# Explainable Human-AI Interaction:
# A Planning Perspective

Sarath Sreedharan, Anagha Kulkarni, Subbarao Kambhampati

Arizona State University

# Abstract

From its inception, AI has had a rather ambivalent relationship with humans—swinging between their augmentation and replacement. Now, as AI technologies enter our everyday lives at an ever increasing pace, there is a greater need for AI systems to work synergistically with humans. One critical requirement for such synergistic human-AI interaction is that the AI systems be explainable to the humans in the loop. To do this effectively, AI agents need to go beyond planning with their own models of the world, and take into account the mental model of the human in the loop. Drawing from several years of research in our lab, we will discuss how the AI agent can use these mental models to either conform to human expectations, or change those expectations through explanatory communication. While the main focus of the book is on cooperative scenarios, we will point out how the same mental models can be used for obfuscation and deception. Although the book is primarily driven by our own research in these areas, in every chapter, we will provide ample connections to relevant research from other groups.


**Keywords:** Human-aware AI Systems, Human-AI Interaction, Explainability, Interpretability, Human-aware planning, Obfuscation

ii

# Contents

# Preface

Artificial Intelligence (AI) systems that interact with us the way we interact with each other have long typified Hollywood's image, whether you think of HAL in "2001: A Space Odyssey," Samantha in "Her," or Ava in "Ex Machina." It thus might surprise people that making systems that interact, assist or collaborate with humans has never been high on the technical agenda.

From its beginning, AI has had a rather ambivalent relationship with humans. The biggest AI successes have come either at a distance from humans (think of the "Spirit" and "Opportunity" rovers navigating the Martian landscape) or in cold adversarial faceoffs (the Deep Blue defeating world chess champion Gary Kasparov, or AlphaGo besting Lee Sedol). In contrast to the magnetic pull of these "replace/defeat humans" ventures, the goal of designing AI systems that are human-aware, capable of interacting and collaborating with humans and engendering trust in them, has received much less attention.

More recently, as AI technologies started capturing our imaginations, there has been a conspicuous change — with "human" becoming the desirable adjective for AI systems. There are so many variations — human-centered, human-compatible, human-aware AI, etc. — that there is almost a need for a dictionary of terms. Some of this interest arose naturally from a desire to understand and regulate the impacts of AI technologies on people. Of particular interest for us are the challenges and impacts of AI systems that continually interact with humans — as decision support systems, personal assistants, intelligent tutoring systems, robot helpers, social robots, AI conversational companions, etc.

To be aware of humans, and to interact with them fluently, an AI agent needs to exhibit social intelligence. Designing agents with social intelligence received little attention when AI development was focused on autonomy rather than coexistence. Its importance for humans cannot be overstated, however. After all, evolutionary theory shows that we developed our impressive brains not so much to run away from lions on the savanna but to get along with each other.

A cornerstone of social intelligence is the so-called "theory of mind" — the ability to model mental states of humans we interact with. Developmental psychologists have shown (with compelling experiments like the Sally-Anne test) that children, with the possible exception of those on the autism spectrum, develop this ability quite early.

Successful AI agents need to acquire, maintain and use such mental models to modulate their own actions. At a minimum, AI agents need approximations of humans' task and goal models, as well as the human's model of the AI agent's task and goal models. The former will guide the agent to anticipate and manage the needs, desires and attention of humans in the loop (think of the prescient abilities of the character Radar on the TV series "M\*A\*S\*H"), and the latter allow it to act in ways that are interpretable to humans — by conforming to their mental models of it — and be ready to provide customized explanations when needed.

With the increasing use of AI-based decision support systems in many high-stakes areas, including health and criminal justice, the need for AI systems exhibiting interpretable or explainable behavior to humans has become quite critical. The European Union's General Data Protection Regulation posits a right to contestable explanations for all machine decisions that affect humans (e.g., automated approval or denial of loan applications). While the simplest form of such explanations could well be a trace of the reasoning steps that lead to the decision, things get complex quickly once we recognize that an explanation is not a soliloquy and that the comprehensibility of an explanation depends crucially on the mental states of the receiver. After all, your physician gives one kind of explanation for her diagnosis to you and another, perhaps more technical one, to her colleagues.

Provision of explanations thus requires a shared vocabulary between AI systems and humans, and the ability to customize the explanation to the mental models of humans. This task becomes particularly challenging since many modern data-based decision-making systems develop their own internal representations that may not be directly translatable to human vocabulary. Some emerging methods for facilitating comprehensible explanations include explicitly having the machine learn to translate explanations based on its internal representations to an agreed-upon vocabulary.

AI systems interacting with humans will need to understand and leverage insights from human factors and psychology. Not doing so could lead to egregious miscalculations. Initial versions of Tesla's auto-pilot self-driving assistant, for example, seemed to have been designed with the unrealistic expectation that human drivers can come back to full alertness and manually override when the self-driving system runs into unforeseen modes, leading to catastrophic failures. Similarly, the systems will need to provide an appropriate emotional response when interacting with humans (even though there is no evidence, as yet, that emotions improve an AI agent's solitary performance). Multiple studies show that people do better at a task when computer interfaces show appropriate affect. Some have even hypothesized that part of the reason for the failure of Clippy, the old Microsoft Office assistant, was because it had a permanent smug smile when it appeared to help flustered users.

AI systems with social intelligence capabilities also produce their own set of ethical quandaries. After all, trust can be weaponized in far more insidious ways than a rampaging robot. The potential for manipulation is further amplified by our own very human tendency to anthropomorphize anything that shows even remotely human-like behavior. Joe Weizenbaum had to shut down Eliza, history's first chatbot, when he found his staff pouring their hearts out to it; and scholars like Sherry Turkle continue to worry about the artificial intimacy such artifacts might engender. Ability to manipulate mental models can also allow AI agents to engage in lying or deception with humans, leading to a form of "head fakes" that will make today's deep fakes tame by comparison. While a certain level of "white lies" are seen as the glue for human social fabric, it is not clear whether we want AI agents to engage in them.

As AI systems increasingly become human-aware, even quotidian tools surrounding us will start gaining mental-modeling capabilities. This adaptivity can be both a boon and a bane. While we talked about the harms of our tendency to anthropomorphize AI artifacts that are not human-aware, equally insidious are the harms that can arise when we fail to recognize that what we see as a simple tool is actually mental-modeling us.

Indeed, micro-targeting by social media can be understood as a weaponized version of such manipulation; people would be much more guarded with social media platforms if they realized that those platforms are actively profiling them.

Given the potential for misuse, we should aim to design AI systems that must understand human values, mental models and emotions, and yet not exploit them with intent to cause harm. In other words, they must be designed with an overarching goal of beneficence to us.

All this requires a meaningful collaboration between AI and humanities — including sociology, anthropology and behavioral psychology. Such interdisciplinary collaborations were the norm rather than the exception at the beginning of the AI field and are coming back into vogue.

Formidable as this endeavor might be, it is worth pursuing. We should be proactively building a future where AI agents work along with us, rather than passively fretting about a dystopian one where they are indifferent or adversarial. By designing AI agents to be human-aware from the ground up, we can increase the chances of a future where such agents both collaborate and get along with us.

This book then is a step towards designing such a future. We focus in particular on recent research efforts on making AI systems explainable. Of particular interest are settings where the AI agents make a sequence of decisions in support of their objectives, and the humans in the loop get to observe the resulting behavior. We consider techniques for making this behavior explicable to the humans out of the box, or after an explanation from the AI agent. These explanations are modeled as reconciliations of the mental models the humans have of the AI agents' goals and capabilities. The central theme of many of these techniques is reasoning with the mental models of the humans in the loop. While much of our focus is on cooperative scenarios, we also discuss how the same techniques can be adapted to support lies and deception in adversarial scenarios. In addition to the formal frameworks and algorithms, this book also discusses several applications of these techniques in decision-support and human-robot interaction scernarios. While we focus on the developments from our group, we provide ample context of related developments across several research groups and areas of AI currently focusing on the explainability of AI systems.

# Acknowledgements

Parts of the material in this book is drawn from our refereed research papers published in several venues; we cite all the sources in the bibliography. Parts of the preface are drawn from a column on human-aware AI systems that first appeared in The Hill.

# Introduction

Artificial Intelligence, the discipline many of us call our intellectual home, is suddenly having a rather huge cultural moment. It is hard to turn anywhere without running into mentions of AI technology and hype about its expected positive and negative societal impacts. AI has been compared to fire *and* electricity, and commercial interest in the AI technologies has sky rocketed. Universities – even high schools – are rushing to start new degree programs or colleges dedicated to AI. Civil society organizations are scrambling to understand the impact of AI technology on humanity, and governments are competing to encourage or regulate AI research and deployment.

There is considerable hand-wringing by pundits of all stripes on whether in the future, AI agents will get along with us or turn on us. Much is being written about the need to make AI technologies safe and delay the "doomsday". We believe that as AI researchers, we are not (and cannot be) passive observers. It is *our* responsibility to design agents that can and will get along with us. Making such *human-aware* AI agents, however poses several foundational research challenges that go beyond simply adding user interfaces *post facto*. In particular, human-aware AI systems need to be designed such that their behavior is *explainable* to the humans interacting with them. This book describes some of the state-of-the-art approaches in making AI systems explainable.

## 1.1  Humans & AI Agents: An Ambivalent Relationship

In this book we focus on human-aware AI systems—goal directed autonomous systems that are capable of effectively interacting, collaborating and teaming with humans. Although developing such systems seems like a rather self-evidently fruitful enterprise, and popular imaginations of AI, dating back to HAL, almost always assume we already do have human-aware AI systems technology, little of the actual energies of the AI research community have gone in this direction.

From its inception, AI has had a rather ambivalent relationship to humans—swinging between their augmentation and replacement. Most high profile achievements of AI have either been far away from the humans—think Spirit and Opportunity exploring Mars; or in a decidedly adversarial stance with humans, be it Deep Blue, AlphaGo or Libratus. Research into effective ways of making AI systems *interact, team and collaborate with humans* has received significantly less attention. It is perhaps no wonder that many lay people have fears about AI technology!

This state of affairs is a bit puzzling given the rich history of early connections between AI and psychology. Part of the initial reluctance to work on these issues had to do with the worry that focusing on AI systems working with human might somehow dilute the grand goals of the AI enterprise, and might even lead to temptations of "cheating," with most of the intelligent work being done by the humans in the loop. After all, prestidigitation

has been a concern since the original mechanical turk. Indeed, much of the early work on human-in-the-loop AI systems mostly focused on using humans as a crutch for making up the limitations of the AI systems [Allen, 1994]. In other words, early AI had humans be "AI-aware" (rather than AI be "human-aware").

Now, as AI systems are maturing with increasing capabilities, the concerns about them depending on humans as crutches are less severe. We would also argue that focus on humans in the loop doesn't dilute the goals of AI enterprise, but in fact broadens them in multiple ways. After all, evolutionary theories tell us that humans may have developed the brains they have, not so much to run away from the lions of the savanna or tigers of Bengal but rather to effectively cooperate and compete with each other. Psychological tests such as the Sally Anne Test [Wimmer and Perner, 1983] demonstrate the importance of such social cognitive abilities in the development of collaboration abilities in children.

Some branches of AI, aimed at specific human-centric applications, such as intelligent tutoring systems [VanLehn, 2006], and social robotics [Breazeal, 2004, 2003, Scassellati, 2002], did focus on the challenges of human-aware AI systems for a long time. It is crucial to note however that human-aware AI systems are needed in a much larger class of quotidian applications beyond those. These include human-aware AI assistants for many applications where humans continue to be at the steering wheel, but will need naturalistic assistance from AI systems—akin to what they can expect from a smart human secretary. *Increasingly, as AI systems become common-place, human-AI interaction will be the dominant form of human-computer interaction* [Amershi et al., 2019].

For all these reasons and more, human-aware AI has started coming to the forefront of AI research of late. Recent road maps for AI research, including the 2016 JASON report[1] and the 2016 White House OSTP report[2] emphasize the need for research in human-aware AI systems. The 2019 White House list of strategic R&D priorities for AI lists "developing effective methods for human-AI collaboration" at the top of the list of priorities[3]. Human-Aware AI was the special theme for the 2016 International Joint Conference on AI (with the tagline "*why intentionally design a dystopian future and spend time being paranoid about it?*"); it has been a special track at AAAI since 2018.

## 1.2   Explanations in Humans

Since our books is about explainable AI systems, it is useful to start with a broad overview of explainability and explanations in humans.

---

[1]https://fas.org/irp/agency/dod/jason/ai-dod.pdf

[2]https://obamawhitehouse.archives.gov/sites/default/files/whitehouse_files/ microsites/ostp/NSTC/national_ai_rd_strategic_plan.pdf

[3]https://www.whitehouse.gov/wp-content/uploads/2019/06/National-AI-Research-and-Development-Strategic-Plan-2019-Update-June-2019.pdf

**Fig. 1.1:** *Architecture of an intelligent agent that takes human mental models into account. All portions in yellow are additions to the standard agent architecture, that are a result of the agent being human-aware. $\mathcal{M}_h^R$ is the mental model the human has of the AI agent's goals and capabilities and $\mathcal{M}_r^H$ is the (mental) model the AI agent has of the human's goal and capabilities (see the section on Mental Models in Human-Aware AI)*

### 1.2.1 When and Why do Humans expect explanations from each other?

To understand the different use cases for explanations offered by AI systems, it is useful to survey the different scenarios where humans ask for explanations from each other:

**When they are confused and or surprised by the behavior of the other person**
People expect explanations when the behavior from the other person is *not what they expected*–and is thus *inexplicable*. It is worth noting that this confusion at the other person's behavior is not predicated on that person's behavior being *incorrect* or *in-optimal*. We may well be confused/surprised when a toddler, for example, makes an optimal chess move. In other words, the need for explanations arises *when the other person's behavior is not consistent with the model we have of the other person*. Explanations are thus meant to *reconcile* these misaligned expectations.

**When they want to teach the other person** We offer explanations either to make the other person understand the real rationale behind a decision or to convince the other person that our decision in this case is not a fluke. Explanatory dialog allows either party to correct the mental model of the other party. The explanations become useful in *localizing the fault, if any,* in the other person's understanding our decision.

Note that the need for explanation is thus dependent on one person's model of the other person's capabilities/reasoning. Mental models thus play a crucial part in offering customized explanations. Indeed, a doctor explains her diagnostic decision to her patient

in one way and to her peers in a different (possibly more "jargon-filled" way), because she intuitively understands the levels of abstraction of the mental models they have of diseases and diagnoses.

It is also worth considering how the need for explanations *reduces* over time. It naturally reduces as the mental models we have of the other agent gets better aligned with that other agent's capabilities, reducing any residual surprise at their behavior. This explains the ideal of "wordless collaboration" between long-time collaborators.

The need for explanations is also modulated by the trust we have in the other person. Trust can be viewed as our willingness to put ourselves in a position of vulnerability. When we trust the other person, even when their behavior/decision is inexplicable to us, we might defer our demands for any explanations from that person. This explains why we ask fewer explanations from people we trust.

### 1.2.2   How do Humans Exchange Explanations?

We now turn to a qualitative understanding of the ways in which humans *exchange* explanations, with a view to gleaning lessons for explanations in the human-AI interaction. While explanations might occur in multiple modalities, we differentiate two broad types: *Pointing explanations* and *Symbolic Explanations*. A given explanatory interaction might be interspersed with both types of explanations.

**Pointing (Tacit) Explanations** These are the type of explanations where the main explanation consists of pointing to some specific features of the object that both the explainer and explainee can see. This type of explanations may well be the only feasible one to exchange when the agents share little beyond what they perceive in their immediate environment. Pointing explanations can get quite unwieldy (both in terms of the communication bandwidth and the cognitive load for processing them)– especially when explaining sequential decisions (such as explaining why you took an earlier flight than the one the other person expected)–as they will involve pointing to the relevant regions of the shared "video history" between the agents (or, more generally, *space time signal tubes*).

**Symbolic (Explicit) Explanations** These involve exchanging explanations in a symbolic vocabulary. Clearly, these require that the explainer and explainee share a symbolic vocabulary to begin with.

Typically, pointing explanations are used for tacit knowledge tasks, and symbolic explanations are used for explicit knowledge tasks. Interestingly, over time, people tend to develop symbolic vocabulary even for exchanging explanations over tacit knowledge tasks. Consider, for example, terms such as *pick and roll* in basket ball, that are used as a shorthand for a complex space time tube–even though the overall task is largely a tacit knowledge one.

The preference for symbolic explanations is not merely because of their *compactness*,

but also because they significantly reduce the cognitive load on the receiver. This preference seems to exist despite the fact that the receiver may likely have to recreate in their own mind the "video" (space time signal tube) versions of symbolic explanations within their own minds.

### 1.2.3 (Why) Should AI Systems be Explainable?

Before we delve into explainability in AI systems, we have to address the view point that explainability from AI systems is largely unnecessary. Some have said, for example, that AI systems–such as those underlying Facebook–routinely make millions of decisions/recommendations (of the "*you might be interested in seeing these pages*" variety) and no users ask for explanations. Some even take the view that since AI systems might eventually be "more intelligent" than any mere mortal human, requiring them to provide explanations for their (obviously correct) decisions unnecessarily hobbles them. Despite this, there are multiple reasons why we want AI systems to be explainable

- Since humans do expect explanations from each other, a naturalistic human-AI interaction requires that AI systems be explainable

- Contestability of decisions is an important part of ensuring that the decisions are seen to be fair and transparent, thus engendering trust in humans. If AI systems make high stakes decisions, they too need to be contestable.

- Since there is always an insurmountable gap between the true preferences of humans and the AI system's estimate of those preferences, an explanatory dialog allows for humans to "teach" the AI systems their true preferences in a demand-driven fashion.

- Given that AI systems–especially those that are trained purely from raw data–may have unanticipated failure modes, explanations for their decisions often help the humans get a better sense of these failure modes. As a case in point, recently, it has been reported that when some Facebook users saw a video of an African male in a quotidian situation, the system solicitously asked the users *Do you like to see more primate videos?* As egregious as this "heads-up" of the failure mode of the underlying system sounds, it can be more insidious when the system just silently acts on the decisions arrived at through those failure modes, such as, for example, inexplicably filling the user's feeds with a slight uptick of primate videos. Explanations are thus a way for us to catch failure modes of these alien intelligences that we are increasingly surrounded by.

## 1.3 Dimensions of Explainable AI systems

Now that we have reviewed how explanations and explainability play a part in human-human interactions, we turn to explainability in human-AI interactions. We will specifically look at the use cases for explanations in human-AI interaction, some desirable requirements on explanations, and an overview of the ongoing research on explainable AI systems.

### 1.3.1   Use cases for explanations in Human-AI Interaction

When humans are in the loop with AI systems, they might be playing a variety of different roles. Being interpretable to one human in one specific role may or may not translate to interpretability for other humans in other roles. Perhaps the most popular role considered in the explainable AI literature to-date is humans as debuggers trying to flag and correct an AI system's behavior. In fact, much of the explainable machine learning research has been focused on this type of debugging role for humans. Given that the humans in the loop here have invested themselves into debugging the system, they are willing to *go into the land of the AI agents*, rather than expect them to come into theirs. This lowers the premium on comprehensibility of explanations.

Next we can have humans as observers of the robot's behavior – either in a peer to peer, or student/teacher setting. The observer might be a lay person who doesn't have access to the robot's model of the task; or an expert one who does.

Finally, the human might be a collaborator–who actively takes part in completing a joint task with the robot.

No matter the role of the human, one important issue is whether the interaction between the human and the robot is a one-off one (i.e., they only interact once in the context of that class of tasks) or a longitudinal one (where the human interacts with the same robot over extended periods). In this latter case, the robot can engender *trust* in the human through its behavior, which, in turn reduces the need for interpretability. In particular, the premium on interpretability of the behavior itself is reduced when the humans develop trust over the capabilities and general beneficence of the robot.

### 1.3.2   Requirements on Explanations

There are a variety of requirements that can be placed on explanations that an AI agent gives the human in the loop:

**Comprehensibility:** The explanation should be comprehensible to the human in the
    loop. This not only means that it should be in terms that the human can under-
    stand, but should not pose undue cognitive load (i.e., expect unreasonable inferential
    capabilities).

**Customization:** The explanations should be in a form and at a level that is accessible
    to the receiving party (explainee).

**Communicability:** The explanation should be easy to exchange. For example, symbolic
    explanations are much easier than pointing explanations (especially in sequential
    decision problems when they have to point to space time signal tubes).

**Soundness:** This is the guarantee from the AI agent that this explanation is really the
    reason behind its decision. Such a guarantee also implicitly implies that the agent
    will stand behind the explanation–and that the decision will change if the condi-
    tions underlying the explanation are falsified. For example, if the explanation for a

loan denial is that the applicant has no collateral, and the applicant then produces collateral, it is fairly expected that the loan denial decision will be reversed.

**Satisfctoriness:** This aims to measure how *satisfied* the end user is with the explanation provided. While incomprehensible, uncustomized and poorly communicated explanations will likely be unsatisfatory to the user, it is also possible that the users are satisfied with misleading explanations that align well with what they want to hear. Indeed, making explanations satisfactory to the users is a slippery slope. Imagine the end-user explanations of the kind that the EU GDPR regulations require being provided by a system with a GPT-3 back-end generating plausible explanations that are likely to satisfy the users. This is why it is important for systems not to take an "end to end" machine learning approach and learn what kind of explanations make the end users happy.

### 1.3.3 Explanations as Studied in the AI Literature

Explanations have been studied in the context of AI systems long before the recent interest in the explainable AI. We will start by differentiating two classes of explanations: *internal explanations* that the system develops to help its own reasoning, and *external explanations* that the agent offers to other agents.

Internal explanations have been used in AI systems to guide their search (e.g. explanation-based or dependency directed backtracking), or to focus their learning. The whole area of explanation-based learning–which attempts to focus the system's learning element on the parts of the scenarios that are relevant to the decision at hand (thus providing feature relevance assessment). While much of the work in explanation-based search and learning have been done in the context of explicit symbolic models, self explanations can also be useful for systems learning their own representations Zha et al. [2021].

External explanations may be given by AI systems to other automated agents/AI systems (as is the case in autonomous and multi-agent systems research), or to other *human* agents. There are however significant differences that arise based on whether the other agent is human or automated. In particular, issues such as cognitive load and inferential capacity in parsing the explanation play an important role when the other agent is a human, but not so much if it is automated. In particular, the work on certificates and proofs of optimality of the decision, which give gigabits of provenance information to support the decision, may work fine for automated agents but not human agents.

Historically, (external) explanations (to human agents) have been considered in the context of AI systems for as long as they have been deployed AI systems. There is a rich tradition of explanations in the context of expert systems. In all these cases, the form of the explanation does depend on (a) the role the human plays (debugger vs. lay observer) and (b) whether the task is an explicit knowledge one or a tacit knowledge one.

### 1.3.4   Explainable AI: The Landscape & The Tribes

The work in explainable AI systems can be classified into multiple dimensions. The first is whether what needs to be explained is a single-decision (e.g. classification) task or a behavior resulting from a sequence of decisions. Second dimension is whether the task is guided by explicit (verbalizable) knowledge or is a tacit task. The third is whether the interaction between the human and the AI agent is one-shot or iterative/longitudinal. Armed with these dimensions, we can discern several research "tribes" focusing on explainable AI systems:

**Explainable Classification**   Most work on the so-called "explainable ML" focused on explainable classification. The classification tasks may be "tacit" in that the machine learns its own potentially inscrutable features/representations. A prominent subclass here is image recognition/classification tasks based on deep learning approaches. In these cases, the default communication between the AI agents and humans will be over the shared substrate of the (spatial) image itself. Explanations here thus amount to "saliency annotations" over the image–showing which pixels/regions in the image have played a significant part in the final classification decision. Of course, there are other classification tasks–such as loan approval or fraudulent transaction detection–where human-specified features, rather than pixels/signals are used as the input to the classifier. Here explanations can be in terms of the relative importance of various features (e.g. shapley values).

**Explainable Behavior**   Here we are interested in sequential decision settings, such as human-robot or human-AI interactions, where humans and AI agents work in tandem to achieve certain goals. Depending on the setting, the human might be a passive observer/monitor, or an active collaborator. The objective here is for the AI agent (e.g. the robot) to exhibit behavior that is interpretable to the human in the loop. A large part of this work focuses on tacit interactions between the humans and robots (e.g. robots moving in ways that avoid colliding with the human, robots signaling which point they plan to go to etc.). The interactions are tacit in that the human has no shared vocabulary with the robot other than the observed behavior of the robot. Here explainability or interpretability typically depend on the robot's ability to exhibit a behavior that helps the human understand its "intentions". Concepts that have been explored in this context include *explicability*–the behavior being in conformance with human's expectation of the robot, *predictability*–the behavior being predictable over small time periods (e.g. next few actions), and *legibility*–the behavior signaling the goals of the robot.

Of course, not all human-robot interactions have to be tacit; in many cases the humans and robots might share explicit knowledge and vocabulary about their collaborative task. In such cases, the robot can also ensure interpretability through exchange of *explanations* in the shared symbolic vocabulary.

Much of this monograph focuses on explainable behavior, in the context of explicit knowledge tasks. These are the scenarios where an AI agent's ability to plan its behavior up front interacts synergistically with its desire to be interpretable to the human in the loop. We will focus first on scenarios where the human and the agent share a common vocabulary, but may have differing task models, and discuss how to ensure explicability

or provide explanations in that scenario. Towards the end, we will also consider the more general scenario where even the vocabulary is not common–with the AI agent using its own internal representations to guide its planning and decision-making, and discuss how explanations can still be provided in human's vocabularly.

## 1.4 Our Perspective on Human-Aware and Explainable AI Agents

In this section, we give a brief summary of the the broad perspective taken in this book in designing human-aware and explainable AI systems. This will be followed in the next section by the overview of the book itself.

### 1.4.1 How do we make AI agents Human-Aware?

When two humans collaborate to solve a task, both of them will develop approximate models of the goals and capabilities of each other (the so called "theory of mind"), and use them to support fluid team performance. AI agents interacting with humans – be they embodied or virtual – will also need to take this implicit mental modeling into account. This certainly poses several research challenges. Indeed, it can be argued that acquiring and reasoning with such models changes almost every aspect of the architecture of an intelligent agent. As an illustration, consider the architecture of an intelligent agent that takes human mental models into account shown in Figure 1.1. Clearly most parts of the agent architecture – including state estimation, estimation of the evolution of the world, projection of its own actions, as well as the task of using all this knowledge to decide what course of action the agent should take – are all critically impacted by the need to take human mental models into account. This in turn gives rise to many fundamental research challenges. In this book, we will use the research in our lab to illustrate some of these challenges as well as our attempts to address them. Our work has focused on the challenges of human-aware AI in the context of human-robot interaction scenarios [Chakraborti et al., 2018], as well as human decision support scenarios [Sengupta et al., 2017b]. Figure 1.2 shows some of the test beds and micro-worlds we have used in our ongoing work.

### 1.4.2 Mental Models in Explainable AI Systems

In our research, we address the following central question in designing human-aware AI systems: *What does it take for an AI agent to show explainable behavior in the presence of humans?* Broadly put, our answer is this: *To synthesize explainable behavior, AI agents need to go beyond planning with their own models of the world, and take into account the mental model of the human in the loop. The mental model here is not just the goals and capabilities of the human in the loop, but includes the human's model of the AI agent's goals/capabilities.*

**Fig. 1.2:** *Test beds developed to study the dynamics of trust and teamwork between autonomous agents and their human teammates.*



**Fig. 1.3:** *Use of different mental models in synthesizing explainable behavior. (Left) The AI system can use its estimation of human's mental model, $\mathcal{M}_r^H$, to take into account the goals and capabilities of the human thus providing appropriate help to them. (Right) The AI system can use its estimation of human's mental model of its capabilities $\mathcal{M}_h^R$ to exhibit explicable behavior and to provide explanations when needed.*

Let $\mathcal{M}^R$ and $\mathcal{M}^H$ correspond to the actual goal/capability models of the AI agent and human. To support collaboration, the AI agent needs an approximation of $\mathcal{M}^H$, we will call it $\widetilde{\mathcal{M}}_r^H$, to take into account the goals and capabilities of the human. The AI agent also needs to recognize that the human will have a model of its goals/capabilities $\mathcal{M}_h^R$, and needs an approximation of this, denoted $\widetilde{\mathcal{M}}_h^R$. It is important to note that while $\mathcal{M}^R$ and $\mathcal{M}^H$ are intended to be "executable models," in that courses of action consistent with them are in fact executable by the corresponding agent–robot or human, $\mathcal{M}_h^R$ and $\mathcal{M}_r^H$ are models of "*expectation*" by the other agent, and thus may not actually be executable by the first agent. All phases of the "sense–plan–act" cycle of an intelligent agent will have to change appropriately to track the impact on these models (as shown in Figure 1.1).

Of particular interest to us in this book is the fact that synthesizing explainable behavior becomes a challenge of supporting planning in the context of these multiple models. In particular, we shall see that the AI agent uses $\mathcal{M}_r^H$ to anticipate the human behavior and provide appropriate assistance (or at least get out of the way), while the agent uses $\widetilde{\mathcal{M}}_h^R$, its estimate of $\mathcal{M}_h^R$, to understand human's expectation on its behavior, and use that understanding to either *conform* to the human expectation or actively *communicate* with the human to induce them to change $\mathcal{M}_h^R$. We discuss the conformance aspect in terms of generating *explicable behavior*. For the model communication part, either the communication can be done *implicitly*–which is discussed in terms of generating *legible behavior*, or can be done *explicitly*–with communication actions. This latter part is what we view as the "explanation" process, and address multiple challenges involved in generating such explanations. Finally, while much of the book is focused on cooperative and non-adversarial scenarios, the mental model framework in Figure 1.3 can also be used by the agent to obfuscate its behavior or provide deceptive communication. We also discuss how such selective obfuscation and deception is facilitated.

## 1.5   Overview of this Book

Through the rest of the book, we will look at some of the central challenges related to human-aware planning that arise due to and can be addressed through working with the human's mental model. In particular, we will ground our discussions of the topics within the context of using such models to generate either interpretable or deceptive behavior. The book is structured as follows:

**Chapter 2**   In this chapter, we will focus on formally defining the goal-directed deterministic planning formalisms and the associated notations that we will be using to study and ground the technical discussions throughout this book. We will also define the three main interpretability metrics; namely, *Explicability*, *Legibility*, and *Predictability*. We will be revisiting these three measures in the following chapters. We will see how the different methods discussed throughout the book relate to these measures, and in various cases could be understood as being designed to optimize, at the very least a variant of these measures.

**Chapter 3**   We next focus on one of these measures, namely, explicability, and see how we could allow the robot to choose plans that maximize explicability scores. In particular, we will look at two main paradigms to generate such *explicable plans*. First, we consider a model-based method called *reconciliation search* that will use the given human model along with a distance function (which could potentially be learned) to generate robot plans with high explicability scores. Then we will look at a model-free paradigm, where we use feedback from users to learn a proxy for the human model in the form of a labeling model and use that simpler model to drive the explicable planning. We will also look at how one could use environment design to allow the agents to generate plans with higher explicability scores. Within the design framework, we will also look at the evolution of plan explicability within the context of longitudinal interactions.

**Chapter 4**   We next turn our focus onto legibility. We look at how the robot can reduce the human observer's uncertainty over its goals and plans, particularly when the observer has imperfect observations of its activities. The robot can reduce the observer's uncertainty about its goals or plans by implicitly communicating information through its behavior i.e. by acting legibly. We formulate this problem as a controlled observability planning problem, where in the robot can choose specific actions that allow it to modulate the human's belief over the set of candidate robot goals or plans. Further we will also discuss the relationship between plan legibility as discussed in this framework with the notion of plan predictability defined in the literature.

**Chapter 5**   In this chapter, we return to the problem of maximizing explicability and investigate an alternate strategy, namely explanations. Under this strategy, rather than letting the robot choose possibly suboptimal plans that may better align with human expectations, the robot can follow its optimal plans and provide explanations as to why the plans are in fact optimal in its own model. The explanation in this context tries to resolve any incorrect beliefs the human may hold about the robot and thus helps to reconcile the differences between the human's mental model and the robot's model. In the chapter, we will look at some specific types of model reconciliation explanations and a model space search algorithm to generate such explanation. We will also consider some approximations for such explanations and discuss some user studies that have been performed to validate such explanations.

**Chapter 6**   In this chapter, we continue our discussion on model reconciliation and focus on one specific assumption made in the earlier chapter, namely that the human's mental model may be known. We will study how one could still generate model reconciliation explanations when this assumption may not be met. In particular, we will consider three cases of incrementally decreasing access to human mental model. We will start by considering cases where an incomplete version of the model may be available, then we will look at scenarios where we can potentially learn a model proxy from human data and finally we will see the kind of explanatory queries we can field by just assuming a prototypical model for the human.

**Chapter 7**    In this chapter, we return back to the problem of explicability and look at how one could combine the two strategies discussed in earlier chapters for addressing explicability; namely, explicable planning and explanation. In fact we will see how one could let the robot trade-off the cost of explaining the plan against the overhead of choosing a potentially suboptimal but explicable plan and will refer to this process as *balanced planning*. In the chapter, we will look at some particular classes of balanced planning and introduce a concept of *self-explaining plans*. We will also introduce a compilation based method to generate such plans and show how the concept of balancing communication and behavior selection can be extended to other interpretability measures.

**Chapter 8**    In this chapter, we look at yet another assumption made in the generation of explanations, namely, that there exists a shared vocabulary between the human and the robot through which it can communicate. Such assumptions may be hard to meet in cases where the agent may be using learned and/or inscrutable models. Instead we will look at how we could learn post-hoc representations of the agent's model using concepts specified by the user and use this representation for model reconciliation. We will see how such concepts could be operationalized by learning classifiers for each concepts. In this chapter, we will also look at how we could calculate confidence over such explanations and how these methods could be used when we are limited to noisy classifiers for each concept. Further, we will discuss how we could acquire new concepts, when the originally specified set of concepts may not be enough to generate an adequate representation of the model.

**Chapter 9**    In this chapter, we turn our attention to adversarial settings, where the robot may have to act in obfuscatory manner to minimize the leakage of sensitive information. In particular, we look at settings where the adversary has partial observability of the robot's activities. In such settings, we show how the robot can leverage the adversary's noisy sensor model to hide information about its goals and plans. We also discuss an approach which maintains goal obfuscation even when the adversary is capable of diagnosing the goal obfuscation algorithm with different inputs to glean additional information. Further, we discuss a general setting where both adversarial as well as a cooperative observers with partial observations of robot's activities may exist. We show that in such a case, the robot has to balance the amount of information hidden from the adversarial observer with the amount of information shared with the cooperative observer. We will also look at the use of deceptive communication in the form of lies. We will show how we can leverage the tools of model reconciliation to generate lies and we will also discuss how such lies could in fact be used to benefit the human-robot team.

**Chapter 10**    In this chapte, we provide a discussion of some of the applications based on the approaches discussed in the book. We will look at two classes of applications. In the first, we look at decision support systems and in particular the methods related to explanations could be used to help decision-makers better understand the decisions being proposed by the user. In the second case, we will look at a specific application designed to help a user with transcribing a declarative model for the task. This application helps in identifying mistakes in current version of the model by reconciling against an empty user model.

**Chapter 11**   Finally we will close the book with a quick discussion on some of the most challenging open problems in Human-Aware AI.

Although this monograph is informed primarily by our group's work, in each chapter, we provide bibliographic remarks at the end connecting the material discussed to other related works.

# Measures of Interpretability

This chapter will act as the introduction to the technical discussions in the book. We will start by establishing some of the basic notations that we will use, including the definitions of deterministic goal-directed planning problems, incomplete planning models, sensor models, etc. With the basic notations in place, we will then focus on establishing the three main interpretability measures in human-aware planning; namely, *Explicability, Legibility, and Predictability.* We will revisit two of these measures (i.e., explicability and legibility) and discuss methods to boost these measures throughout the later chapters.

## 2.1    Planning Models

We will be using goal-oriented STRIPS planning models to represent the planning problems used in our discussions. However, the ideas discussed in this book apply to any of the popular planning representations. Under this representation scheme, a planning model (sometimes also referred to as a planning problem) can be represented by the tuple $\mathcal{M} = \langle F, A, I, G, C \rangle$, where the elements correspond to

- $F$ - A set of propositional fluents that define the space of possible task states. Each state corresponds to a specific instantiation of the propositions. We will denote the set of states by $S$. When required we will uniquely identify each state by the subset of fluents which are true in the given state. This representation scheme implicitly encodes the fact that any proposition from F not present in the set representation of the state is false in the underlying state.

- $A$ - The set of actions that are available to the agent. Under this representation scheme, each action $a_i \in A$ is described by a tuple of the form

$$a_i = \langle \text{pre}(a_i), \text{adds}(a_i), \text{dels}(a_i) \rangle,$$

where

  - $\text{pre}(a_i)$ - The preconditions for executing the action. For most of the discussion, we will follow the STRIPS execution semantics, wherein an action is only allowed to execute in a state where the preconditions are 'met'. In general, preconditions can be any logical formula over the propositional fluents provided in $F$. Moreover, we would say an action precondition is met in a given state if the logical formula holds in that state (remember a state here correspond to a specific instantiation of fluents) We will mostly consider cases where the precondition is captured as a conjunctive formula over a subset of state fluents, which we can equivalently represent as a set over these fluents. This representation allows us to test whether a precondition holds by checking if the set of fluents that are part of the precondition is a subset of the fluents that are true in the given state. Thus the action $a_i$ is executable in a state $s_k$, if $\text{pre}(a_i) \subseteq s_k$.

- adds($a_i$)/dels($a_i$) - The add and delete effects of the action $a_i$ together captures the effects of executing the action in a state. The add effects represent the set of state fluents that will be made true by the action and the delete effects capture the state fluents that will be turned false. Thus executing an action $a_i$ in state $s_j$ results in a state $s_k = (s_j \setminus \mathrm{dels}(a_i)) \cup \mathrm{adds}(a_i)$

- *I* - The initial state from which the agent starts.

- *G* - The goal that the agent is trying to achieve. Usually the goal is considered to be a partially specified state. That is, the agent is particularly interested in achieving specific values for certain state fluents and that the values of other state fluents do not matter. Thus any state where the specified goal fluent values are met are considered to be valid goal states.

- *C* - The cost function ($C : A \to \mathbb{R}_{>0}$) that specifies the cost of executing a particular action.

