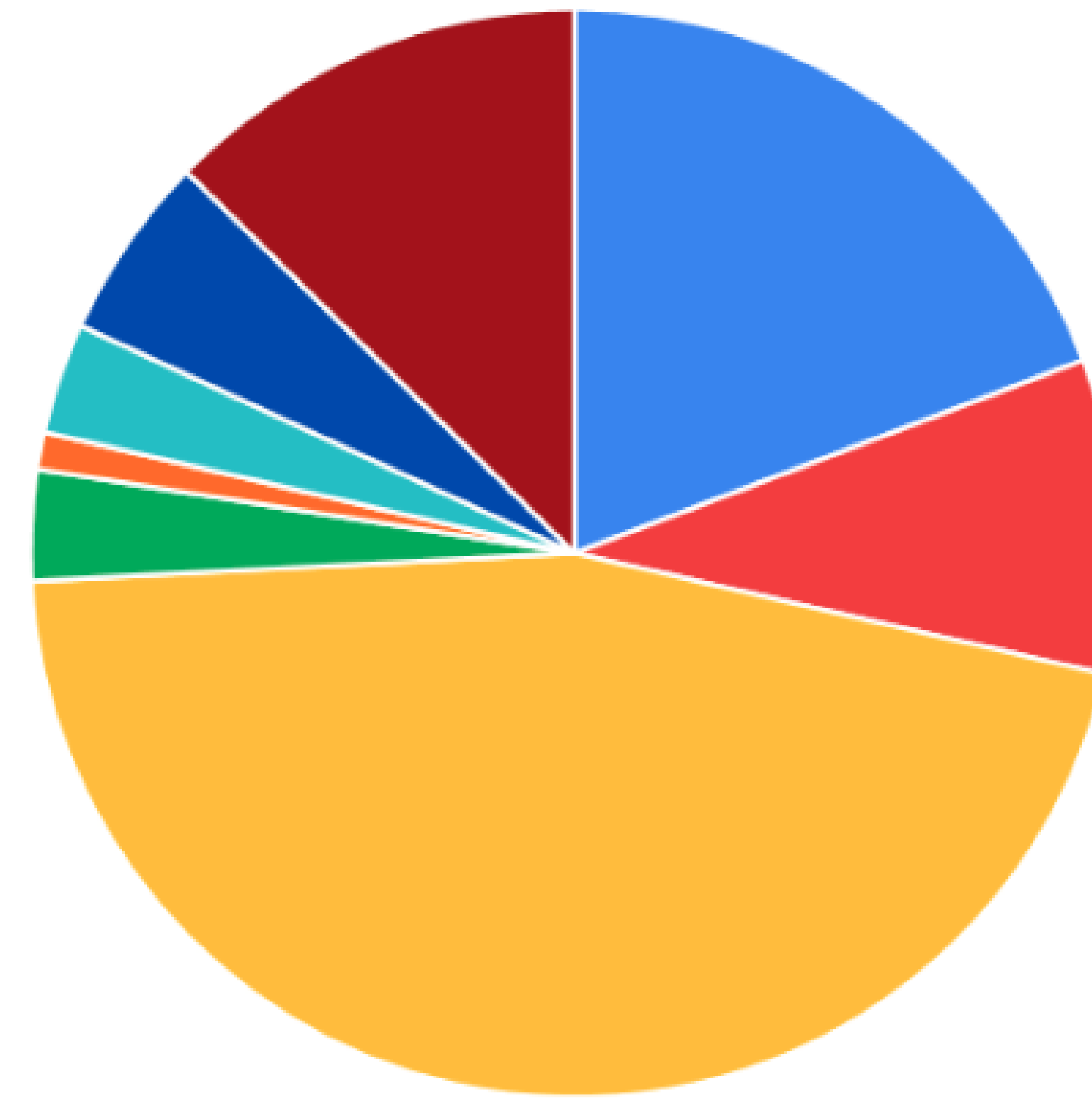


In automated planning, the need for explanations arises when there is a mismatch between a proposed plan and the user's expectation. We frame Explainable AI Planning as an iterative plan exploration process, in which the user asks a succession of contrastive questions that lead to the generation and solution of hypothetical planning problems that are restrictions of the original problem. The object of the exploration is for the user to understand the constraints that govern the original plan and, ultimately, to arrive at a satisfactory plan. We present the results of a user study that demonstrates that when users ask questions about plans, those questions are usually contrastive, i.e. "why A rather than B?". We use the data from this study to construct a taxonomy of user questions that often arise during plan exploration. Our approach to iterative plan exploration is a process of successive model restriction. Each contrastive user question imposes a set of constraints on the planning problem, leading to the construction of a new hypothetical planning problem as a restriction of the original. Solving this restricted problem results in a plan that can be compared with the original plan, admitting a contrastive explanation. We formally define model-based compilations in PDDL2.1 for each type of constraint derived from a contrastive user question in the taxonomy, and empirically evaluate the compilations in terms of computational complexity. The compilations were implemented as part of an explanation framework supporting iterative model restriction. We demonstrate its benefits in a second user study.

Why Contrastive Explanations?

