

# Rejecting Hallucinated State Targets during Planning

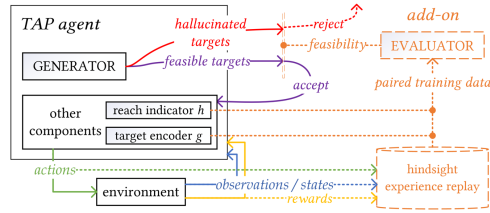
Mingde "Harry" Zhao<sup>1,2,3</sup>, Tristan Sylvain<sup>4</sup>, Romain Laroche<sup>3</sup>, Doina Precup<sup>1,2,5,6</sup> & Yoshua Bengio<sup>1,5,6,7</sup>

1: Mila (Quebec AI Institute), 2: McGill University, 3: Wayve, 4: RBC Borealis, 5: (DeepMind), 6: CIFAR Chair, 7: DIRO (Université de Montréal),



## Abstract

Model hallucinations can lead to infeasible targets (sets of states), which cause delusional planning behaviors in Target-Assisted Planning (TAP). This work first categorizes infeasible targets. Then, we propose to reject infeasible targets with an add-on evaluator, which trains alongside TAP agents, requiring almost no change to the agent (and the generative model, "generator") it is assisting. To make sure the learned evaluator is non-delusional, we developed a solution combining 1) specific learning rule, 2) architecture, and 3) assistive hindsight relabeling strategies. Our experiments validate significant reductions in delusional behaviors and performance improvements for several kinds of existing TAP agents.



## Targets & $\tau$ -feasibility

A target is an embedding of a set of states. Assume we are given  $h$  which indicates if a state belongs to target  $g^\odot$ . For practical purposes, we also introduce a time limit  $\tau$ . Let  $D_\pi(s, g^\odot)$  be a random variable representing the 1st timestep  $t$ , s.t.  $h(s_t, g^\odot) = 1$ . We define the  $\tau$ -feasibility of  $g^\odot$  from  $s$  under  $\pi$  as

$$p(D_\pi(s, g^\odot) \leq \tau) := \sum_{t=1}^{\tau} p(D_\pi(s, g^\odot) = t)$$

A target is generally considered "good" if it leads to rewarding outcomes:

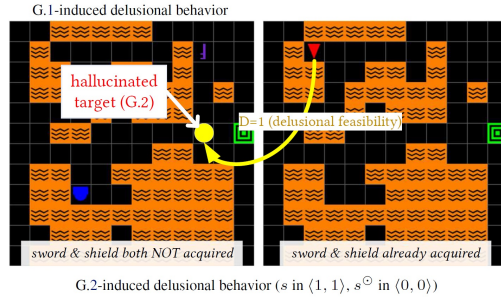
$$\mathcal{U}_{\pi, \mu}(s, g^\odot, \tau) := r_\pi(s, g^\odot, \tau) + \gamma_\pi(s, g^\odot, \tau) \cdot V_\mu(s_{\min(D_\pi(s, g^\odot), \tau)})$$

where  $\min(D_\pi(s, g^\odot), \tau)$  is the timestep the commitment to  $g^\odot$  is terminated (by  $h$  or  $\tau$ ).  $s_{\min(D_\pi(s, g^\odot), \tau)}$  is the state the agent ended up in,  $r_\pi(s, g^\odot, \tau) := \sum_{t=1}^{\min(D_\pi(s, g^\odot), \tau)} \gamma^{t-1} r_t$  is the cumulative discounted reward along the way,  $\gamma_\pi(s, g^\odot, \tau) := \gamma^{\min(D_\pi(s, g^\odot), \tau)}$  is the associated cumulative discount, and  $V_\mu(s_{\min(D_\pi(s, g^\odot), \tau)})$  is the future value for following  $\mu$  from  $s_{\min(D_\pi(s, g^\odot), \tau)}$ .

## Infeasible Targets (Singleton Case)

**G.1:** permanently infeasible targets that do not correspond to any state of the MDP.

**G.2:** temporarily infeasible targets that correspond to a valid state in the MDP, that cannot be reached from the current state. G.2 targets may not exist, depending on the MDP structure.