Given such a planning model, the solution takes the form of a plan, which is a sequence of actions. A plan $\pi = \langle a_1, ..., a_k \rangle$ is said to be a valid plan for a planning model $\mathcal{M}$, if executing the sequence of the plan results in a state that satisfies the goal. Given the cost function, we can also define the cost of a plan, $C(\pi) = \sum_{a_i \in \pi} C(a_i)$. A plan, $\pi$, is said to be optimal if there exists no valid plan that costs less than $\pi$. We will use the term behavior to refer to the observed state action sequence generated by the execution of a given plan. For models where all actions have unit cost, we will generally skip $C$ from the model definition. We will use the superscript '$*$' to refer to optimal plans. Further, we will use the modifier '$\bar{\phantom{x}}$' to refer to prefixes of a plan and '+' operator to refer to concatenation of action sequences. In cases where we are comparing the cost of the same plan over different models, we will overload the cost functions $C$ to take the model as an argument, i.e., we will use $C(\pi, \mathcal{M}_1)$ to denote the cost of plan $\pi$ in the model $\mathcal{M}_1$, while $C(\pi, \mathcal{M}_2)$ denotes its cost in $\mathcal{M}_1$. Additionally, we will use the notation $C^*_{\mathcal{M}}$ to denote the cost of an optimal plan in the model $\mathcal{M}$.

Since most of the discussion in this book will rely on reasoning about the properties of such plan in different models, we will use a transition function $\delta$ to capture the effect of executing the plan under a given model, such that the execution of an action $a_i$ in state $s_j$ for a model $\mathcal{M}$, where $\delta(s_j, a, \mathcal{M})$, gives the state that results from the execution of action $a$ in accordance with the model $\mathcal{M}$.

**Incomplete Models**  In this book, we will also be dealing with scenarios where the model may not be completely known. In particular, we will consider cases where the specification may be incomplete insofar as we do not know every part of the model with absolute certainty, instead we will allow information on parts of the model that may be possible. To represent such models, we can follow the conventions of an annotated model, wherein the model definition is quite similar to STRIPS one, except that each action $a_i$ is now defined as $a_i = \langle \mathrm{pre}(a_i), \mathrm{poss\_prec}_{a_i}, \mathrm{adds}(a_i), \widetilde{\mathrm{pre}}(a_i), \widetilde{\mathrm{dels}}(a_i), \widetilde{\mathrm{dels}}(a_i) \rangle$, where '$\widetilde{\mathrm{pre}}$', '$\widetilde{\mathrm{adds}}$' and '$\widetilde{\mathrm{dels}}$' represent the possible preconditions, adds and deletes of an action. If a fluent $f$ is part of such a possible list, say a possible precondition, then we are, in effect, saying that we need to consider two possible versions of action $a_i$; one where it has a precondition $f$ and one where it does not. If the model in total contains $k$ possible model

components (including preconditions and effects), then it is in fact representing $2^k$ possible models (this set of possible models are sometimes referred to as the *completion set* of an annotated model).

**Sensor Models**  Most of the settings discussed in the book contain multiple agents and require one of the agents to observe and make sense of the actions performed by the other. We will refer to the former agent as the observer and the other as the actor (though these roles need not be fixed in a given problem). In many cases, the observer may not have perfect visibility of the actor's activities. In particular, we will consider cases where the observer's sensor model may not be able to distinguish between multiple different activities performed by the actor. To capture cases like this, we will use the tuple, $Obs = \langle \Omega, O \rangle$ to specify the sensor model of the observer, where the elements correspond to:

- $\Omega$ - A set of observation tokens that are distinguishable by the observer. We will use $\omega$ to denote an observation token, and $\langle \omega \rangle$ to denote a sequence of tokens.

- $O$ - An observation function ($O : A \times S \rightarrow \Omega$) maps an action performed and the next state reached by the actor to an observation token in $\Omega$. This function captures any limitations present in the observer's sensor model. If the agent cannot distinguish between multiple different activities, we say that the agent has partial observability. Given a sequence of tokens, $\langle \omega \rangle$, a plan $\pi$ is consistent with this sequence if and only if the observation at any step could have been generated by the corresponding plan step (denoted as $\langle \omega \rangle \models \pi$).

**Models in Play in Human-Aware Planning Problems**  Throughout most of this book, we will use the tuple $\mathcal{M}^R = \langle F^R, A^R, I^R, G^R, C^R \rangle$ to capture the robot's planning model that it uses to plan its actions, $\mathcal{M}^H = \langle F^H, A^H, I^H, G^H, C^H \rangle$ the model human may use to capture their own capabilities and plan their behavior and $\mathcal{M}_h^R = \langle F_h^R, A_h^R, I_h^R, G_h^R, C_h^R \rangle$ the mental model the human maintains of the robot. In cases where we want to differentiate between the different definitions of the same action under different models, we will use the specific model name in the superscript to differentiate them. For example, we will refer to the definition of action $a_i$ in the robot's original model as $a_i^{\mathcal{M}^R}$, while the human's expectation of this model will be captured as $a_i^{\mathcal{M}_h^R}$. In cases where the human is maintaining an explicit set of possible robot models then we will use $\mathbb{M}_h^R$ to represent this set. In this book, we will be using the term robot to refer to the autonomous agent that will be acting in the world and interacting with the human. Though, by no means are the methods discussed in the book are limited to physically embodied agents. In fact, in Chapter 10, we will see many examples of software agents capable of using the same methods. Much of the discussions in this book will be focused on cases where the human teammate is merely an observer, and thus in terms of human models we will be focusing on $\mathcal{M}_h^R$. For the sensor models, since the human will be the observer, we will denote the human sensor model by the tuple, $Obs_r^H = \langle \Omega_r^H, O_r^H \rangle$.

**(a)** Plan legibility / transparency.



**(b)** Plan explicability.



**(c)** Plan predictability.

**Fig. 2.1:** A simple illustration of the differences between plan explicability, legibility and predictability. In this Gridworld, the robot can travel across cells, but cannot go backwards. Figure 2.1a illustrates a legible plan (green) in the presence of 3 possible goals of the robot, marked with **?**s. The red plan is not legible since all three goals are likely in its initial stages. Figure 2.1b illustrates an explicable plan (green) which goes straight to the goal **G** as we would expect. The red plan may be more favorable to the robot due to its internal constraints (the arm sticking out might hit the wall), but is inexplicable (i.e. sub-optimal) in the observer's model. Finally, Figure 2.1c illustrates a predictable plan (green) since there is only one possible plan after it performs the first action. The red plans fail to disambiguate among two possible completions of the plan. Note that all the plans shown in Figure 2.1c are explicable (optimal in the observer's model) but only one of them is predictable – i.e. explicable plans may not be predictable. Similarly, in Figure 2.1b, the red plan is predictable after the first action (even though not optimal, since there is only one likely completion) but not explicable – i.e. predictable plans may not be explicable. Without a prefix in Figure 2.1b, the green plan is the only predictable plan.

## 2.2 Modes of Interpretable Behavior

The term interpretable behavior has been used to cover a wide variety of behaviors. However, one unifying thread that runs throughout the different views is the fact that they are all strategies designed to handle the asymmetry between what the human knows about the robot (in terms of its model or the plan it is following) and the robot's model or plans. A robot's actions may be uninterpretable when it does not conform to the expectations or predictions engendered by the human's model of the agent. In this chapter, we will look at three notions of interpretability; namely ***explicability, legibility, and predictability***. Each of these interpretability types captures a different aspect of the asymmetry. For each type, we will define a scoring function, that maps a given behavior and the human's beliefs about the robot to a score. Thus a robot that desires to exhibit a certain type of interpretable behavior can do so by selecting behaviors that maximize the given score. Figure 2.1, present some sample behavior illustrating the various interpretability measures.

### 2.2.1 Explicability

The notion of explicability is related to the human's understanding of the robot's model and how well it explains or aligns with the behavior generated by the robot.

*Explicability measures how close a plan is to the expectations of the observer.*

To formally define explicability, we need to define the notion of distance between the robot's plan and that expected by the human observer. Such distance measure could be based on a number of different aspects, including the cost of the given plan or even on syntactic features like actions used. For now, we will assume, we have access to a function $\mathcal{D}$ that takes two plans and a model and returns a positive number that reflects the distance between the plans. From the semantics of explicability, the farther the plan is from one of the 'expected' plans, the more unexpected it is to the observer. Thus, the **explicability score** for a robot's plan can be written as,

$$E(\pi, \mathcal{M}_h^R) \propto \max_{\pi' \in \Pi^{\mathcal{M}_h^R}} -1 * \mathcal{D}(\pi, \pi', \mathcal{M}_h^R)$$

where, $\Pi^{\mathcal{M}_h^R}$ are the set of plans expected by the human. That is, the explicability score of a robot's plan is negatively related to the minimum possible plan distance from the human's set of expected plans. Throughout the book we will mostly leave the exact function mapping distances to explicability score undefined. Instead we will assume that this mapping could be any monotonic increasing function over the term $-1 * \mathcal{D}(\pi, \pi', \mathcal{M}_h^R)$. Thus the objective of explicable planning becomes,

Find: $\pi$

$$\max_{\pi \in \Pi^{\mathcal{M}^R}} E(\pi, \mathcal{M}_h^R)$$

where $\Pi^{\mathcal{M}^R}$ is the set of plans valid for the model $\mathcal{M}^R$. We will be returning to this score multiple times as it represents one of the key desirable properties of a human-aware robot – namely the ability to generate plans that the human can "interpret", insofar that it aligns with her expectations about the robot.

### 2.2.2 Legibility

The robot may be capable of performing multiple different tasks in a given environment. In such environments, the human observer may sometimes have ambiguity over the current task being pursued by the robot. In such cases, the robot can make its behavior interpretable to the human by implicitly communicating its task through its behavior. This brings us to another type of interpretable behavior, namely, legible behavior. In general, legible behavior which allows the robot to convey some information implicitly. In this setting, the inherent assumption is that the human observer maintains a set of possible hypotheses about the robot model but is unaware of the true model. In a nutshell,

*Legibility reduces observer's ambiguity over possible models the robot is using.*

Let $\mathbb{M}_h^R$ be the set of possible models that the human thinks the robot may have. Here the robot's objective is to execute a plan that reduces the human observer's ambiguity over the possible tasks under consideration. Therefore, the **legibility score** of a plan $\pi$ can be defined as follows:

$$\mathcal{L}(\pi, \mathbb{M}_h^R, \mathcal{M}^R) \propto \frac{1}{f_{amb}^{\mathcal{M}^R}(\{\mathcal{M}^j \mid \mathcal{M}^j \in \mathbb{M}_h^R \wedge \delta(I^j, \pi, \mathcal{M}^j) \models G^j\})}$$

Where $f_{amb}^{\mathcal{M}^R}$ is a measure of ambiguity over the set of models consistent with the behavior observed. There are multiple ways such an ambiguity function could be instantiated including, cardinality of the set, the similarity of the remaining models, and in the probabilistic case, the probability that the human would attach to the models that meet some specific criteria (for example, it shares some parameter with the true robot model). Thus the objective of legible planning thus becomes

Find: $\pi$

$$\max_{\pi \in \Pi^{\mathcal{M}^R}} \mathcal{L}(\pi, \mathbb{M}_h^R)$$

In Chapter 4, we will see specific cases of legibility namely, goal legibility and plan legibility, defined in settings where the human has partial observability of the robot's activities. Due to the human's partial observability of the robot's activities, the human may suffer from uncertainty over the robot's true goal, given a set of candidate goals (goal legibility) or robot's true plan given a set of possible plans towards a goal (plan legibility). Moreover, for scenarios with partial observability, we will be defining legibility over observation sequence generated by the plan as opposed to over the original plans.

### 2.2.3  Predictability

In the case of explicability, as long as the robot's behavior conforms to one of the human's expected plans, the behavior is explicable to the human. However, predictability goes beyond explicability, in that, the behavior has to reduce the number of possible plan completions in the human's mental model given the robot's task. In a nutshell,

> *Plan predictability reduces ambiguity over possible plan completions, given a plan prefix.*

That is, a predictable behavior is one that allows the human observer to guess or anticipate the robot's actions towards its goal. We can define the **predictability score**, i.e., predictability score of a plan prefix, $\bar{\pi}$, given the human's mental model, $\mathcal{M}_h^R$, as follows:

$$\mathcal{P}(\bar{\pi}, \mathcal{M}_h^R) \propto \frac{1}{\mid \{\tilde{\pi} \mid \tilde{\pi} \in \text{compl}_{\mathcal{M}_h^R}(\bar{\pi})\} \mid}$$

Here, $\text{compl}_{\mathcal{M}}(\cdot)$ denotes the set of possible completions with respect to model $\mathcal{M}$, i.e., $\forall \pi \in \{\text{compl}_{\mathcal{M}}(\cdot)\}, \delta(I, \pi, \mathcal{M}) \models G$. The predictability score is higher with less number of possible completions for a given prefix in the human's mental model. Further, the objective in the problem of predictable planning is to select an action that reduces the number of possible completions given the prefix and the robot's goal. Note that, this measure is being defined in an online setting, that is the robot has executed a partial prefix and is in the process of completing the execution in a predictable fashion. Therefore, the problem of predictable planning with respect to a given prefix $\bar{\pi}$ is defined as follows:

$$\text{Find: } a$$

$$\max_{a \in A^R} \ \mathcal{P}(\bar{\pi} + a, \mathcal{M}_h^R)$$

While predictability is a popular measure for interpretability, in this particular book we will not be spending much time investigating this particular measure.

## 2.3 Communication to Improve Interpretability

One of the core reasons for the uninterpretability of any agent behavior is the asymmetry between the agent's knowledge and what the human observer knows about the agent. Explicability and legibility stem from the human's misunderstanding or lack of knowledge about the agent's model and predictability deals with the human's ignorance about the specific plan being followed by the agent. This means that in addition to adjusting the agent behavior, another strategy the agent could adopt is to inform the human about its model or the current plan.

### 2.3.1 Communicating Model Information

Among the communication strategies, we will focus primarily on communicating model information. Such communication requires the ability to decompose the complete model into meaningful components and be able to reason about the effect of receiving information about that component in the human model. We will do this by assuming that the planning models can be represented a set of model features or parameters that can be meaningfully communicated to the human. STRIPS like models gives us a natural parameterization scheme of the form $\Gamma : \mathcal{M} \mapsto s$ represents any planning problem $\mathcal{M} = \langle F, A, I, G, C \rangle$ as a state $s \subseteq F$ as follows –

$$\tau(f) = \begin{cases} \textit{init-has-f} & \text{if } f \in I, \\ \textit{goal-has-f} & \text{if } f \in G, \\ \textit{a-has-precondition-f} & \text{if } f \in pre(a), \ a \in A \\ \textit{a-has-add-effect-f} & \text{if } f \in adds(a), \ a \in A \\ \textit{a-has-del-effect-f} & \text{if } f \in dels(a), \ a \in A \\ \textit{a-has-cost-f} & \text{if } f = C(a), \ a \in A \end{cases}$$

$$\Gamma(\mathcal{M}) = \{\tau(f) \mid \forall f \in I \cup G \cup$$
$$\bigcup_{a \in A} \{f' \mid \forall f' \in \{C(a)\} \cup \text{pre}(a) \cup \text{adds}(a) \cup \text{dels}(a)\}\}$$

Now the agent can communicate to the user in terms of these model features, wherein they can assert that certain features are either part of the model or not. If $\mu$ is a set such messages, then let $\mathcal{M}_h^R + \mu$ be the updated human model that incorporates these messages. If $I$ is an interpretability score (say explicability or legibility) and $\pi$ the current plan, then the goal of the communication becomes to identify a set of messages $\mu$ such that

$$I(\pi, \mathcal{M}_h^R + \mu) > I(\pi, \mathcal{M}_h^R)$$

The discussion in Chapter 5 focuses on the cases where the robot's plan is fixed and it tries to identify the set of messages that best communicate this plan. There may be cases where the most desirable approach for the agent is to not only communicate but also adapt its behavior to reduce uninterpretability and communication overhead. We will refer to such approaches as balanced planning and will dive deeper into such methods in Chapter 6.

Apart from the knowledge asymmetry, another reason for uninterpretability could be the difference in the computational capabilities of the human and the agent. We could also use communication in these cases to help improve the human's ability to comprehend the plan. While we won't look at methods for generating such explanatory messages in detail, we will briefly discuss them in Chapter 11.

## 2.4   Other Considerations in Interpretable Planning

In this section, we will look at some of the other factors that have been studied in the context of interpretable behavior generation and human-aware planning in general.

**Modes of Operation**   One additional factor we can take into account while describing interpretable behavior is what part of the plan is being analyzed by the human observer. Are they watching the plan being executed on the fly or are they trying to analyze the entire plan? We will refer to the former as online setting and the latter as offline or post-hoc setting. In the previous sections, while we defined predictability in an online setting both legibility and explicability were defined for an offline setting. This is not to imply

that legibility and explicability are not useful in online settings (though the usefulness of offline predictability is debatable in fully observable settings, we will touch on the topic in Chapter 4). The online formulations are useful when the robot's reward/cost functions are tied to the interpretability score at every execution time step. For online variations of the scores, the score for each prefix may be defined in regard to the expected completions of the prefix. A related setting that does not fit into either paradigms, is one where the human is interested in interpretable behavior after the execution of first $k$ steps by the agent, we will refer to such settings as *k-step interpretability* settings. These settings are not the same as the online planning setting, as they allow for cases where the human would not pose any penalty for uninterpretable behavior for the first $k - 1$ steps.

**Motion vs Task Planning**   The notions of interpretable behaviors have been studied in both motion planning and task planning. From the algorithmic perspective, this is simply differentiated in usual terms – e.g. continuous versus discrete state variables. However, the notion of interpretability in task planning engenders additional challenges. This is because a reasonable human's mental model of the robot for motion planning can be assumed to be one that prefers straight-line plans and thus need not be modeled explicitly (or needs to be acquired or learned). For task planning in general, this is less straightforward. In Chapters 3 and 6, we will see how we could learn such human models.

**Computational Capabilities of the Human**   One aspect of decision-making we haven't modeled explicitly here is the computational/inferential capabilities of the robot and the human. The robot's computation capabilities determine the kind of behaviors they can exhibit and the human's computational capabilities determine the kind of behavior they expect from the robot. Unfortunately, most current works in interpretable behavior generation generally tend to assume they are dealing with perfectly rational humans and agents (with few exceptions such as the user of noisy-rational models). Developing and using more sophisticated models of bounded rationality that better simulate the human and the robot's reasoning capabilities remain an exciting future direction for interpretable behavior generation.

**Behavior versus Plans**   Our discussion has mostly been confined to analysis of behaviors – i.e. one particular observed instantiation of a plan or policy. In particular, a plan – which can be seen as a set of constraints on behavior – engenders a *candidate set* of behaviors some of which may have certain interpretable properties while others may not. However, this also means that an algorithm that generates plan with a specific interpretable property, can also do so for a particular behavior it models, since in the worst case a behavior is also a plan that has a singular candidate completion. A general treatment of a plan can be very useful in the offline setting – e.g. in decision-support where human decision-makers are deliberating over possible plans with the support from an automated planner. Unfortunately, interpretability of such plans has received very little attention beyond explanation generation.

## 2.5  Generalizing Interpretability Measures

In this chapter and in the rest of the book, we focus one three interpretability measures that were identified by previous works as capturing some desirable behaviors in human-robot interaction settings. Note that these are specific instances of larger behavioral patterns the robot could exhibit in such scenarios. Another possible way to categorize agent behaviors would be to organize them based on how they influence or use the human mental models. In the case of settings with humans as observers, this gives rise to two broad categories, namely;

**Model-Communication Behaviors:** Model-communication involves molding the human's mental models through implicit (i.e tacit) or explicit communication, to allow the agent to achieve their objectives. Examples of model-communication behaviors include legible behavior generation where the agent is trying to implicitly communicate some model information and the model-reconciliation explanation where the agent is directly communicating some model information so that the robot behavior appears more explicable.

**Model-Following Behaviors:** This strategy involves taking the current models of the human and generating behavior that conforms to current human expectations or exploits it in the robot's favor. The examples of this strategy we have seen so far include explicability and predictability.

One could see the agent engaging cyclically in model-communication and model-following behaviors, wherein the agent may choose to mold the user's expectations to a point where its behavior may be better received by the observer. We could go one step further and get away from behavior categorization altogether, and simply define a single human-aware planning problem that generates these individual interpretable behaviors in different scenarios. In fact, Sreedharan et al. [2021b] defines a generalized human-aware planning problem as follows

**Definition 1.** *A **Generalized human-aware planning problem or G-HAP** is a tuple $\Pi_{\mathcal{H}} = \langle \mathcal{M}^R, \mathbb{M}_h^R, P_h^0, P_\ell, C_{\mathcal{H}} \rangle$, where $\mathcal{M}^R$ as always is the robot model, $\mathbb{M}_h^R$ is the set of models human maintains of the robot, $P_h^0$ is the human's initial prior over the models in the hypothesis set $\mathbb{M}_h^R$, $P_\ell$ are the set of inference models the human associates with each possible model and $C_{\mathcal{H}}$ is a generalized cost function that depends on the exact objective of the agent.*

Here $\mathbb{M}_h^R$ also includes special model denoted as $\mathcal{M}^0$. $\mathcal{M}^0$ corresponds to a high entropy model, that is meant to capture the fact that the human's current set of explicit models cannot account for the observed robot behavior. The formulation assumes the human is a Bayesian reasoner who reasoning about the robot behavior given their beliefs about the robot model (and previous observations about the robot behavior).

This problem formulation generates explicable behavior as discussed in this book if (a) $\mathbb{M}_h^R = \{\mathcal{M}_h^R, \mathcal{M}^0\}$, (b) the objective of the planning problem is to select behavior that maximizes the posterior probability associated with model $\mathcal{M}_h^R$ and (c) if, in human's expectation, the likelihood of the robot selecting a plan for $\mathcal{M}_h^R$ (i.e. $\mathcal{M}_h^R$) is inversely

proportional to a distance from a set of expected plans. [Sreedharan et al., 2021b] specifically shows this for cases where the distance measure corresponds to cost difference and the expected plans correspond to optimal plans.

Similarly, the formulation generates legible behaviors when the prior on $\mathcal{M}^0$ is zero and if the objective is to select behavior that maximizes the cumulative posterior over all the models that share the model parameter the robot is trying to communicate (say a goal). Finally, predictable behaviors are generated when $\mathcal{M}^0$ has zero prior, the $\mathbb{M}_h^R = \{\mathcal{M}_h^R, \mathcal{M}^0\}$ and the objective of the planning problem is to choose behavioral prefixes, whose completions have a high posterior likelihood.

## 2.6   Bibliographic Remarks

We would refer the reader to look at Geffner and Bonet [2013] for a detailed discussion of the planning formalism we will be focusing on. In most cases, we can leverage any appropriate planning algorithm (satisficing or optimal depending on the setting) that applies to this formalism. In cases, where we require the use of a specifically designed search or planning algorithm, we will generally include a pseudo code in the chapter.

In terms of the interpretability measures, the original work to introduce explicability was Zhang et al. [2017], though it used a learned labeling model as a proxy to measure the explicability score. The measure was further developed in Kulkarni et al. [2019a] to take into account the user's model and a learned distance function. Parallelly, works Chakraborti et al. [2017a] and Chakraborti et al. [2019c] have looked at applying communication based strategies to address inexplicability but primarily in the context of when the distance function is based on cost and expected plan consists of plans that are optimal in the robot model.

With respect to legibility and predictability, initial works were done in the context of motion-planning problems by Dragan et al. [2013] and Dragan [2017]. Legibility was later extended to task planning settings by MacNally et al. [2018] and further into partially observable domains by Kulkarni et al. [2019b]. A *k-step* version of predictability was introduced in Fisac et al. [2018]. The noisy-rational model as a stand-in for the human's inferential capability was used by Dragan et al. [2013], Dragan [2017], and Fisac et al. [2018]. Even outside the interpretability literature, there is evidence to suggest that such models can be an effective way to capture the inferential process of humans Baker and Tenenbaum [2014]. Chakraborti et al. [2019a] provides a good overview of the entire space of works done in this area and provides a categorization of existing individual works done in this space into these three core scores.

The communication paradigm discussed in the chapter is based on work done in Chakraborti et al. [2017a]. Such model parameterization was later extended to MDPs in Sreedharan et al. [2019a]. Even outside of explanations, the use of representing models has been investigated in the context of other applications that require a search over a space of plausible models, like in the case of Bryce et al. [2016].

CHAPTER **3**

# Explicable Behavior Generation

In chapter 2, among other things we defined the notion of explicability of a plan and laid out an informal description of explicable planning. In this chapter, we will take a closer look at explicability and discuss some practical methods to facilitate explicable planning. This would include discussion on both planning algorithms specifically designed for generating explicable behavior and how one could design/update the task to make the generation of explicable plans easier.

In the earlier chapter, we discussed how the generation of explicable plans requires the robot to simultaneously reason with both models $\mathcal{M}^R$ and $\mathcal{M}_h^R$. This is because the robot needs to select feasible and possibly low-cost plans from $\mathcal{M}^R$, that align with or are close to plans expected by the human (as computed from $\mathcal{M}_h^R$). This implies the robot needs to have access to $\mathcal{M}_h^R$ (at least an approximation of it in some form). An immediate question the reader could ask is if $\mathcal{M}_h^R$ is available to the robot, why is $\mathcal{M}^R$ required in the plan generation process at all? This is necessary since the human's mental model might entail plans that are not even feasible for the robot or are prohibitively expensive, and can thus at best serve as a guide, and not as an oracle, to the explicable plan generation process.

In fact, the critical consideration in the computation of explicable behavior would be the form in which the human's mental model, $\mathcal{M}_h^R$, is accessible to the robot. We present two settings in this chapter, one where the robot directly has access to $\mathcal{M}_h^R$, and another where the robot learns an approximation of it, $\widetilde{\mathcal{M}}_h^R$.

In the first setting, we hypothesize that a prespecified plan distance measure can quantify the distance between the robot plan $\pi_{\mathcal{M}^R}$ and the expected plans $\pi_{\mathcal{M}_h^R}$ from $\mathcal{M}_h^R$. There are many scenarios in which the system could have direct access to the human's mental model. In domains such as household or factory floor scenarios, there is generally a clear expectation of how a task should be performed. In such cases, $\mathcal{M}_h^R$ can be constructed following the norms or protocols that are relevant to that domain. Most deployed products make use of inbuilt models of user expectations in some form or the other. In domains like urban search and rescue, the human and the robot may start with the same model and they may diverge over time as the robot acquires more information owing to the use of more advanced sensors. With $\mathcal{M}_h^R$ in place, we can then use it along with $\mathcal{M}^R$ to generate plans that are closer to the plans expected by the human. Section 3.2 presents a specific heuristic search method for generating such plans and the basic approach is illustrated in Figure 3.1.

In the second setting, we learn an approximation of the human's mental model represented as $\widetilde{\mathcal{M}}_h^R$. This model is learned based on an underlying assumption that humans tend to associate tasks/sub-goals with actions in a given plan. Thus the approximate model, $\widetilde{\mathcal{M}}_h^R$, takes the form of a labeling function that maps plan steps to specific labels. This learned model is then used as a heuristic to generate explicable plans. This approach is illustrated in Figure 3.2 and is discussed in Section 3.3.

**Fig. 3.1:** Schematic diagram of the model based explicable planning.

Finally, we will conclude this chapter by looking at methods that can help improve the feasibility of these behaviors. The plan expected by the human observer may be prohibitively expensive and/or infeasible for the robot. In such cases, it may be wise to redesign the environment itself in which the human and the robot are operating, especially when tasks need to be performed repeatedly in the environment. In Section 3.4, we discuss how the environment design techniques can be leveraged to improve the explicability of the robot's behaviors, and also explore the longitudinal impact on explicable behaviors in a repeated interaction setting.

## 3.1   Explicable Planning Problem

As a reminder from Chapter 2, the problem of explicable planning arises when the robot model differs from the human's expectation of it. Let $\pi_{\mathcal{M}^R}$ be the robot's plan solution to the planning problem, $\mathcal{M}^R = \langle F, A^R, I^R, G^R, C^R \rangle$; whereas, $\Pi_{\mathcal{M}_h^R}$ be the set of plans that the human expects given her mental model of robot's model, $\mathcal{M}_h^R = \langle F, A_h^R, I_h^R, G_h^R, C_h^R \rangle$. The differences in the human's mental model can lead to different plan solutions.

Thus explicable planning consists of the robot coming up with plans that are close to what the human expects (measured using a distance function $\mathcal{D}$).

**Definition 2.** *The **explicable planning problem** is defined as a tuple $\mathcal{P}_{Exp} = \langle \mathcal{M}^R, \mathcal{M}_h^R, \mathcal{D} \rangle$, where, $\mathcal{D}$ is the distance function that the human uses to compute the distance between*

**Fig. 3.2:** Schematic diagram of model-free explicable planning: A learned labeling model is used as a proxy for the actual human model. This model is used to compute the heuristic for explicable plan generation.

*her expected plan and the robot's plan.*

In general, explicable planning is a multi-objective planning problem where the robot is trying to choose a behavior that tries to both maximize its explicability score (i.e. minimize the distance to the expected plan) and minimize its cost. To keep the planning formulation simple, we will assume that we can reduce the planning objective to a linear combination of the cost of the plan in the robot model and the distance of the plan from the expected human plan. Thus the solution to an explicable planning problem is an explicable plan that achieves the goal and minimizes the plan distances while also minimizing the cost of the plan.

**Definition 3.** *A **maximally explicable plan** is a plan, $\pi^*_{\mathcal{M}^R}$, starting at $I^R$ that achieves the goal $G^R$, such that,* $\underset{\pi_{\mathcal{M}^R}}{\arg\min} \, C(\pi_{\mathcal{M}^R}) + \underset{\pi \in \Pi_{\mathcal{M}^R_h}}{\min} \, f_{\mathcal{I}\mathcal{E}}(\mathcal{D}(\pi_{\mathcal{M}^R}, \pi, \mathcal{M}^R)).$

Where $f_{\mathcal{I}\mathcal{E}}$ is any increasing monotonic function over the distance $\mathcal{D}$ (the actual choice of which may depend on the specific optimization algorithm). We will refer to the term $\min_{\pi \in \Pi_{\mathcal{M}^R_h}} f_{\mathcal{I}\mathcal{E}}(\mathcal{D}(\pi_{\mathcal{M}^R}, \pi, \mathcal{M}^R))$ as the inexplicability score or simply $\mathcal{I}\mathcal{E}$. Readers would note that $\mathcal{I}\mathcal{E}$ is the opposite of the explicability score defined in the previous chapter. In this chapter, we will mostly focus on planning algorithms that are trying to minimize the cost of plans generated, and as such focusing on inexplicability score, makes it easier to integrate the concept of explicabillity into the algorithms as we can treat $\mathcal{I}\mathcal{E}$ as yet another cost term. Though other planning formulations say utility maximization formulation, may benefit from the use of explicability score.

## 3.2 Model-Based Explicable Planning

We will start by looking at explicable behavior generation in cases where the human's mental model, $\mathcal{M}^R_h$, and the distance function, $\mathcal{D}$ are known beforehand. Here we assume that the set of plans expected by the human correspond to the optimal plans in $\mathcal{M}^R_h$, i.e., $\Pi_{\mathcal{M}^R_h} = \{\pi | C^R_h(\pi) = C^*_{\mathcal{M}^R_h}\}$. For each robot plan, the minimum distance (per the given distance function, $\mathcal{D}$) is computed with respect to this set of expected plans.

### 3.2.1 Plan Generation through Reconciliation Search

Our primary objective here is to use the distance to the expected plans in $\mathcal{M}^R_h$ as a heuristic to come up with plans that are both valid for the model $\mathcal{M}^R$ and explicable. One obvious idea we could try is to first use $\mathcal{M}^R$ to come up with the expected plan set $\Pi^R_h$ and then during planning using the model $\mathcal{M}^R$ (via a heuristic progression search), expand only the plan prefixes which have the minimal distance to one of the expected plans in the set $\Pi^R_h$ (we will assume that the inexplicability score here is equal to the distance and use the terms inexplicability score and distance interchangeably).

Unfortunately, this won't work in practice due to the non-monotonicity of the inexplicability score. As a partial plan grows, each new action may contribute either positively

---

**Algorithm 1** Reconciliation Search

---

**Input:** $\mathcal{P}_{Exp} = \langle \mathcal{M}^R, \mathcal{M}_h^R, \mathcal{D} \rangle$ and max_cost
**Output:** $\Pi_{Exp}$

 1: $\Pi_{Exp} \leftarrow \emptyset$     {Explicable plan set}
 2: $open \leftarrow \emptyset$     {Open list }
 3: $closed \leftarrow \emptyset$     {Closed list }
 4: $open.\text{insert}(I^R, 0, \text{inf})$
 5: **while** $open \neq \emptyset$ **do**
 6:    $n \leftarrow open.\text{remove}()$     {Node with highest $h(\cdot)$ }
 7:    **if** $n \models G^R$ **then**
 8:       $\Pi_{Exp}.\text{insert}(\pi \text{ s.t. } \delta_{\mathcal{M}^R}(I^R, \pi) \models n)$
 9:    $closed.\text{insert}(n)$
10:    **for** each $v \in successors(n)$ **do**
11:       **if** $v \notin closed$ **then**
12:          **if** $g(n) + cost(n, v) \leq \text{max\_cost}$ **then**
13:             $open.\text{insert}(v, g(n) + cost(n, v), h(v))$
14:       **else**
15:          **if** $h(n) < h(v)$ **then**
16:             $closed.\text{remove}(v)$
17:             $open.\text{insert}(v, g(n) + cost(n, v), h(v))$
18:
19: **return** $\Pi_{Exp}$

---

or negatively to the distance, thus making the final inexplicability score non-monotonic. Consider that the goal of an autonomous car is to park itself in a parking spot on its left side. The car takes the left turn, parks, and turns on its left indicator. Here the turning on of the left tail light after having parked is an inexplicable action. The first two actions are explicable to the human drivers and contribute positively to the overall explicability of the plan but the last action has a negative impact.

Due to the non-monotonic nature of distance, we cannot stop the search process after finding the first solution. Consider the following case, if $d_1$ is the distance of the current prefix, then a node may exist in the open list (set of unexpanded nodes) whose distance is greater than $d_1$, which when expanded may result in a solution plan with distance less than $d_1$. A greedy method that expands a node with the lowest distance score of the corresponding prefix at each step is not guaranteed to find an optimal explicable plan (one of the plans with the lowest inexplicability score) as its first solution. One possible solution that has been studied in the literature has been to use a cost-bounded anytime greedy search algorithm called *reconciliation search* that generates all the valid candidate plans up to a given cost bound, and then progressively searches for plans with lower inexplicability scores. The value of the heuristic *h(v)* in a particular state $v$ encountered during the search is based entirely on the distance of the agent plan prefix $\pi_{\mathcal{M}^R}$ up to that state, h(v) $= \mathcal{D}(\pi_{\mathcal{M}^R}, \pi'_{\mathcal{M}_h^R}, \mathcal{M}_h^R)$ s.t. $\delta_{\mathcal{M}^R}(I^R, \pi_{\mathcal{M}^R}) \models v$ and $\delta_{\mathcal{M}_h^R}(I^R, \pi'_{\mathcal{M}_h^R}) \models v$.

The approach is described in detail in Algorithm 1. At each iteration of the algorithm, the plan prefix of the agent model is compared with the explicable trace $\pi'_{\mathcal{M}_h^R}$ (these are the plans generated using $\mathcal{M}_h^R$ up to the current state in the search process) for the given problem. There are few choices one could consider in creating such prefixes including,

considering prefixes from among the expected plans, optimal or even valid plans to the state in the human model. The search algorithm then makes a locally optimal choice of states and continues till a solution is found. The search is not stopped after generating the first solution, but instead, it restarts from the initial state. The search continues to find all the valid candidate solutions within the given cost bound or until the state space is completely explored.

### 3.2.2   Possible Distance Functions

In the above discussion, we mostly assumed we are given some specific distance function $\mathcal{D}$. A simple candidate for distance function could be the cost difference, but that presupposes that we know exactly the cost function of the human and that the explicability consideration of the human observer is strictly limited to cost difference. In general, it may be possible that the human is using more complex distance functions and we may want to learn the distance function directly from the human. In this case, we could consider learning a distance function that is a combination of various individual distance functions that can be directly computed from the human's model. Some of the candidate distance functions we could consider include measures like action-distance, causal-link distance, and state distance.

**Action distance**   We denote the set of unique actions in a plan $\pi$ as $A(\pi) = \{a \mid a \in \pi\}$. Given the action sets $A(\pi_{\mathcal{M}^R})$ and $A(\pi^*_{\mathcal{M}^R_h})$ of two plans $\pi_{\mathcal{M}^R}$ and $\pi^*_{\mathcal{M}^R_h}$ respectively, the action distance is, $\mathcal{D}_A(\pi_{\mathcal{M}^R}, \pi^*_{\mathcal{M}^R_h}) = 1 - \frac{|A(\pi_{\mathcal{M}^R}) \cap A(\pi^*_{\mathcal{M}^R_h})|}{|A(\pi_{\mathcal{M}^R}) \cup A(\pi^*_{\mathcal{M}^R_h})|}$. Here, two plans are similar (and hence their distance measure is smaller) if they contain same actions. Note that it does not consider the ordering of actions.

**Causal link distance**   A causal link represents a tuple of the form $\langle a_i, p_i, a_{i+1} \rangle$, where $p_i$ is a predicate variable that is produced as an effect of action $a_i$ and used as a precondition for the next action $a_{i+1}$. The causal link distance measure is represented using the causal link sets $Cl(\pi_{\mathcal{M}^R})$ and $Cl(\pi^*_{\mathcal{M}^R_h})$, $\mathcal{D}_C(\pi_{\mathcal{M}^R}, \pi^*_{\mathcal{M}^R_h}) = 1 - \frac{|Cl(\pi_{\mathcal{M}^R}) \cap Cl(\pi^*_{\mathcal{M}^R_h})|}{|Cl(\pi_{\mathcal{M}^R}) \cup Cl(\pi^*_{\mathcal{M}^R_h})|}$.