Frequency of Questions Categorised by the Contrastive Taxonomy (Overall)



- Why is action A used in state S, rather than action B?
- Why is action A not used in the plan, rather than being used?
- Why is action A used in the plan, rather than not being used?
- Why is action A used outside of time window W, rather than only being allowed within W?
- Why is action A not used in time window W, rather than being used within W?
- Why is action A used at time T, rather than at least some time T' after/before T?
- Why is action A not performed before (after) action B, rather than A being performed after (before) B?
- Non-Contrastive

Example Compilation

$$\Pi' = \langle \langle Ps', Vs, As', arity' \rangle, \langle Os, I', G', W' \rangle \rangle$$

where:

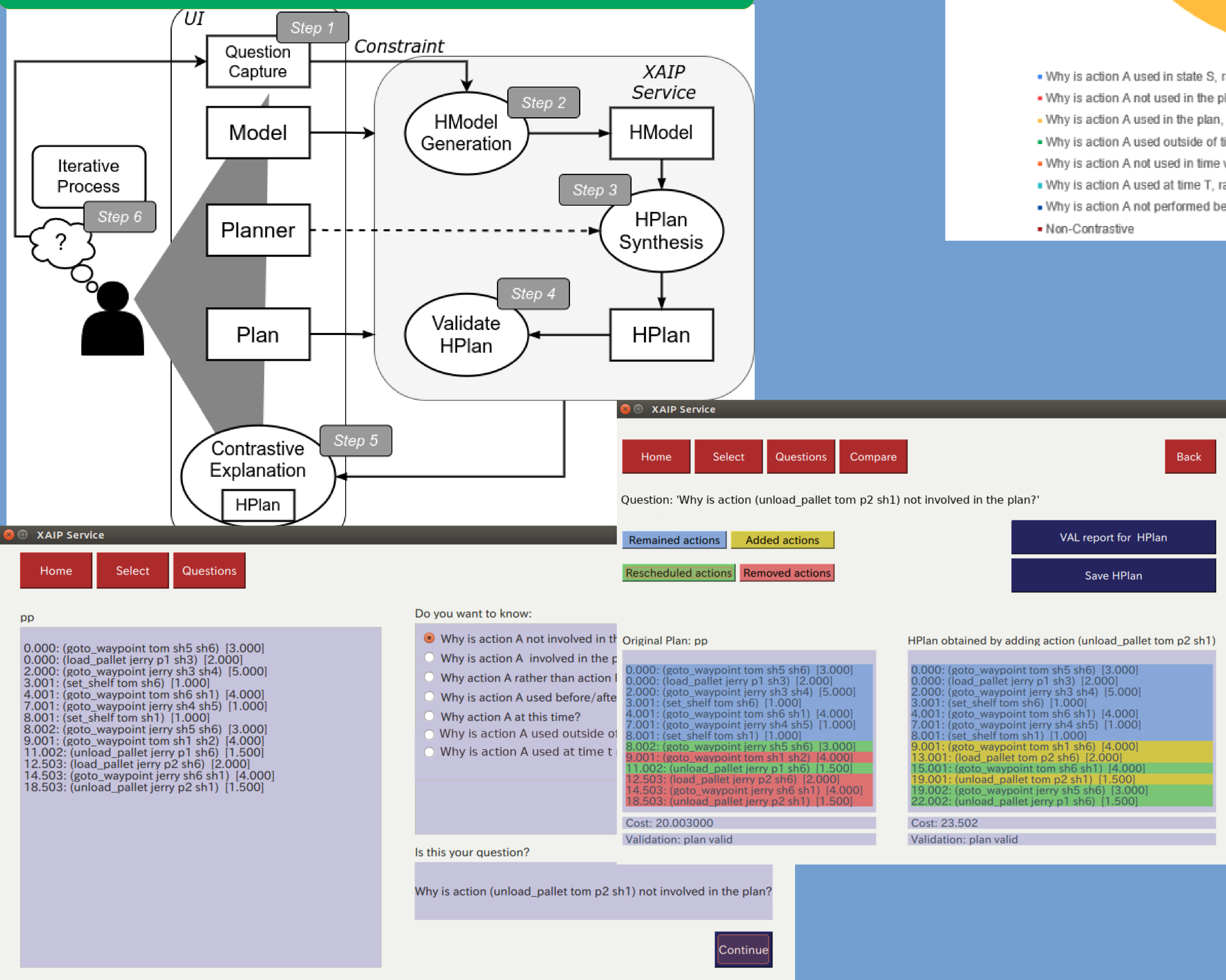
- $Ps' = Ps \cup \{can_do_a, not_done_a\}$
- $As' = \{o_a, o_{\neg a}\} \cup As \setminus \{o\}$
- $arity'(x) = arity(x), \forall x \in arity$
- $arity'(can_do_a) = arity'(not_done_a) = arity'(o_a) = arity'(o_{\neg a}) = arity(o)$
- $I' = I \cup \{ground(not_done_a, \chi)\}$
- $G' = G \cup \{ground(not_done_a, \chi)\}$
- $W' = W \cup \{\langle lb, ub, ground(can_do_a, \chi) \rangle\}$

where the new operators $o_{\neg a}$ and o_a extend o with the delete effect not_done_a and the precondition can_do_a , respectively. i.e:

$$Eff_{\neg}^-(o_{\neg a}) = Eff_{\neg}^-(o) \cup \{not_done_a\}$$

$$Pre_{\neg}^-(o_a) = Pre_{\neg}^-(o) \cup \{can_do_a\}$$

Framework



Iterative Process