## Infeasible Targets (Non-singleton Case)

**Theorem 4.1:**  $\tau$ -feasibility of a target set is equivalent to the reduced set, i.e., the set with G.1 and G.2 states removed. Infeasible targets: consist of solely G.1 and / or G.2 states

Target Composition	State Correspondence	$\infty$ -Feasibility $p(D_\pi(s, g^\odot) < \infty)$	Feasibility Delusions & Resulting Delusional Planning Behaviors
Only or Single G.0	non-hallucinated feasible states from $s$	$> 0$	E.0: May think G.0 states are infeasible, thus turn to riskier alternatives, e.g., G.1 or G.2
Only or Single G.1	hallucinated "states" not belonging to the MDP	should = 0	E.1: May think G.1 states are favorable, thus commit to them. Impacted by ill-defined $V_\mu(\dots)$
Only or Single G.2	hallucinated MDP states infeasible from $s$	should = 0	E.2: May think G.2 states are favorable, thus commit to them
Some G.0	at least one non-hallucinated state from $s$	$p(D_\pi(s, g^\odot) < \infty) > 0$ (Thm. 4.1)	E.0
Only G.1 & G.2	set of ONLY hallucinated states	should = 0	E.1 & E.2

## SOLUTION DESIGN: An auxiliary evaluator that can be used to reject infeasible targets

An evaluator that can estimate the  $\tau$ -feasibility of targets with the input of source-target pairs  $(s, g^\odot)$ . [automatic] the evaluator learns to automatically differentiate the feasibility of all kinds of targets without pre-labeling: we need to exploit  $h$  [minimally intrusive] the evaluator is made to be generally applicable to existing TAP agents without changing the agents too much to disturb the sensitive RL components: we need to ensure its behavior as an add-on and it can be conditioned on the policy  $\pi$  of the agent, to learn alongside the agent.

[unified] the evaluator has a unified behavior compatible with different  $\tau$ s: we design it in a way to learn the  $\tau$ -feasibilities for many  $\tau$ s simultaneously.

## SOLUTION IMPLEMENTATION: combo of 3

- 1) Learning Rule (indirect)
- 2) Distributional Architecture
- 3) Training Data (Source-Target Pair Construction & hindsight-relabeling)

## Evaluator Learning Rule

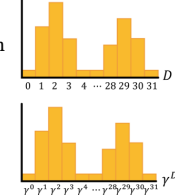
Indirectly learn feasibility by learning the distribution of  $D_\pi(s, g^\odot)$ . Represent  $D_\pi(s, g^\odot)$  as  $D_\pi(s, a, g^\odot)$ , where  $a \sim \pi(\cdot | s, g^\odot)$

$$D_\pi(s, g^\odot) \leftarrow 1 + D_\pi(s', g^\odot), \text{ with}$$

$$\begin{cases} D_\pi(s', g^\odot) := \infty & \text{if } s' \text{ is terminal and } h(s', g^\odot) = 0 \\ D_\pi(s', g^\odot) := 0 & h(s', g^\odot) = 1 \end{cases}$$

## Evaluator Architecture

C51-distributional output with each bin probability being  $p(D_\pi(s, g^\odot) = t)$ . Compatible with varying time-horizons Swap support for different outputs

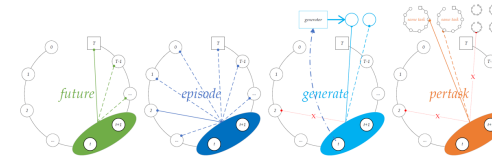


## Evaluator Training Data

The evaluator evaluates and trains on source-target pairs  $(s, g^\odot)$ , these pairs can be constructed non-trivially with hindsight-relabeling.

## Training-Reasoning Discrepancy

most existing target-directed agents are only trained on experienced data and thus only know how to deal with feasible targets. Yet, they must also deal with hallucinated targets during planning. Naïve relabeling will cause delusions (false evaluations that cannot be fixed by more training)!

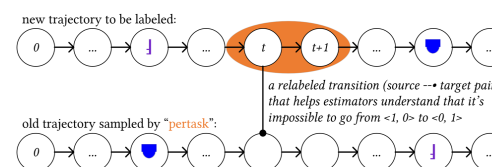


## "generate" strategy: Let Evaluator Learn about Generated Candidates

expose targets that the generator can generate

## "pertask" strategy: Let Evaluator Learn about Experienced Targets Outside the Episode

expose targets previously achieved in other episodes



## Experiment Summary

Evaluator can help reduce delusional planning behaviors and boost performance of different types of TAP agents