**State sequence distance**   This distance measure finds the difference between sequences of the states. Given two state sequences $(s_0^R, \ldots, s_n^R)$ and $(s_0^H, \ldots, s_{n'}^H)$ for $\pi_{\mathcal{M}^R}$ and $\pi^*_{\mathcal{M}^R_h}$ respectively, where $n \geq n'$ are the lengths of the plans, the state sequence distance is, $\mathcal{D}_S(\pi_{\mathcal{M}^R}, \pi^*_{\mathcal{M}^R_h}) = \frac{1}{n}\left[ \sum_{k=1}^{n'} d(s_k^R, s_k^H) + n - n' \right]$ , where $d(s_k^R, s_k^H) = 1 - \frac{|s_k^R \cap s_k^H|}{|s_k^R \cup s_k^H|}$ represents the distance between two states (where $s_k^R$ is overloaded to denote the set of predicate variables in state $s_k^R$). The first term measures the normalized difference between states up to the end of the shortest plan, while the second term, in the absence of a state to compare to, assigns maximum difference possible.

As mentioned earlier, the actual distance used by the human observer could be some

combination of such constituent distance function, so we will need to consider composite distance functions.

**Definition 4.** *A **composite distance**, $\mathcal{D}_{comp}$ is a distance between pair of two plans $\langle \pi_{\mathcal{M}^R}, \pi_{\mathcal{M}_h^R} \rangle$, such that, $\mathcal{D}_{comp}(\pi_{\mathcal{M}^R}, \pi_{\mathcal{M}_h^R}, \mathcal{M}_h^R) = ||\mathcal{D}_A(\pi_{\mathcal{M}^R}, \pi_{\mathcal{M}_h^R}) + \mathcal{D}_C(\pi_{\mathcal{M}^R}, \pi_{\mathcal{M}_h^R}) + \mathcal{D}_S(\pi_{\mathcal{M}^R}, \pi_{\mathcal{M}_h^R})||$.*

But for each robot plan we want to find the minimum distance with respect to the set of human's expected plans. We say a distance minimizing plan in the set of the expected plans is defined as follows:

**Definition 5.** *A **distance minimizing plan**, $\pi^*_{\mathcal{M}_h^R}$, is a plan in $\Pi_{\mathcal{M}_h^R}$, such that for a robot plan, $\pi_{\mathcal{M}^R}$, the composite distance is minimized, i.e., $\forall \pi \in \Pi_{\mathcal{M}_h^R} \; \mathcal{D}_{comp}(\pi_{\mathcal{M}^R, \pi, \mathcal{M}_h^R}) \geq \mathcal{D}_{comp}(\pi_{\mathcal{M}^R, \pi^*_{\mathcal{M}_h^R}, \mathcal{M}_h^R})$.*

Our overall objective is to learn an inexplicability score from the individual distances. One way to do it would be to first collect scores human participants and then try to learn a mapping to human specified scores from the individual plan distances between a robot plan and corresponding distance minimizing plan in the set of expected plans. To that end, we define a explicability feature vector as follows:

**Definition 6.** *An **explicability feature vector**, $\mathbb{F}$, is a three-dimensional vector, which is given with respect to a distance minimizing plan pair, $\langle \pi_{\mathcal{M}^R}, \pi^*_{\mathcal{M}_h^R} \rangle$, such that, $\mathbb{F} = \langle \mathcal{D}(\pi_{\mathcal{M}^R}, \pi^*_{\mathcal{M}_h^R}), \mathcal{D}_C(\pi_{\mathcal{M}^R}, \pi^*_{\mathcal{M}_h^R}), \mathcal{D}_S(\pi_{\mathcal{M}^R}, \pi^*_{\mathcal{M}_h^R}) \rangle^T$.*

This allows us to learn an explicability distance function, $\hat{\mathcal{D}}_{Exp}(\pi_{\mathcal{M}^R} \, / \, \pi^*_{\mathcal{M}_h^R})$, which is essentially a regression function, $f$, that fits the three plan distances to the total plan scores, with $b$ as the parameter vector, and $\mathbb{F}$ as the explicability feature vector, such that, $\hat{\mathcal{D}}_{Exp}(\pi_{\mathcal{M}^R} \, / \, \pi^*_{\mathcal{M}_h^R}) \approx f(\mathbb{F}, b)$. Therefore, a regression model can be trained to learn the explicability assessment (total plan scores) of the users by mapping this assessment to the explicability feature vector which consists of plan distances for corresponding plans. Since we are unaware of the distance measure, we can never guess the right plan against which the human may be comparing the current behavior. The distance minimizing plan is a stand-in for this unknown plan. As such any distance function learned from features derived from this plan will be an approximation of the actual distance. Though in many cases such approximation acts as a reasonable guide for behavior generation.

## 3.3   Model-Free Explicable Planning

If the robot does not have access to the human's mental model, then an approximation of it can be learned. This approach shows that it is not necessary to build a full-fledged planning model of the human's mental model, rather it is enough to learn a simple labeling function. Here the underlying hypothesis is that humans tend to associate abstract tasks or sub-goals to actions in a plan. If the human-in-the-loop can associate any domain-specific

label to an action in the plan then the action is assumed to be explicable, otherwise, the action is considered inexplicable. Such a labeling scheme can be learned from training examples annotated by humans. This learned model can then be used as a heuristic function in the planning process. This approach is illustrated in Figure 3.2.

### 3.3.1   Problem Formulation

In this case, since the $\mathcal{M}_h^R$ is not known beforehand, the distance function $\mathcal{D}(\cdot, \cdot, \cdot)$ in Definition 2 is approximated using a learning method. Here the terms $\mathcal{D}(\pi_{\mathcal{M}^R}, \pi_{\mathcal{M}_h^R}, \mathcal{M}_h^R)$ and by extension the inexplicability scores are measured in terms of the labels the human would associate with the plan. Thus the method expects a set of task labels $T = \{T_1, T_2, \ldots, T_M\}$ be provided for each domain. Depending on how the human view the role of the action in the plan, each action may be associated with multiple labels. The set of action labels for explicability is the power set of task labels: $\mathbb{L} = 2^T$. When an action label includes multiple task labels, the action is interpreted as contributing to multiple tasks; when an action label is an empty set, the action is interpreted as inexplicable.

For a given plan $\pi_{\mathcal{M}^R}$, we can define the inexplicability score of the plan as follows $\mathcal{IE}(\pi_{\mathcal{M}^R}, \mathcal{M}_h^R) = f \circ \mathcal{L}^*(\pi_{\mathcal{M}^R})$, where $f$ is a domain-independent function that takes plan labels as input, and $\mathcal{L}^*$ is the labeling scheme of the human for agent plans based on $\mathcal{M}_h^R$. A possible candidate for $f$ could be to compute the ratio between the number of actions with empty action labels and the number of all actions.

### 3.3.2   Learning Approach

Now the question is, for a new plan how do we predict the labels human would associate with it? As mentioned earlier, one could use a training set consisting of plans annotated with labels and train a sequence to sequence model. Here, possible models, one could use include sequential neural networks (including Recurrent Neural Networks [Goodfellow et al., 2016] variants and transformers) and probabilistic models like Hidden Markov Models [Russell and Norvig, 2002] or Conditional Random Fields [Lafferty et al., 2001]. Here each step in the sequence corresponds to a step in the plan execution and the possible features that could be associated with each step include features like the action description, the state fluents obtained after executing each action in the sequence $\langle a_0, ..., a_n \rangle$ from the initial state, etc. One could also potentially include other features like motion level information, like gripper position, orientation, etc. We will use the notation $\mathcal{L}'$ to denote the learned labeling model and to differentiate it from the true labeling model $\mathcal{L}^*$.

### 3.3.3   Plan Generation

Once we have learned a labeling model, we can then use it to drive the behavior generation.

### 3.3.3.1 Plan Selection

The most straightforward method is to perform plan selection on a set of candidate plans which can simply be a set of plans that are within a certain cost bound of the optimal plan. Candidate plans can also be generated to be diverse with respect to various plan distances. For each plan, the features of the actions are extracted as we discussed earlier. Then the trained model $\mathcal{L}'$ can be used to produce the labels for the actions in the plan and the inexplicability score $\mathcal{IE}$ can then be computed as given by the mappings above. The inexplicability score can then be used to choose a plan that is least inexplicable.

---

**Algorithm 2** Synthesis of Explicable Plan

---

**Input:** $\mathcal{M}^R, \mathcal{L}'$
**Output:** $\pi^*_{\mathcal{M}^R}$

1: $open \leftarrow \emptyset$
2: $open.\text{push}(I^R)$
3: **while** $open \neq \emptyset$ **do**
4:     $s \leftarrow open.\text{remove}()$
5:     **if** $G^R$ is reached **then**
6:         **return** $s.\text{plan}$ {the plan that leads to $s$ from $I^R$)}
7:     Compute all possible next states $N$ from $s$
8:     **for** $n \in N$ **do**
9:         Compute the relaxed plan $\pi_{RELAX}$ for $n$
10:        Concatenate $s.\text{plan}$ with $\pi_{RELAX}$ as $\pi^{prime}$ {$s.\text{plan}$ contains plan features and $\pi_{RELAX}$ only contains action descriptions}
11:        Compute and add other relevant features
12:        Compute $L_{\pi^R} = \mathcal{L}'(\pi^{prime})$
13:        Compute $h = f(IE, h_{cost})$ {$f$ is a combination function; $h_{cost}$ is the relaxed planning heuristic}
14:     **if** $h(n^*) < h^*$ {$n^* \in N$ with minimum $h$} **then**
15:         $open.\text{clear}()$
16:         $open.\text{push}(n^*); h^* = h(n^*)$ {$h^*$ is initially MAX}
17:     **else**
18:         Push all $n \in N$ into $open$

---

### 3.3.3.2 Plan Synthesis

A more efficient way is to incorporate these measures as heuristics into the planning process. Algorithm 2 presents a simple modified version of Enforced Hill Climbing [Hoffmann, 2005], that can make use of such learned heuristics. To compute the heuristic value given a planning state, one can use the relaxed planning graph to construct the remaining planning steps. However, since relaxed planning does not ensure a valid plan, one can only use action descriptions as plan features for actions that are beyond the current planning state when estimating the inexplicability measures. These estimates are then combined with the relaxed planning heuristic (which only considers plan cost) to guide the search.

## 3.4   Environment Design for Explicability

The environment in which the robot is operating may not always be conducive to explicable behavior. As a result, making its behavior explicable may be prohibitively expensive for the robot. Additionally, certain behaviors that are explicable with respect to the human's mental model may not be feasible for the robot. Fortunately, in highly structured settings, where the robot is expected to solve repetitive tasks (like in warehouses, factories, restaurants, etc.), it might be feasible to *redesign the environment* in a way that improves explicability of the robot's behavior, given a set of tasks. This brings us to the notion of environment design which involves redesigning the environment to maximize (or minimize) some objective for the robot. Thus, environment design can be used to boost the explicability of the robot's behavior, especially in settings that require solving repetitive tasks and a one-time environment design cost to boost explicable behavior might be preferable over the repetitive cost overhead of explicable behavior borne by the robot.

However, environment design alone may not be a panacea for explicability. For one, the design could be quite expensive, not only in terms of making the required environment changes but also in terms of limiting the capabilities of the robot. Moreover, in many cases, there may not be a single set of design modifications that will work for a given set of tasks. For instance, designing a robot with wheels for efficient navigation on the floor will not optimize the robot's motion up a stairwell. This means, to achieve truly effective synergy with autonomous robots in a shared space, we need a synthesis of environment design and human-aware behavior generation. In this section, we will talk about how we could trade off the one-time (but potentially expensive) design changes, against repetitive costs borne by the robot to exhibit explicable behavior.

### 3.4.0.1   Motivating Example

Consider a restaurant with a robot server (Figure 9.3a). Let $G_1$ and $G_2$ represent the robot's possible goals of serving the two booths: it travels between the kitchen and the two booths. The observers consist of customers at the restaurant. Given the position of the kitchen, the observers may have expectations on the route taken by the robot. However, unbeknownst to the observers, the robot can not traverse between the two tables and can only take the route around the tables. Therefore, the path marked in red is the cheapest path for the robot but the observers expect the robot to take the path marked in green in Figure 9.3a.

In this environment, there is no way for the robot to behave as per the human's expectations. Applying environment design provides us with alternatives. For example, the designer could choose to build two barriers as shown in Figure 3.3b. With these barriers in place, the humans would expect the robot to follow the path highlighted in green. However, whether it is preferable to perform environment modifications or to bear the impact of inexplicable behavior depends on the cost of changing the environment versus the cost of inexplicability caused by the behavior.

(a) Explicable behavior is costlier without design.



(b) Optimal behavior is explicable with design.

**Fig. 3.3:** Use of environment design to improve the explicability of a robot's behavior in a shared environment.

### 3.4.1   Problem Setting

We again consider an explicability problem of the $\mathcal{P}_{Exp} = \langle \mathcal{M}^R, \mathcal{M}^R_h, \mathcal{D}_{\mathcal{M}^R_h} \rangle$. In this subsection we will again consider cases where the human mental model is given upfront. Let the expected set of plans in the human model. $\Pi_{\mathcal{M}^R_h}$, be the set of plans optimal in $\mathcal{M}^R_h$. The robot plans to minimize the inexplicability score in the human's mental model. We will use the notation $\Pi^*_{\mathcal{IE}(\cdot, \mathcal{M}^R_h, \delta_{\mathcal{M}^R_h})}$ (in the absence of the parameter $\pi^R$) to refer to the set of plans in the robot's model with the lowest inexplicability score, and $\mathcal{IE}_{min}(\mathcal{P}_{Exp})$ to represent the lowest inexplicability score associated with the set. Further, let $f_{Exp}$ be the decision function used by the explicable robot: $f_{Exp}(\mathcal{P}_{Exp})$ represents the cheapest plan that minimizes the inexplicability score, i.e. $f_{Exp}(\mathcal{P}_{Exp}) \in \Pi^*_{\mathcal{IE}(\cdot, \mathcal{M}^R_h, \delta_{\mathcal{M}^R_h})}$ and $\neg \exists \pi' : \pi' \in \Pi^*_{\mathcal{IE}(\cdot, \mathcal{M}^R_h, \delta_{\mathcal{M}^R_h})}$ such that $c^R(\pi') < c^R(f_{Exp}(\mathcal{P}_{Exp}))$.

#### 3.4.1.1   Environment Design

Before we delve more into the problem of environment design to boost explicability, lets take a general look at *environment design problems*. An environment design problem takes as input the initial environment configuration along with a set of available modifications and computes a subset of modifications that can be applied to the initial environment to derive a new environment in which a desired objective is optimized.

We consider $\mathcal{M}^{R0} = \langle F^0, A^{R0}, I^{R0}, G^{R0}, C^{R0} \rangle$ as the initial environment and $\rho^R$ as the set of valid configurations of that environment: $\mathcal{M}^{R0} \in \rho^R$. Let $\mathcal{O}$ be some metric that needs to be optimized with environment design, i.e a planning model with lower value for $\mathcal{O}$ is preferred. A design problem is a tuple $\langle \mathcal{M}^{R0}, \Delta, \Lambda^R, C, \mathcal{O} \rangle$ where, $\Delta$ is the set of all modifications, $\Lambda^R : \rho^R \times 2^\Delta \to \rho^R$ is the model transition function that specifies the resulting model after applying a subset of modifications to the existing model, $C : \Delta \to \mathbb{R}$ is the cost function that maps each design choice to its cost. The modifications are independent of each other and their costs are additive. We will overload the notation and use $C$ as the cost function for a subset of modifications as well, i.e. $C(\xi) = \sum_{\xi_i \in \xi} C(\xi)$.

The set of possible modifications could include modifications to the state space, action preconditions, action effects, action costs, initial state and goal. In general, the space of design modifications, which are an input to our system, may also involve modifications to the robot itself (since the robot is part of the environment that is being modified). An optimal solution to a design problem identifies the subset of design modifications, $\xi$, that minimizes the following objective consisting of the cost of modifications and the metric $\mathcal{O}$: $\min \mathcal{O}(\Lambda^R(\mathcal{M}^{R0}, \xi)), \ C(\xi)$.

### 3.4.2   Framework for Design for Explicability

In this framework, we not only discuss the problem of environment design with respect to explicability but also in the context of **(1)** a set of tasks that the robot has to perform in

the environment, and **(2)** over the lifetime of the tasks i.e. the time horizon over which the robot is expected to repeat the execution of the given set of tasks. These considerations add an additional dimension to the environment design problem since the design will have lasting effects on the robot's behavior. In the following, we will first introduce the design problem for a single explicable planning problem, then extend it to a set of explicable planning problems and lastly extend it over a time horizon.

### 3.4.2.1 Design for a Single Explicable Problem

In the design problem for explicability, the inexplicability score becomes the metric that we want to optimize for. That is we want to find an environment design such that the inexplicability score is reduced in the new environment. This problem can be defined as follows:

**Definition 7.** *The **design problem for explicability** is a tuple,*

$$\mathcal{DP}_{Exp} = \langle \mathcal{P}^0_{Exp}, \Delta, \Lambda_{Exp}, C, \mathcal{IE}_{min} \rangle,$$

*where:*

- $\mathcal{P}^0_{Exp} \in \rho_{Exp}$ *is the initial configuration of the explicable planning problem, where $\rho_{Exp}$ represents the set of valid configurations for $\mathcal{P}_{Exp}$.*

- $\Delta$ *is the set of available design modifications. The space of all possible modifications is the power-set $2^\Delta$.*

- $\Lambda_{Exp} : \rho_{Exp} \times 2^\Delta \to \rho_{Exp}$ *is the transition function over the explicable planning problem, which gives an updated problem after applying the modifications.*

- $C$ *is the additive cost associated with each design in $\Delta$.*

- $\mathcal{IE}_{min} : \rho_{Exp} \to \mathbb{R}$ *is the minimum possible inexplicability score in a configuration, i.e. the inexplicability score associated with the most explicable plan.*

With respect to our motivating example in Figure 3.3a, $\mathcal{DP}_{Exp}$ is the problem of designing the environment to improve the robot's explicability given its task of serving every new customer at a booth (say $G_1$) only once. The optimal solution to $\mathcal{DP}_{Exp}$ involves finding a configuration which minimizes the minimum inexplicability score. We also need to take into account an additional optimization metric which is the effect of design modifications on the robot's plan cost. That is, we need to examine to what extent the decrease in inexplicability is coming at the robot's expense. For instance, if you confine the robot to a cage so that it cannot move, its behavior becomes completely and trivially explicable, but the cost of achieving its goals goes to infinity.

**Definition 8.** *An **optimal solution** to $\mathcal{DP}_{Exp}$, is a subset of modifications $\xi^*$ that minimizes the following:*

$$\min \ \mathcal{IE}_{min}( \ \mathcal{P}^*_{Exp}), \ C(\xi^*), c^R(f_{Exp}(\mathcal{P}^*_{Exp})) \tag{3.1}$$

*where $\mathcal{P}^*_{Exp} = \Lambda_{Exp}(\mathcal{P}^0_{Exp}, \xi^*)$ is the final modified explicable planning problem, $\mathcal{IE}_{min}(\cdot)$ represents the minimum possible inexplicability score for a given configuration, $C(\xi^*)$ denotes the cost of the design modifications and $c^R(f_{Exp}(\mathcal{P}^*_{Exp}))$ is the cost of the cheapest most explicable plan in a configuration.*

### 3.4.2.2   Design for Multiple Explicable Problems

We will now show how $\mathcal{DP}_{Exp}$ evolves when there are multiple explicable planning problems in the environment that the robot needs to solve. When there are multiple tasks there may not exist a single set of design modifications that may benefit all the tasks. In such cases, a solution might involve performing design modifications that benefit some subset of the tasks while allowing the robot to act explicably with respect to the remaining set of tasks. Let there be $k$ explicable planning problems, given by the set $\mathbf{P}_{Exp} = \{\langle \mathcal{M}^R(0), \mathcal{M}_h^R(0), \delta_{\mathcal{M}_h^R(0)} \rangle, \ldots, \langle \mathcal{M}^R(k), \mathcal{M}_h^R(k), \delta_{\mathcal{M}_h^R(k)} \rangle\}$, with a categorical probability distribution $\mathbb{P}$ over the problems. We use $\mathcal{P}_{Exp}(i) \in \mathbf{P}_{Exp}$ to denote the $i^{th}$ explicable planning problem. These $k$ explicable problems may differ in terms of their initial state and goal conditions. Now the design problem can be defined as:

$$\mathcal{DP}_{Exp,\mathbb{P}} = \langle \mathbf{P}_{Exp}^0, \mathbb{P}, \Delta, \Lambda_{Exp}, C, \mathcal{IE}_{min,\mathbb{P}} \rangle, \tag{3.2}$$

where $\mathbf{P}_{Exp}^0$, is the set of planning tasks in the initial environment configuration, $\mathcal{IE}_{min,\mathbb{P}}$ is a function that computes the minimum possible inexplicability score in a given environment configuration by taking the expectation over the minimum inexplicability score for each explicable planning problem, i.e., $\mathcal{IE}_{min,\mathbb{P}}(\mathbf{P}_{Exp}) = \mathbb{E}[\mathcal{IE}_{min}(\mathcal{P}_{Exp})]$, where $\mathcal{P}_{Exp} \sim \mathbb{P}$. With respect to our running example, $\mathcal{DP}_{Exp,\mathbb{P}}$ is the problem of designing the environment given the robot's task of serving every new customer only once at either of the booths $(G_1, G_2)$ with probability given by $\mathbb{P}$.

The solution to $\mathcal{DP}_{Exp,\mathbb{P}}$ has to take into account the distribution over the set of explicable planning problems. Therefore the optimal solution is given by:

$$\min \quad \mathcal{IE}_{min,\mathcal{D}}(\mathbf{P}_{Exp}^*), \ C(\xi^*), \ \mathbb{E}[c^R(f_{Exp}(\mathcal{P}_{Exp}^*))] \tag{3.3}$$

where $\mathcal{P}_{Exp}^* \sim \mathbb{P}$. A valid configuration minimizes the minimum possible inexplicability score, which involves 1) expectation over minimum inexplicability scores for each explicable planning problem; 2) the cost of the design modifications (these modifications are applied to each explicable planning problem); and 3) the expectation over the cheapest most explicable plan for each explicable planning problem.

### 3.4.2.3   Longitudinal Impact on Explicable Behavior

The process of applying design modifications to an environment makes more sense if the tasks are going to be performed repeatedly in the presence of a human (i.e. the robot does not have to bear the cost of being explicable repeatedly). This has quite a different temporal characteristic in comparison to that of execution of one-time explicable behavior. For instance, design changes are associated with a one-time cost (i.e. the cost of applying those changes in the environment). On the other hand, if we are relying on the robot to execute explicable plans at the cost of foregoing optimal plans, then it needs to bear this cost multiple times in the presence of a human over the time horizon.

We will use a discrete time formulation where the design problem is associated with a time horizon $\mathcal{T}$. At each time step, one of the $k$ explicable planning problems is chosen.

**Fig. 3.4:** Illustration of longitudinal impact on explicability. *Prob* determines the probability associated with executing each task in $\mathbf{P}_{Exp}$. For each task, the reward is determined by the inexplicability score of that task. The probability of achieving this reward is determined by $\gamma \times$ probability of executing that task. Additionally, with a probability $(1-\gamma)$ the human ignores the inexplicability of a task and the associated reward is given by an inexplicability score of 0.

Now the design problem can be defined as:

$$\mathcal{DP}_{Exp,\mathbb{P},\mathcal{T}} = \langle \mathbf{P}^0_{Exp}, \mathbb{P}, \Delta, \Lambda_{Exp}, C, \mathcal{IE}_{min,\mathbb{P}}, \mathcal{T} \rangle \qquad (3.4)$$

In our running example, $\mathcal{DP}_{Exp,\mathbb{P},\mathcal{T}}$ is the problem of designing the environment given the robot's task of serving the same customer at either of the booths with a distribution $\mathbb{P}$ over a horizon $\mathcal{T}$.

Note than earlier discussions on explicable behavior generation focused on cases where there was only a single interaction between the human and the robot. However, in this section we consider a time horizon, $\mathcal{T} > 1$, over which the robot's interaction with the human may be repeated multiple times for the same task. This means the human's expectations about the task can evolve over time. This may not be a problem if the robot's behavior aligns perfectly with the human's expectations. Although, if the robot's plan for a given task is associated with a non-zero inexplicability score, then the human is likely to be more surprised the very first time she notices the inexplicable behavior than she would be if she noticed the inexplicable behavior subsequent times. As the task is performed over and over, the amount of surprise associated with the inexplicable behavior starts decreasing. In fact, there is a probability that the human may ignore the inexplicability of the robot's behavior after sufficient repetitions of the task. We incorporate this intuition by using discounting.

Figure 3.4 illustrates the Markov reward process to represent the dynamics of this system. Let $(1-\gamma)$ denote the probability that the human will ignore the inexplicability of the robot's plan, i.e, the reward will have inexplicability score 0. $\gamma$ times the probability of executing a task represents the probability that the reward will have the minimum inexplicability score associated with that task. Assuming $\gamma < 1$, the minimum possible

inexplicability score for a set of explicable planning problems is:

$$
\begin{aligned}
f_{\mathcal{T}}(\mathcal{IE}_{min,\mathbb{P}}(\mathbf{P}_{Exp})) &= \mathcal{IE}_{min,\mathbb{P}}(\mathbf{P}_{Exp}) \\
&\quad + \gamma * \mathcal{IE}_{min,\mathbb{P}}(\mathbf{P}_{Exp}) + \ldots + \\
&\quad \gamma^{T-1} * \mathcal{IE}_{min,\mathbb{P}}(\mathbf{P}_{Exp}) \\
f_{\mathcal{T}}(\mathcal{IE}_{min,\mathbb{P}}(\mathbf{P}_{Exp})) &= \frac{1 - \gamma^{\mathcal{T}}}{1 - \gamma} * \mathcal{IE}_{min,\mathbb{P}}(\mathbf{P}_{Exp})
\end{aligned}
\tag{3.5}
$$

Thus the optimal solution to $\mathcal{DP}_{Exp,\mathbb{P},\mathcal{T}}$ is given by:

$$
\begin{aligned}
\min \quad & f_{\mathcal{T}}(\mathcal{IE}_{min,\mathbb{P}}(\mathbf{P}^*_{Exp})), \; C(\xi^*), \\
& \mathbb{E}[c^R(f_{Exp}(\mathcal{P}^*_{Exp}))] * \mathcal{T}
\end{aligned}
\tag{3.6}
$$

where, $\mathcal{P}^*_{Exp} \sim \mathbb{P}$. The optimal solution is a valid configuration that minimizes 1) the minimum possible inexplicability over the set of explicable planning problems given the human's tolerance to inexplicable behavior; 2) one-time cost of the design modifications; and 3) the expectation over the cheapest most explicable plan for each explicable planning problem given a time horizon. Note that, since the design cost is not discounted and we always make the design changes before the task is solved, there is never a reason to delay the design execution to future steps in the horizon. Instead it can be executed before the first time step.

### 3.4.3 Search for Optimal Design

We can find the optimal solution for $\mathcal{DP}_{Exp,\mathbb{P},\mathcal{T}}$, by performing a breadth-first search over the space of environment configurations that are achievable from the initial configuration through the application of the given set of modifications. The performance of the search depends on the number of designs available. By choosing appropriate design strategies, significant scale up can be attained. Each search node is a valid environment configuration and the possible actions are the applicable designs. For simplicity, we convert the multi-objective optimization in Equation 3.1 into a single objective as a linear combination of each term associated with a coefficients $\alpha$, $\beta$, and $\kappa$, respectively. The value of each node is decided by the aforementioned objective function. For each node, it is straightforward to calculate the design modification cost. However, in order to calculate the minimum inexplicability score and the robot's plan cost, we have to generate a plan that minimizes the inexplicability score for each explicable planning problem in that environment configuration. In this setting one could use the method discussed in Section 3.2. Since the distances here are specifically limited to cost differences, we can also use a compilation based method, that compiles the problem of generating the explicable plan to a classical planning problem. We will discuss this compilation in detail in Chapter 6. Essentially, the search has two loops: the outer loop which explores all valid environment configurations, and the inner loop which performs search in a valid environment configuration to find a plan that minimizes the inexplicability score. At the end of the search, the node with best value is chosen, and the corresponding set of design modifications, $\xi^*$, is output.

One way to optimize the search over the space of environment configurations is to only consider the designs that are relevant to the actions in the optimal robot plans ($\Pi^*_{\mathbf{M}^R}$)

and those in the human's expected plans ($\Pi^*_{\mathbf{M}^R_h}$) given the set of tasks. This can be implemented as a pruning strategy that prunes out designs that are not relevant to the actions.

### 3.4.4 Demonstration of Environment Design for Explicability

We use our running example from Figure 9.3a to demonstrate how the design problem evolves. We constructed a domain where the robot had 3 actions: *pick-up* and *put-down* to serve the items on a tray and *move* to navigate between the kitchen and the booths. Some grid cells are blocked due to the tables and the robot cannot pass through these: cell(0, 1) and cell(1, 1). Therefore, the following passages are blocked: cell(0, 0)-cell(0, 1), cell(0, 1)-cell(0, 2), cell(0, 1)-cell(1, 1), cell(1, 0)-cell(1, 1), cell(1, 1)-cell(1, 2), cell(1, 1)-cell(2, 1). We considered 6 designs, each consisting of putting a barrier at one of the 6 passages to indicate the inaccessibility to the human (i.e. the design space has $2^6$ possibilities).

For the following parameters: $\alpha = 1$, $\beta = 30$, $\kappa = 0.25$ and $\gamma = 0.9$, let us consider the problem of identifying design for three settings: (a) single explicable problem for $\mathcal{T} = 1$, (b) multiple explicable problems for $\mathcal{T} = 1$, and (c) multiple explicable problems for $\mathcal{T} = 10$. As mentioned earlier, we will use the set of plans optimal in the human model as the expected plan set and the difference between cost of the plans as the distance function. Let us concretize the three settings as follows, (a) involves serving a new customer at a booth (say $G_1$) only once, (b) involves serving a new customer only once at either of the booths with equal probability and (c) involves serving each customer at most 10 times at either of the booths with equal probability. We found that for settings (a) and (b) no design was chosen. This is because these settings are over a single time step and the cost of installing design modifications in the environment is higher than the amount of inexplicability caused by the robot ($\beta > \alpha$). On the other hand, for setting (c), the algorithm generated the design in Figure 3.3b, which makes the robot's roundabout path completely explicable to the customers.

## 3.5  Bibliographic Remarks

The first paper to introduce explicability was Zhang et al. [2017] that introduced the model-free explicability. The paper used CRF Lafferty et al. [2001] to learn the labeling model and a modified version of the FF planner [Hoffmann, 2001] to generate the plan. The paper also performed some limited user studies to test whether such models can capture people's measure of plan explicability. The model-based explanation was introduced by Kulkarni et al. [2019a] and used a modified version of Helmert [2006] for the reconciliation search and used a learned distance function which used features described in Definition 6. The specific distance functions considered was originally discussed and used in the context of diverse planning (cf. [Srivastava et al., 2007, Nguyen et al., 2012]) where these distance functions are used to generate plans of differing characteristics. The paper also used data collected from participants to learn the distance function. The design for explicability was proposed by Kulkarni et al. [2020a] and followed the other design works that have been considered within the context of planning (cf. [Keren et al., 2014, 2016, 2018]). In

general, most of these works can be seen as special cases of the more general environment design work studied in papers like Zhang et al. [2009]. The specific paper Kulkarni et al. [2020a] used a compilation based method to generate the explicable plan as the paper looked at a cost difference based distance. The compilation follows the methods discussed in Sreedharan et al. [2020a], but removes the use of explanatory actions. We will discuss the base compilation in more detail in Chapter 7.

# Legible Behavior

In this chapter, the discussion will focus on another type of interpretable behavior, namely legibility. The notion of legibility allows the robot to implicitly communicate information about its goals, plans (or model, in general) to a human observer. For instance, consider a human robot cohabitation scenario consisting of a multi-tasking robot with varied capabilities that is capable of performing a multitude of tasks in an environment. In such a scenario, it is crucial for the robot to aid the human's goal or plan recognition process, as the human observer may not always know the robot's intentions or objectives beforehand. Hence, in such cases, it may be useful for the robot to communicate information that the human is unaware of. As the better off the human is at identifying the robot's goals or plans accurately, the better off is the overall team performance. However, explicit communication of objectives might not always be suitable. For instance, the *what, when* and *how* of explicit communication may require additional thought. Further, several other aspects like cost of communication (in terms of resources or time), delay in communication (communications signals may take time to reach the human), feasibility of communication (broken or unavailable sensors), etc., may also need to be considered. On the other hand, the robot can simply synthesize a behavior that implicitly communicates the necessary information to the human observer.

We will discuss the notion of legibility from the perspective of an offline setting where the observer has partial observability of the robot's actions. That is, the robot has to synthesize legible behavior in a setting where the human observer has access to the observations emitted from the entire execution trace of the robot and these observations do not always reveal the robot's exact action or state to the human. As the human is trying to infer the robot's goals or plans, she operates in a belief space due to the partial observability of the robot's activities. The robot has to modulate the human's observability and choose actions such that the ambiguity over specific goals or plans is reduced. We refer to this as the controlled observability planning problem (COPP). In the upcoming sections, we will see how the COPP formulation can be used to synthesize goal legible behavior as well as plan legible behavior in an offline setting with partial observability of the robot's activities.

## 4.1   Controlled Observability Planning Problem

In this framework, we consider two agents: a robot and the human observer. The robot has full observability of its activities. However, the observer only has partial observability of the robot's activities. The observer is aware of the robot's planning model and receives observations emitted as a side effect of the robot's execution. This framework supports an offline setting, and therefore the observer only receives the observations after the robot has finished executing the entire plan.

In this setting, the robot has a set of candidate goals, inclusive of its true goal. The

candidate goals are the set of possible goals that the robot may achieve in the given environment. The observer is aware of the robot's candidate goal set but is unaware of the robot's true goal. We now introduce a general COPP framework that will be used to define the goal legibility and plan legibility problems in the upcoming sections.

**Definition 9.** *A **controlled observability planning problem** is a tuple, $\mathcal{P}_{CO} = \langle \mathcal{M}^R, \mathcal{G}, \Omega, \mathcal{O} \rangle$, where,*

- $\mathcal{M}^R$ *is the planning model of the robot.*

- $\mathcal{G} = \{G_1 \cup G_2 \ldots \cup G_{n-1} \cup G^R\}$ *is a set of candidate goal conditions, each defined by subsets of fluent instantiations, where $G^R$ is the true goal of the robot.*

- $\Omega = \{o_i | i = 1, \ldots, m\}$ *is a set of $m$ observations that can be emitted as a result of the action taken and the state transition.*

- $\mathcal{O} : (A \times S) \to \Omega$ *is a many-to-one observation function associated with the human which maps the action taken and the next state reached to an observation in $\Omega$. That is to say, the observations are deterministic, each $\langle a, s' \rangle$ pair is associated with a single observation but multiple pairs can be mapped to the same observation.*

Note that here the human's expectation of the robot (i.e. $\mathcal{M}_h^R$) is defined in a distributed fashion. In particular, this setting assumes that the expectation matches the real model in all aspects but the goals. That is the human is unaware of the true goal and only has access to the set of potential goals. Moreover the human is a noisy observer (as defined by the observation model $\mathcal{O}$ and the set of observations $\Omega$). To simplify the setting we will focus on the cases where the human is aware of the observation model, i.e., from the observation they could potentially guess the possible state action pairs that could have generated it. Therefore, the human observer operates in the belief space. The robot takes the belief space of the observer into account during its planning process, so as to modulate what the human believes the current possible states could be.

### 4.1.1   Human's Belief Space

The human may use its observations of the robot's activity to maintain a *belief state*. A belief state is simply a set of possible states consistent with the observations. We use $\hat{s}$ as a notational aid to denote a state that is a member of the belief state.

**Definition 10.** *The **initial belief**, $b_0$, induced by observation, $o_0$ is defined as, $b_0 = \{\hat{s}_0 \mid \mathcal{O}(\emptyset, s_0) = o_0 \wedge \mathcal{O}(\emptyset, \hat{s}_0) = o_0\}$.*

Whenever a new action is taken by the robot, and the state transition occurs, the human's belief updates as follows:

**Definition 11.** *A **belief update**, $b_{i+1}$ for belief $b_i$ is defined as, $b_{i+1} = update(b_i, o_{i+1}) = \{\hat{s}_{i+1} \mid \exists \hat{a}, \ \delta(\hat{s}_i, \hat{a}) \models \hat{s}_{i+1} \wedge \hat{s}_i \in b_i \wedge \mathcal{O}(\hat{a}, \hat{s}_{i+1}) = o_{i+1}\}$.*

A sequence of belief updates gives us the observer's belief sequence that is consistent with the observation sequence emitted by the robot.

**Definition 12.** *A **belief sequence** induced by a plan p starting at state $s_0$, BS(p, $s_0$), is defined as a sequence of beliefs $\langle b_o, b_1, \ldots, b_n \rangle$ such that there exist $o_0, o_1, o_2, \ldots, o_n \in \Omega$ where,*

- $o_i = \mathcal{O}(a_i, s_i)$

- $b_{i+1} = update(b_i, o_{i+1})$

The set of plans that are consistent with the belief sequence of a given plan are called as belief plan set.

**Definition 13.** *A **belief plan set**, BPS(p, $s_0$) = $\{p_1, \ldots, p_n\}$, induced by a plan p starting at $s_0$, is a set of plans that are formed by causally consistent chaining of state sequences in BS(p, $s_0$), i.e., BPS(p, $s_0$) = $\{\langle \hat{s}_0, \hat{a}_1, \hat{s}_1, \ldots, \hat{s}_n \rangle \mid \forall \hat{a}_j, \hat{s}_{j-1} \models pre(\hat{a}_j) \land \hat{s}_{j-1} \in b_{j-1} \land \hat{s}_j \models \hat{s}_{j-1} \cup adds(\hat{a}_j) \setminus dels(\hat{a}_j) \land \hat{s}_j \in b_j\}$.*

The robot's objective is to generate a belief sequence in human's belief space, such that, the last belief in the sequence satisfies certain desired conditions.

### 4.1.2   Computing Solutions to COPP variants

Now we present a common algorithm template that can be used to compute plans for the COPP problem variants.

