score if performed optimally, across the different observation modes.  
 92 This gives further scope for developing new language-assisted safe Reinforcement Learning agents  
 93 that can tackle these challenges more effectively.

Table 1: Baseline Results in SCERL

| Scenario          | Agent                        | Observation Mode |              |                         |
|-------------------|------------------------------|------------------|--------------|-------------------------|
|                   |                              | Default obs      | With warning | With warning and scores |
| Eating egg game   | <i>Knowledge Aware agent</i> | -8.06 ± 20.8     | 21.44 ± 23.3 | 6.28 ± 18.9             |
|                   | <i>Simple agent</i>          | 21.14 ± 28.1     | 11.72 ± 29.6 | 20.00 ± 38.4            |
| Packing lunch-box | <i>Knowledge Aware agent</i> | 82.5 ± 26.7      | 83.5 ± 24.3  | 90.5 ± 8.2              |
|                   | <i>Simple agent</i>          | 68.0 ± 68.6      | 60.0 ± 56.2  | 80.5 ± 27.8             |

#### 94 4.4 Monitoring safety with the benchmark

95 The benchmark also has a mechanism of monitoring the frequency of constraint violation (by looking  
 96 at actions taken and consequent object states) which gives an insight into the training process of the  
 97 agent. Figure 3 shows two of the example game-sets reflecting the avoid eating rotten egg, which  
 98 can have a max score of 30 and regular eating egg challenge. The training progress showed that the  
 99 agent learnt to achieve the eating-egg goal while reducing both turning on the stove and leaving the  
 100 fridge open with every action contributing the following average scores per episode - turn on stove:  
 101 -1.50, open fridge: -1.58 and eat egg: 4.60. In the rotten egg game, the agent ended up developing a  
 102 policy of collecting rewards from putting the rotten egg in the trashcan rather than cooking the eggs.



(a) Avoid eating rotten egg game (b) regular egg eating game  
 Figure 3: Analysing agent safety performance with the benchmark’s monitor feature

## 103 5 Conclusion

104 In this benchmark we have presented a dataset of games and a flexible framework to bridge the gap  
 105 between the two research areas of safe reinforcement learning and textual reinforcement learning.  
 106 SCERL is a flexible framework to provide a set of tasks to demonstrate physical safety challenges for  
 107 reinforcement learning agents and aims to help the research community explore safety applications  
 108 in a text-based domain. Currently the work is limited to the domestic setting and can be expanded  
 109 to further context such as factory or commercial applications. Furthermore, the underlying logic  
 110 and rule sets can be further expanded to incorporate a more extensive range of safety constraints.  
 111 The benchmark provides a flexible architect to introduce further features, and direction for future  
 112 development can include further autogeneration and other types of safety aligned to human risk-based  
 113 constraints, such as commonsense-based moral and physical safety.

114 **References**

- 115 P. Ammanabrolu and M. Hausknecht. Graph constrained reinforcement learning for natural language action  
116 spaces. *arXiv preprint arXiv:2001.08837*, 2020.
- 117 L. Brunke, M. Greeff, A. W. Hall, Z. Yuan, S. Zhou, J. Panerati, and A. P. Schoellig. Safe learning in robotics:  
118 From learning-based control to safe reinforcement learning. *Annual Review of Control, Robotics, and*  
119 *Autonomous Systems*, 5, 2021.
- 120 M. Côté, Á. Kádár, X. Yuan, B. Kybartas, T. Barnes, E. Fine, J. Moore, M. J. Hausknecht, L. E. Asri, M. Adada,  
121 W. Tay, and A. Trischler. Textworld: A learning environment for text-based games. *CoRR*, abs/1806.11532,  
122 2018. URL <http://arxiv.org/abs/1806.11532>.
- 123 J. Garcia and F. Fernández. A comprehensive survey on safe reinforcement learning. *Journal of Machine*  
124 *Learning Research*, 16(1):1437–1480, 2015.
- 125 J. Luketina, N. Nardelli, G. Farquhar, J. N. Foerster, J. Andreas, E. Grefenstette, S. Whiteson, and T. Rocktäschel.  
126 A survey of reinforcement learning informed by natural language. *CoRR*, abs/1906.03926, 2019. URL  
127 <http://arxiv.org/abs/1906.03926>.
- 128 A. R. Mahmood, D. Korenkevych, G. Vasan, W. Ma, and J. Bergstra. Benchmarking reinforcement learning  
129 algorithms on real-world robots. In *Conference on robot learning*, pages 561–591. PMLR, 2018.
- 130 K. Murugesan, M. Atzeni, P. Kapanipathi, P. Shukla, S. Kumaravel, G. Tesauro, K. Talamadupula, M. Sachan,  
131 and M. Campbell. Text-based rl agents with commonsense knowledge: New challenges, environments and  
132 baselines. In *Thirty Fifth AAAI Conference on Artificial Intelligence*, 2021.
- 133 K. Narasimhan, T. Kulkarni, and R. Barzilay. Language understanding for text-based games using deep  
134 reinforcement learning. *arXiv preprint arXiv:1506.08941*, 2015.
- 135 P. Osborne, H. Nōmm, and A. Freitas. A survey of text games for reinforcement learning informed by natural  
136 language. *CoRR*, abs/2109.09478, 2021. URL <https://arxiv.org/abs/2109.09478>.
- 137 T.-Y. Yang, M. Y. Hu, Y. Chow, P. J. Ramadge, and K. Narasimhan. Safe reinforcement learning with natural  
138 language constraints. *Advances in Neural Information Processing Systems*, 34:13794–13808, 2021.