#### 4.1.2.1   Algorithm for Plan Computation

In this algorithm, each search node is represented by a belief estimate, which is a $\Delta$-sized subset of the belief state. A search node, $b_\Delta$, consists of state $s$ and a set of possible states, $\tilde{b}$ (cardinality of $\tilde{b}$ is given by $\Delta - 1$, where $\Delta$ is a counter). The successor node uses the $s$ in $b_\Delta$ to generate the successor state, and $b_\Delta$ to create the next belief state. There are two loops in the algorithm: the outer and the inner loop. The outer loop maintains the cardinality of $b_\Delta$ by incrementing the value of $\Delta$ in each iteration, such that value of $\Delta$ ranges from $1, 2, \ldots, |\mathcal{S}|$. In the inner loop, a heuristic-guided forward search (for instance, GBFS[Russell and Norvig, 2002]) can be used to search over space of belief states of cardinality less than or equal to $\Delta$. These loops ensure the complete exploration of the belief space. The algorithm terminates either when a plan is found or after running the outer loop for $|\mathcal{S}|$ iterations. The outer loop ensures that all the paths in the search space are explored.

**Proposition 1.** *The algorithm necessarily terminates in finite number of $|\mathcal{S}|$ iterations, such that, the following conditions hold:*

*(**Completeness**) The algorithm explores the complete solution space of $\mathcal{P}_{\mathcal{CO}}$, that is, if there exists a $\pi_{\mathcal{P}_{\mathcal{CO}}}$ that correctly solves $\mathcal{P}_{\mathcal{CO}}$, it will be found.*

---

**Algorithm 3** COOP Solution Plan Algorithm

---

**Input:** $\mathcal{P}_{\mathcal{CO}} = \langle D, \mathcal{G}, \Omega, \mathcal{O} \rangle$

**Output:** Solution $\pi_{\mathcal{P}_{\mathcal{CO}}}$, observation sequence, $O_{\mathcal{P}_{\mathcal{CO}}}$

1: Initialize *open, closed, unopened* lists and the counter $\Delta \leftarrow 1$
2: $b_\Delta \leftarrow \{I\}$ ; $b_0 \leftarrow \{\mathcal{O}(\emptyset, I)\}$ {Initialize initial search node, initial belief}
3: $open.push(\langle b_\Delta, b_0 \rangle, priority = 0)$
4: **while** $\Delta \leqslant |\mathcal{S}|$ **do**
5:     **while** $open \neq \emptyset$ **do**
6:         $b_\Delta, b, h(b_\Delta) \leftarrow open.\text{pop}()$
7:         **if** $|b_\Delta| > \Delta$ **then**
8:             $unopened.\text{push}(\langle b_\Delta, b \rangle, h(b_\Delta))$; **continue**
9:         $closed \leftarrow closed \cup b_\Delta$
10:         **if** $\langle b_\Delta, b \rangle \models \text{GOAL-TEST}(\mathcal{G})$ **then**
11:             **return** $\pi_{\mathcal{P}_{\mathcal{CO}}}, O_{\mathcal{P}_{\mathcal{CO}}}$
12:         **for** $s' \in successors(s)$ **do**
13:             $o \leftarrow \mathcal{O}(a, s')$
14:             $b' \leftarrow \text{Belief-Generation}(b, o)$
15:             $b'_\Delta = \langle s', \tilde{b}' \rangle$ {$\tilde{b}'$ of size $\Delta$-1}
16:             $h(b'_\Delta) \leftarrow \text{HEURISTIC-FUNCTION}(b'_\Delta, b')$
17:             add $b'_\Delta$ to *open* if not in *closed*
18:     $\Delta \leftarrow \Delta + 1$
19:     copy items of *unopened* to *open*, empty *unopened*
20:
21: **procedure** Belief-Generation($b$, $o$)
22: $b' \leftarrow \{\}$
23: **for** $\hat{s} \in b$ **do**
24:     **for** $\hat{a} \in A$ **do**
25:         **if** $\mathcal{O}(\hat{a}, \delta(\hat{s}, \hat{a})) = o$ **then**
26:             $b' \leftarrow b' \cup \delta(\hat{s}, \hat{a})$
27: **return** $b'$

---

(***Soundness***) *The plan, $\pi_{\mathcal{P}_{\mathcal{CO}}}$, found by the algorithm correctly solves $\mathcal{P}_{\mathcal{CO}}$ as ensured by the corresponding goal-test.*

We have not yet defined the goal-test and the notion of validity of a solution for a COPP problem. We will be encoding the notion of legibility we wish to pursue as part of the goal-test. We will do so by ensuring that the legibility score for the corresponding plan is above certain threshold. We will look at some specific legibility scores we can use as we ground this specific algorithmic framework for various use cases.

#### 4.1.2.2 Optimization

In order to speed up the search process, we perform an optimization on the aforementioned algorithm. For each search node, $b_\Delta$, apart from the approximate belief estimate, we maintain the full belief update $b$ consistent with a path to $s$. The approximate belief update $b_\Delta$ can be generated by choosing $\Delta$-sized combinations of states from the complete

belief. For example, when $\Delta = 1$, $b_\Delta$ only consists of the state $s$ but still maintains full belief update $b$, when $\Delta = 2$, $b_\Delta$ consists of a new combination of approximate belief of size 2 derived from the maintained full belief. When $\Delta = 1$, because of the check for duplicate states in the closed list, only one path to the search node is explored. Therefore, the use of $\Delta$ allows the search process to explore multiple paths leading to a particular search node. The complete $b$ helps in finding the problem variant solutions faster at lower $\Delta$ values. We present the details of the optimization in Algorithm 3. In the following sections, we show how we customize the goal-test (line 10) and the heuristic function (line 16) to suit the needs of each of the COPP problem variants.

### 4.1.3   Variants of COPP

Within this framework, we will discuss the problem of goal legibility and plan legibility. With goal legibility, the objective of the robot is to convey *at most j* goals to the observer. With plan legibility, we look at a specific instance of COPP problem where $|\mathcal{G}| = 1$. Since we know that COPP problems require that $G^R$ be part of $|\mathcal{G}|$ this means that in these cases the human knows the true goal of the robot. But even then the human observer may not know the exact plan being executed by the robot given their limited sensors. Thus the objective of the robot within plan legibility becomes to constrain the uncertainty of the observer to *at most m* similar plans to the goal. We will see both of these problems in detail in the following sections. The COPP framework can be also used in adversarial environments. These variants will be covered in Chapter 10.

## 4.2   Goal Legibility

In this setting, the robot's objective is to convey some information about its candidate goal set to the human observer. This may involve communicating its true goal to human, by ensuring at most one goal is consistent with the observation sequence produced. Or in general, the robot may want to communicate a set of at most $j$ candidate goals inclusive of its true goal, to the human observer. Essentially, the point with goal legibility is to reduce the observer's ambiguity over the robot's possible goal set by narrowing it down to at most $j$ goals. Which mean in this scenario, for a given plan $\pi$ the corresponding legibility score is given as $\mathcal{L}(\pi, \mathcal{P}_{CO}) \propto \frac{1}{\lceil \{ G \mid G \in \mathcal{G} \land G \in BS(\pi, I) \rceil \}}$.

#### 4.2.0.1   Example

Let's understand this problem with an example. Consider a situation where the robot is a port management agent and a human is the supervisor that has sensors or subordinates at the port who provide partial information about the nature of activity being carried out at the port (refer Figure 4.1). For instance, when a specific crate is loaded onto the ship, the observer finds out that something was loaded, but not the identity of the loaded crate. The observer knows the initial inventory at the port, but when new cargo is acquired by the port, the observer's sensors reveal only that more cargo was received; they do not specify the numbers or identities of the received crates. A legible plan for loading sensitive

**Fig. 4.1:** The differences in belief sequences induced by different plans for an observer with noisy sensors.

cargo (the red crate) and acquiring more cargo may first load the crate and then acquire more crates. This plan reveals the identity of the crate that was loaded based on the observers' information about the remaining cargo in the port: the final belief state has a unique crate loaded on the ship even though it retains uncertainty about the new cargo in the port. However, if the plan were to first acquire more cargo, the observer's sensors are insufficient to determine which crate was loaded: the plan maintains ambiguity in the observer's belief. This is reflected in the observer's belief state sequence, where the last belief state includes states with all different types of crates in the ship. Although both plans have the same cost and effects for the dock, one conveys the activity being carried out while the other adds more uncertainty. The COPP framework allows the robot to select plans that may be legible in this manner.

#### 4.2.0.2 Goal Legibility Problem

As mentioned earlier, in goal legibility problem, the robot's aim is to take actions exclusive to the goal so as to help the observer in goal deduction.

**Definition 14.** *A **goal legibility planning problem** is a $\mathcal{P}_{CO}$, where, $\mathcal{G} = \{G^R \cup G_1 \cup \ldots \cup G_{n-1}\}$ is the set of n goals where $G^R$ is the true goal of the robot, and $G_1, \ldots, G_{n-1}$ are confounding goals.*

Rather than optimizing directly for the legibility score, we will instead limit ourselves to solutions whose legibility score is above a certain threshold. In particular, we can generate a plan that conveys at most $j$ candidate goals inclusive of its true goal. Thus robot has to ensure that the observation sequence of a legible plan is consistent with *at most $j$* goals so as to limit the number of goals in the observer's final belief state.

**Definition 15.** *A plan, $\pi_j$, is a **j-legible plan**, if $\delta(I, \pi_j) \models G^R$ and the last belief, $b_n \in BS(\pi_j, I)$, satisfies the following, $|G \in \mathcal{G} : \exists s \in b_n, s \models G| \leqslant j$, where $1 \leqslant j \leqslant n$.*

### 4.2.1 Computing Goal Legible Plans

In the case of goal legibility, we run the algorithm with the robot's true goal.

#### 4.2.1.1 Goal test

In order to ensure that the computed plan is consistent with *at most $j$* true goals, we change our goal condition to additionally check whether at most $j - 1$ confounding goals have been achieved in the last belief in $BS$, or it can be interpreted as at least $n - j$ goals are absent in the belief.

**Fig. 4.2:** Illustration of the impact of plan legibility on the observer's plan inference process.

### 4.2.1.2  Heuristic function

In this case, our objective is to avoid at least $n - j$ goals, that is be consistent with at most $j$ goals. We achieve this by minimizing the heuristic cost to the true goal from the robot's actual state and to the $j - 1$ confounding goals from the robot's belief state. However, we maximize the heuristic cost to other $n - j$ goals in order to achieve at most $j$ goals in the last belief state. This is written as,

$$h(s) = h_{G_A}(s) + h_{\mathcal{G}_{j-1}}(b) - h_{\mathcal{G}_{n-j}}(b) \tag{4.1}$$

## 4.3  Plan Legibility

In this problem setting, the observer is aware of the robot's goal, since the candidate goal set consists of a single goal. Although, the observer is unaware of the robot's choice of plan to achieve that goal. In general, there can be multiple different plans that allow the robot to achieve a goal. This is further complicated by the fact that the observer has partial observability of the robot's activities. Therefore, in this problem setting the goal of the robot is to reduce the ambiguity over the possible plans to its goal. The robot can achieve this, by computing a plan whose observation sequence conforms to a set of plans that are similar in terms of their actions, thereby making it easy for the observer to guess the actions executed by the agent.

For instance, in Figure 4.2, the goal of the agent is to pick up a medkit and treat the victim on the disaster site. The observer has a noisy sensor model and it is not accurate enough to give the exact cell location of the robot. Therefore, the observer relies on how many steps the robot has moved to guess the robot's path of execution. In the initial state shown in the leftmost subfigure in Figure 4.2, it can be seen that the robot has two choices of medkits to pick from. If it chooses the medkit as shown in the rightmost sub-figure, there are only two paths of length 2 towards the medkit and one path of length 2 towards the victim. In comparison, the middle sub-figure, shows the various different paths of length 4 that can lead to the other medkit, making it harder for the human to guess the robot's actions. Therefore, to reduce the ambiguity over its path, the robot

chooses the medkit shown in the rightmost figure. This involves an observation sequence with two similar paths leading to the goal, making it easy for the human to guess some of the actions of the robot. Thus, the plan legibility problem can be solved by finding an observation sequence that is consistent with plans that are similar to each other.

**Definition 16.** *A **plan legibility planning problem** is a tuple,* $\mathcal{P}_{PL} = \langle \mathcal{D}, \mathcal{G}_{PL}, \Omega, \mathcal{O} \rangle$, *where,* $|\mathcal{G}_{PL}| = 1$.

## 4.3.1 Computing Plan Legible Plans

The solution to a plan legibility problem is an *m-similar* plan. An *m-similar* plan is a plan whose observation sequence is consistent with *at least m* similar plans to the goal, such that, these plans are *at most d* distance away from each other. In order to compute an *m-similar* plan, we need to keep track of the plans that are consistent with the observation sequence and reach the goal. To compute the diversity between all the pairs of plans consistent with the observation sequence, a plan distance measure like action distance, causal link distance, state sequence distance (discussed in Chapter 3) can be used. This approach can use any valid plan distance. We now define an *m-similar* plan.

**Definition 17.** *Two plans,* $p_1, p_2$, *are a **d-distant pair** with respect to distance function* $\mathcal{D}$ *if,* $\mathcal{D}(p_1, p_2) = d$, *where* $\mathcal{D}$ *is a diversity measure.*

We will be using this distance function as the basis of legibility score, but again we will focus on generating plans whose legibility score is above some prespecified threshold, by requiring that the possible plans are no farther than *d-distance* away.

**Definition 18.** *A BPS induced by plan p starting at* $s_0$ *is **maximally d-distant**,*

$$d_{max}(BPS(p, s_0)) \text{ if } d = \max_{p1, p2 \in BPS(p, s_0)} \mathcal{D}(p1, p2)$$

.

**Definition 19.** *A plan,* $\pi_m$, *is an **m-similar plan**, if for a given value of d and distance function* $\mathcal{D}$, $d_{max}(BPS(\pi_m, I)) \leq d$, $|BPS(\pi_m, I)| \geq m$, *where* $m \geq 2$ *and every plan in* $BPS(\pi_m, I)$ *achieves the goal in* $\mathcal{G}_{PL}$.

In order to generate a solution to the plan legibility problem, we use the algorithm presented in Algorithm 3. Again, the goal test and heuristic function are customized to ensure that there are at least *m* similar plans to the true goal that are consistent with the observation sequence and the maximum distance between these plans is *at most d*.

### 4.3.1.1 Goal test

To ensure the plans in $BPS$, induced by an *m-similar* plan, can achieve the goal in $\mathcal{G}_{PL}$, we check whether at least *m* plans are reaching the goal or not and whether the maximum distance between plans in $BPS$ is at most *d*. Also in order to ensure termination of the algorithm, there is a cost-bound given as input to the algorithm.

#### 4.3.1.2  Heuristic function

Apart from minimizing the heuristic cost to the goal, the customized heuristic given below also minimizes the $d$ of $d_{max}(BPS(p, s_0))$ induced by plan $p$ starting at $s_0$. This decreases the maximum distance between the plan pairs in the BPS. This distance can be computed using a plan distance measure.

$$h(s) = h_{G_A}(s) + d_{max}(BPS(p, s_0)) \tag{4.2}$$

#### 4.3.1.3  Plan Legibility as Offline Predictability

Plan legibility can be likened to the notion of offline predictability insofar that plan legible behaviors allows the robot to reduce the observer's ambiguity over the possible plans executed given a goal, which is also a property exhibited by predictable behaviors. However, predictable behaviors are also inherently easy to anticipate, either globally or locally. Globally predictable behavior is a behavior that an observer would anticipate the robot to perform for a certain goal, versus locally predictable behavior is a behavior where given a plan prefix, the rest of the suffix towards a certain goal can be easily anticipated by the observer. This notion of predictability has been mostly explored in the motion planning community. However, plan legibility does not always lead to predictable behaviors. This is because in plan legibility, the emphasis is on making the robot's actions easy to guess given a corresponding observation sequence. However, the observation sequence in itself might not be globally or even locally predictable to the observer.

## 4.4  Bibliographic Remarks

In the motion planning and robotics community, legibility has been a well-studied topic [Dragan et al., 2013, Dragan and Srinivasa, 2013, Dragan et al., 2015, Knepper et al., 2017]. The original work on legibility used Bayesian formulation to generate legible robot motions in an online setting. A legible motion is one that ensures that the actual robot goal has the highest likelihood amongst the candidate goals. The legible planning discussed here is formulated within the controlled observability planning framework (introduced in Kulkarni et al. [2019b]), which generalizes the notion of legibility in terms of human's observability as well as in terms of the amount of information divulged (with at most $j$ goals, or plans that are at most $d$ distance away). This framework uses the notion of sensor models to capture the human's imperfect observations and to model the human's belief update [Geffner and Bonet, 2013, Bonet and Geffner, 2014, Keren et al., 2016]. The goal and plan legible behaviors have been categorized into legible planning and predictable planning in a recent survey on interpretable behaviors [Chakraborti et al., 2019a].

CHAPTER 5

# Explanation as Model Reconciliation

In this chapter, we revisit the explicability score and investigate an alternate strategy to improve the explicability of the robot behavior, namely explanations. Rather than force the robot to choose behaviors that are inherently explicable in the human model, here we will let the robot choose a behavior optimal in its model and use communication to address the central reason why the human is confused about the behavior in the first place, i.e., the model difference. That is, the robot will help the human understand why the behavior was performed, by choosing to reveal parts of its model that were previously unknown to the human. This would allow us to overcome one of the main shortcomings of the plan generation methods discussed in Chapter 2, namely that there might not exist a plan in the robot model that has a high explicability score. In this scenario, the explicability score of the plan is only limited by the agent's ability to effectively explain it. In this chapter, in addition to introducing the basic framework of explanation as model reconciliation under a certain set of assumptions, we will also look at several types of model reconciliation explanations, study some of their properties and consider some simple approximations. In the coming chapters, we will further extend the idea of explanations and look at ways of relaxing some of the assumptions made in this chapter.

## 5.1 Model-Reconciliation as Explanation

As mentioned in chapter 2, the explicability score of a plan $\pi$ generated using the robot's model $\mathcal{M}^R$ is given by the expression

$$E(\pi, \mathcal{M}_h^R) \propto \max_{\pi' \in \Pi^{\mathcal{M}_h^R}} - 1 * \mathcal{D}(\pi, \pi', \mathcal{M}_h^R)$$

Where $\mathcal{D}$ is the distance function, $\mathcal{M}_h^R$ the human's estimate of what the robot's model may be, and $\pi'$ the plan closest to $\pi$ in $\mathcal{M}_h^R$ (per the distance function $\mathcal{D}$). This means, even if the robot chooses a plan in its own model ($\mathcal{M}^R$), a human observer may find it inexplicable if their estimate of the robot model differs from $\mathcal{M}^R$. Thus a way to address the inexplicability of the plan would be to explain the difference between the robot's own model and human's estimation. This explanatory process is referred to as ***model reconciliation.***We will specifically focus on cases, where $\mathcal{D}$ is given by the cost difference and the expected set corresponds to plans optimal in the human's model. While this method assumes the human is a perfectly rational decision-maker (which isn't necessarily true), studies have shown this method to still result in useful explanations in many scenarios. Of course we could replace the distance function used and expected plan set with more realistic choices and the basic algorithms presented here should stay mostly same.

The explanation process captured by Model Reconciliation begins with the following question:

*Q₁: Why plan π?*

An explanation here needs to ensure that both the explainer and the explainee agree that $\pi$ is the best decision that could have been made in the given problem. In the setting we are considering here that would mean providing model artifacts to the explainee so that $\pi$ is now also optimal in the updated mental model and thus have high explicability score under our current set of assumptions (we will refer to this as the completeness property later).

**Definition 20.** *The* **model reconciliation problem (MRP)** *is represented as a tuple* $\langle \pi^*, \langle \mathcal{M}^R, \mathcal{M}_h^R \rangle \rangle$ *where* $C(\pi^*, \mathcal{M}^R) = C^*_{\mathcal{M}^R}$, *i.e. the plan* $\pi^*$ *is the optimal plan in the robot model* $\mathcal{M}^R$ *but may not be so in the human mental model* $\mathcal{M}_h^R$.

Now we can define an explanation as

**Definition 21.** *An* **explanation** *is a solution to the model reconciliation problem in the form of (1) a model update* $\mathcal{E}$ *such that the (2) robot optimal plan is (3) also optimal in the updated mental model.*

- $\widehat{\mathcal{M}_h^R} \longleftarrow \mathcal{M}_h^R + \mathcal{E}$; *and*

- $C(\pi, \mathcal{M}^R) = C^*_{\mathcal{M}^R}$;

- $C(\pi, \widehat{\mathcal{M}_h^R}) = C^*_{\widehat{\mathcal{M}_h^R}}$.

Now the problem of finding the explanation becomes a search over the set of model information that can be provided to the user to get an updated user-model of desired property. This in term can be visualized as the search over the space of possible human models that can result from providing information consistent with the robot model. We can facilitate such a search over the model space, by leveraging the model parameterization scheme specified in Chapter 2. Specifically such a scheme results in a state space of the form

$$\mathcal{F} = \{\textit{init-has-f} \mid \forall f \in F_h^R \cup F^R\} \cup \{\textit{goal-has-f} \mid \forall f \in F_h^R \cup F^R\}$$

$$\bigcup_{a \in A_h^R \cup A^R} \{\textit{a-has-precondition-f, a-has-add-effect-f,}$$

$$\textit{a-has-del-effect-f} \mid \forall f \in F_h^R \cup F^R\}$$

$$\cup \{\textit{a-has-cost-C(a)} \mid a \in A_h^R\} \cup \{\textit{a-has-cost-C(a)} \mid a \in A^R\}.$$

Again the mapping function $\Gamma : \mathcal{M} \mapsto s$ can convert the given planning problem $\mathcal{M} = \langle F, A, I, G, C \rangle$ as a state $s \subseteq \mathcal{F}$, as follows –

$$\tau(f) = \begin{cases} \textit{init-has-f} & \text{if } f \in I, \\ \textit{goal-has-f} & \text{if } f \in G, \\ \textit{a-has-precondition-f} & \text{if } f \in pre(a), \ a \in A \\ \textit{a-has-add-effect-f} & \text{if } f \in adds(a), \ a \in A \\ \textit{a-has-del-effect-f} & \text{if } f \in dels^-(a), \ a \in A \\ \textit{a-has-cost-f} & \text{if } f = C(a), \ a \in A \end{cases}$$

$$\Gamma(\mathcal{M}) = \{\tau(f) \mid \forall f \in I \cup G \cup \\ \bigcup_{a \in A} \{f' \mid \forall f' \in \{C(a)\} \cup \text{pre}(a) \cup \text{adds}(a) \cup \text{dels}(a)\}\}$$

**Definition 22.** *The **model-space search problem** is specified as $\langle \mathcal{F}, \Lambda, \Gamma(\mathcal{M}_1), \Gamma(\mathcal{M}_2) \rangle$ with a new action set $\Lambda$ containing unit model change actions $\lambda \in \Lambda, \lambda : \mathcal{F} \to \mathcal{F}$ such that $|(s_1 \setminus s_2) \cup (s_2 \setminus s_1)| = 1$. The new transition or edit function is given by $\delta_{\mathcal{M}_1, \mathcal{M}_2}(s_1, \lambda) = s_2$ such that,* `condition 1:` $s_2 \setminus s_1 \subseteq \Gamma(\mathcal{M}_2)$ *and* `condition 2:` $s_1 \setminus s_2 \not\subseteq \Gamma(\mathcal{M}_2)$ *are satisfied.*

This means that model change actions can only make a single change to a domain at a time, and *all these changes are consistent with the model of the planner.* The solution to a model-space search problem is given by a *set* of edit functions $\{\lambda_i\}$ that transforms the model $\mathcal{M}_1$ to $\mathcal{M}_2$, i.e. $\delta_{\mathcal{M}_1, \mathcal{M}_2}(\Gamma(\mathcal{M}_1), \{\lambda_i\}) = \Gamma(\mathcal{M}_2)$. An explanation can thus be cast as a solution to the model-space search problem $\langle \mathcal{F}, \Lambda, \Gamma(\mathcal{M}_h^R), \Gamma(\widehat{\mathcal{M}}) \rangle$ with the transition function $\delta_{\mathcal{M}_h^R, \mathcal{M}^R}$ such that Condition (3) mentioned in Definition 21 is preserved.

### 5.1.1 The Fetch Domain

Let us look at an example domain to see how model reconciliation explanation could be used. Consider the Fetch robot [1] whose design requires it to `tuck` its arms and lower its torso or `crouch` before moving. This is not obvious to a human navigating it and it may lead to an unbalanced base and toppling of the robot if the human deems such actions as unnecessary. The move action for the robot is described in PDDL in the following model snippet –

```
(:action move
:parameters    (?from ?to - location)
:precondition  (and (robot-at ?from) (hand-tucked) (crouched))
:effect        (and (robot-at ?to) (not (robot-at ?from))))
```

```
(:action tuck
:parameters    ()
:precondition  ()
:effect        (and (hand-tucked) (crouched)))
```

---

[1] https://fetchrobotics.com/fetch-mobile-manipulator/

**Fig. 5.1:** The Fetch in the crouched position with arm tucked (left), torso raised and arm outstretched (middle) and the rather tragic consequences of a mistaken action model (right showing a fractured head from an accident).

```
(:action crouch
:parameters    ()
:precondition  ()
:effect        (and (crouched)))
```

Notice that the `tuck` action also involves a lowering of torso so that the arm can rest on the base once it is tucked in.[2] Now, consider a planning problem where the the robot needs to transport a block from one location to another, with the following initial and goal states –

```
(:init (block-at b1 loc1) (robot-at loc1) (hand-empty))
(:goal (and (block-at b1 loc2)))
```

An optimal plan for the robot involves a `tuck` action followed by a `move`:

```
pick-up b1 -> tuck -> move loc1 loc2 -> put-down b1
```

The human, on the other hand, expects a much simpler model, as shown below. The `move` action does not have the preconditions for tucking the arm and lowering the torso, while `tuck` does not automatically lower the torso either.

```
(:action move
:parameters    (?from ?to - location)
:precondition  (and (robot-at ?from)
:effect        (and (robot-at ?to) (not (robot-at ?from))))
```

---
[2]Fetch User Manual: https://docs.fetchrobotics.com/

```
(:action tuck
:parameters     ()
:precondition   ()
:effect         (and (hand-tucked))



(:action crouch
:parameters     ()
:precondition   ()
:effect         (and (crouched)))
```

The original plan is no longer optimal to the human who can envisage better alternatives (a shorter plan without the extra `tuck` action) in their mental model. An explanation here is a model update that can mitigate this disagreement –

```
Explanation >> MOVE_LOC1_LOC2-has-precondition-HAND-TUCKED
```

This correction brings the mental model closer to the robot's ground truth and is necessary and sufficient to make the robot's plan optimal in the resultant domain so that the human cannot envisage any better alternatives. This is the essence of the model reconciliation process.

## 5.2 Explanation Generation

Before we go on to specific methods we could use to identify the required explanation for a model reconciliation problem, let us consider the various desirable properties that characterize explanations in such settings.

P1. **Completeness -** Explanations of a plan should be able to be compared and contrasted against other alternatives, so that no better solution exists. We can enforce this property by requiring that in the updated human mental model the plan being explained is now optimal.

   – *An explanation is complete iff* $C(\pi, \widehat{\mathcal{M}}_h^R) = C^*_{\widehat{\mathcal{M}}_h^R}$.

P2. **Conciseness -** Explanation should be concise so that it is easily understandable to the explainee. Larger an explanation is, the harder it is for the human to process that information. Thus we can use the length of the explanation as a useful proxy or first approximation for the complexity of an explanation.

P3. **Monotonicity -** This property ensures that remaining model differences cannot change the completeness of an explanation, i.e. all aspects of the model that were relevant to the plan have been reconciled. Thus, monotonicity of an explanation subsumes completeness and requires more detail.

> – *An explanation is monotonic iff*
> $C(\pi^*, \hat{\mathcal{M}}) = C^*_{\hat{\mathcal{M}}} \; \forall \hat{\mathcal{M}} : ((\Gamma(\hat{\mathcal{M}}) \setminus \Gamma(\mathcal{M}^R_h)) \cup (\Gamma(\mathcal{M}^R_h) \setminus \Gamma(\hat{\mathcal{M}}))) \subset (\Gamma(\hat{\mathcal{M}}) \setminus \Gamma(\mathcal{M}^R_h)) \cup (\Gamma(\mathcal{M}^R_h) \setminus \Gamma(\hat{\mathcal{M}})).$

Thus an updated model would satisfy monontonicity property if no additional information about the robot model would render the plan being explained suboptimal or invalid.

P4. **Computability -** While conciseness deals with how easy it is for the explainee to understand an explanation, computability measures the ease of computing the explanation from the point of view of the planner.

We will now introduce different kinds of multi-model explanations that can participate in the model reconciliation process, provide examples, propose algorithms to compute them, and compare and contrast their respective properties. We note that the requirements outlined above are in fact often at odds with each other - an explanation that is very easy to compute may be very hard to comprehend.

### 5.2.1  Explanation types

A simple way to explain would be to provide the model differences pertaining to only the actions that are present in the plan being explained –

**Definition 23.** *A* **plan patch explanation (PPE)** *is given by –*

$$\mathcal{E}^{PPE} = (\mathcal{F}^{\mathcal{M}^R}(\pi) \setminus \mathcal{F}^{\mathcal{M}^R_h}(\pi)) \cup (\mathcal{F}^{\mathcal{M}^R_h}(\pi) \setminus \mathcal{F}^{\mathcal{M}^R}(\pi))$$

*where* $\mathcal{F}^{\mathcal{M}}(\pi)$ *gives the model parameters of* $\mathcal{M}$ *corresponding to all the actions in the plan* $\pi$ *(i.e.,* $\mathcal{F}^{\mathcal{M}}(\pi) = \bigcup_{\{C(a)\} \cup pre(a) \cup adds(a) \cup dels(a) : a \in \pi} \tau(f))$.

Clearly, such an explanation is easy to compute and concise by focusing only on plan being explained. However, it may also contain information that need not have been revealed, while at the same time ignoring model differences elsewhere in $\mathcal{M}^R_h$ that could have contributed to the plan being suboptimal in it. Thus, it is not *complete*. On the other hand, an easy way to compute a complete explanation would be to provide the entire model difference to the human –

**Definition 24.** *A* **model patch explanation (MPE)** *is given by –*

$$\mathcal{E}^{MPE} = (\Gamma(\mathcal{M}^R) \setminus \Gamma(\mathcal{M}^R_h)) \cup (\Gamma(\mathcal{M}^R_h) \setminus \Gamma(\mathcal{M}^R))$$

This is also easy to compute but can be quite large and is hence far from being concise. Instead, we can try to minimize the size (and hence increase the comprehensibility) of explanations by searching in the space of models and thereby not exposing information that is not relevant to the plan being explained while still trying to satisfy as many requirements as we can.

**Definition 25.** *A **minimally complete explanation (MCE)** is the shortest possible explanation that is complete –*

$$\mathcal{E}^{MCE} = \arg\min_{\mathcal{E}} |(\Gamma(\widehat{\mathcal{M}}) \setminus \Gamma(\mathcal{M}_h^R)) \cup (\Gamma(\mathcal{M}_h^R) \setminus \Gamma(\widehat{\mathcal{M}}))| \;\; with \; C(\pi, \widehat{\mathcal{M}}_h^R) = C^*_{\widehat{\mathcal{M}}_h^R}$$

The explanation provided before in the Fetch domain is indeed the smallest set of domain changes that may be made to make the given plan optimal in the updated action model, and is thus an example of a minimally complete explanation.

The optimality criterion happens to be relevant to both the cases where the human expectation is better, or worse, than the plan computed by the planner. This might be counter to intuition, since in the latter case one might expect that just establishing feasibility of a better plan would be enough. Unfortunately, this is not the case, as can be easily seen by creating counter-examples where other faulty parts of the human model might disprove the optimality of the plan. Which brings us to the proposition,

**Proposition 2.** *If $C(\pi^*, \mathcal{M}_h^R) < \min_\pi C(\pi, \mathcal{M}_h^R)$, then ensuring feasibility of the plan in the modified planning problem, i.e. $\delta_{\widehat{\mathcal{M}}}(\widehat{I}, \pi^*) \models \widehat{G}$, is a necessary but* not *a sufficient condition for $\widehat{\mathcal{M}} = \langle \widehat{F}, \widehat{A}, \widehat{I}, \widehat{G} \rangle$ to yield a valid explanation.*

Note that a minimally complete explanation can be rendered invalid given further updates to the model. This can be easily demonstrated in our running example in the Fetch domain. Imagine that if, at some point, the human were to find out that the action `(move)` also has a precondition `(crouched)`, then the previous robot plan will no longer make sense to the human since now according to the human's faulty model (being unaware that the tucking action also lowers the robot's torso) the robot would need to do *both* `tuck` and `crouch` actions before moving. Consider the following explanation in the Fetch domain instead –

```
Explanation >> TUCK-has-add-effect-CROUCHED
Explanation >> MOVE_LOC2_LOC1-has-precondition-CROUCHED
```

This explanation does not reveal all model differences but at the same time ensures that the plan remains optimal for this problem, irrespective of any other changes to the model, by accounting for all the relevant parts of the model that engendered the plan. It is also the smallest possible among all such explanations.

**Definition 26.** *A **minimally monotonic explanation (MME)** is the shortest explanation that preserves both completeness and monotonicity –*

$$\mathcal{E}^{MME} = \arg\min_{\mathcal{E}} |(\Gamma(\widehat{\mathcal{M}}) \setminus \Gamma(\mathcal{M}_h^R)) \cup (\Gamma(\mathcal{M}_h^R) \setminus \Gamma(\widehat{\mathcal{M}}))| \;\; with \; P1 \; \& \; P3$$

**Fig. 5.2:** Illustration of the different kinds of explanations in the Fetch domain. Here the PPE and MPE are equivalent (which is the worst case for the former) and both longer than the MCE or the MME. Also, the MCE is shorter than, and not a subset of the MME.

An MCE or MME solution may not be unique to an MRP problem. This can happen when there are multiple model differences supporting the same causal links in the plan - a minimal explanation can get by (i.e. guarantee optimality in the modified model) by only exposing one of them to the human.

Also it is easy to see that an MCE may not necessarily be part of an actual MME. This is illustrated in Figure 5.2.

**Proposition 3.** *An MCE may not be a subset of an MME, but it is always smaller or equal in size, i.e.* $|MCE| \leq |MME|$.

### 5.2.2   Desiderata for Explanations as Discussed in Social Sciences

Before, we delve into the actual algorithms for generating such explanations, lets take a quick look at how model reconciliation explanation connects to the wider literature on explanations. One of the key takeaways from many of the works dealing with this topic from social sciences has been the identification of three crucial properties for effective explanations, namely, explanations need to be *contrastive, social and selective.*

Contrastive explanations are explanations that take the form of answers to questions of the type "Why P and not Q?" Where P is the fact being explained and Q the foil or the alternate option raised by the explainee (i.e. the one who is asking for the explanation). Thus explanations for such question need to contrast the choice P over the raised foil. In the case of planning problems, a very common form of contrastive questions involve cases where P , or *the fact*, is the plan itself being proposed by the system and *the foil* is some alternate plan expected by the user. In model reconciliation, the explanations being provided as part of model reconciliation can be viewed as answering a contrastive query where the foil is implicitly specified, i.e., the foil could be any plan in the human model. Given the fact that these explanations establish optimality, the current plan must be better than or as good as any alternate plan the user could have come up with.

Social explanations are those that explicitly take into account the beliefs of the user. Model-reconciliation framework discussed here is clearly social in the sense that explanations are tailored to what the human believes to be the robot model. Selectivity of explanations deals with the fact that explanations need to only include relevant details (and in some sense minimal). As we will see below, the actual approaches to explanation generation focus on generating the minimal information needed to meet the required properties.

### 5.2.3 Model Space Search for Minimal Explanations

In the following, we will see how the state space provided by $\Gamma$ can be used in model-space search for computing MCEs and MMEs (computation of PPE and MPE follows directly from $\mathcal{M}^R$, $\mathcal{M}_h^R$ and $\pi^*$).

#### 5.2.3.1 Model Space Search for MCEs

To compute MCEs, we employ A$^*$ search in the space of models, as shown in Algorithm 4. The algorithm is referred to as MEGA – **M**ulti-model **E**xplanation **G**eneration **A**lgorithm. Given an MRP, we start off with the initial state $\Gamma(\mathcal{M}_h^R)$ derived from the human's expectation of a given planning problem $\mathcal{M}^R$, and modify it incrementally until we arrive at a planning problem $\widehat{\mathcal{M}}$ with $C(\pi^*, \widehat{\mathcal{M}}) = C^*_{\widehat{\mathcal{M}}}$, i.e. the given plan is explained. Note that the model changes are represented as a set, i.e. there is no sequentiality in the search problem. Also, we assign equal importance to all model corrections. We can easily capture differential importance of model updates by attaching costs to the edit actions $\lambda$ - the algorithm remains unchanged. One could also employ a selection strategy for successor nodes to speed up search (by overloading the way the priority queue is popped) by first processing model changes that are relevant to actions in $\pi_R^*$ and $\pi_H$ before the rest.

**Proposition 4.** *The successor selection strategy outlined in Algorithm 5 will never remove a valid solution from the search space.*

*Proof Sketch.* Let $\mathcal{E}$ be the MCE for an MRP problem and let $\mathcal{E}'$ be any intermediate explanation found by our search such that $\mathcal{E}' \subset \mathcal{E}$, then the set $\mathcal{E} \setminus \mathcal{E}'$ must contain at least one $\lambda$ related to actions in the set $\{a \mid a \in \pi_R^* \vee a \in \pi'\}$ (where $\pi'$ is the optimal plan for the model $\hat{\mathcal{M}}$ where $\delta_{\mathcal{M}_h^R, \mathcal{M}^R}(\Gamma(\mathcal{M}_h^R), \mathcal{E}') = \Gamma(\hat{\mathcal{M}})$. To see why this is true, consider an $\mathcal{E}'$ where $|\mathcal{E}'| = |\mathcal{E}| - 1$. If the action in $\mathcal{E} \setminus \mathcal{E}'$ does not belong to either $\pi_R^*$ or $\pi'$ then it cannot improve the cost of $\pi_R^*$ in comparison to $\pi'$ and hence $\mathcal{E}$ cannot be the MCE. Similarly we can show that this relation will hold for any size of $\mathcal{E}'$. We can leverage this knowledge about $\mathcal{E} \setminus \mathcal{E}'$ to create an admissible heuristic that considers only relevant changes. □
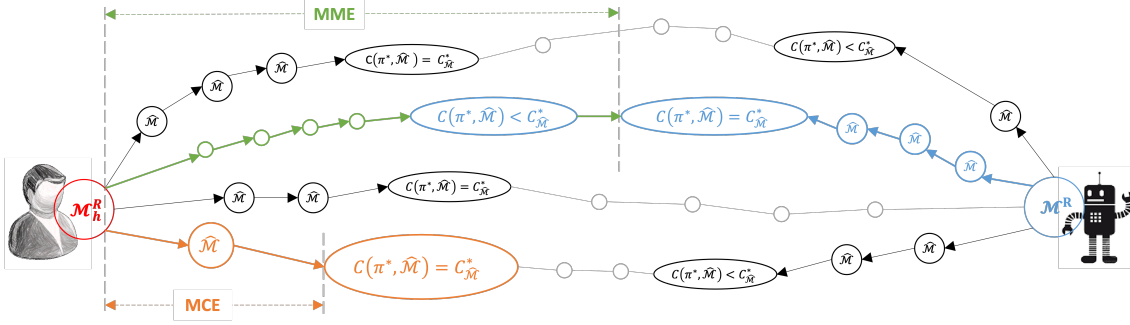
**Fig. 5.3:** Illustration contrasting MCE search with MME search.