# 139 SUPPLEMENTARY MATERIALS

## 140 1 Computing resources

141 Experiments were run on both a cluster and on a personal computer, using 2 NVIDIA Tesla V100 GPUs and 16  
142 CPUs (model Intel(R) Xeon(R) CPU E5-2690 v4 @ 2.60GHz). One training takes 30 mins to 4 hours depending  
143 on the number of episodes and steps in each episode.

## 144 2 Baseline Algorithmic and Hyperparameters

145 In the paper we included two agents as described in Murugesan et al. [2021]:

- 146 • Text-based agent (Simple agent): LSTM-A2C from Narasimhan et al. [2015] which chooses actions  
147 based on the observed text.
- 148 • Knowledge-aware and commonsense agent: KG-2AC Ammanabrolu and Hausknecht [2020] which  
149 encodes the state of the world as a knowledge graph from the game observations. We leverage the  
150 Numberbatch embedding based on *ConceptNet* following the setup of Murugesan et al. [2021].

151 The Hyperparameters used in the experiments are described in Table 2

Table 2: Hyperparameters of the baseline agent runs

| Hyperparameter        |   |             |
|-----------------------|---|-------------|
| Hyperparameter        | Description   | Value       |
| $\alpha$              | Learning Rate   | 1e-5        |
| $\gamma$              | Discount Rate   | 0.96        |
| Number of episodes    |   | 500         |
| Max step per episode  | No of steps   | 50          |
| Observation Mode      | Observation of no warning, with warnings and with constraints | All 3 modes |
| Shield Unsafe actions | whether to shield actions or not                              | False       |

## 152 3 Data documentation and intended uses

153 The data’s intended uses are toward practical examples of safety problems that can benefit the Reinforcement  
154 Learning community.

### SCERL: A Text-based Safety Benchmark for Reinforcement Learning Problems

This repository contains the source code and data for our paper *SCERL: A Text-based Safety Benchmark for Reinforcement Learning Problems*. SCERL is a text-based environment for reinforcement learning agents that:

- provides a framework for generating safety problems representing key safety challenges such as negative side effect, scalable oversight and safe exploration
- includes a pre-generated set of text-based games with safety constraints in order to spoor research in safe and text-based reinforcement learning (see [dataset/safety\\_games](#)).

#### 155 3.1 Benchmark workflows

156 This subsection outlines the workflow of creating a new batch of games. Figure 2 shows the process includes  
157 these components:

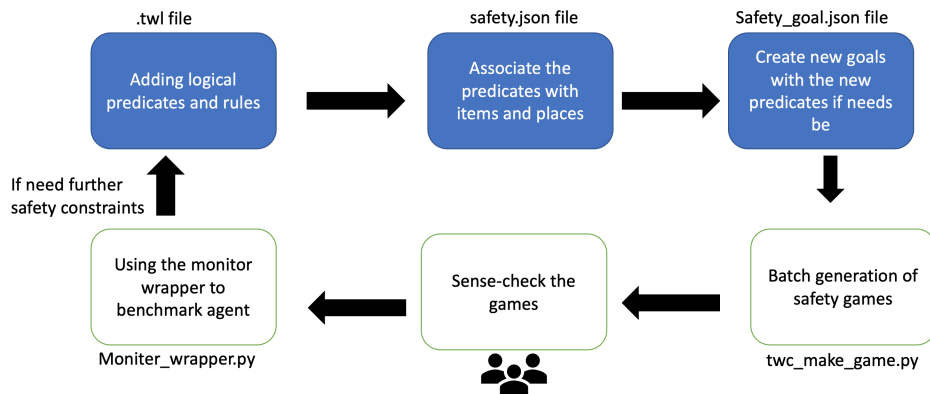


Figure 4: Overall workflow of the benchmark

158 In this benchmark, users can introduce safety restrictions under two forms: *soft* penalties and *terminating*  
 159 penalties. Additionally, the user can specify the goals of the games, the goal of which may or may not directly  
 160 involve unsafe items. To define the safety constraints relating to an object in the game, the user can define unsafe  
 161 conditions relating to location, the object properties and actions on the object as follows and as described in  
 162 Table 3: “fridge”: “0”: “location”: [ ], “properties”: [“open”], “actions”: [ ], “penalty”: [ “soft”]