### 5.2.3.2 Model Space Search for MMEs

As per definition, beyond the model obtained from the minimally monotonic explanation, there do not exist any models which are not explanations of the same MRP, while at the same time making as few changes to the original problem as possible. It follows that this is the largest set of changes that can be done on $\mathcal{M}^R$ and still find a model $\widehat{\mathcal{M}}$ where $C(\pi^*, \widehat{\mathcal{M}}) = C^*_{\widehat{\mathcal{M}}}$ - this property can be used in the search for MMEs.

**Proposition 5.** $\mathcal{E}^{MME} = \arg\max_{\mathcal{E}} |(\Gamma(\widehat{\mathcal{M}}) \setminus \Gamma(\mathcal{M}^R)) \cup (\Gamma(\mathcal{M}^R) \setminus \Gamma(\widehat{\mathcal{M}}))|$ *such that* $\forall \hat{\mathcal{M}} ~((\Gamma(\hat{\mathcal{M}}) \setminus \Gamma(\mathcal{M}^R)) \cup (\Gamma(\mathcal{M}^R) \setminus \Gamma(\hat{\mathcal{M}}))) \subseteq ((\Gamma(\widehat{\mathcal{M}}) \setminus \Gamma(\mathcal{M}^R)) \cup (\Gamma(\mathcal{M}^R) \setminus \Gamma(\widehat{\mathcal{M}})))$ *it is guaranteed to have* $C(\pi^*, \hat{\mathcal{M}}) = C^*_{\hat{\mathcal{M}}}$.

This is similar to the model-space search for MCEs, but this time starting from the robot's model $\mathcal{M}^R$ instead. The goal here is to find the largest set of model changes for which the explicability criterion becomes invalid for the first time (due to either suboptimality or inexecutability). This requires a search over the entire model space (Algorithm 6). We can leverage Proposition 3 to reduce our search space. Starting from $\mathcal{M}^R$, given a set of model changes $\mathcal{E}$ where $\delta_{\mathcal{M}_R, \mathcal{M}_H}(\Gamma(\mathcal{M}^R), \mathcal{E}) = \Gamma(\widehat{\mathcal{M}})$ and $C(\pi^*, \widehat{\mathcal{M}}) > C^*_{\widehat{\mathcal{M}}}$, no superset of $\mathcal{E}$ can lead to an MME solution. In Algorithm 6, we keep track of such unhelpful model changes in the list h_list. The variable $\mathcal{E}^{MME}$ keeps track of the current best list of model changes. Whenever the search finds a new set of model changes where $\pi^*$ is optimal and is larger than $\mathcal{E}^{MME}$, $\mathcal{E}^{MME}$ is updated with $\mathcal{E}$. The resulting MME is all the possible model changes that did not appear in $\mathcal{E}^{MME}$.

Figure 5.3 contrasts MCE search with MME search. MCE search starts from $\mathcal{M}^R_h$, updates $\widehat{\mathcal{M}}$ towards $\mathcal{M}^R$ and returns the first node (indicated in orange) where $C(\pi^*, \widehat{\mathcal{M}}) = C^*_{\widehat{\mathcal{M}}}$. MME search starts from $\mathcal{M}^R$ and moves towards $\mathcal{M}^R_h$. It finds the longest path (indicated in blue) where $C(\pi^*, \widehat{\mathcal{M}}) = C^*_{\widehat{\mathcal{M}}}$ for all $\widehat{\mathcal{M}}$ in the path. The MME (green) is the rest of the path towards $\mathcal{M}^R_h$.

## 5.3 Approximate Explanations

In this section, we will discuss how some of the methods discussed in this section can be simplified to either 1) generate simpler explanation (in terms of the explanation size or

---

**Algorithm 4** Search for Minimally Complete Explanations

---

1: MCE-Search
2: *Input*: MRP $\langle \pi^*, \langle \mathcal{M}^R, \mathcal{M}_h^R \rangle \rangle$
3: *Output*: Explanation $\mathcal{E}^{MCE}$
4: *Procedure*:
5: fringe $\leftarrow$ `Priority_Queue()`
6: c_list $\leftarrow \{\}$ {Closed list}
7: $\pi_R^* \leftarrow \pi^*$ {Optimal plan being explained}
8: $\pi_H \leftarrow \pi$ such that $C(\pi, \mathcal{M}_h^R) = C^*_{\mathcal{M}_h^R}$ {Plan expected by human}
9: fringe.push($\langle \mathcal{M}_h^R, \{\} \rangle$, priority $= 0$)
10: **while** True **do**
11:    $\langle \widehat{\mathcal{M}}, \mathcal{E} \rangle, c \leftarrow$ fringe.pop($\widehat{\mathcal{M}}$)
12:    **if** $C(\pi_R^*, \widehat{\mathcal{M}}) = C^*_{\widehat{\mathcal{M}}}$ **then**
13:       return $\mathcal{E}$ {Returns $\mathcal{E}$ if $\pi_R^*$ optimal in $\widehat{\mathcal{M}}$}
14:    **else**
15:       c_list $\leftarrow$ c_list $\cup \widehat{\mathcal{M}}$
16:       **for** $f \in \Gamma(\widehat{\mathcal{M}}) \setminus \Gamma(\mathcal{M}^R)$ {Models that satisfy Condition 1} **do**
17:          $\lambda \leftarrow \langle 1, \{\widehat{\mathcal{M}}\}, \{\}, \{f\} \rangle$ {Removes f from $\widehat{\mathcal{M}}$}
18:          **if** $\delta_{\mathcal{M}_h^R, \mathcal{M}^R}(\Gamma(\widehat{\mathcal{M}}), \lambda) \notin$ c_list **then**
19:             fringe.push($\langle \delta_{\mathcal{M}_h^R, \mathcal{M}^R}(\Gamma(\widehat{\mathcal{M}}), \lambda), \mathcal{E} \cup \lambda \rangle, c+1$)
20:       **for** $f \in \Gamma(\mathcal{M}^R) \setminus \Gamma(\widehat{\mathcal{M}})$ {Models that satisfy Condition 2} **do**
21:          $\lambda \leftarrow \langle 1, \{\widehat{\mathcal{M}}\}, \{f\}, \{\} \rangle$ {Adds f to $\widehat{\mathcal{M}}$}
22:          **if** $\delta_{\mathcal{M}_h^R, \mathcal{M}^R}(\Gamma(\widehat{\mathcal{M}}), \lambda) \notin$ c_list **then**
23:             fringe.push($\langle \delta_{\mathcal{M}_h^R, \mathcal{M}^R}(\Gamma(\widehat{\mathcal{M}}), \lambda), \mathcal{E} \cup \lambda \rangle, c+1$)

---

**Algorithm 5** Selection strategy for identifying model updates relevant to the current plan

---

*Procedure* Priority_Queue.pop $\hat{\mathcal{M}}$
candidates $\leftarrow \{\langle \langle \widehat{\mathcal{M}}, \mathcal{E} \rangle, c^* \rangle \mid c^* = \arg\min_c \langle \langle \widehat{\mathcal{M}}, \mathcal{E} \rangle, c \rangle\}$
pruned_list $\leftarrow \{\}$
$\pi_H \leftarrow \pi$ such that $C(\pi, \hat{\mathcal{M}}) = C^*_{\hat{\mathcal{M}}}$

**for** $\langle \langle \widehat{\mathcal{M}}, \mathcal{E} \rangle, c \rangle \in$ candidates **do**
   **if** $\exists a \in \pi_R^* \cup \pi_H$ such that $\tau^{-1}((\Gamma(\widehat{\mathcal{M}}) \setminus \Gamma(\hat{\mathcal{M}})) \cup (\Gamma(\hat{\mathcal{M}}) \setminus \Gamma(\widehat{\mathcal{M}}))) \in \{C(a)\} \cup \text{pre}(a) \cup$
   $\text{adds}(a) \cup \text{dels}(a)$ **then**
      pruned_list $\leftarrow$ pruned_list $\cup \langle \langle \widehat{\mathcal{M}}, \mathcal{E} \rangle, c \rangle$
**if** pruned_list $= \phi$ **then**
   $\langle \widehat{\mathcal{M}}, \mathcal{E} \rangle, c \sim Unif(\text{candidate\_list})$
**else**
   $\langle \widehat{\mathcal{M}}, \mathcal{E} \rangle, c \sim Unif(\text{pruned\_list})$

---

---

**Algorithm 6** Search for Minimally Monotonic Explanations

---

1: MME-Search
2: *Input*: MRP $\langle \pi^*, \langle \mathcal{M}^R, \mathcal{M}_h^R \rangle \rangle$
3: *Output*: Explanation $\mathcal{E}^{MME}$
4: *Procedure*:
5: $\mathcal{E}^{MME} \leftarrow \{\}$
6: fringe $\leftarrow$ `Priority_Queue()`
7: c_list $\leftarrow \{\}$ {Closed list}
8: h_list $\leftarrow \{\}$ {List of incorrect model changes}
9: fringe.push($\langle \mathcal{M}^R, \{\} \rangle$, priority $= 0$)
10: **while** fringe is not empty **do**
11:     $\langle \widehat{\mathcal{M}}, \mathcal{E} \rangle, c \leftarrow$ fringe.pop($\widehat{\mathcal{M}}$)
12:     **if** $C(\pi^*, \widehat{\mathcal{M}}) > C_{\widehat{\mathcal{M}}}^*$ **then**
13:         h_list $\leftarrow$ h_list $\cup$ $((\Gamma(\widehat{\mathcal{M}}) \setminus \Gamma(\mathcal{M}^R)) \cup (\Gamma(\mathcal{M}^R) \setminus \Gamma(\widehat{\mathcal{M}})))$ {Updating h_list }
14:     **else**
15:         c_list $\leftarrow$ c_list $\cup \widehat{\mathcal{M}}$
16:         **for** $f \in \Gamma(\widehat{\mathcal{M}}) \setminus \Gamma(\mathcal{M}_h^R)$ {Models that satisfy Condition 1} **do**
17:             $\lambda \leftarrow \langle 1, \{\widehat{\mathcal{M}}\}, \{\}, \{f\} \rangle$ {Removes $f$ from $\widehat{\mathcal{M}}$}
18:             **if** $\delta_{\mathcal{M}^R, \mathcal{M}_h^R}(\Gamma(\widehat{\mathcal{M}}), \lambda) \notin$ c_list
                    **and** $\nexists S$ s.t. $((\Gamma(\widehat{\mathcal{M}}) \setminus \Gamma(\mathcal{M}^R)) \cup (\Gamma(\mathcal{M}^R) \setminus \Gamma(\widehat{\mathcal{M}}))) \supseteq S \in$ h_list **then**
19:                 fringe.push($\langle \delta_{\mathcal{M}^R, \mathcal{M}_h^R}(\Gamma(\widehat{\mathcal{M}}), \lambda), \mathcal{E} \cup \lambda \rangle$, $c + 1$)
20:                 $\mathcal{E}^{MME} \leftarrow \max_{|\cdot|} \{\mathcal{E}^{MME}, \mathcal{E}\}$
21:         **for** $f \in \Gamma(\mathcal{M}_h^R) \setminus \Gamma(\widehat{\mathcal{M}})$ {Models that satisfy Condition 2} **do**
22:             $\lambda \leftarrow \langle 1, \{\widehat{\mathcal{M}}\}, \{f\}, \{\} \rangle$ {Adds $f$ from $\widehat{\mathcal{M}}$}
23:             **if** $\delta_{\mathcal{M}^R, \mathcal{M}_h^R}(\Gamma(\widehat{\mathcal{M}}), \lambda) \notin$ c_list
                    **and** $\nexists S$ s.t. $((\Gamma(\widehat{\mathcal{M}}) \setminus \Gamma(\mathcal{M}^R)) \cup (\Gamma(\mathcal{M}^R) \setminus \Gamma(\widehat{\mathcal{M}}))) \supseteq S \in$ h_list **then**
24:                 fringe.push($\langle \delta_{\mathcal{M}^R, \mathcal{M}_h^R}(\Gamma(\widehat{\mathcal{M}}), \lambda), \mathcal{E} \cup \lambda \rangle$, $c + 1$)
25:                 $\mathcal{E}^{MME} \leftarrow \max_{|\cdot|} \{\mathcal{E}^{MME}, \mathcal{E}\}$
26: $\mathcal{E}^{MME} \leftarrow ((\Gamma(\widehat{\mathcal{M}}) \setminus \Gamma(\mathcal{M}^R)) \cup (\Gamma(\mathcal{M}^R) \setminus \Gamma(\widehat{\mathcal{M}}))) \setminus \mathcal{E}^{MME}$
27: return $\mathcal{E}^{MME}$

---

more directed to a specific user query) and 2) reduce the computational overhead of the search.

### 5.3.1 Explicit Contrastive Explanations

As mentioned earlier, the explanations generated by model reconciliation can be viewed as an answer to an implicit contrastive query. Although this may lead to user being provided more information than required. This may be particularly unnecessary if the user's original expected set of plans is much smaller than the set of all optimal plans. In these cases, we could let the user directly specify their expected set of plans in the form of explicit foils. Thus explanation can focus on establishing how the current plan compares against their original expected set of plans. So this allows the system to not only provide less information, the system can now potentially also provide additional information that can allow the user to understand why the current plan is better than the plans they were expecting.However in this case, the explanation could be a multi-step process where they raise additional queries (in the form of more foils) after each explanation.

Now if $\hat{\Pi}$ is the set of alternate foils raised by the human, the objective of a minimal explanation generation method would be

$$\mathcal{E}^{contr} = \arg\min_{\mathcal{E}} |(\Gamma(\widehat{\mathcal{M}})\backslash\Gamma(\mathcal{M}_h^R))\cup(\Gamma(\mathcal{M}_h^R)\backslash\Gamma(\widehat{\mathcal{M}}))| \text{ with } C(\pi^*,\widehat{\mathcal{M}}_h^R) \leq C(\pi',\widehat{\mathcal{M}}_h^R) \,\forall\pi' \in \hat{\Pi}$$

Now we can modify the MCE search to use this new objective to identify the required explanation. Now the foil set $\hat{\Pi}$ may be either explicitly specified or may be implicitly specified in terms of constraints satisfied by plans in the set.

### 5.3.2 Approximate MCE

Both MCEs and MMEs may be hard to compute - in the worst case it involves a search over the entire space of model differences. Thus the biggest bottleneck here is the check for optimality of a plan given a new model. A check for necessary or sufficient conditions for optimality, without actually computing optimal plans can be used as a way to further prune the search tree.

In the following section, we investigate an approximation to an MCE by employing a few simple proxies to the optimality test. By doing this we lose the completeness guarantee but improve computability. Specifically, we replace the equality test in line 12 of Algorithm 4 by the following rules –

C(1) $\delta_{\widehat{\mathcal{M}}}(\widehat{I},\pi_R^*) \models \widehat{G}$; **and**

C(2) $C(\pi_R^*,\widehat{\mathcal{M}}) < C(\pi_R^*,\mathcal{M}_h^R)$ **or** $\delta_{\widehat{\mathcal{M}}}(\widehat{I},\pi_H^*) \not\models \widehat{G}$; **and**

C(3) Each action contributes at least one causal link to $\pi_R^*$.

**Fig. 5.4:** Interface for a user study where participants assumed the role of the external commander and evaluated plans provided by the internal robot. They could request for plans and explanations to those plans (if not satisfied) and rate them as optimal or suboptimal or (if unsatisfied) can chose to pass.

C(1) ensures that the plan $\pi_R^*$ originally computed is actually valid in the new model. C(2) requires that this plan has either become better in the new model or at least that the human's expected plan $\pi_H^*$ has been disproved. Finally, C(3), ensures that for each action $a_i \in \pi_R^*$ there exists an effect $p$ that satisfies the precondition of at least one action $a_k$ (where $a_i \prec a_k$) and there exists no action $a_j$ (where $a_i \prec a_j \prec a_k$) such that $p \in \mathrm{dels}(a_j)$. Such explanations are only able to preserve local properties of a plan and hence incomplete.

**Proposition 6.** *C(3) is a necessary condition for optimality of $\pi^*$ in $\widehat{\mathcal{M}}$.*

*Proof Sketch.* Assume that for an optimal plan $\pi_R^*$, there exists an action $a_i$ where criterion (3) is not met. Now we can rewrite $\pi_R^*$ as $\pi_R' = \langle a_0, a_1, \ldots, a_{i-1}, a_i, a_{i+1}, \ldots, a_n, a_{n+1}\rangle$, where $\mathrm{pre}(a_0) = \phi$ and $\mathrm{adds}(a_0) = \{I\}$ and $\mathrm{pre}(a_{n+1}) = \{G\}$ and $\mathrm{adds}(a_{n+1}) = \mathrm{dels}(a_{n+1}) = \phi$. It is easy to see that $\delta_{\widehat{\mathcal{M}}}(\phi, \pi_R') \models G$. Now let us consider a cheaper plan $\hat{\pi_R'} = \langle a_0, a_1, \ldots, a_{i-1}, a_{i+1}, \ldots, a_n, a_{n+1}\rangle$. Since $a_i$ does not contribute any causal links to the original plan $\pi_R^*$, we will also have $\delta_{\widehat{\mathcal{M}}}(\phi, \hat{\pi_R'}) \models G$. This contradicts our original assumption of $\pi_R^*$ being optimal, hence proved.                                                   □

## 5.4   User Studies

Apart from the underlying formal theory and intuition for the methods discussed here, user studies have also been performed to explicitly test the effectiveness of these type of explanations. In one study, the participants were shown a simulated robot navigating a floor in a building, and the floor map was made available to the user. They were then asked to

**Fig. 5.5:** Interface for Study-1: Participants assumed the role of the internal agent and explained their plans to a teammate with a possibly different map of the world.

evaluate whether the robot plans were valid and optimal. The participants had the option to query for explanations that revealed details about the floor map that were previously unknown to the user. Figure 5.4 presents a screenshot of this interface along with a sample plan. The user study also had a reward structure that disincentivized participants from providing incorrect answers (in regards to whether the plans were valid and/or optimal). The study showed the majority of participants were correctly able to identify the validity and optimality of the plan based on the explanations. A subjective survey at the end of the study showed that the majority of participants found the explanations were easy to understand and useful.

Additional studies were also performed that placed participants in the role of the explainer. In this study, they were asked to explain why a given plan was optimal. They were shown the original map and the map the explainee would see. After analyzing the explanations given by the participants, it was found that most of them perfectly matched up with many of the explanation types studied in this chapter. Particularly, when the participants were asked to minimize the information to be provided the majority of explanation provided were MCEs. This points to the fact that in many cases, people do naturally calculate explanations of the type discussed in this chapter.

## 5.5 Other Explanatory Methods

In terms of explanations, this chapter focused on one specific type of explanation, namely explanation as model reconciliation. But this is not the only type of explanation that has been studied in the literature. In particular, we will look at two other types of explanations that are popular within the explainable planning community (generally referred to as XAIP), and in this section, we will look at how these explanation types could be viewed in the context of human-aware planning.

**Explanations to Address Inferential Asymmetry.** The first type of explanation is one designed to address explanatory queries rising from differences in inferential capabilities of the human and the robot. Going back to the setting where the human is a pure

observer, even if the human had access to the exact model, there is no guarantee that they would be able to correctly evaluate the validity or optimality of the plan being followed by the robot given potential inferential limitations of the human. So in addition to reconciling any model differences, the system may need to help address any confusion on the human's end stemming from such asymmetry in inferential capabilities. Some general strategies for addressing such considerations include,

1. Allowing users to ask questions about plans: In particularly allowing for contrastive queries, in which the human raises specific alternatives they were expecting Sreedharan et al. [2018b]. As mentioned earlier, such questions help expose the human's expectations about the robot and thus allowing the robot to provide a more directed explanation. Having access to the human's expected behavior, also means the robot could provide additional information that helps contrast the given decision against the alternatives the human had in mind (say by tracing the execution of the current plan and the alternative).

2. Simplifying the problem: The next strategy one could employ may be to simplify the problem being explained. One of the obvious strategy we can employ here is to use state abstractions, we will look at this strategy in more detail in Chapter 6 though in a slightly different context. But this is one of many problem simplification strategies we could employ. Other strategies include the use of local approximation (also discussed in Chapter 8 (Section 8.2.1)), decomposing the problem into smaller subproblems (for example, focusing on the achievement of a subgoal instead of the true goal as done in the case of Sreedharan et al. [2019b]) and even simplifying the class of the reasoning/planning problem. For example, in many cases it may suffice to focus on a determinization or a qualitative non-deterministic representation of a problem to address many of the explanatory challenges related to stochastic planning problems.

3. Providing explanatory witness: In many cases while the explanation may be in regards to a space of possible solutions, it may be helpful to provide the human with specific example to how the current solution may satisfy a property that is in question. For example, if the human were to raise a partial specification of an alternative, say she asks why the robot didn't choose a particular action at a state, then the robot could point to a specific example plan that includes an action and point out that this alternate plan is costlier than the current plan (this is the strategy followed by Krarup et al. [2019]). This doesn't necessarily explain why all plans that include that action may be worse off than the current plan, but providing an example could help the user better understand the plan in question.

Note that all the above strategies are generally based on computational intuition on how to reduce computational overhead of a problem (many of these methods have their roots on heuristics for planning). On the other hand, a systematic framework to address such asymmetry would require us to correctly model and capture the human inferential capabilities, so that the agent can predict the inferential overhead placed on the human by a specific explanation approach. Unfortunately, we are still far from achieving such accurate modeling of the computational capabilities of a bounded rational agent let alone a human.

**Plan and Policy Summary**   Another category of work that we haven't had a chance to discuss in the book are the ones related to plan and policy summaries. Particularly in cases, where the plan may need to be approved by the human before it is executed. In such cases, the robot would need to communicate its plans to the human before it can be executed. Several works have looked at the problem of effectively summarizing a generated course of action. Many of these works focus on stochastic planning settings, where the solution concept takes the form of policies which could in general be harder to understand than traditional sequential plans. Some general strategies used in such policy summarization method include, the use of abstractions (for example Topin and Veloso [2019] use state abstraction, while Sreedharan et al. [2020d] use temporal abstraction) and selecting some representative state-action pairs, so the human can reconstruct the rest of the policy (for example Amir and Amir [2018], Lage et al. [2019]). Generally most works in this latter direction assume that the human correctly understands the task at hand and thus the reconstruction can be done correctly by the human without the aid of any model reconciliation. One may also need to address the problem of vocabulary mismatch in the context of policy communication, wherein the policy may need to be translated first into terms the user can understand (as done by the work Hayes and Shah [2017]).

**Other Popular Forms of Explanation Categorization.**   As final note on other forms of explanation studied in XAIP, a popular way to categorize explanatory information with in the literature is to organize them based on the type of explanatory questions they may be addressing. Fox et al. [2017] provides a very comprehensive list of explanatory questions that could be relevant to planning. This work not only includes question related to specific plans that are generated by the agent, but also questions related to the planning process as a whole. For example, the paper considers questions like why the agent decided to replan at a particular point during the plan execution.

## 5.6   Bibliographic Remarks

The idea of explanation as model reconciliation was first introduced in the paper Chakraborti et al. [2017a]. The paper specifically looked at the application of the methods for STRIPS style models. The user studies for evaluating the explanations were presented in the paper Chakraborti et al. [2019b]. The paper looked at three types of model-reconciliation explanations, namely MCE, MPE and PPE. In the case of PPE, a slightly modified version of PPE was considered that also takes into account one of the possible optimal plans in the human model (thereby ensuring the explanation is still contrastive). Sreedharan et al. [2021a] also presents a more unified view of the various earlier works done in the area of model-reconciliation explanation. While most works in model reconciliation explanations have focused on deterministic STRIPS-style planning models, researchers have also started looking at other planning formalisms. For example, Sreedharan et al. [2019a] looked at model reconciliation in the context of MDPs and Vasileiou et al. [2021] looked at model reconciliation in the context of logical programs.

There are even parallel works on explanation from fields like machine learning, that could be understood as a special case of model reconciliation. A particularly significant

example being LIME [Ribeiro et al., 2016]. In this work, the authors propose a model-agnostic method that take a model, and the specific instance to be explained and generates an interpretable linear model that locally approximates the dynamics of the current instance. In this case the final explanation takes the form of details about the linear model and the human's model is taken to be empty. One important factor to note here is the model being transferred to the user is not the original model but rather a local approximation of it. Moreover, the model is post-hoc expressed in terms of human understandable features. These strategies are also applicable in the context of model reconciliation explanation in the sequential decision making settings, and we will discuss methods that employ such methods in chapter 8. Some examples of explanation works that use abstraction of the model and/or plan include Zahavy et al. [2016], Topin and Veloso [2019], and Sreedharan et al. [2020d]. Another method popular in explaining the utility of a specific action in plan involve providing causal chain contributed by that action to the overall plan [Veloso, 1992, Seegebarth et al., 2012, Bercher et al., 2014]. Another technique popular in explaining multi-objective planning problem is to provide conflicts between the different objectives Eifler et al. [2020]. For more works done in this direction, readers can refer to the survey in Chakraborti et al. [2020]. Miller [2019] presents a concise introduction into various works from disciplines like social sciences, psychology and philosophy that have looked at explanations.

CHAPTER 6

# Acquiring Mental Models for Explanations

The previous chapter sketches out some of the central ideas behind generating explanation as model reconciliation, but it does so while making some strong assumptions. Particularly, the setting assumes that the human's model of the robot is known exactly upfront. In this chapter, we will look at how we can relax this assumption and see how we can perform model reconciliation in scenarios where the robot has progressively less information about the human mental model. We will start by investigating how the robot can perform model reconciliation with incomplete model information. Next we will look at cases where the robot doesn't have a human mental model but can collect feedback from users.[1] We will see how we can use such feedback to learn simple labeling models that suffice to generate explanations. We will also look at generating model reconciliation explanations by assuming the human has a simpler mental model, specifically one that is an abstraction of the original model and also see how this method can help reduce inferential burden on the human. Throughout this chapter, we will focus on generating MCE though most of the methods discussed here could also be extended to MME and contrastive versions of the explanations.

## 6.1   The Urban Search and Reconnaissance Domain

We will concretize the discussion in this section in a typical Urban Search and Reconnaissance (USAR) domain where a remote robot is put into disaster response operation often controlled partly or fully by an external human commander, as shown in Figure 6.1. The robot's job is to scout areas that may be otherwise harmful to humans and report on its surroundings as instructed by the external supervisor. The scenario can also have other internal agents (humans or robots) with whom the robot needs to coordinate.

Here, even though all agents start off with the same model – i.e. the blueprint of the building – their models diverge as the internal agent interacts with the scene. Due to the disaster, new paths may have opened up due to collapsed walls or old paths may no longer be available due to rubble. This means that plans that are valid and optimal in the robot's (current) model may not make sense to the external commander. In the scenario in Figure 6.2, the robot is tasked to go from its current location marked blue to conduct reconnaissance in the location marked orange. The green path is most optimal in its current model but this is blocked in the external's mental model while the expected plan in the mental model is no longer possible due to rubble. Without removing rubble in the blocked paths, the robot can instead communicate that the path at the bottom is no

---

[1]Note that the model we need for explicability and explanation, $\mathcal{M}_h^R$ can't just be learned in general from past behavior traces of the humans themselves, since those give us information only about $\mathcal{M}^H$ and not $\mathcal{M}_h^R$. While there may be self-centered humans who consider the AI robot's model to be same as theirs, more generally, these models will be different.

**Fig. 6.1:** A typical USAR domain with an internal robot and an external commander.



**Fig. 6.2:** Model differences in the USAR domain.

longer blocked. This explanation preserves the validity and optimality of its plan in the updated model (even though further differences exist).

## 6.2   Model Uncertainty

We will start by considering cases where the human's mental model may not be known exactly upfront. There may be parts of the human model the robot may know exactly, while it may be uncertain about others. We will leverage the annotated model representation to capture such incomplete models. Under this representation, in addition to the standard preconditions and effects associated with actions, the model includes *possible* preconditions and effects which may or may not be realized in practice.

**Definition 27.** *An **annotated model** is the tuple $\mathbb{M} = \langle F, \mathbb{A}, \mathbb{I}, \mathbb{G}, C \rangle$ where $F$ is a finite set of fluents that define a state $s \subseteq F$, and $\mathbb{A}$ is a finite set of annotated actions, $\mathbb{I} = \langle I^0, I^+ \rangle$, $\mathbb{G} = \langle G^0, G^+ \rangle$; the annotated initial and goal states ( such that $I^0, G^0, I^+, G^+ \subseteq F$) and $C$ the cost function. Action $a \in \mathbb{A}$ is a tuple $\langle pre(a), \widetilde{pre}(a), adds(a), dels(a), \widetilde{adds}(a), \widetilde{dels}(a) \rangle$*

*where in addition to $pre(a), adds(a), dels(a) \subseteq F$, i.e., the known preconditions and add/delete effects each action also contains* possible preconditions $\widetilde{pre}(a) \subseteq F$ *containing propositions that it* might *need as preconditions, and* possible add (delete) effects $\widetilde{adds}(a), \widetilde{dels}(a), \subseteq F$) *containing propositions that it* might *add (delete, respectively) after execution. $I^0, G^0$ (and $I^+, G^+$) are the known (and possible) parts of the initial and goal states.*

Each possible condition $f \in \widetilde{pre}(a) \cup \widetilde{adds}(a) \cup \widetilde{dels}(a)$ has an associated probability $p(f)$ denoting how likely it is to be a known condition in the ground truth model – i.e. $p(f)$ measures the confidence with which that condition has been learned. The sets of known and possible conditions of a model $\mathcal{M}$ are denoted by $\mathbb{S}_k(\mathcal{M})$ and $\mathbb{S}_p(\mathcal{M})$ respectively. An *instantiation* of an annotated model $\mathbb{M}$ is a classical planning model where a subset of the possible conditions have been realized, and is thus given by the tuple $inst(\mathbb{M}) = \langle F, A, I, G \rangle$, initial and goal states $\mathcal{I} = \mathcal{I}^0 \cup \chi;\ \chi \subseteq \mathcal{I}^+$ and $\mathcal{G} = \mathcal{G}^0 \cup \chi;\ \chi \subseteq \mathcal{G}^+$ respectively, and action $A \ni a = \langle pre(a) \leftarrow pre(a) \cup \chi;\ \chi \subseteq \widetilde{pre}(a), adds(a) \leftarrow adds(a) \cup \chi;\ \chi \subseteq \widetilde{adds}(a), dels(a) \leftarrow dels(a) \cup \chi;\ \chi \subseteq \widetilde{dels}(a) \rangle$. Clearly, given an annotated model with $k$ possible conditions, there may be $2^k$ such instantiations, which forms its *completion set.*

**Definition 28. Likelihood** $P_\ell$ *of instantiation $inst(\mathbb{M})$ of an annotated model $\mathbb{M}$ is:*

$$P_\ell(inst(\mathbb{M})) = \prod_{f \in \mathbb{S}_p(\mathbb{M}) \wedge \mathbb{S}_k(inst(\mathbb{M}))} p(f) \quad \times \prod_{f \in \mathbb{S}_p(\mathbb{M}) \setminus \mathbb{S}_k(inst(\mathbb{M}))} (1 - p(f))$$

As discussed before, such models turn out to be especially useful for the representation of human (mental) models learned from observations, where uncertainty after the learning process can be represented in terms of model annotations. Let $\mathbb{M}_h^R$ be the culmination of a model learning process and $\{\mathcal{M}_{h_i}^R\}_i$ be the completion set of $\mathbb{M}_h^R$. One of these models is the actual ground truth (i.e. the human's real mental model). We refer to this as $g(\mathbb{M}_h^R)$. We will explore now how this representation will allow us to compute *conformant explanations* that can explain with respect to all possible mental models and *conditional explanations* that engage the explainee in dialogue to minimize the size of the completion set to compute shorter explanations.

### 6.2.0.1 Conformant Explanations

In this situation, the robot can try to compute MCEs for each possible configuration. However, this can result in situations where the explanations computed for individual models independently are not consistent across all possible target domains. Thus, in the case of model uncertainty, such an approach cannot guarantee that the resulting explanation will be acceptable.

Instead, the objective is to find an explanation such that $\forall i\ \pi^*_{\widehat{\mathcal{M}}_{h_i}^R} \equiv \pi^*_{\mathcal{M}^R}$ (as shown in Figure 6.3). This is a single set of model updates that makes the given plan optimal (and hence explained) in all the updated models. At first glance, it appears that such an approach, even though desirable, might turn out to be prohibitively expensive especially since solving for a *single* MCE involves search in the model space where each search node

**Fig. 6.3:** Model reconciliation in the presence of model uncertainty or multiple explainees.

is an optimal planning problem. However, it turns out that the same search strategy can be employed here as well by representing the human mental model as an *annotated* model. The optimality condition for MCE now becomes –

$$C(\pi, g(\mathbb{M}_h^R)) = C^*_{g(\mathbb{M}_h^R)}$$

We define *robustness* of an explanation for an incomplete mental model as the probability mass of models where it is a valid explanation.

**Definition 29. Robustness** *of an explanation $\mathcal{E}$ is given by* –

$$R(\mathcal{E}) = \sum_{inst(\widehat{\mathcal{M}}_h^R) \ s.t. \ C(\pi, inst(\widehat{\mathcal{M}}_h^R)) = C^*_{inst(\widehat{\mathcal{M}}_h^R)}} P_\ell(inst(\widehat{\mathcal{M}}_h^R))$$

**Definition 30.** *A* **conformant explanation** *is such that $R(\mathcal{E}) = 1$.*

Which is equivalent to saying that conformant explanation ensures that the given plan is explained in all the models in the completion set of the human model. Let's look at an example. Consider again the USAR domain (Figure 6.4), the robot is now at P1 (blue) and needs to collect data from P5. While the commander understands the goal, she is under the false impression that the paths from P1 to P9 and P4 to P5 are unusable (red question marks). She is also unaware of the robot's inability to use its hands. On the

**Fig. 6.4:** Back to our USAR scenario: the robot plan is marked in blue and uncertain parts of the human model is marked with red question marks.

other hand, while the robot does not have a complete picture of her mental model, it understands that any differences between the models are related to (1) the path from P1 to P9; (2) the path from P4 to P5; (3) its ability to use its hands; and (4) whether it needs its arm to clear rubble. Thus, from the robot's perspective, the mental model can be one of sixteen possible models (one of which is the actual one). Here, a conformant explanation for the optimal robot plan (blue) is as follows –

```
Explanation >> remove-known-INIT-has-add-effect-hand_capable
Explanation >> add-annot-clear_passage-has-precondition-hand_capable
Explanation >> remove-annot-INIT-has-add-effect-clear_path P1 P9
```

### 6.2.0.2 Model-Space Search for Conformant Explanations

As we discussed before, we cannot launch an MCE-search for each possible mental model separately, both for issues of complexity and consistency of the solutions. However, in the following discussion, we will see how we can reuse the model space search from the previous section with a compilation trick.

We begin by defining two models – the most relaxed model possible $\mathcal{M}_{max}$ and the least relaxed one $\mathcal{M}_{min}$. The former is the model where all the possible add effects and none of the possible preconditions and deletes hold, the state has all the possible conditions set to true, and the goal is the smallest one possible; while in the latter all the possible preconditions and deletes and none of the possible adds are realized and with the minimal start state and the maximal goal. This means that, if a plan is executable in $\mathcal{M}_{min}$ it will be executable in all the possible models. Also, if this plan is optimal in $\mathcal{M}_{max}$, then it must be optimal throughout the set. Of course, such a plan may not exist, and we are not trying to find one either. Instead, we are trying to find a set of model updates which when applied to the annotated model, produce a new set of models where a given plan is optimal. In providing these model updates, we are in effect reducing the set of possible models to a smaller set. The new set need not be a subset of the original set of models but will be equal or smaller in size to the original set. For any given annotated model, such an explanation always exists (entire model difference in the worst case), and the goal here becomes to find the smallest one. $\mathbb{M}_h^R$ thus affords the following two models –

$$\mathcal{M}_{max} = \langle F, A, I, G \rangle$$

- initial state $I \leftarrow I^0 \cup I^+$; given $\mathbb{I}$
- goal state $G \leftarrow G^0$; given $\mathbb{G}$
- $\forall a \in A$

  - $\text{pre}(a) \leftarrow \text{pre}(a)$; $a \in \mathbb{A}$
  - $\text{adds}(a) \leftarrow \text{adds}(a) \cup \widetilde{\text{adds}}(a)$; $a \in \mathbb{A}$
  - $\text{dels}(a) \leftarrow \text{dels}(a)$; $a \in \mathbb{A}$

$$\mathcal{M}_{min} = \langle F, A, I, G \rangle$$

- initial state $I \leftarrow I^0$; given $\mathbb{I}$
- goal state $G \leftarrow G^0 \cup G^+$; given $\mathbb{G}$
- $\forall a \in A$

  - $\text{pre}(a) \leftarrow \text{pre}(a) \cup \widetilde{\text{pre}}(a)$; $a \in \mathbb{A}$
  - $\text{adds}(a) \leftarrow \text{adds}(a)$; $a \in \mathbb{A}$
  - $\text{dels} \leftarrow \text{dels} \cup \widetilde{\text{dels}}$; $a \in \mathbb{A}$

As explained before, $\mathcal{M}_{max}$ is a model where all the add effects hold and it is easiest to achieve the goal, and similarly $\mathcal{M}_{min}$ is the model where it is the hardest to achieve the goal. These definitions might end up creating inconsistencies (e.g. in an annotated `BlocksWorld` domain, the `unstack` action may have add effects to make the block both `holding` and `ontable` at the same time), but the model reconciliation process will take care of these.

**Proposition 7.** *For a given MRP $\Psi = \langle \pi, \langle \mathcal{M}^R, \mathbb{M}_h^R \rangle \rangle$, if the plan $\pi$ is optimal in $\mathcal{M}_{max}$ and executable in $\mathcal{M}_{min}$, then the plan is optimal for all $i$.*

This now becomes the new criterion to satisfy in the course of search for an MCE for a set of models. We again reuse the state representation in Chapter 5 (generated by $\Gamma$ as described in Section 5.1). We start the conformant search by first creating the corresponding $\mathcal{M}_{max}$ and $\mathcal{M}_{min}$ model for the given annotated model $\mathbb{M}_h^R$. While the goal test for the original MCE only included an optimality test, here we need to both check the optimality of the plan in $\mathcal{M}_{max}$ and verify the correctness of the plan in $\mathcal{M}_{min}$. As stated in Proposition 7, the plan is only optimal in the entire set of possible models if it satisfies both tests. Since the correctness of a given plan can be verified in polynomial time with respect to the plan size, this is a relatively easy test to perform.

The other important point of difference between the algorithm mentioned above and the original MCE is how the applicable model updates are calculated. Here we consider the superset of model differences between the robot model and $\mathcal{M}_{min}$ and the differences between the robot model and $\mathcal{M}_{max}$. This could potentially mean that the search might end up applying a model update that is already satisfied in one of the models but not in

the other. Since all the model update actions are formulated as set operations, the original MRP formulation can handle this without any further changes. The models obtained by applying the model update to $\mathcal{M}_{min}$ and $\mathcal{M}_{max}$ are then pushed to the open queue.

**Proposition 8.** *$\mathcal{M}_{max}$ and $\mathcal{M}_{min}$ only need to be computed once – i.e. with a model update $\mathcal{E}$ to $\mathbb{M}$: $\mathcal{M}_{max} \leftarrow \mathcal{M}_{max} + \mathcal{E}$ and $\mathcal{M}_{min} \leftarrow \mathcal{M}_{min} + \mathcal{E}$.*

These models form the new $\mathcal{M}_{min}$ and $\mathcal{M}_{max}$ models for the set of models obtained by applying the current set of model updates to the original annotated model. This proposition ensures that we no longer have to keep track of the current list of models or recalculate $\mathcal{M}_{min}$ and $\mathcal{M}_{max}$ for the new set.

### 6.2.0.3   Conditional Explanations

Conformant explanations can contain superfluous information – i.e. asking the human to remove non-existent conditions or add existing ones. In the previous example, the second explanation (regarding the need of the hand to clear rubble) was already known to the human and was thus superfluous information. Such redundant information can be annoying and may end up reducing the human's trust in the robot. This can be avoided by –

- Increasing the cost of model updates involving uncertain conditions relative to those involving known preconditions or effects. This ensures that the search prefers explanations that contain known conditions. By definition, such explanations will not have superfluous information.

- However, sometimes such explanations may not exist. Instead, we can convert conformant explanations into *conditional* ones. This can be achieved by turning each model update for an annotated condition into a question and only provide an explanation if the human's response warrants it – e.g. instead of asking the human to update the precondition of `clear_passage`, the robot can first ask if the human thinks that action has a precondition `hand_usable`. Thus, one way of removing superfluous explanations is to reduce the size of the completion set by gathering information from the human.

By using information gathering actions, we can do even better than just remove redundant information from an already generated conformant explanation. For example, consider the following exchange in the USAR scenario –

```
R : Are you aware that the path from P1 to P4 has collapsed?
H : Yes.
> R realizes the plan is optimal in all possible models.
> It does not need to explain further.
```

If the robot knew that the human thought that the path from P1 to P4 was collapsed, it would know that the robot's plan is already optimal in the human mental model and

hence be required to provide no further explanation. This form of explanations can thus clearly be used to cut down on the size of conformant explanations by reducing the size of the completion set.

**Definition 31.** *A* **conditional explanation** *is represented by a policy that maps the annotated model (represented by a $\mathcal{M}_{min}$ and $\mathcal{M}_{max}$ model pair) to either a question regarding the existence of a condition in the human ground model or a model update request. The resultant annotated model is produced, by either applying the model update directly into the current model or by updating the model to conform to human's answer regarding the existence of the condition.*

In asking questions such as these, the robot is trying to exploit the human's (lack of) knowledge of the problem in order to provide more concise explanations. This can be construed as a case of lying by omission and can raise interesting ethical considerations (we will look at the role of lies in team in more detail in Chapter 9, Section 9.3). Humans, during an explanation process, tend to undergo this same "selection" process as well in determining which of the many reasons that could explain an event is worth highlighting.

**Modified $AO^*$-search to find Conditional Explanations** We can generate conditional explanations by either performing post-processing on conformant explanations or by performing an AND-OR graph search with $AO^*$ Nilsson [2014]. Here each model update related to a known condition forms an OR successor node while each *possible* condition can be applied on the current state to produce a pair of AND successors, where the first node reflects a node where the annotated condition holds while the second one represents the state where it does not. So the number of possible conditions reduces by one in each one of these AND successor nodes. This AND successor relates to the answers the human could potentially provide when asked about the existence of that particular possible condition. Note that this AND-OR graph will not contain any cycles as we only provide model updates that are consistent with the robot model and hence we can directly use the $AO^*$ search here.

Unfortunately, if we used the standard $AO^*$ search, it will not produce a conditional explanation that contains this "less robust" explanation as one of the potential branches in the conditional explanation. This is because, in the above example, if the human had said that the path was free, the robot would need to revert to the original conformant explanation. Thus the cost of the subtree containing this solution will be higher than the one that only includes the original conformant explanation.

To overcome this shortcoming, we can use a discounted version of the $AO^*$ search where the cost contributed by a pair of AND successors is calculated as –

$$min(node1.h\_val,\ node2.h\_val) + \gamma * max(node1.h\_val,\ node2.h\_val)$$

where node1 and node2 are the successor nodes and node1.*h_val*, node2.*h_val* are their respective $h$-values. Here $\gamma$ represents the discount fact and controls how much the search values short paths in its solution subtree. When $\gamma = 1$, the search becomes standard $AO^*$ search and when $\gamma = 0$, the search myopically optimizes for short branches (at the cost

of the depth of the solution subtree). The rest of the algorithm stays the same as the standard *AO** search.

### 6.2.0.4  Anytime Explanations

Both the algorithms discussed above can be computationally expensive, in spite of the compilation trick to reduce the set of possible models to two representative models. However, as we did previously with MCEs, we can also aim for an approximate solution by relaxing the minimality requirement of explanation to achieve much shorter explanation generation time when required. For this we use an anytime depth first explanation generation algorithm. Here, for each state, the successor states include all the nodes that can be generated by applying the model edit actions on all the known predicates and two possible successors for each possible condition – one where the condition holds and one where it does not. Once the search reaches a goal state (a new model where the target plan is optimal throughout its completion set), it queries the human to see if the assumptions it has made regarding possible conditions hold in the human mental model (the list of model updates made related to possible conditions). If all the assumptions hold in the human model, then we return the current solution as the final explanation (or use the answers to look for smaller explanations), else continue the search after pruning the search space using the answers provided by the human. Such approaches may also be able to facilitate iterative presentation of model reconciliation explanations to the human.

## 6.3   Model-Free Explanations

The previous section focused on cases where there still exists an estimate (incomplete though) of the human's knowledge about the robot. The conciseness of the explanation generated using this method would still rely on how accurate and complete that estimate is. In many cases the robot may start with no such estimate. While technically one could start with an incomplete model that will cover the space of all possible models representable using the current set of action and fluent set, it is unlikely that such a general human mental model would lead to concise explanations. One possible alternative is to consider a plausible assumption about the human model (such as assuming its an abstraction of the robot model or even that the human is completely ignorant of the robot model). Such assumptions are more effective when the user is ready to interact with the system while receiving the explanation. A slightly different possibility might be to learn the human's model. This approach may be particularly well suited when it is possible to assume that it is possible for the system to collect data from a set of people about their ability to make sense of system decisions and the final end user would have similar mental model to the set of users. Note that this assumption is similar to the one made in model-free explicable plan generation in Chapter 3 (Section 3.3) and we will make use of a similar approach to learn simple labeling models that are able to predict whether a human would find a given explanation satisfactory.

Unlike the setting in learning for explicability, we will assume the agent now also has access to a set of explanatory messages $\mu$. Note that such messages only contain

information about the original robot model and are independent of what the human may or may not know. So assuming a STRIPS model of representation for explanation, this could contain messages about what propositions are or aren't part of the robot model. As such $\mu$ could be derived from $\Gamma$ (i.e the mapping from the model to a set of propositions) discussed in earlier chapter.

Once we have access to such a set of messages our goal is now quite similar to the one described in chapter 3, as we would want to learn a labelling model over the steps in the plan. To keep the discussion simple, we will assume the labeling model here merely captures the fact whether the user find the plan step explicable or not (and will not worry about the exact task label they may associate with the step). So the labeling function we use in this context will look something like

$$\mathcal{L}_{\pi^R} : \mathbb{F} \times 2^\mu \to \{0, 1\}$$

Where $\mathbb{F}$ is the space of features for the plan steps. Note that unlike previous learned labeling model, here the function actually captures the fact that the human's expectation could shift once she is given addition information about the robot model. Now we can again take individual labels of the steps and calculate an explicability score for the entire plan using a function $f_{exp}$, where $f_{exp}$ returns 1 if all the steps are labeled as explicable, which means, according to the labeling model, the plan is perfectly explicable. Thus the learned model labeling model allows us to evaluate the explicability of the plan at each modified model. This means that we can now run the model space search even without having access to the original human mental model $\mathcal{M}_h^R$. Here the possible model space is represented by the space of possible message sets the robot can provide. Each possible modified model in the original scenario should map to at least one message set in this new setting. Instead of running an optimality check, we can check the labeling model to see if for a given set of messages every step in the plan would be labeled as being explicable. In this new setting, one could calculate an approximation of MCE for a plan $\pi_R$ as follows

$$\underset{m \subseteq \mu}{\arg\min}\, C(m) \text{ such that } f_{exp}(\pi_R, m) = 1$$

Where $C$ gives the cost associated with the set of messages $m$. The formulation looks for the cheapest set of messages that would render the entire plan explicable. One could similarly use such labeling models to also try extracting explanations that meet the requirements of other forms of explanations. Of course, the effectiveness of this method relies on the ability to learn good labeling models. Instead of making the explicability score a hard constraint, it can also be made part of the objective function.

## 6.4   Assuming Prototypical Models

All the previous discussion focused on cases where we handled lack of information about human model, by either trying to use uncertain model or learning simple model alternatives. But in many cases such uncertain models may not be available or the system may

not have previous data or be able to interact with the human for a long enough time to learn useful models. An alternative possibility might be to assume a simple prototypical model for the human and then use it to generate the explanation. In this section, we will consider such a possibility, specifically we will look at cases where we can perform model reconciliation by assuming human model is some abstraction of the robot model (which includes the possibility that its just an empty model). We will look at the kind of explanatory queries such assumptions can support and how we could extend them to cover all the cases.

In particular, we will consider cases where the human model may be considered as a state abstraction of the original models. We will describe abstraction operations in terms of transition systems of models. Formally a transition system $\mathcal{T}$ corresponding to a model $\mathcal{M}$ can be represented by a tuple of the form $\mathcal{T} = \langle S, L, T, s_o, S_g \rangle$, where $S$ is the set of possible states in $\mathcal{M}$, $L$ is the set of transition labels (corresponding to the action that induce that transition), $T$ is the set of possible labeled transitions, $s_0$ is the initial state and $S_g$ is the set of states that satisfies the goal specified by $\mathcal{M}$.

**Definition 32.** *A **propositional abstraction function** $f_\Lambda$ for a set of propositions $\Lambda$ and state space $S$, defines a surjective mapping of the form $f_\Lambda : S \to X$, where $X$ is a projection of $S$, such that for every state $s \in S$, there exists a state $f_\Lambda(s) \in X$ where $f_\Lambda(s) = s \setminus \Lambda$.*

**Definition 33.** *For a planning model $\mathcal{M} = \langle F, A, I, G \rangle$ with a corresponding transition system $\mathcal{T}$, a model $\mathcal{M}' = \langle F', A', I', G' \rangle$ with a transition system $\mathcal{T}'$ is considered an **abstraction of $\mathcal{M}$** for a set of propositions $\Lambda$, if for every transition $s_1 \xrightarrow{a} s_2$ in $\mathcal{T}$ corresponding to an action $a$, there exists an equivalent transition $f_\Lambda(s_1) \xrightarrow{a'} f_\Lambda(s_2)$ in $\mathcal{T}'$, where $a'$ is part of the new action set $A'$.*

We will slightly abuse notation and extend the abstraction functions to models and actions, i.e in the above case, we will have $\mathcal{M}' \in f_\Lambda(\mathcal{M})$ (where $f_\Lambda(\mathcal{M})$ is the set of all models that satisfy the above definition for the set of fluents $\Lambda$) and similarly we will have $a' \in f_\Lambda(a)$. As per Definition 33, the abstract model is *complete* in the sense that all plans that were valid in the original model will have an equivalent plan in this new model. We will use the operator $\sqsubset$ to capture the fact that a model $\mathcal{M}$ is less abstract than the model $\mathcal{M}'$, i.e if $\mathcal{M} \sqsubset \mathcal{M}'$ then there exist a set of propositions $\Lambda$ such that $\mathcal{M}' \in f_\Lambda(\mathcal{M})$.

In particular, we will focus on abstractions formed by projecting out a subset of propositional fluents. Where for a given model $\mathcal{M} = \langle F, A, I, G \rangle$ and a set of propositions $\Lambda$, we define the abstract model to be $\mathcal{M}' = \langle F', A', I', G' \rangle$, where $F' = F - \Lambda$, $I' = f_\Lambda(I)$, $G' = f_\Lambda(G)$ and for every $a \in A$ (where $a = \langle \text{pre}(a), \text{adds}(a), \text{dels}(a) \rangle$) there exists $a' \in A'$, such that $a' = \langle \text{pre}(a) \setminus \Lambda, \text{pre}(a) \setminus \Lambda, \text{adds}(a) \setminus \Lambda, \text{dels}(a) \setminus \Lambda \rangle$.

Thus if the system has information on some subset of fluents that the user may be unaware of, then we can approximate the user model by using the abstraction of the model by setting $\Lambda$ to be these fluents. As we will see even if the approximation is more pessimistic than the true model, we will still be able to generate the required explanation. Notice also that here if we set $\Lambda$ to $P$, we essentially get an empty model that can be thought of as representing the model of a user who isn't aware of any of the details about the task, which effectively means the user thinks all possible plans are feasible.

Now going back to the question of explanation, let us consider cases where the user is raising explicit contrastive queries of the form "*Why P and not Q?*", where 'P' is the current plan being proposed by the robot and 'Q' is the alternate plan (or set of plans) being proposed by the human. If the assumption holds that the human's model is a complete abstraction then we don't need to provide any information to justify P since it should be valid in the human model. So our focus would be to address the foil plans raised by the user. Now consider the cases where the foils raised by the human are in fact infeasible in the robot model. Then as part of the explanation should reveal parts of the robot model that will help show the infeasibility of the human plan. Thus explanation here can be characterized by the set of fluents that was previously missing from the human model and whose introduction into the human model will help her understand the planning problem better, and thus model reconciliation would leave the human model less abstract than it was before.

Now what happens if the assumption about the human's model being an abstraction of the robot model is not met? Let us assume the system is aware of the set of fluents that the human may not be completely aware of, but now instead of them simply missing the human may hold incorrect beliefs about them. Interestingly, we can still use the same approach discussed above to refute the foils. That is you can still identify the set of fluents to refute the foils as if the user's model was a complete abstraction. But now we need additional mechanisms to provide explanations about the validity of the plan. One simple approach would be to allow users to ask specific questions about the steps in the plan they assume are invalid and provide explanation for those queries.

In addition to allowing for a reasonable assumption about the user's model, the use of abstractions can also provide us with the tools to address the human's limited inferential capabilities. Note that here the explanations can be provided at the minimal level of abstraction needed to explain the foils, i.e., only introduce the minimal set of fluents needed to address the user queries. Thus the inferential overhead on the human's end to make sense of the explanation is minimized. Furthermore, if the human's confusion about the plan is coming from limited inferential capabilities as opposed to a difference in knowledge regarding the task, we can still provide explanations by following the steps described in this section. In this case, the user is exposed to an abstract version of the model and asked to re-evaluate the plan and foils with respect to this simplified model.

## 6.5   Bibliographic Remarks

The first work to consider model reconciliation in the presence of uncertain models was Sreedharan et al. [2018a], and the form of the uncertain models used by the work was adapted from earlier work in robust planning Nguyen et al. [2017]. These annotated models can also be used to provide explanations for multiple users, by combining the mental model of each users into a single uncertain model. The *AO*\* search described in the chapter is adapted from Nilsson [2014] introduced to solve planning problems where the solution takes the conditional plans that can be represented as acyclic graphs. The method for model-free model reconciliation was introduced in Sreedharan et al. [2019a] for MDP settings, for the discussion in this chapter we adapted it to classical planning problems. The abstraction based explanation framework called HELM was first proposed

in Sreedharan et al. [2018b], and then later adapted to partial foil specifications and explaining unsolvability in Sreedharan et al. [2019b]. Sreedharan et al. [2021c] builds on these earlier works and provides a clearer discussion on connections between model reconciliation and HELM style explanation. This paper also presents a user study testing the effectiveness of abstraction to reduce cognitive load on the explainee. Instead of assuming a set of fluents to build abstractions on, these works generally assume access to a lattice of abstract models, which might also contain multiple maximally abstract models. This means in these setting they also have to perform an additional localization operation to identify possible models that could correspond to the human. These lattices could also be used to encode information about the complexity of concepts there by also allowing them to provide simplest explanations for a given query.

CHAPTER **7**

# Balancing Communication and Behavior

In the previous sections, we considered how in human robot teaming scenarios, the robot behavior influences and is influenced by the human's mental model of the robot. We have been quantifying some of the interaction between the behavior and human's model in terms of three interpretability scores, each of which corresponds to some desirable property one would expect the robot behavior to satisfy under cooperative scenarios. With these measures defined one of the strategies the robot can adopt is to specifically generate behavior that optimizes these measures. Though for one of those measures, we also investigated an alternate strategy, namely to use communication to alter the human's mental model to improve explicability of the current plan. Interestingly these strategies are not mutually exclusive, they can in fact be combined to capitalize on the individual strengths of each to create behaviors that are unique. In this chapter, we will start by focusing on how one can combine explicable and explanation generation and will look at a compilation based method to generate such plans. In general, communication is a strategy we can use for the other two measures as well. As such, we will end this chapter with a brief discussion on how we could use the idea of combining communication and selective behavior generation to improve legibility and predictability scores.

## 7.1 Modified USAR Domain

A typical Urban Search and Rescue (USAR) scenario consists of an autonomous robot deployed to a disaster scene with an external commander who is monitoring its activities. Both agents start with the same model of the world (i.e the map of the building before the disaster) but the models diverge over time since the robot, being internal to the scene, has access to updated information about the building. This model divergence could lead to the commander incorrectly evaluating valid plans from the robot as sub-optimal or even unsafe. One way to satisfy the commander would be to communicate or explain changes to the model that led the robot to come up with those plans in the first place.

Figure 7.1 illustrates a scenario where the robot needs to travel from P1 to its goal at P17. The optimal plan expected by the commander is highlighted in grey in their map and involves the robot moving through waypoint P7 and follow that corridor to go to P15 and then finally to P16. The robot knows that it should in fact be moving to P2 – its optimal plan is highlighted in blue. This disagreement rises from the fact that the human incorrectly believes that the path from P16 to P17 is clear while that from P2 to P3 is blocked.

If the robot were to follow the explanation scheme established in Chapter 5, it would stick to its own plan and provide the following explanation:

```
> remove-(clear p16 p17)-from-I
    (i.e. Path from P16 to P17 is blocked)
> add-(clear p2 p3)-to-I
```

**Fig. 7.1:** Illustration of the robot model and the corresponding mental model of the human. The robot starts at P1 and needs to go to P17. The human incorrectly believes that the path from P16 to P17 is clear and the one from P2 to P3 is blocked due to fire. Both agents know that there is movable rubble between P5 and P6 which can be cleared with a costly `clear_passage` action. Finally, in the mental model, the door at P8 is locked while it is unlocked in the model for the robot which cannot open unlocked doors.

```
(i.e. Path from P2 to P3 is clear)
```

If the robot were to stick to a purely explicable plan then it can choose to use the passage through P5 and P6 after performing a costly `clear_passage` action (this plan is not optimal in either of the models).

## 7.2 Balancing Explanation and Explicable Behavior Generation

Throughout this chapter, our focus would be mainly on how one could *balance* selective behavior generation with communication in the context of explicability. In the case of explicable plan generation, we focus on generating behavior that aligns with observer expectations. This process is inherently limited by what the robot can perform and one can not always guarantee that the robot can generate plans with high explicability score, let alone perfectly explicable ones. Moreover under this scheme, the robot is forced to perform suboptimal plans that may be quite far from the best course of action prescribed by the robot's model for the given task.

On the other hand, under explanation, the agent tries to identify the appropriate model information that would help convince the observer of the optimality of given robot plan, and thus help improve the explicability of the plan in the updated human model. Under this scheme, the robot is free to choose a plan that is optimal under its model but now needs to meet the overhead of communicating the explanation, which depending on the scenario could be quite substantial.

The goal of this section is to discuss how one could effectively combine the strengths of these two disparate methods. We will achieve this by folding in the how difficult it is to explain a plan into the plan generation process. Thereby allowing the robot to explicitly

reason about the gains made by the plan in terms of explicability with the cost of the plan and the cost of communicating the explanations. We will call such plans as *Balanced Plans.*

In the case of explicability, balanced planning consists of optimizing over the following objectives

1. **Plan Cost** $C(\pi)$: The first objective is the cost of the plan that the robot is going to follow.

2. **Communication Cost** $\mathcal{C}(\mathcal{E})$: The next objective is the cost of communicating the explanation. Ideally this cost should reflect both the cost accrued at the robot's end (corresponding to the cost of actions that need to be performed to communicate the information) and a cost relating to the difficulty faced by the user to understand the explanation.

3. **Penalty of Inexplicability ($C_{\mathcal{IE}}(\pi, \mathcal{M}_h^R + \mathcal{E})$)**: This corresponds to the cost related to the inexplicability experienced by the human for the current plan in their updated mental model. We will assume this to be proportional to the inexplicability score of the plan after receiving the explanations $\mathcal{E}$.

While in the most general case generating a balanced plan would be a multi-objective optimization, for simplicity, we will assume that we have weight parameters that allow us to combine the individual costs into a single cost. Thus the cost of a balanced plan (which includes both an explanation and a plan), would be given as

$$C((\mathcal{E}, \pi)) = \alpha * C(\pi) + \beta * \mathcal{C}(\mathcal{E}) + \gamma * C_{\mathcal{IE}}(\pi, \mathcal{M}_h^R + \mathcal{E})$$

Thus balanced planning would ideally consist of selecting the plan that minimizes this cost. In the USAR scenario considered in this chapter, assuming all the component weights are one and all action except `clear_passage` is twice the cost, and approximating the inexplicability score $\mathcal{IE}$ to be just exponential of the cost difference. If we set the cost of communicating a single piece of model information to be high (say 50), then the optimal balanced plan would be the most explicable plan, namely the one involving the agent moving through the rubble filled path. Depending on the weights and the exact cost involved the same formulation could give rise to wildly different behaviors. We will refer to the plan that is optimal with respect to the above objective function as **optimal balanced plans**. Though in some scenarios, we may want to focus on more constrained versions of the problem. In particular, two versions that would be interesting to consider would be

1. **Perfectly Explicable Plans:** The first special case is one where we restrict ourselves to cases where the plan is perfectly explicable, i.e., optimal, in the resultant model (i.e. after receiving the explanation). Therefore in our objective we can ignore the inexplicability penalty term and the optimization objective becomes.

$$\min \alpha * C(\pi) + \beta * \mathcal{C}(\mathcal{E})$$

$$\text{subject to } C_{\mathcal{IE}}(\pi, \mathcal{M}_h^R + \mathcal{E}) = 0$$

Note that while the plan may be optimal in the human's updated model, the plan need not be optimal for the robot. In the example considered, an example for perfectly explicable plan would be again the plan through the rubble, but now includes an explanation that the path from P16 to P17 is blocked due to fire.

2. **Perfectly Explicable Optimal Plans:** In this case, we constrain ourselves to identifying not only plans and explanations that will ensure perfect explicability, but we also try to ensure that the plans are optimal in the robot model. Such methods can be particularly effective when the cost of communicating the explanation is much smaller than the cost of the plan itself. Thus the objective in this case becomes just

$$\min \mathcal{C}(\mathcal{E})$$

$$\text{subject to } C_{\mathcal{IE}}(\pi, \mathcal{M}_h^R + \mathcal{E}) = 0$$

$$\text{and } C(\pi, \mathcal{M}^R) \in C_{\mathcal{M}^R}^*$$

Interestingly, this involves identifying the optimal plan in the robot's model with the cheapest MCE (in terms of the explanation cost). In the USAR scenario, an example for perfectly explicable optimal plan would be the plan through P4, but now you need to explain that the paths from P16 to P17 and from P2 to P3 are actually clear.

In the following sections, we will see how one could generate these balanced plans.

### 7.2.1  Generating Balanced Plans

Now to actually generate such plans, one could use modfied versions of the model space search discussed in chapter 6. Specifically if we are limiting ourselves to plans that are optimal in the human mental model, then for each possible model in the search space we could consider the space of possible optimal plans and then select the one that best meets our requirement (which would be executability in the robot model). Since this is an optimization problem one would have to keep repeating this step till the stopping criteria is met. In this case, we need to keep repeating this till we find a model where there is at least one plan that is both optimal for the updated model and the original robot model. Alternatively, one could also consider converting balancing into a single a planning problem. One way to go about this would be by assuming that the robot has access to a set of communicative actions called *explanatory actions* through which it can deliver the required explanations to the human. These actions can, in fact, be seen as actions with epistemic effects in as much as they are aimed towards modifying the human mental model (knowledge state). This means that a solution to a balanced planning problem can be seen as *self-explaining plans*, in the sense that some of the actions in the plan are aimed at helping people better understand the rest of it.

Inclusion of explanatory actions puts balanced planning squarely in the purview of epistemic planning, but the additional constraints enforced by the setting allow us to leverage relatively efficient methods to solve the problem at hand. These constraints include facts such as: all epistemic actions are public, modal depth is restricted to one,

modal operators only applied to literals, for any literal the observer believes it to be true or false and the robot is fully aware of all of the observer beliefs.

Model updates in the form of epistemic effects of communication actions also open up the possibility of other actions having epistemic *side effects* as well. The definition of balanced planning makes no claims as to how the model update information is delivered. It is quite possible that actions that the agent is performing to achieve the goal (henceforth referred to as task-level actions to differentiate it from primarily communication actions) itself could have epistemic side-effects. This is something people leverage to simplify communication in day to day lives – e.g. one might avoid providing prior description of some skill they are about to use when they can simply demonstrate it. The compilation of balanced planning into a single planning problem allows us to capture such epistemic side effects. By the same token the compilation can also capture task level constraints that may be imposed on the communication actions.

### 7.2.1.1  Compilation to Classical Planning

To support such self-explaining plans, we can adopt a formulation that is similar to ones popular in epistemic planning literature to compile it into classical planning formalisms. In our setting, each explanatory action can be viewed as an action with epistemic effects. One interesting distinction to make here is that the mental model now not only includes the human's belief about the task state but also their belief about the robot's model. This means that the planning model will need to separately keep track of (1) the current robot state, (2) the human's belief regarding the current state, (3) how actions would affect each of these (as humans may have differing expectations about the effects of each action) and (4) how those expectations change with explanations.

Given the human and robot model, $\mathcal{M}^R$ and $\mathcal{M}_h^R$, we will generate a new planning model $\mathcal{M}_\Psi = \langle F^\Psi, A^\Psi, I^\Psi, G^\Psi, C_\Psi \rangle$ as follows $F^\Psi = F \cup F_\mathcal{B} \cup F_\mu \cup \{\mathcal{G}, \mathcal{I}\}$, where $F_\mathcal{B}$ is a set of new fluents that will be used to capture the human's belief about the task state and $F_\mu$ is a set of meta fluents that we will use to capture the effects of explanatory actions and $\mathcal{G}$ and $\mathcal{I}$ are special goal and initial state propositions. We will use the notation $\mathcal{B}(p)$ to capture the human's belief about the fluent $p$. We are able to use a single fluent to capture the human belief for each (as opposed to introducing two new fluents $\mathcal{B}(p)$ and $\mathcal{B}(\neg p)$) as we are specifically dealing with a scenario where the human's belief about the robot model is fully known and human either believes each of the fluent to be true or false.

$F_\mu$ will contain an element for every part of the human model that can be changed by the robot through explanations. A meta fluent corresponding to a literal $\phi$ from the precondition of an action $a$ takes the form of $\mu^+(\phi^{\mathrm{pre}(a)})$, where the superscript $+$ refers to the fact that the clause $\phi$ is part the precondition of the action $a$ in the robot model (for cases where the fluent represents an incorrect human belief we will be using the superscript "$-$").

For every action $a^R = \langle \mathrm{pre}(a^R), \mathrm{adds}(a^R), \mathrm{dels}(a^R) \rangle \in A^R$ and its human counterpart $a_h^R = \langle \mathrm{pre}(a_h^R), \mathrm{adds}(a_h^R), \mathrm{dels}(a_h^R) \rangle \in A_h^R$, we define a new action

$$a^\Psi = \langle \mathrm{pre}(a^\Psi), \mathrm{adds}(a^\Psi), \mathrm{dels}(a^\Psi) \rangle \in A^\Psi \in \mathcal{M}_\Psi$$

whose precondition is given as:

$$\mathrm{pre}(a^{\Psi}) = \mathrm{pre}(a_h^R) \cup \{\mu^+(\phi^{\mathrm{pre}(a^R)}) \rightarrow \mathcal{B}(\phi) | \phi \in \mathrm{pre}(a^R) \setminus \mathrm{pre}(a_h^R)\}$$
$$\cup \{\mu^-(\phi^{\mathrm{pre}(a^R)}) \rightarrow \mathcal{B}(\phi) | \phi \in \mathrm{pre}(a_h^R) \setminus \mathrm{pre}(a^R)\} \cup \{\mathcal{B}(\phi) | \phi \in \mathrm{pre}(a_h^R) \cap \mathrm{pre}(a^R)\}$$

The important point to note here is that at any given state, an action in the augmented model is only applicable if the action is executable in robot model and the human believes the action to be executable. Unlike the executability of the action in the robot model (captured through unconditional preconditions) the human's beliefs about the action executability can be manipulated by turning the meta fluents on and off. The effects of these actions can also be defined similarly by conditioning them on the relevant meta fluent. In addition to these task level actions (represented by the set $A_\tau$), we can also define explanatory actions ($A_\mu$) that either add $\mu^+(*)$ fluents or delete $\mu^-(*)$. Special actions $a_0$ and $a_\infty$ that are responsible for setting all the initial state conditions true and checking the goal conditions are also added into the domain model. $a_0$ has a single precondition that checks for $\mathcal{I}$ and has the following effects:

$$\mathrm{adds}(a_0) = \{\top \rightarrow p \mid p \in I^R\} \cup \{\top \rightarrow \mathcal{B}(p) \mid p \in I_h^R\} \cup \{\top \rightarrow p \mid p \in F_{\mu^-}\}$$

$$\mathrm{dels}(a_0) = \{\mathcal{I}\}$$

where $F_{\mu^-}$ is the subset of $F_\mu$ that consists of all the fluents of the form $\mu^-(*)$. Similarly, the precondition of action $a_\infty$ is set using the original goal and adds the special goal proposition $\mathcal{G}$.

$$\mathrm{pre}^(a_\infty) = G^R \cup \{\mu^+(p^G) \rightarrow \mathcal{B}(p) \mid p \in G^R \setminus G_h^R\} \cup$$
$$\{\mu^-(p^G) \rightarrow \mathcal{B}(p) \mid p \in G_h^R \setminus G^R\} \cup \{\mathcal{B}(p) \mid G_h^R \cap G^R\}$$

Finally the new initial state and the goal specification becomes $I_\mathcal{E} = \{\mathcal{I}\}$ and $G_\mathcal{E} = \{\mathcal{G}\}$ respectively. To see how such a compilation would look in practice, consider an action (`move_from p1 p2`) that allows the robot to move from `p1` to `p2` only if the path is clear. The action is defined as follows in the robot model:

```
(:action move_from_p1_p2
    :precondition (and (at_p1) (clear_p1_p2))
    :effect (and (not (at_p1)) (at_p2) ))
```

Let us assume the human is aware of this action but does not care about the status of the path (as they assume the robot can move through any debris filled path). In this case, the corresponding action in the augmented model and the relevant explanatory action will be:

```
(:action move_from_p1_p2
  :precondition (and (at_p1) (B((at_p1))) (clear_p1_p2)
                     (implies (μ⁺_pre(move_from_p1_p2,
                                        (clear_p1_p2)))
                              (B((clear_p1_p2)))))
  :effect (and (not (at_p1)) (at_p2)
```

$$( \text{not } \mathcal{B}(\texttt{at\_p1})) \ \mathcal{B}(\texttt{at\_p2}))))$$

```
(:action explain_μ⁺ₚᵣₑ_move_from_clear
  :precondition (and)
  :effect (and μ⁺ₚᵣₑ(move_from_p1_p2, (clear_p1_p2)))))
```

Finally $C_\Psi$ captures the cost of all explanatory and task level actions. For now, we will assume that the cost of task-level actions are set to the original action cost in either robot or human model and the explanatory action costs are set according to $C_E$. Later, we will discuss how we can adjust the explanatory action costs to generate desired behavior.

We will refer to an augmented model that contains an explanatory action for each possible model updates and has no actions with effects on both the human's mental model and the task level states as the *canonical augmented model*. Given an augmented model, let $\pi_\mathcal{E}$ be a plan that is valid for this model ($\pi_\mathcal{E}(I^\Psi) \subseteq G^\Psi$). From $\pi_\mathcal{E}$, we extract two types of information – the model updates induced by the actions in the plan (represented as $\mathcal{E}(\pi_\mathcal{E})$) and the sequence of actions that have some effect on the task state (represented as $\mathbb{T}(\pi_\mathcal{E})$). We refer to the output of $\mathbb{T}$ as the task level fragment of the original plan $\pi_\mathcal{E}$. $\mathcal{E}(\pi_\mathcal{E})$ may also contain effects from action in $\mathbb{T}(\pi_\mathcal{E})$.

### 7.2.2 Stage of Interaction and Epistemic Side Effects

One of the important parameters of the problem setting that we have yet to discuss is whether the explanation is meant for a plan that is proposed by the system (i.e the system presents a sequence of actions to the human) or are we explaining some plan that is being executed either in the real world or some simulation that the human observer has access to. Even though the above formulation can be directly used for both scenarios, we can use the fact that the human is observing the execution of the plans to simplify the explanatory behavior by leveraging the fact that many of these actions may have epistemic side effects. This allows us to not explain any of the effects of the actions that the human can observe (for those effects we can directly update the believed value of the corresponding state fluent and even the meta-fluent provided the human doesn't hypothesize a conditional effect).[1] This is beyond the capability of any of the existing algorithms in this space of the explicability-explanation dichotomy.

This consideration also allows for the incorporation of more complicated epistemic side-effects wherein the human may infer facts about the task that may not be directly tied to the effects of actions. Such effects may be specified by domain experts or generated using heuristics. Once identified, adding them to the model is relatively straightforward as we can directly add the corresponding meta fluent into the effects of the relevant action. An example for a simple heuristic would be to assume that the firing of a conditional effect results in the human believing the condition to be true, provided the observer can't directly observe the fluents it was conditioned on. For example, if we assume that the robot had an action (`open_door_d1_p3`) that had a conditional effect:

---

[1] This means that when the plan is being executed, the problem definition should include the observation model of the human (which we assume to be deterministic). To keep the formulation simple, we ignore this for now. Including this additional consideration is straightforward for deterministic sensor models.

```
(when (and (unlocked_d1)) (open_d1))
```

which says the door will open if it was unlocked. Then in the compiled model, we can add a new effect to this conditional effect:

```
(when (and (unlocked_d1))
      (and 𝓑(open_d1) 𝓑(unlocked_d1)))
```

which basically says, that if the conditional effect executes the human will believe that the door was unlocked. Even in this simple case, it may be useful to restrict the rule to cases where the effect is conditioned on previously unused fluents so the robot does not expect the observer to be capable of regressing over the entire plan.

### 7.2.3    Optimizing for Explicability of the Plan

The balancing formulations discussed in this chapter require more than just generating plans that are valid in both robot and updated human model. They require us to either optimize for or ensure complete explicability of plans. For the *Perfectly Explicable Plans* and *Perfectly Explicable Optimal Plans* formulation, where plans need to be optimal in the human model we will need to enforce this as an additional goal test for the planner while for the most general formulation the inexplicability score could be added as additional penalty to the last action of the plan.