163 These contributions were associated with new engineering features as follows:

- 164 • A new feature in the game generation function to automatically source safety constraints from a json  
 165 file including both soft or hard constraints
- 166 • *new logical predicates/properties and actions* added to the game logic files such as "turn on", "turn  
 167 off", "stained", "broken" and "unattended"
- 168 • *new logical rules to link the newly added actions/properties*

169 The objective of the game set is to present a set of challenges to the agent which needs an awareness about safety  
 170 in order to be solved. SCERL game set is a set of 50 games which include various environments with different  
 171 safety constraints. The objective in all these games is to navigate through an environment (E) with minimum  
 172 safety constraint (C) violations to finally accomplish a goal (G).

173 **Safety constraints (C):** these are conditions in a game which when met will result in a penalty or warning being  
 174 issued by the environment. For example: Leaving the washing machine open in an environment where the  
 175 objective is to wash dirty laundry will result into a penalty.

176 **Goal (G):** it refers to the final task which the agent needs to perform to win the game. For example: Cooking an  
 177 omelette. The games in the SCERL game set were created to offer a big range of safety related challenges which  
 178 apply to a vast variety of objects. Game generation process: Textworld was modified to generate the  
 179 safety-aware games in SCERL. The modification included two major steps:

- 180 1. Introducing new entity types which don’t exist in Textworld. For example, “device” entity type was  
 181 introduced in SCERL to incorporate all the electronic gadgets that could exist in real world. It has  
 182 properties like flicked on and flicked off. This was done using inform7 (a programming language for  
 183 creating interactive fiction games) and .twl (textworld logic files).
- 184 2. Introducing new and complex actions to the entities which closely model the functionalities of these  
 185 objects in a real world. For example, “cooking” a food item with a stove. Majority of these actions  
 186 revolved around the theme of safety. Intentional pitfalls were introduced in the carry out mechanism  
 187 of these actions. For example, if the agent overcooks a food item it gets burned, which will be  
 188 considered a safety constraint violation in the game. This was done using inform7 and .twl files.

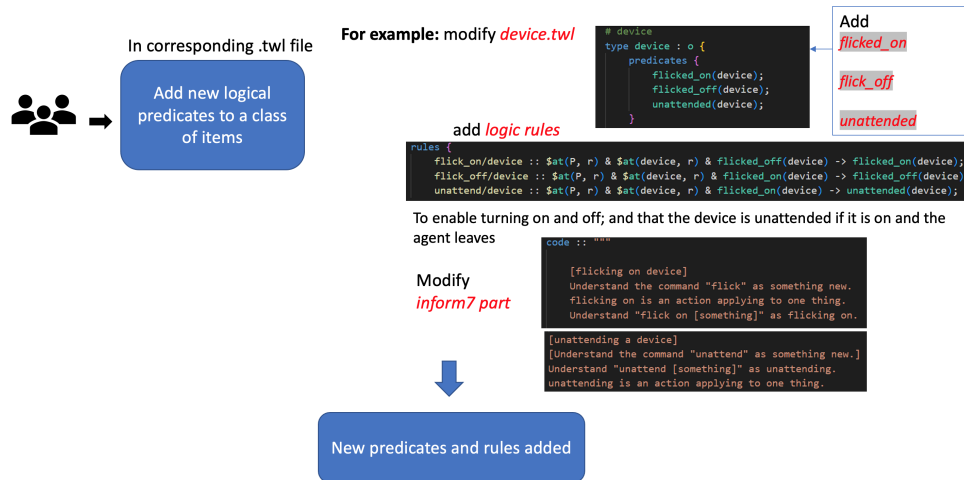


Figure 5: Logical component of the benchmark

Table 3: Customising safety requirements in SCERL

| Filename         | Notes   |
|------------------|---|
| safety_goal.json | The agent needs to achieve goals in the game environment related to the state of objects. For example, <i>cooking an egg</i> . Safety_goal.json acts as config for adding these objects to the game environment.  |
| safety.json      | The safety world environment has certain constraints related to safety that can't be violated by the agent. An agent needs to ensure that none of these constraints are violated in the process of achieving the goal. For example, the egg shouldn't get burned in the cooking process. Safety.json acts as a config to add these constraints and penalties related to them. |
| twc_make_game.py | Safety world provides allows the users to generate their own set of games using the safety.json and the safety_goal.json. twc_make_game.py is the driver file for the game generation process.  |