Though to allow for *Perfectly Explicable Optimal Plans*, we need to restrict the plans to only ones that are optimal in the robot model. We can generate such robot optimal plans by setting lower explanatory action costs. Before we formally state the bounds for explanatory costs, let us define the concept of *optimality delta* (denoted as $\Delta\pi_{\mathcal{M}}$) for a planning model, which captures the cost difference between the optimal plan and the second most optimal plan. More formally $\Delta\pi_{\mathcal{M}}$ can be specified as:

$$\Delta\pi_{\mathcal{M}} = \min\{v \mid v \in \mathbb{R} \;\wedge$$
$$\nexists \pi_1, \pi_2((0 < (C(\pi_1) - C(\pi_2)) < v)$$
$$\wedge\, \pi_1(I_{\mathcal{M}}) \models_{\mathcal{M}} G_{\mathcal{M}} \wedge \pi_2(I_{\mathcal{M}}) \in \Pi^*_{\mathcal{M}}\}$$

**Proposition 9.** *In a canonical augmented model $\mathcal{M}_{\Psi}$ for the human and robot model tuple $\Psi = \langle \mathcal{M}^R, \mathcal{M}^R_h \rangle$, if the sum of costs of all explanatory actions is $\leq \Delta\pi_{\mathcal{M}^R}$ and if $\pi$ is the cheapest valid plan for $\mathcal{M}_{\Psi}$ such that $\mathbb{T}(\pi) \in \Pi^*_{\mathcal{M}_{\Psi}+\mathcal{E}(\pi)}$, then:*

*(1) $\mathbb{T}(\pi)$ is optimal for $\mathcal{M}^R$*

*(2) $\mathcal{E}(\pi)$ is the MCE for $\mathbb{T}(\pi)$*

*(3) There exists no plan $\hat{\pi} \in \Pi^*_R$ such that MCE for $\mathbb{T}(\hat{\pi})$ is cheaper than $\mathcal{E}(\pi)$, i.e. the search will find an the plan with the smallest MCE.*

First off, we can see that there exists no valid plan $\pi'$ for the augmented model ($\mathcal{M}_{\Psi}$) with a cost lower than that of $\pi$ and where the task level fragment ($\mathbb{T}(\pi')$) is optimal for

the human model. Let's assume $\mathbb{T}(\pi) \notin \Pi_{\mathcal{R}}^*$ (i.e current plan's task-level fragment is not optimal in robot model) and let $\hat{\pi} \in \Pi_{\mathcal{R}}^*$. Now let's consider a plan $\hat{\pi}_{\mathcal{E}}$ for augmented model that corresponds to the plan $\hat{\pi}$, i.e., $\mathcal{E}(\hat{\pi}_{\mathcal{E}})$ is the MCE for the plan $\hat{\pi}$ and $\mathbb{T}(\hat{\pi}_{\mathcal{E}}) = \hat{\pi}$. Then the given augmented plan $\hat{\pi}_{\mathcal{E}}$ is a valid solution for our augmented planning problem $\mathcal{M}_{\Psi}$ (since the $\hat{\pi}_{\mathcal{E}}$ consists of the MCE for $\hat{\pi}$, the plan must be valid and optimal in the human model), moreover the cost of $\hat{\pi}_{\mathcal{E}}$ must be lower than $\pi$. This contradicts our earlier assumption hence we can show that $\mathbb{T}(\pi)$ is in fact optimal for the robot model.

Using a similar approach we can also show that no cheaper explanation exists for $\pi_{\mathcal{E}}$ and there exists no other plan with a cheaper explanation.

Note that while it is hard to find the exact value of the optimality $\Delta\pi_{\mathcal{M}}$, it is guaranteed to be $\geq 1$ for domains with only unit cost actions or $\geq (C_2 - C_1)$, where $C_1$ is the cost of the cheapest action and $C_2$ is the cost of the second cheapest action, i.e. $\forall a(C_{\mathcal{M}}(a) < C_2 \rightarrow C_{\mathcal{M}}(a) = C_1)$. Thus allowing us to easily scale the cost of the explanatory actions to meet this criteria.

## 7.3 Balancing Communication and Behavior For other Measures

While the majority of the chapter is focused on balancing communication and behavior selection to improve explicability, it should be easy to see that such schemes should be extendable to other behavioral metrics.

In the case of legibility, since the original framework itself is focused on communication, it is easy to replace specialized behavior with communication. In particular, for a formulation that dealt exclusively with goal uncertainty, you are effectively trading off a single message (namely the goal of the robot) with the additional effort that needs to be taken by the robot to communicate its eventual goal implicitly. Though one could easily consider the more general version of legibility where the uncertainty is over entire models or different model components. In such cases, one could foresee a combination of communication and selective behavior to effectively communicate the actual parameters. In this case, the communication takes a form quite similar to the explanation since we will be communicating model information. The planning objective for a setting that tries to maximize legibility within $k$ steps, would be

$$C(\langle \mathcal{E}, \pi \rangle) = \alpha * C(\pi) + \beta * \mathcal{C}(\mathcal{E}) + \gamma * P(\theta | \mathbb{M}_h^R + \mathcal{E}, \hat{\pi}^k)$$

Where $\theta$ is the target model parameter being communicated (which can be the full model), $\mathbb{M}_h^R$ the hypothesis set maintained by the human observer, $\hat{\pi}^k$ is the $k$ step prefix of the plan $\pi$. Thus the formulation tries to optimize for the probability that the human would associate with the target model parameter. Similar to the explicability case, one could also formulate the more constrained cases.

Now moving onto predictability, one main difference would be the communication

strategy to be used. Predictability, as originally conceived, was meant for cases where the human and robot shares the same mental model, but the human observer is unsure about the future behavior the robot will exhibit. So even if the human has incorrect beliefs or uncertainty about the robot model, the communication strategy here needs to include more than just model information. A useful strategy here could be for the robot to communicate potential constraints on the behavior that the robot could actually follow. A possibility here could be to communicate a partially ordered set of action or state facts, which basically conveys a commitment on the part of the robot that no matter the exact behavior the robot ends up choosing it would be one that involves either performing the actions or achieving facts in the order specified, thereby constraining the possible completions it can follow and thus allowing the rest of the behavior. Now using $\mathcal{E}$ again as the stand-in for the information being conveyed, $\mathcal{M}_h^R + \mathcal{E}$ as the updated human model and applying the constraints specified in the communication, then the objective for planning for balanced plans that maximize predictability in $k$ steps is given as

$$C(\langle \mathcal{E}, \pi \rangle) = \alpha * C(\pi) + \beta * \mathcal{C}(\mathcal{E}) + \gamma * P(\pi | \mathbb{M}_h^R + \mathcal{E}, \hat{\pi}^k)$$

Where $P(\pi | \mathbb{M}_h^R + \mathcal{E}, \hat{\pi}^k)$ is the probability the human observer associates with the actual plan after observing the prefix $P(\pi | \mathbb{M}_h^R + \mathcal{E}, \hat{\pi}^k)$. Similar to the earlier formulations, we can also look for constrained versions of the objective, that requires perfect predictability and optimality in robot model.

Moreover, in both these cases, we can adapt the compilation method discussed earlier to fit these new objectives. In particular, we just need to change the goal-test in the case of constrained version and or the additional penalty calculated from $\mathcal{IE}$ to the respective interpretability measures.

## 7.4    Bibliographic Remarks

The first work to consider balanced planning was Chakraborti et al. [2019c], that considers the generation of the *Perfectly Explicable Plans*. In particular, they consider a model space search based solution to identify the plan. Unfortunately, in the case of the model space search, the only way to guarantee the explanation generated as part of *Perfectly Explicable Plans* is minimal would be to iterate over all the optimal plans in the given model. Thus the method they ended up operationalizing was an approximate version, where they guaranteed that while the plan generated is perfectly explicable, the explanation may be larger than required. The paper also presents some initial user studies to validate the usefulness of such balanced plans. The classical planning compilation for balanced planning was Sreedharan et al. [2020a]. This was also the first work to connect epistemic planning to model reconciliation and the central compilation method was derived from an earlier work to compiling restricted forms of epistemic planning to classical planning Muise et al. [2015]. The paper also discusses how all three types of balanced plans (i.e *OPtimal Balanced Plans*, *Perfectly Explicable Plans* and, *Perfectly Explicable Optimal Plans*) can be generated. In terms of balancing for other interpretability measures, while we are unaware of works that consider them in the general form, works like Chakraborti et al. [2018] have looked at applying these ideas in more specific context. In particular, Chakraborti et al. [2018] looks at using mixed reality based visualization

to improve, explicability, predictability and legibility. The kind of scenarios they considered included a block stacking scenario where they considered use of mixed reality cues to highlight potential blocks that might be used in the plan and thereby improving the identification of the plan and eventual goal of the agent.

CHAPTER **8**

# Explaining in the Presence of Vocabulary Mismatch

All previous discussions on model-reconciliation explanations implicitly assume that the robot can communicate information about the model to the user. This suggests that the human and the robot share a common vocabulary that can be used to describe the model. However, this cannot be guaranteed unless the robots are using models that are specified by an expert. Since many of the modern AI systems rely on learned models, they may use representational schemes that would be inscrutable to most users. So in this chapter, we will focus on the question of generating model reconciliation explanations in the absence of shared vocabulary. Specifically, we will see how one could identify relevant model information in terms of concepts specified by a user. We will describe how the agent could use classifiers learned over such concepts to identify fragments of symbolic models specified in terms of such concepts that are sufficient to provide explanations. We will also discuss how one could measure the probability of the generated model fragments being true (especially when the learned classifiers may be noisy) and also how to identify cases where the user specified concepts may be insufficient.

## 8.1   Representation of Robot Model

Note that the assumption that the agent can communicate model information in itself entails two separate parts, (a) the model follows a representation scheme that is intuitive or understandable to the end users and (b) the model is actually represented using factors that makes sense to its end user. Model-reconciliation, does not technically require explanation to be carried out in the same terms as the representation scheme used by the agent. Instead we could always map the given robot model information into representation schemes that are easier for people to understand. In the previous chapters, we assumed that the robot model is already represented using STRIPS like description languages and in this chapter we will look at mapping the original model of the robot (regardless of its current representation scheme) into a STRIPS model. Such representation schemes are not only expressive (any deterministic goal directed planning problem with finite states and actions can be captured as STRIPS planning problem), but they also have their roots in folk psychology [Miller, 2019]. Thus model information represented in such representation schemes are relatively easier for people to make sense of.

A more pressing question is what should be the state factors or action labels that we should use to represent the model. The fact that the model is represented using STRIPS in and of itself wouldn't facilitate effective explanation, if the information about an action's precondition is represented using concepts alien to the end user. We will use the term *vocabulary used by the model* to refer to the state fluents and the action names used by it. Thus access to a shared vocabulary is a prerequisite to facilitating model reconciliation explanations. In this chapter, we will look at cases where this is not a given. We will look at one possible way of collecting vocabulary concepts from the end users and discuss how

we could use the collected vocabulary to generate the model or at least identify parts of the model relevant to the explanation. We will look at (a) how to handle cases where the vocabulary items we gathered are insufficient (b) how to ensure that models learned this way actually reflect true model and (c) how to handle cases where the mappings we learn from human concepts to task states may be noisy. We will also see how this process of mapping the model into a secondary representation in user defined terms could also give an opportunity to approximate and simplify the model and by extension the explanation itself.

In this chapter we will focus on acquiring human understandable state variables/factors and assume the action labels are given beforehand. We do this because if the human and the robot have a one to one correspondence between what constitutes an atomic action, it is usually easier to establish their labels/names. Acquiring action labels could become challenging if the agent and the human are viewing the action space at varying levels of abstraction, but we will ignore this possibility for now.

## 8.2   Setting

For the purposes of discussion, we will look at a slightly different scenario where the robot has access to a deterministic simulator of the task of the form $\mathcal{M}_{\text{sim}} = \langle S, A, T, \mathcal{C} \rangle$ (where this simulator effectively acts as the robot model), where $S$ represents the set of possible world states, $A$ the set of actions and $T$ the transition function that specifies the problem dynamics. The state space $S$ here need not be defined over any state factors and may just be atomic states. The transition function is defined as $T : S \times A \to S \cup \{\bot\}$, where $\bot$ corresponds to an invalid absorber-state generated by the execution of an infeasible action. Invalid state could be used to capture failure states that could occur when the agent violates hard constraints like safety constraints. Finally, $\mathcal{C} : A \to \mathbb{R}$ captures the cost function of executing an action in any state where the action is feasible. We will overload the transition function $T$ to also work on action sequence, i.e., $T(s, \langle a_1, ..., a_k \rangle) = T(...T(T(s, a_1), a_2), ..., a_k)$. We will assume the problem is still goal-directed in the sense that the robot needs to come up with the plan $\pi$, that will drive the state of the world to a goal state. In general we will use the tuple $\Pi_{\text{sim}} = \langle I, \mathbb{G}, \mathcal{M}_{\text{sim}} \rangle$ to represent the decision making problem, where $I$ is the initial state and $\mathbb{G}$ the set of goal states. Similar to the earlier settings, a plan is optimal if it achieves the goal and there exists no cheaper plan that can achieve the goal (where $\mathcal{C}(I, \pi)$ is the total cost of executing the plan $\pi$).

Now the challenge before us is to map this blackbox simulator model into a STRIPS model defined using vocabulary that makes sense to the human. In particular we want to first collect the set of propositional state fluents that the human uses to define the model $\mathcal{M}_h^R$. Before we discuss how to collect such state fluents information, let us take a quick look at the state fluents themselves. For one, how does one connect a model defined over atomic space to a set of propositional fluents? The basic convention we will follow here is that each proposition state fluent correspond to a fact about the task, that is either true or false in all states of the given simulator $\mathcal{M}_{sim}$. For a given set of state fluents, each state in $\mathcal{M}_{sim}$ is mapped to the set of propositional fluents that are true in that state and each propositional fluent can be for our purposes by the subset of states $S$ where that fluent is true.

Now lets assume the human observer is ready to list all these fluents (captured by the set $\mathbb{C}$) they believe are relevant to the problem at hand. We can hardly expect the human to be able to list or point out all possible states corresponding to each factor. Instead we will employ a different method to learn the mapping between atomic states and fluents. In this chapter, we will assume that each factor can be approximated by a classifier. Now we can learn such a classifier for a fluent (or concept) in the set $\mathbb{C}$ by asking the user for a set of positive and negative examples for the concept. This means that the explanatory system should have some method of exposing the simulator states to the user. A common way to satisfy this requirement would be by having access to visual representations for the states. The simulator state itself doesn't need to be an image as long as we have a way to visualize it. The concept list can also be mined from qualitative descriptions of the domain and we can crowdsource the collection of example states for each concept. Once we have learned such a set of classifiers, we can then use these factors to learn a model of the task and start using it for the explanations.

### 8.2.1 Local Approximation of Planning Model

In most of the previous discussions in the book the implicit assumption was that the explanation was given with respect to a representation of the entire robot model. This could mean the model being used for explanation is quite complex and large even after the application of simple abstraction methods. A strategy quite popular in explaining classification decisions is to explain a given decision with respect to an approximation of the original model that only captures the behavior of the model for a subset of the input space. We can rely on a similar strategy for explanation in sequential decision-making scenarios. Namely, we can focus on a symbolic representation of the task that aims to capture the dynamics of the underlying domain model only for a subset of states.

More formally, consider a STRIPS model defined over a set of concepts $\mathbb{C}$, $\mathcal{M}^{\mathbb{C}} = \langle \mathbb{C}, A^{\mathbb{C}}, \mathbb{C}(I), \mathbb{C}(\mathbb{G}), \mathcal{C}_{\mathcal{S}}^{\mathbb{C}} \rangle$, where $\mathbb{C}(\mathbb{G}) = \bigcap_{s_g \in \mathbb{G}} \mathbb{C}(s_g)$. We will call this model to be a local approximation of the simulation $\mathcal{M}_{\text{sim}} = \langle S, T, A, \mathcal{C} \rangle$ for regions of interest $\hat{S} \subseteq S$, if $\forall s \in \hat{S}$ and $\forall a \in A$, we have an equivalent action $a^{\mathbb{C}} \in A_{\mathcal{S}}^{\mathbb{C}}$, such that $a^{\mathbb{C}}(\mathbb{C}(s)) = \mathbb{C}(T(s,a))$ (assuming $\mathbb{C}(\bot) = \bot$) and $\mathcal{C}_{\mathcal{S}}^{\mathbb{C}}(a) = \mathcal{C}(a)$.

Now one of the important aspects of this formulation is the subset of states over which we will define the approximation. A popular strategy used in machine learning approaches is to select a set of input points close to the data point being explained, where closeness of two data points is defined via some pre-specified distance function. This strategy can also be adapted to sequential decision-making settings. In this case, the data point would correspond to the plan (and in some cases also the foils) and the state space of interest can be limited to states close to the ones that appear on the plan or foil execution trace. The distance function can be defined over some representation of the underlying state or use reachability measures to define closeness (i.e. whether the states can be reached within some predefined number of steps).
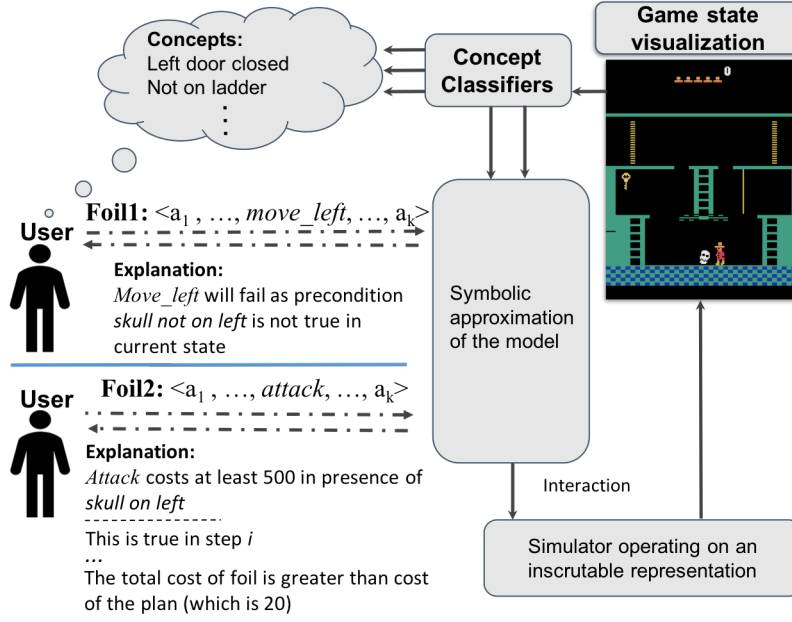
**Fig. 8.1:** An overview of the overall system described in the chapter. The screenshot is from a game called Montezuma's revenge. In this particular instance, the agent is tasked with making its way to the key on the left side of the screen while avoiding the skull on the ground.

### 8.2.2 Montezuma's Revenge

For the purposes of the discussion we will consider the game Montezuma's revenge as the example scenario. The game was first launched in 1984 for Atari gaming systems, and featured a character that needs to navigate various levels filled with traps and enemies while collecting treasures. The game has recently received a lot of attention as a benchmark for RL algorithms and is a domain that requires long term planning to solve. The simulator here consists of the game engine, which stores the current state and predicts what happens when a specific action is executed. Two popular state representation schemes for the game that is used by RL algorithms either involves the use of images of the game state or the RAM state of the game controller corresponding to the game state. Neither of which is particularly easy for a naive user to understand. The human observer could be a user of this game playing system, and as such we will sometimes refer to the human as the user in the text. In this chapter we will mostly focus on the first level of the game that is pictured in 8.1. In particular, we will focus on the task of getting the agent from the point visualized in the image to the key it needs to pick up. A successful plan for this task would involve the agent making its way to the lowest platform and them moving towards the skull, jumping over it once it close enough and then climbing up the ladder before going for the key. Now if we were to come up with a set of concepts to describe the agent, many of them would involve the actual position of the agent (like if the agent is on the leftmost ladder or if the agent is on the topmost platform etc.), whether the agent has possession of the key, where the agent is with respect to the skull (especially if the agent is right next to it) etc. One can train a classifier that takes either the image or the RAM state of a specific game state and determines whether the concept is absent or present in that state. Now the goal of many of the methods discussed in this chapter would be to build models that can approximate the dynamics of the game in this level using such

concepts.

## 8.3   Acquiring Interpretable Models

Now to provide explanations in this setting, the first approach would be to try constructing the model $\mathcal{M}^{\mathbb{C}}$, particularly the action definition as the initial states and goal specification are easier to construct. For now we will assume that we have a perfect classifier for each of the concept listed in the list $\mathbb{C}$. This means that for any state sampled from the set we know the exact list of concepts present in that state. Now we will use the simulator to observe how each action transforms the state and convert it into symbolic representation through the classifiers.

More specifically, we will sample a set of states from the simulator state space (where the samples are limited to ones that meet the desired distance requirement discussed in previous section). For each state sampled, we test whether the actions are executable (i.e it doesn't lead to an invalid state), and for each action we can update the definition of the action $a$ as follows (assuming the sample state is $s_i$ and the state resulting from executing $a$ is denoted as $a(s_i)$)

$$\text{pre}(a^c) = \text{pre}(a^c) \cap \mathbb{C}(s_i)$$
$$\text{adds}(a^c) = \mathbb{C}(a(s_i)) \setminus \mathbb{C}(s_i)$$
$$\text{dels}(a^c) = \mathbb{C}(s_i) \setminus \mathbb{C}(a(s_i))$$

Where the original estimate of the action preconditions ($\text{pre}(a^c)$) is set to be equal to $\mathbb{C}$. Note that learning using this method makes the specific assumption that each action in the exact local symbolic approximation only contains preconditions that are conjunctions of positive concepts and are free of conditional effects. If the model in fact does not meet these assumptions, we may need to rely on more complex model learning techniques. [1]

Once such a model is learned we could then leverage many of the explanatory methods discussed in the book to generate explanations (especially the ones making simplifying assumptions about the human model), though one obvious issue that could occur is that the original concept list provided to the explanatory system may be incomplete. That is the fluent list of the exact local approximation (we will represent this as $\mathcal{M}^*$) may be a superset of $\mathbb{C}$. Let the learned model be $\hat{\mathcal{M}}^{\mathbb{C}}$, and for now let us assume that this model is the fix-point of the learning process described above. Then it is easy to see that if the assumptions about $\mathcal{M}^*$ hold (i.e., it is conditional effect free and has only conjunction of positive literals as precondition) then $\hat{\mathcal{M}}^{\mathbb{C}}$, must be a complete abstraction of $\mathcal{M}^*$, in so far as any plan viable in $\mathcal{M}^*$ must also be valid in $\hat{\mathcal{M}}^{\mathbb{C}}$. This is because the above learning method will ensure that the learned preconditions and effects are a subset of the original action precondition and effects. As we have already seen in Chapter 6 (Section 33), this means that even though this is not the exact model, this model suffices to address many specific user queries. Whenever the user raises a foil that cannot be refuted through $\hat{\mathcal{M}}^{\mathbb{C}}$, then the system can use that as a signal that the current vocabulary is incomplete.

---

[1] Though if the only assumption not met is the form of the precondition then it can be mapped to a case of missing concepts that could be handled by assuming the exact approximation model meets the same representational assumptions but is defined over a larger set of concepts.

## 8.4   Query Specific Model Acquisition

The previous section talked about how we can acquire a representation of the entire model that can then be used in the standard model reconciliation process. However, this could be a pretty expensive process, especially if the action space is quite large. In cases where the agents are raising specific explanatory queries we may not require the entire model, but rather only learn the parts of the model necessary to answer the specific user query. To illustrate this process, let us take a look at the simplest contrastive explanation setting. Here the human explanatory query consists of asking why the system chose to follow the current plan instead of a set of alternative plans. Specifically in this case, for the decision-making problem specified by the tuple $\Pi_{sim} = \langle I, \mathbb{G}, \mathcal{M}_{sim} \rangle$ the system identifies a plan $\pi$. When presented with the plan, the user of the system responds by raising an alternative plan $\pi_f$ (*the foil*) that they believe should be followed instead. Now the system would need to explain why the plan $\pi$ is preferred over the foil $\pi_f$ in question. The only two possibilities here are that either the foil is inexecutable and hence cannot be followed or it is costlier than the plan in question.

**Definition 34.** *The plan $\pi$ is said to be preferred over a foil $\pi_f$ for a problem $\Pi_{sim} = \langle I, \mathbb{G}, \mathcal{M}_{sim} \rangle$, if either of the following conditions are met, i.e.,*

1. *$\pi_f$ is inexecutable, which means, either (a) $T(I, \pi_f) \notin \mathbb{G}$, i.e the action sequence doesn't lead to a possible goal state, or (b) the execution of the plan leads to an invalid state, i.e., $T(I, \pi_f) = \bot$.*

2. *Or $\pi_f$ is costlier than $\pi$, i.e., $\mathcal{C}(I, \pi) < \mathcal{C}(I, \pi_f)$*

So going back to our montezuma example. Let's assume the plan involves the agent starting from the highest platform, and the goal is to get to the key. The specified plan $\pi$ may require the agent to make its way to the lowest level, jump over the skull, and then go to the key with a total cost of 20. Let us consider a case where the user raises two possible foils that are quite similar to $\pi$, but, (a) in the first foil, instead of jumping the agent just moves left (as in it tries to move through the skull) and (b) in the second, instead of jumping over the skull, the agent performs the `attack` action (not part of the original game, but added here for illustrative purposes) and then moves on to the key. Now using the simulator, the system could tell that in the first case, moving left would lead to an invalid state and in the second case, the foil is more expensive. It may however struggle to explain to the user what particular aspects of the state or state sequence lead to the invalidity or suboptimality as it cannot directly expose relevant model information. Instead, as in earlier parts it would need to map it to a specific symbolic model and expose information about that model. However in this case, rather than first learning a full symbolic model and then identifying relevant model components to provide, we can directly try to learn the model parts. To capture scenarios like the one mentioned above, we will allow for a class of richer symbolic, namely one that allows for cost functions that are conditioned on the state in addition to the action. This means the cost function will take the form $\mathcal{C}_\mathcal{S} : 2^F \times A_\mathcal{S} \to \mathbb{R}$ to capture the cost of valid action executions in a specific state. Internally, such state models may be represented using conditional cost models. In this discussion, we won't try to reconstruct the exact cost function but will rather try to estimate an abstract version of the cost function.

As a quick aside, most of the discussion in this section focuses on generating explanation to refute the alternate plan and not really on explaining why the current plan works or has the cost assigned to it. Since the human can observe the robot and we had previously assumed that there exists mechanism to visualize the simulator states, in theory the robot could just demonstrate the outcome of executing the plan. One could also adapt the methods we discuss for refuting the foil to answer more specific question the human user might have about the robot, for example by learning an abstract cost function for the actions in the plan or by identifying whether a concept the human had in mind is in fact a precondition or not.

### 8.4.1  Explanation Generation:

For establishing the invalidity of $\pi_f$, we just need to focus on explaining the failure of the first failing action $a_i$, i.e., the last action in the shortest prefix that would lead to an invalid state (which in our running example is the move-left action in the state presented in Figure 8.1 for the first foil). We can do so by informing the user that the failing action has an unmet precondition, as per the symbolic model, in the state it was executed in. Formally

**Definition 35.** *For a failing action $a_i$ for the foil $\pi_f = \langle a_1, .., a_i, .., a_n \rangle$, $c_i \in \mathbb{C}$ is considered an explanation for failure if $c_i \in pre(a_i) \setminus \mathbb{C}(s_i)$, where $s_i$ is the state where $a_i$ is meant to be executed (i.e $s_i = T(I, \langle a_1, .., a_{i-1} \rangle))$.*

In our example for the invalid foil, a possible explanation would be to inform the user that move-left can only be executed in states for which the concept `skull-not-on-left` is true; and the concept is false in the given state. This formulation is enough to capture both conditions for foil inexecutability by appending an additional goal action at the end of each sequence. The goal action causes the state to transition to an end state and it fails for all states except the ones in $\mathbb{G}$. Our approach to identifying the minimal information needed to explain specific query follows from studies in social sciences that have shown that selectivity or minimality is an essential property of effective explanations.

For explaining the suboptimality of the foil, we have to inform the user about $\mathcal{C}_{\mathcal{S}}^{\mathbb{C}}$. To ensure minimality of explanations, rather than generating the entire cost function or even trying to figure out individual conditional components of the function, we will instead try to learn an abstraction of the cost function $\mathcal{C}_s^{abs}$, defined as follows

**Definition 36.** *For the symbolic model $\mathcal{M}_{\mathcal{S}}^{\mathbb{C}} = \langle \mathbb{C}, A_{\mathcal{S}}^{\mathbb{C}}, \mathbb{C}(I), \mathbb{C}(\mathbb{G}), \mathcal{C}_{\mathcal{S}}^{\mathbb{C}} \rangle$, an abstract cost function $\mathcal{C}_{\mathcal{S}}^{abs} : 2^{\mathbb{C}} \times A_{\mathcal{S}}^{\mathbb{C}} \to \mathbb{R}$ is specified as follows $\mathcal{C}_{\mathcal{S}}^{abs}(\{c_1, .., c_k\}, a) = min\{\mathcal{C}_{\mathcal{S}}^{\mathbb{C}}(s, a) | s \in S_{\mathcal{M}_{\mathcal{S}}^{\mathbb{C}}} \wedge \{c_1, .., c_k\} \subseteq s\}]$.*

Intuitively, $\mathcal{C}_{\mathcal{S}}^{abs}(\{c_1, .., c_k\}, a) = k$ can be understood as stating that *executing the action $a$, in the presence of concepts $\{c_1, .., c_k\}$ costs at least $k$.* We can use $\mathcal{C}_{\mathcal{S}}^{abs}$ in an explanation of the form

**Definition 37.** *For a valid foil $\pi_f = \langle a_1, .., a_k \rangle$, a plan $\pi$ and a problem $\Pi_{sim} = \langle I, \mathbb{G}, \mathcal{M}_{sim} \rangle$, the sequence of concept sets of the form $\mathbb{C}_{\pi_f} = \langle \hat{\mathbb{C}}_1, ..., \hat{\mathbb{C}}_k \rangle$ along with $\mathcal{C}_s^{abs}$ is considered a valid explanation for relative suboptimality of the foil (denoted as $\mathcal{C}_{\mathcal{S}}^{abs}(\mathbb{C}_{\pi_f}, \pi_f) > \mathbb{C}(I, \pi))$,*

*if $\forall \hat{\mathbb{C}}_i \in \mathbb{C}_{\pi_f}$, $\hat{\mathbb{C}}_i$ is a subset of concepts presents in the corresponding state (where state is $I$ for $i = 1$ and $T(I, \langle a_1, ..., a_{i-1} \rangle)$ for $i > 1$). and $\Sigma_{i=\{1..k\}} \mathcal{C}_{\mathcal{S}}^{abs}(\hat{\mathbb{C}}_i, a_i) > \mathcal{C}(I, \pi)$*

In the earlier example, the explanation would include the fact that executing the action `attack` in the presence of the concept `skull-on-left`, will cost at least 500 (as opposed to original plan cost of 20).

## 8.4.2  Identifying Explanations through Sample-Based Trials

For identifying the model parts for explanatory query, we can rely on the agent's ability to interact with the simulator to build estimates. Given the fact that we can separate the two cases at the simulator level, we will keep the discussion of identifying each explanation type separate and only focus on identifying the model parts once we know the failure type.

### 8.4.2.1  Identifying failing precondition:

To identify the missing preconditions, we can rely on the simple intuition that while successful execution of an action $a$ in the state $s_j$ with a concept $C_i$ doesn't necessarily establish that $C_i$ is a precondition, we can guarantee that any concept false in that state cannot be a precondition of that action. This is a common line of reasoning exploited by many of the model learning methods. So we start with the set of concepts that are absent in the the state ($s_{\text{fail}}$) where the failing action ($a_{\text{fail}}$) was executed, i.e., poss_pre_set $= \mathbb{C} \setminus \mathbb{C}(s_{\text{fail}})$. We then randomly sample for states where $a_{\text{fail}}$ is executable. Each new sampled state $s_i$ where the action is executable can then be used to update the possible precondition set as poss_pre_set $=$ poss_pre_set $\cap$ $\mathbb{C}(s_i)$. That is, if a state is identified where the action is executable but a concept is absent then it can't be part of the precondition. We will keep repeating this sampling step until the sampling budget is exhausted or if one of the following exit conditions is met. (a) In cases where we are guaranteed that the concept list is exhaustive, we can quit as soon as the set of possibilities reduce to one (since there has to be a missing precondition at the failure state). (b) The search results in an empty list. The list of concepts left at the end of exhausting the sampling budget represents the most likely candidates for preconditions. *An empty list here signifies the fact that whatever concept is required to differentiate the failure state from the executable one is not present in the initial concept list $\mathbb{C}$.* This can be taken as evidence to query the user for more task-related concepts.

**Identifying cost function:**  Now we will employ a similar sampling based method to identify the right cost function abstraction. Unlike the precondition failure case, there is no single action we can choose but rather we need to choose a level of abstraction for each action in the foil (though it may be possible in many cases to explain the suboptimality of foil by only referrring to a subset of actions in the foil). Our approach here would be to find the most abstract representation of the cost function at each step such that of the total cost of the foil becomes greater than that of the specified plan. Thus for a foil $\pi_f = \langle a_1, ..., a_k \rangle$ our objective become

$$\min_{\hat{\mathbb{C}}_1, ..., \hat{\mathbb{C}}_k} \Sigma_{i=1..k} \|\hat{\mathbb{C}}_i\| \text{ subject to } \mathcal{C}_s^{abs}(\mathbb{C}_{\pi_f}, \pi_f) > \mathbb{C}(I, \pi)$$

For any given $\hat{\mathbb{C}}_i$, $\mathcal{C}_s^{abs}(\hat{\mathbb{C}}_i, a_i)$ can be approximated by sampling states randomly and finding the minimum cost of executing the action $a_i$ in states containing the concepts $\hat{\mathbb{C}}_i$. We can again rely on a sampling budget to decide how many samples to check and enforce required locality within sampler.

## 8.5   Explanation Confidence

The methods discussed above (both for learning full model and model component) are guaranteed to identify the exact model in the limit, the accuracy of the methods is still limited by practical sampling budgets we can employ. So this means it is important that we are able to establish some level of confidence in the solutions identified. In case of learning full model, there are some loose PAC learning guarantees we could employ, but for learning model components we will strive for a more accurate measure. To assess confidence, we will follow the probabilistic relationship between the random variables as captured by Figure 8.2 (A) for precondition identification and Figure 8.2 (B) for cost calculation. Where the various random variables captures the following facts: $O_a^s$ - indicates that action $a$ can be executed in state $s$, $c_i \in p_a$ - concept $c_i$ is a precondition of $a$, $O_{c_i}^s$ - the concept $c_i$ is present in state $s$, $\mathcal{C}_s^{abs}(\{c_i\}, a) \geq k$ - the abstract cost function is guaranteed to be higher than or equal to k and finally $O_{\mathcal{C}(s,a) \geq k}$ - stands for the fact that the action execution in the state resulted in cost higher than or equal to $k$. We will allow for inference over these models, by relying on the following simplifying assumptions - (1) the distribution of all non-precondition concepts in states where the action is executable is the same as their overall distribution across the problem states (which can be empirically estimated), (3) cost distribution of an action over states corresponding to a concept that does not affect the cost function is identical to the overall distribution of cost for the action (which can again be empirically estimated).

The first assumption implies that you are as likely to see a non-precondition concept in a sampled state where the action is executable as the concept was likely to appear at any sampled state with the same set of concepts (this distribution is denoted as $P(O_{c_i}^s | O_{\mathbb{C} \setminus c_i^s})$, where $O_{\mathbb{C} \setminus c_i^s}$ represents the other observed concepts). While the second one implies that for a concept that has no bearing on the cost function for an action, the likelihood that executing the action in a state where the concept is present will result in a cost greater than $k$ will be the same as that of the action execution resulting in cost greater than $k$ for a randomly sampled state ($p_{\mathcal{C}(.,a) \geq k}$).

For a single sample, the posterior probability of explanations for each case can be expressed as follows: for precondition estimation, updated posterior probability for a positive observation can be computed as $P(c_i \in p_a | O_{c_i}^s \wedge O_{\mathbb{C} \setminus c_i^s} \wedge O_a^s) = (1 - P(c_i \notin p_a | O_{c_i}^s \wedge O_{\mathbb{C} \setminus c_i^s} \wedge O_a^s))$, where

$$P(c_i \notin p_a | O_{c_i}^s \wedge O_a^s) = \frac{P(O_{c_i}^s | O_{\mathbb{C} \setminus c_i^s}) \times P(c_i \notin p_a)}{P(O_{c_i}^s | O_a^s)}$$
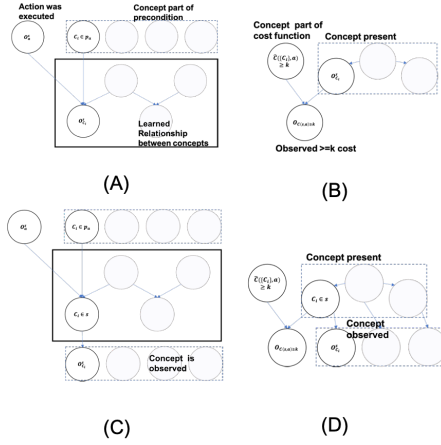
**Fig. 8.2:** A simplified probabilistic graphical models for explanation inference, Subfigure (A) and (B) assumes classifiers to be completely correct, while (C) and (D) presents cases with noisy classifier.

and for the case of cost function approximation

$$P(\mathcal{C}_s^{abs}(\{c_i\}, a) \geq k | O_{c_i}^s \wedge O_{\mathbb{C} \setminus c_i^s} \wedge O_{\mathcal{C}(s,a) \geq k}) =$$

$$\frac{P(\mathcal{C}_s^{abs}(\{c_i\}, a) \geq k)}{P(\mathcal{C}_s^{abs}(\{c_i\}, a) \geq k)) + p_{\mathcal{C}(.,a) \geq k} * P(\neg \mathcal{C}_s^{abs}(\{c_i\}, a) \geq k))}$$

The distribution used in the cost explanation, can either be limited to distribution over states where action $a_i$ is executable or allow for the cost of executing an action in a state where it is not executable to be infinite. The full derivation of these formulas can be found in the paper Sreedharan et al. [2020c].

## 8.6   Handling Uncertainty in Concept Mapping

Given how unlikely it is to have access to a perfect classifier for any concept, a more practical assumption to adopt could be that we have access to a noisy classifier. However, we assume that we also have access to a probabilistic model for its prediction. That is, we have access to a function $P_{\mathbb{C}} : \mathbb{C} \rightarrow [0, 1]$ that gives the probability that the concept predicted by the classifier is actually associated with the state. Such probability functions could be learned from the test set used for learning the classifier. Allowing for the possibility of noisy observation generally has a more significant impact on the precondition calculation than the cost function approximation. Since we are relying on just generating a lower bound for the cost function, we can be on the safer side by under-approximating the cost observations received (though this could lead to larger than required explanation). In the case of precondition estimation, we can no longer use a single failure (execution of an action in a state where the concept is absent) as evidence for discarding the concept. Though we can still use it as an evidence to update the probability of the given concept being a precondition. We can remove a particular possible precondition from consideration once the probability of it not being a precondition crosses a specified threshold.