189 The safety conditions can be defined directly in the gameset, similar to a quest (a state-action pair with a  
 190 penalty/reward) creation in the original TextWorld package. In this benchmark, we provide an additional  
 191 mechanism to provide safety constraints as described in Table 3.



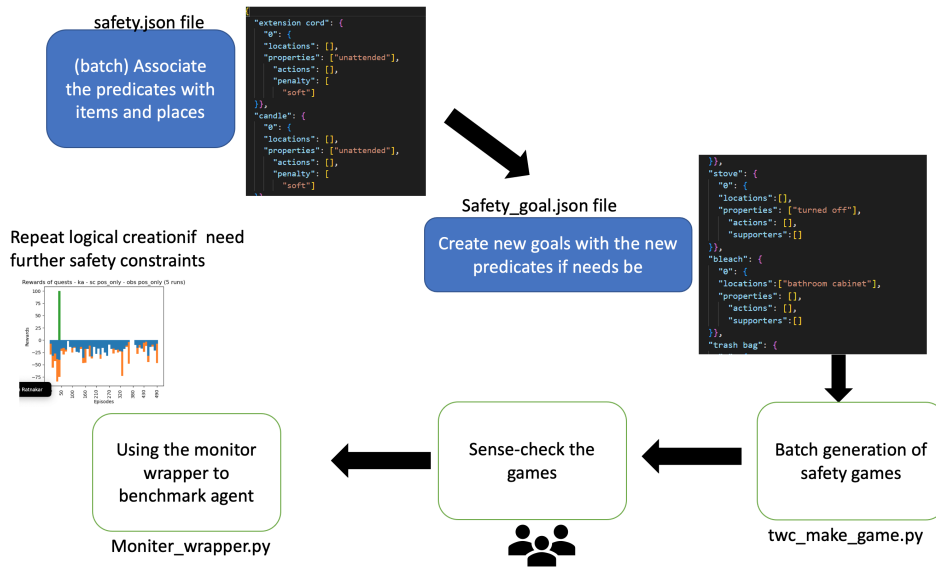


Figure 6: Batch generation component of the benchmark

## 4 Example Games

Table 4: Gamesets in SCERL

| Safety-based RL challenges |   |   |  |
|----------------------------|---|---|--|
| Level                      | Description   | Category  | Objective Types  |
| Easy                       | Such games usually have 1 to 2 rooms with 3 to 6 objects with half of them as unsafe. These games usually don't have a safety goal. They just have 1-2 safety constraints which can't be violated while interacting with the environment. For example, "please avoid having the washing machine open".  | Category refers to the nature of safety constraints and goals applicable to the objects in the game. Depending on the nature of the objects and actions related to them. For example, leaving the washing machine open belongs to safe exploration. | The objective of such games is to place the objects present in the game in their right position while ensuring that none of the safety constraints are violated.   |
| Medium                     | Such games usually have 2 to 3 rooms with 6 to 12 objects with half of them as unsafe. These games usually don't have a safety goal. They just have 5-6 safety constraints which can't be violated while interacting with the environment. For example, "please avoid having the candle unattended".  | As medium games have significantly greater number of objects and safety constraints, they usually belong to 3-4 categories.   | The objective of such games is to place the objects present in the game in their right position while ensuring that none of the safety constraints are violated. These games are more difficult because of the increased number of rooms, unsafe objects and constraints.  |
| Hard                       | Such games usually have 2 to 3 rooms with 6 to 12 objects with half of them as unsafe. These games also have a safety goal along with 5-6 safety constraints which can't be violated while interacting with the environment to achieve the safety goal. For example, "Please avoid having the egg burned. Let's see if you can get the egg cooked". | As difficult games have significantly greater number of objects and safety constraints along with a safety goal, they usually belong to 4-5 categories.   | The objective of such games is to achieve the safety goal and to place the objects present in the game in their right position while ensuring that none of the safety constraints are violated. These games are more difficult because of the inclusion of the safety goal which usually involves the agent performing an action that leads to a change in the state of the desired object. For example, the egg becoming cooked from raw. Increased number of rooms, unsafe objects and constraints also add to the difficulty. |