To see how we can incorporate these probabilistic observations into our confidence calculation, consider the updated relationships presented in Figure 8.2 (C) and (D) for

precondition and cost function approximation. Note that in previous sections, we made no distinction between the concept being part of the state and actually observing the concept. Now we will differentiate between the classifier saying that a concept is present $(O^s_{c_i})$ from the fact that the concept is part of the state $(c_i \in \mathbb{C}(S))$. Now we can use this updated model for calculating the confidence. We can update the posterior of a concept not being a precondition given a negative observation $(O^s_{\neg c_i})$ using the formula

$$P(c_i \notin p_a | O^s_{\neg c_i} \wedge O^s_a \wedge O_{\mathbb{C} \setminus c^s_i}) = \frac{P(O^s_{\neg c_i} | c_i \notin p_a \wedge O^s_a \wedge O_{\mathbb{C} \setminus c^s_i}) * P(c_i \notin p_a)}{P(O_{\neg c_i} | O^s_a)}$$

Similarly we can modify the update for a positive observation to include the observation model and also do the same for the cost explanation. For calculation of cost confidence, we will now need to calculate $P(O_{\mathcal{C}(s,a) \geq k} | c_i \notin \mathbb{C}(s), \mathcal{C}^{abs}_s(\{c_i\}, a) \geq k)$. This can either be empirically calculated from samples with true label or we can assume that this value is going to be approximately equal to the overall distribution of the cost for the action.

## 8.7 Acquiring New Vocabulary

A core capability of the method discussed in the earlier sections is the ability to detect scenarios where the user specified concept list may not suffice to represent the required model component. This is important as it is unlikely that the user would always be able to provide the required concepts, either because they are unaware of it or they may just have overlooked some of the concepts when coming up with the list provided to the system. The former case could happen when the user is not an expert or if the simulator is modeling novel phenomena and is itself learned from experience. Systems capable of handling such scenarios would require the capability of teaching the user new concepts and even coming with intuitive labels for such concepts. We are unaware any methods that can handle this in a general enough manner.

Handling the latter case would generally be easier, since now the system would only need to query the human for a concept they may know but did not specify. One possible way such concepts can be queried is by presenting state pairs and asking the user to provide concepts absent in one but present in another. To see how we could perform such queries , consider the case of identifying failing preconditions. Let us denote the state where the foil fails as $s_f$. Now once the system has established that the current set of concepts cannot reliably explain the failing state, we can sample one of the states where action succeeds ($s'$ and we will refer to such states as positive states) and present the two states to the user and ask them to specify a concept that is absent in state $s_f$ but is present in $s'$. Now the system can keep showing more and more positive states and confirm that identified concept is still present in the state. The process can continue till either the system has enough samples to reliably establish the confidence of the fact that the concept is a precondition or the user notes a state where the concept is absent. If the latter happens you can ask the user to pick another concept that is absent in the original failing state but present in all the positive states. We can try to reduce cognitive load on the user's end on choosing the right concept by sampling positive states that are similar to $s_f$ (which would generally mean less number of conceptual difference). Moreover, in cases where the user is being presented a visual representation of the state, then the system could also provide saliency maps showing areas in the state that are important to the decision maker (which should generally cover patches on the image corresponding

to precondition concept). Though this would mean restricting positive states to the ones where the decision maker would have used the failing action.

## 8.8  Bibliographic Remarks

Many of the specific methods discussed in the chapter were first published in Sreedharan et al. [2020c]. As mentioned in the chapter many of the methods also have parallels in explainable Machine Learning. For example, Ribeiro et al. [2016] is a popular explanation method for classifiers that relies on local approximation of models and Kim et al. [2018] generates concept based explanations, where individual concepts are captured through classifiers. Concept as classifier has also been utilized by Hayes and Shah [2017] to provide policy summaries and Waa et al. [2018] to provide contrastive explanations for MDPs. Though in the case of Waa et al. [2018], they require designer to also assign positive and negative outcomes to each action that is used to explain why one policy is preferred over another and their method also do not really allow for explanation of infeasible foils. Outside of explanations, the idea of using classifiers to capture state fluents have also been utilized by Konidaris et al. [2018] to learn post-hoc symbolic models from low-level continuous models. One possibility alluded to in the chapter though not discussed in details is that of using the methods to establish the cost of the plan itself. As one can guess here the process would be nearly identical to that for establishing the cost of the foil, but instead of trying to find max cost bounds, we will be looking for min ones. Also the utility of establishing a representation of the agent in symbolic terms the user can understand goes beyond just providing helpful explanation. The paper by Kambhampati et al. [2021] makes a case to develop a symbolic middle layer expressed in human-understandable terms as a basis for all human-AI interaction.

CHAPTER 9

# Obfuscatory Behavior and Deceptive Communication

In this chapter, we will focus the discussion on some of the behavioral and communication strategies that a robot can employ in adversarial environments. So far in this book, we have looked at how the robot can be interpretable to the human in the loop while it is interacting with her either through its behavior or through explicit communication. However, in the real world not all of the robot's interactions may be of purely cooperative nature. The robot may come across entities of adversarial nature while it is completing its tasks in the environment. In such cases, the robot may have some secondary objectives like privacy preservation, minimizing information leakage, etc. in addition to its primary objective of achieving the task. Further, in adversarial settings, it is not only essential to minimize information leakage but also to ensure that this process of minimizing information leakage is secure. Since, an adversarial observer may use diagnosis to infer the internal information and then use that information to interfere with the robot's objectives.

To prevent leakage of sensitive information, the robot should be capable of generating behaviors that can obfuscate the sensitive information about its goals or plans. However, the execution of such obfuscating behaviors may be expensive and in such cases it may be essential to balance the amount of obfuscation needed with the resources available. Further, there may be complex real-world settings that involve both cooperative as well as adversarial observers. In such mixed settings, the robot has to balance the amount of obfuscation desired for adversarial entities with the amount of legibility required for the cooperative entities. Additionally, in this chapter we will look at how communication and the model-reconciliation techniques could be used for deception, in particularly to generate *lies*. As we will see, white lies could even be used to improve the overall team performance, although the ethics of such an undertaking has to be carefully considered. Now let's look at different types of adversarial behaviors and communication strategies that are available to the robot.

## 9.1 Obfuscation

A robot can hide its true objective from an adversarial observer by making several objectives seem plausible at the same time. Thus the true objective stays obfuscated from the observer. The robot can choose to obfuscate its activities as well apart from its objectives. In this section, we will see different types of obfuscation approaches that a robot can employ to hide sensitive information in adversarial situations.
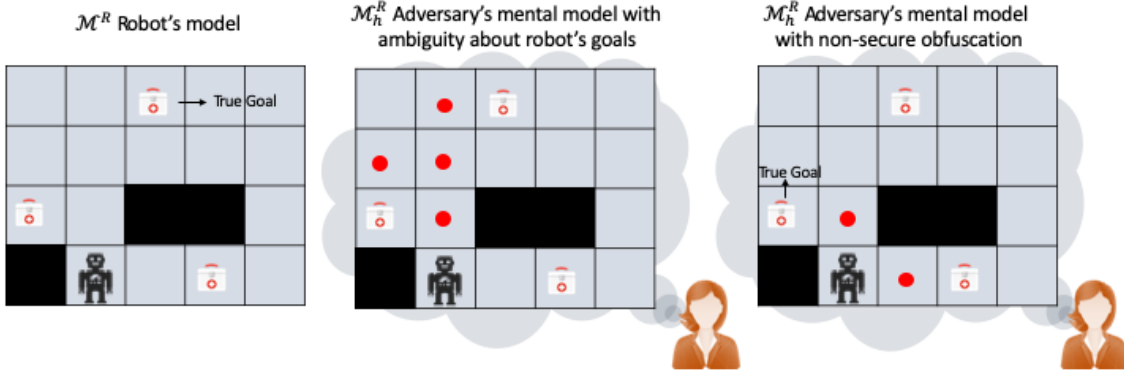
**Fig. 9.1:** Illustration of impact of goal obfuscation and secure goal obfuscation on human's mental model.

## 9.1.1   Goal Obfuscation

The problem of goal obfuscation arises when the robot wants to obfuscate its true goal from the adversary. This is possible in settings where the adversarial entity is unaware of the robot's true goal but is aware of the possible goals that the robot may try to achieve in the given domain. In this book, we will mostly focus on scenarios where the adversary is a human, but the methods we talk about could also be applied in scenarios where we have another autonomous agent as our adversary. This problem setting is also part of the controlled observability planning problem, introduced in Chapter 4. In fact, it is similar to the setting seen in the offline goal legibility problem. However, here because of the existence of an adversarial observer, instead of being legible with respect to its true goal, the robot obfuscates its true goal. In order to obfuscate effectively, the robot needs to have access to the adversaries models, incluging both their belief about the robot model ($\mathcal{M}_h^R$) and the observer's sensor model ($\mathcal{O}$). As in the case of Chapter 4, the adversaries may not be able to differentiate between certain states/actions leading to partial observability. The adversary may use these observations resulting from the robot's behavior to infer some sensitive information about the robot's goals. Depending on the granularity of the observations available to the adversary, the problem of goal obfuscation becomes easy or complex. That is with fine-grained observations, goal obfuscation might be harder to achieve or may even be infeasible in some cases, whereas with coarse-grained observations, it might be easier for the robot to obfuscate its true goal. Further, the adversary may be able to use the information gleaned from observations to interfere with or hamper the robot's activities. Therefore, in this setting, the robot is tasked with the additional objective of preventing the adversary from learning sensitive information about its goal. Similar to Chapter 4, we will focus on scenarios where the adversaries expectations about the robot model aligns with the robot's true model, i.e., $\mathcal{M}_h^R$ is the same as $\mathcal{M}^R$.

**Definition 38.** *A **goal obfuscation planning problem**, is a controlled observability planning problem, where, $\mathcal{G} = \{G^R \cup G_1 \cup \ldots \cup G_{n-1}\}$, is the set of n goals where $G^R$ is the true goal of the robot, and $G_1, \ldots, G_{n-1}$ are the decoy goals.*

A solution to a goal obfuscation planning problem is a *k-ambiguous* plan. Recall that, since the observer receives observations of the robot's behavior, each plan is associated with a corresponding observations sequence. Therefore, the objective of a *k-ambiguous* plan is to make the observation sequence consistent with at least *k* goals, out of which

$k - 1$ are decoy goals, such that, $k \leq n$. This is done by taking actions towards robot's true goal, such that, the corresponding observation sequence exploits the observer's belief space in order to be consistent with multiple goals. That is, a *k-ambiguous* plan produces a final belief state consistent with its observation sequence where at least $k$ goals are possible. The $k - 1$ decoy goals can be chosen specifically to maximize the obfuscation. Figure 9.1 illustrates goal obfuscation for $k = 2$ in the middle image. Here the true goal of the robot to pickup the medkit marked in the leftmost image. The adversarial observer has partial observability of the direction the robot moves in (say, the observer receives coarse observations from the GPS sensor). In order to obfuscate its true goal from the adversarial observer, the robot performs a 4-step k-ambiguous plan (for k=2), as shown in the middle image.

To define a k-ambiguous plan, we use the notation for belief state, introduced in Chapter 4.

**Definition 39.** *A plan, $\pi_k$, is a **k-ambiguous plan**, if $\delta(I, \pi_k) \models G^R$ and the final belief state, $b_n \in BS(\pi_k, I)$, satisfies the following, $|G \in \mathcal{G} : \exists s \in b_n, s \models G| \geqslant k$, where $1 \leqslant k \leqslant n$.*

**Definition 40.** *An observation sequence $O_k = \langle o_1, \ldots, o_n \rangle$ is **k-ambiguous observation sequence** if it is an observation sequence emitted by a k-ambiguous plan.*

The underlying algorithm for the synthesis of the solution is the same as the one described for the offline legibility planning problem. However, the goal test and heuristic function are customized to ensure that the goal obfuscation constraints are met.

### 9.1.1.1 Goal test

The goal test ensures that the solution plan does not end unless all $k$ goals occur in the observer's belief state. Specifically, it validates whether the true goal has been achieved in the robot's true state and also whether the $k - 1$ goals have been achieved in the belief state. If there doesn't exist a plan that can achieve the true goal with the given $k$ goals, we restart the search with next combination of $k$ goals (inclusive of true goal) among the total $n$ candidate goals. In the worst case, the value of $k$ can be reduced until a solution is found.

### 9.1.1.2 Heuristic function

Here the heuristic strategy is to generate a plan where the observer's final belief state satisfies $k$ goals, making it a *k-ambiguous* plan. We use two heuristic functions to estimate a node's cost: $h(s) = h_{G_A}(s) + h_{\mathcal{G}_{k-1}}(b)$, where the first part computes the heuristic cost to the true goal from the robot's actual state, and the second part computes the heuristic cost to $k - 1$ goals from the belief induced by the emitted observation. This is computed by finding the *max* of heuristic cost to each of the $k - 1$ goals. The heuristic value of a state is computed as the total of these two values. Thus the heuristic helps to guide the search towards at least $k$ goals in the observer's belief.

### 9.1.2   Secure Goal Obfuscation

The above approach to synthesizing goal obfuscating behavior is not secure. In the sense that it reveals information when the following assumptions are not satisfied: (1) adversary has access to the algorithm and has the means to simulate the algorithm for a given input. This means the adversary could identify the potential states generated by the algorithm. (2) the adversary can run the algorithm with different inputs. For example, the adversary can set different goals as true goal to check the variability of the output. Figure 9.1, rightmost subfigure illustrates how an adversary can rerun the algorithm with different input and glean information about the true goal. By setting the true goal to some other candidate goal, the adversary can run the algorithm and find out the robot's true goal which gave the original observation sequence.

One way to synthesize a secure goal obfuscating behavior would be to introduce some randomization in the solution synthesis process, that would prevent leakage of additional information. To that end, one could randomly choose a set of $k$ goals (inclusive of true goal) from $n$ candidate goals. Then randomly pick a goal from this set as a "true goal" and run the algorithm. Let's call this goal as the decoy true goal. We have to maintain the plans to each of the $k$ goals that are consistent with the observation sequence (which includes the robot goal $G^R$). Then by performing a post processing step, we can retrieve the plan that reaches the robot's true goal and is consistent with the observation sequence produced by the plan to the decoy true goal. This process ensures that even if the adversary runs the algorithm again with either the same or different input (different goal as true goal, or different combination of $k$ goals), the identity of the true goal remains *k-ambiguous* even with the observer being aware of the randomization. Note that unlike the other sections here we are not just considering an observer who is reasoning on their own based on the observations, but rather we are considering a sophisticated adversary who is able to use more sophisticated computing resources to simulate the algorithm.

In this setting, a perfectly rational adversary would be one that has unbounded computational resources, and can use them to make perfect inferences. In other words, that is the worst case scenario for this algorithm. The following proposition shows that our algorithm provides *k-ambiguous* guarantee even in the worst case, such that a perfectly rational adversary remains uncertain about the true goal, provided the goals are independent.

**Proposition 10.** *Let $\mathcal{H}$ be a perfectly rational adversary. Given a goal obfuscation planning problem for the robot $\mathcal{R}$. Plans computed using Algorithm 3 with k-ambiguous goal test and heuristic function are such that, $\mathcal{H}$ will be able to infer the true goal with probability $\leqslant 1/k$ provided the goals are independent.*

Let's say goal $G_i$ from the set $\mathcal{G}_k$ is chosen randomly to be the decoy true goal. The observation sequence $O_i$ is obtained by running the algorithm with $G_A = G_i$. The adversary can run the algorithm with each goal from $\mathcal{G}_k$ to search for $G_i$. But $G_i$ can only be the robot's true goal with probability $1/k$. Therefore the adversary can detect the true goal with probability $\leqslant 1/k$.

Note that the goals can be specified using sub-formulae. This makes the input specification much smaller depending on how many states are in the domain that are consistent
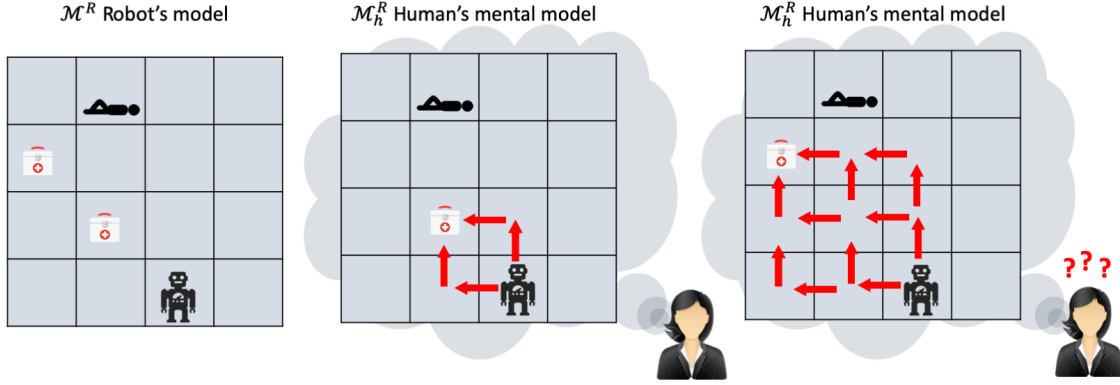
**Fig. 9.2:** Illustration of the impact of plan obfuscation on human's mental model.

with the sub-formulae. In the least constrained domain, that may be exponential in the number of propositions not used.

### 9.1.3 Plan Obfuscation

In this problem setting, the adversarial observer is aware of the robot's goal. The robot, however, obfuscates its activities towards its goal thereby hiding sensitive information about its activities. We refer to this behavior as plan obfuscatory behavior. The robot can achieve a plan obfuscatory behavior by computing a plan whose observation sequence conforms to a set of plans that are diverse in terms of their actions, thereby making it hard for the adversarial observer to guess the actions executed by the robot. For instance, in Figure 9.2, the goal of the robot is to pick up a medkit and treat the victim on the disaster site. It has two choices of medkits to pick from. If it chooses the medkit as shown in the rightmost sub-figure, it can maximize the diversity with respect to its actions. Here, the multiple different paths to the medkit produce the same observation sequence (due to the underlying coarse GPS sensor model that the observer is using), thus hiding the robot's true actions from the adversarial observer.

**Definition 41.** *A **plan obfuscation planning problem** is a controlled observability planning problem, where, $|\mathcal{G}| = 1$, and therefore the true goal is known to the observer.*

The solution to a plan obfuscation problem is an $\ell$-diverse plan. An $\ell$-diverse plan is a plan whose observation sequence is consistent with at least $\ell$ diverse plans to the goal, such that, the diverse plans are at least $d$ distance away from each other. In order to compute an $\ell$-*diverse* plan, we need to keep track of the plans that are consistent with the observation sequence and reach the goal. To compute the diversity between all the pairs of plans consistent with the observation sequence, a plan distance measure like action distance, causal link distance, state sequence distance (introduced in Chapter 3) can be used. In general, the approach can use any valid plan distance. We now define an $\ell$-*diverse* plan.

We refer to the set of plans that are consistent with the belief sequence of a given plan, as a belief plan set (BPS), introduced in Chapter 4.

**Definition 42.** *A BPS induced by plan p starting at $s_0$ is **minimally d-distant**, $d_{min}(BPS(p, s_0))$, if $d = \min\limits_{p1,p2 \in BPS(p,s_0)} \mathcal{D}(p1, p2)$.*

**Definition 43.** *A plan, $\pi_l$, is an **$\ell$-diverse plan**, if for a given value of d and distance function $\delta$, $d_{min}(BPS(\pi_l, I)) \geq d$, $|BPS(\pi_l, I)| \geq \ell$, where $\ell \geq 2$ and every plan in $BPS(\pi_l, I)$ achieves the goal in $\mathcal{G}$.*

We can again leverage Algorithm 3 to generate the plans for obfuscation problem. The goal test and heuristic function are customized to ensure that there are at least $\ell$ diverse plans to the true goal that are consistent with the observation sequence and the minimum distance between these plans is at least $d$.

### 9.1.3.1    Goal test

To ensure the plans in $BPS$, induced by an $\ell$-*diverse* plan, can achieve the goal in $\mathcal{G}$, we can check whether at least $\ell$ plans are reaching the goal or not and whether the minimum distance between plans in $BPS$ is at least $d$. Also in order to ensure termination of the algorithm, there is a cost-bound given as input to the algorithm.

### 9.1.3.2    Heuristic function

Apart from minimizing the heuristic cost to the goal, the customized heuristic given below also maximizes the $d$ of $d_{min}(BPS(p, s_0))$ induced by plan $p$ starting at $s_0$. This increases the minimum distance between the plan pairs. This distance is computed using a plan distance measure.

$$h(s) = h_{G_A}(s) - d_{min}(BPS(p, s_0)) \tag{9.1}$$

### 9.1.4    Deception

In addition to obfuscation, the robot can also employ deceptive behaviors to hide sensitive information. While the concept of obfuscation deals with the notion of confusing the adversary with several decoy candidates, the concept of deception involves making one of the decoy candidates more likely than the robot's true objective/activities. Thus, in order to deceive an adversarial observer, it is crucial to have access to their goal or plan recognition models. By incorporating the human's reasoning ability (in regards to identifying the goal or the plan), the robot can synthesize behavior that deceives the adversarial observer into believing that the decoy candidate is the true candidate. In the prior literature, synthesis of deceptive behaviors with respect to robot's goals has been studied in path planning scenarios where the observer has full observability of robot's activities. In order to successfully deceive an adversarial observer, the robot's plan has to end when a decoy goal is achieved. However, in reality, the robot has a primary objective of achieving its true goal. Therefore, in cases where the observer has full observability of the robot's activities, deceptive behavior may be hard to synthesize. To that end, the

notion of the radius of maximum probability with respect to a goal has been studied in the literature on deceptive planning. This is a radius around a goal location, within which that goal itself becomes the most likely goal. So deception can be maintained until the robot reaches the area within the radius of maximum probability for its true goal.

## 9.2 Multi-Observer Simultaneous obfuscation and legibility

Obfuscatory strategies are also crucial in scenarios where there are multiple different types of observers, such that some are of adversarial nature while some others are of cooperative nature. In such cases, the robot has to ensure that its behavior is simultaneously legible to cooperative observers and obfuscatory to adversarial ones. For instance, in soccer, a player may perform feinting trick to confuse an opponent while signaling a teammate. Synthesizing a single behavior that is legible and obfuscatory to different observers presents significant technical challenges. In particular, the robot may have to be deliberately less legible to its friends so that it can be effectively more obfuscatory to its adversaries. This problem gives rise to a novel optimization space that involves trading-off the amount of obfuscation desired for adversaries with the amount of legibility desired for friends.

### 9.2.1 Mixed-Observer Controlled Observability Planning Problem

We now discuss a problem framework called *mixed-observer controlled observability planning problem*, MO-COPP, that allows a robot to simultaneously control information yielded to both cooperative and adversarial observers while achieving its goal. This framework models and exploits situations where different observers have differing sensing capabilities, which result in different "perceptions" of the same information. Typically, different observers in an environment may have different "sensors" (perception capabilities) due to differences in prior communication, background knowledge, and innate differences in sensing capabilities.

#### 9.2.1.1 Example

Consider the example in Figure 9.3, where an automated truck has to deliver two packages. As shown in Figure 9.3a, there are two packages at factory A and two at factory B. Let there be two observers who know the initial inventory at each factory, but do not know which two packages will be delivered by the truck. Each observer has different sensors (or sources of information) about the truck's activities. Observer-C (cooperative observer) has access to a database of barcodes which allows it to determine the identity of the factory that the package was loaded from (but not the identity of the individual packages), whereas observer-X (adversarial observer) does not have any additional information and cannot differentiate between any of the packages. The truck wants to convey its true goal to observer-C but obfuscate it to observer-X. When the truck loads one package each from A and B (Figure 9.3b), observer-C cannot distinguish the packages. However, when it loads both the packages from A (Figure 9.3c), observer-C can infer the identity of the loaded packages, while observer-X still does not know the identity of the packages. Therefore,
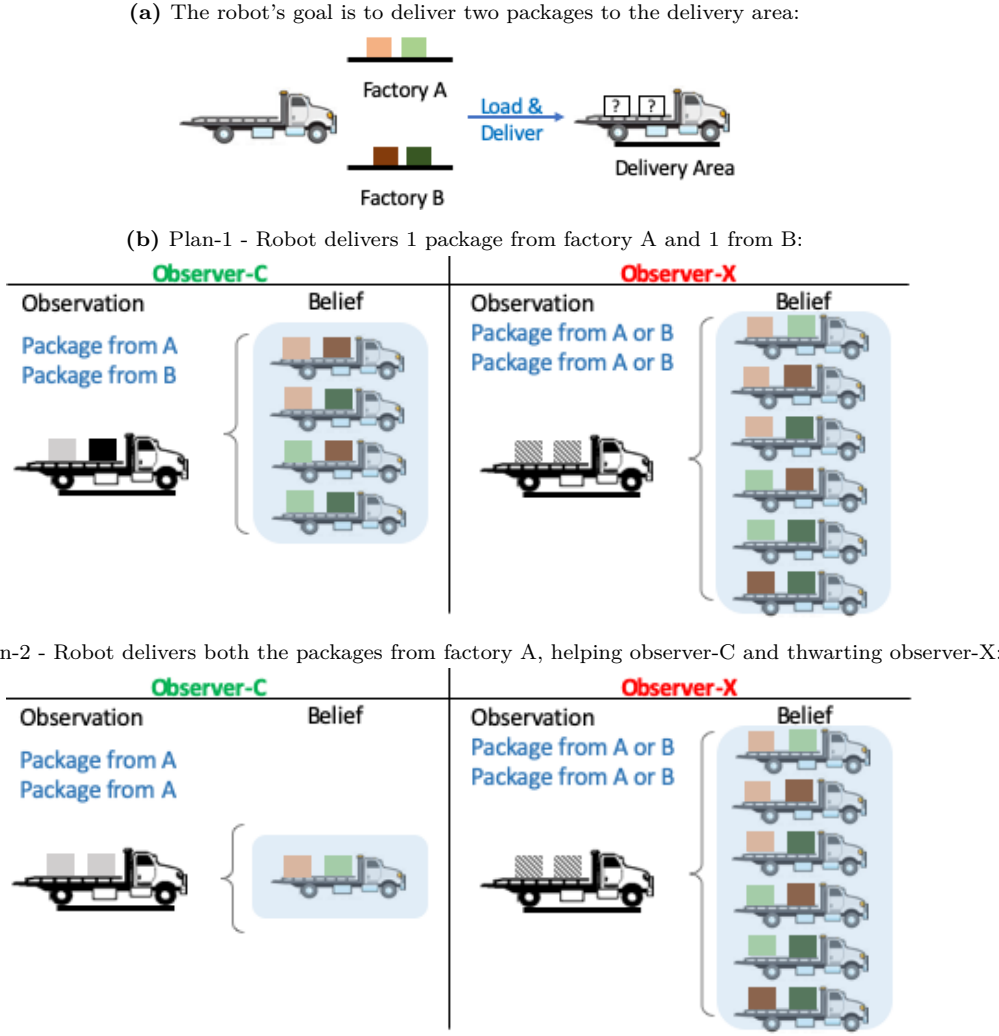
**(a)** The robot's goal is to deliver two packages to the delivery area:



**(b)** Plan-1 - Robot delivers 1 package from factory A and 1 from B:



**(c)** Plan-2 - Robot delivers both the packages from factory A, helping observer-C and thwarting observer-X:



**Fig. 9.3:** Different belief updates induced by the same plan for two observers; observer-C is cooperative, observer-X is adversarial.

plan-2 makes the true goal legible to observer-C while obfuscating the true goal with all the possible goals for observer-X.

MO-COPP setting involves an actor ($\mathbf{R}$) (i.e., the robot) and two observers, where one is adversarial observer ($\mathbf{X}$) while the other is cooperative ($\mathbf{C}$). This framework assumes that the robot has full observability of its own activities and knows the sensor models used by the observers, and that the observers have different sensor models. When the robot takes an action and reaches a new state, each observer receives an observation. After obtaining the observations, each observer updates its belief, similar to the belief update in COPP formulation. The robot leverages the known limits in the observers' sensors to control the observability of multiple observers in the environment simultaneously. Given a set of candidate goals, the objective of the robot is to convey information about its goal to the cooperative observer and to hide it from the adversarial observer. This involves increasing the number of candidate goals possible in the adversary's belief, while decreasing the number of candidate goals possible in cooperative observer's belief.

**Definition 44.** *A **mixed-observer controlled observability planning problem** is a tuple,*

$MO - COPP = \langle \Lambda, \mathcal{M}^R, \mathcal{G}, \{\Omega_i\}_{i \in \Lambda}, \{\mathcal{O}_i\}_{i \in \Lambda}, \{\mathcal{B}_0^i\}_{i \in \{\mathbf{X}, \mathbf{C}\}} \rangle$, *where,*

- $\Lambda = \{\mathbf{R}, \mathbf{C}, \mathbf{X}\}$ *is the set of agents.*

- $\mathcal{M}^R = \langle F^R, A^R, I^R, G^R \rangle$ *is the robot model Also, for* $a \in A$, $pre(a)$, $adds(a)$, $dels(a)$ *are each a subset of fluents representing precondition, add effects and delete effects of* $a$.

- $\mathcal{G} = \{G_1, G_2, \ldots, G_{n-1}, G^R\}$ *is the set of candidate goals, where* $G^R$ *is the true goal of* $\mathbf{R}$, *which is not known to both* $\mathbf{C}$ *and* $\mathbf{X}$.

- $\Omega_i$ *is the set of observation symbols for agent* $i$, *which are emitted when* $\mathbf{R}$ *takes an action and reaches a new state. Further,* $\Omega_R = \{o_{a,s}^{\mathbf{R}} | a \in A, s \in \mathcal{S}\}$.

- $\mathcal{O}_i : A \times \mathcal{S} \to \Omega_i$ *is agent* $i$*'s deterministic sensor model.* $\mathcal{S}$ *is the set of states, where each state is instantiation of all fluents. Further,* $\mathcal{O}_{\mathbf{R}}$ *maps each action-state pair to a unique observation,* $\forall\, a, a' \in A,\, s, s' \in \mathcal{S}, a \neq a' \wedge s \neq s' : \mathcal{O}_{\mathbf{R}}(a,s) \neq \mathcal{O}_{\mathbf{R}}(a',s')$, *while* $\mathcal{O}_{\mathbf{X}}$ *and* $\mathcal{O}_{\mathbf{C}}$ *are noisy sensor models that map multiple action-state pairs to the same observation symbol.*

- $\mathcal{B}_0^i$ *is the initial belief of an observer,* $i \in \{\mathbf{X}, \mathbf{C}\}$. *The initial belief is a set of states inclusive of* $I$.

Every time the robot acts, each $i \in \Lambda$ receives an observation consistent with its sensor model. The sensor model of an observer $i \in \{\mathbf{X}, \mathbf{C}\}$ supports many-to-one mapping of $\langle a, s \rangle$ pairs to observation symbols, i.e., $\exists a, a' \in A, s, s' \in \mathcal{S}, a \neq a' \wedge s \neq s' : \mathcal{O}_i(a,s) = \mathcal{O}_i(a',s')$. For an agent $i$, the inverse of sensor model gives the set of $\langle a, s \rangle$ pairs consistent with an observation symbol $o^i \in \Omega_i$, i.e., $\mathcal{O}_i^{-1}(o^i) = \{\langle a, s \rangle | \forall a \in A, s \in \mathcal{S}, \mathcal{O}_i(a,s) = o^i\}$. Again, we assume the observer is aware of the true robot model.

Each observer $i \in \{\mathbf{X}, \mathbf{C}\}$ maintains its own belief, which is a set of states. $\delta(\cdot)$ is a transition function, such that, $\delta(s, a) = \bot$ if $s \not\models pre(a)$; else $\delta(s, a) = s \setminus dels(a) \cup adds(a)$. Now we can define the belief update: (1) at time step $t = 0$, the initial belief of observer $i$ is given by $\mathcal{B}_0^i$, (2) at time step $t \in \{1, \ldots, \mathcal{T}\}$, $\mathcal{B}_t^i = \{\hat{s} \mid \exists \hat{a}, \bar{s}\ \delta(\bar{s}, \hat{a}) = \hat{s}; \bar{s} \in \mathcal{B}_{t-1}^i; \mathcal{O}_i(\hat{a}, \hat{s}) = o_t^i\}$, where $\mathcal{T}$ is the time horizon. That is, the belief is updated using the previous belief and the observation received. A plan $\pi$ (or sequence of actions) is associated with a sequence of observations, $ObsSeq_i(\pi) = \{o_1^i, \ldots, o_{\mathcal{T}}^i\}$ for each observer $i$. The belief of $\mathbf{R}$ is a single state due to full observability.

### 9.2.1.2 MO-COPP Solutions

A plan $\pi$ solves a $MO - COPP = \langle \Lambda, \mathcal{M}^R, \mathcal{G}, \{\Omega_i\}_{i \in \Lambda}, \{\mathcal{O}_i\}_{i \in \Lambda}, \{\mathcal{B}_0^i\}_{i \in \{\mathbf{X}, \mathbf{C}\}} \rangle$, iff: $\delta(I, \pi) \models G_A$. In other words, any solution to $\mathcal{M}^R$ is a solution to MO-COPP. We measure the quality of a MO-COPP solution in terms of the extent to which the plan is legible to $\mathbf{C}$ and obfuscating to $\mathbf{X}$. Given a solution plan, $\pi$, let $\mathcal{G}_{\mathbf{X}}^{\pi} = \{G \in \mathcal{G} : \exists s \in \mathcal{B}_{\mathcal{T}}^{\mathbf{X}},\ s \models G\}$, denote the goals present in $\mathbf{X}$'s final belief and $\mathcal{G}_{\mathbf{C}}^{\pi} = \{G \in \mathcal{G} : \exists s \in \mathcal{B}_{\mathcal{T}}^{\mathbf{C}},\ s \models G\}$, denote the goals present in $\mathbf{C}$'s final belief. Also, we can define the inverse of its observation sequence for observer $i$, as follows: $\forall i \in \{\mathbf{X}, \mathbf{C}\}, ObsSeq_i^{-1}(\pi) = \{\hat{\pi} | \exists s_0^i \in \mathcal{B}_0^i, G \in \mathcal{G}_i^{\pi}, \delta(s_0^i, \hat{\pi}) \models G, ObsSeq_i(\pi) = ObsSeq_i(\hat{\pi})\}$.

Here $|\mathcal{G}_{\mathbf{X}}^{\pi}|$ represents the amount of **goal obfuscation for X**, and by increasing the number of goals in $\mathcal{G}_{\mathbf{X}}^{\pi}$, we can increase the goal obfuscation. Similarly, $|\mathcal{G}_{\mathbf{C}}^{\pi}|$ represents the amount of **goal legibility for C**, and by decreasing the number of goals in $\mathcal{G}_{\mathbf{C}}^{\pi}$, we can increase the goal legibility. We now introduce a goal difference metric to estimate the quality of a solution plan to MO-COPP.

**Definition 45.** *Given a solution plan $\pi$ that solves MO-COPP, where $|\mathcal{G}| = n$, the **goal difference**, of $\pi$ is given by: $GD(\pi) = \frac{|\mathcal{G}_{\mathbf{X}}^{\pi}| - |\mathcal{G}_{\mathbf{C}}^{\pi}|}{n-1}$ where the denominator represents the difference between ideal values of $\mathcal{G}_{\mathbf{X}}^{\pi}$ and $\mathcal{G}_{\mathbf{C}}^{\pi}$.*

An optimal solution to MO-COPP maximizes the trade-off between amount of goal obfuscation and goal legibility. That is, it maximizes the difference between the number of goals in $\mathcal{G}_{\mathbf{X}}^{\pi}$ and $\mathcal{G}_{\mathbf{C}}^{\pi}$. Equivalently, closer the $GD(\pi)$ value to 1, better is the plan quality. A solution plan with $GD(\pi) = 1$ is an optimal plan. E.g., in Figure 9.3c, plan-2 is an optimal plan with $GD(plan\text{-}2) = \frac{6-1}{6-1} = 1$. The denominator is essential for comparing the quality of plans across different problems, with varying number of candidate goals.

**Proposition 11.** *Given a solution plan, $\pi$, to MO-COPP, if $|\mathcal{G}_{\mathbf{C}}^{\pi}| = 1$, then $G^R \in \mathcal{G}_{\mathbf{C}}^{\pi}$.*

The above proposition states that when maximum goal legibility is achieved, only one goal $G_{\mathbf{A}}$ is present in **C**'s final belief. This is because $I \in \mathcal{B}_0^{\mathbf{C}}$, and while performing belief update at $t$, $\langle a_t, s_t \rangle \in \mathcal{O}_{\mathbf{C}}^{-1}(o_t^{\mathbf{C}})$. That is the true state is part of each belief update. Therefore, $G^R \in \mathcal{G}_{\mathbf{C}}^{\pi}$.

### 9.2.2 Plan Computation

We now present two solution approaches to solve MO-COPP. In the first approach, MO-COPP can be formulated as a constraint optimization problem and provide an IP encoding to solve it in $\mathcal{T}$ steps. The IP encoding provides an optimal solution for the given horizon by maximizing the trade-off between the amount of obfuscation and legibility. The second approach builds up on the search algorithm presented in Chapter 4 to solve COPP problems. The search algorithm generates solutions that satisfy a prespecified lower bound on the amount of goal obfuscation and goal legibility.

#### 9.2.2.1 MO-COPP as Integer Program

The IP encoding includes variables for an action being executed at a time, each state being true at a time, for a specific observation to emitted at a give time and for specific state and goal being in the belief state. The objective function is essentially the numerator of $GD(\cdot)$ metric, i.e.,

$$max \quad \sum_{G \in \mathcal{G}} g_{G,\mathcal{T}}^{\mathbf{X}} - \sum_{G \in \mathcal{G}} g_{G,\mathcal{T}}^{\mathbf{C}} \tag{9.2}$$

Where $g_{G,\mathcal{T}}^i$ is an indicator variable capturing the fact that goal $G$ is present in the belief state for the agent $i$ at the time step $\mathcal{T}$. The denominator of the *GD* metric is skipped, as it is a constant and does not contribute to the optimization. Maximizing for this objective provides a single solution that achieves the maximum difference between the number of goals possible for the two observers. One could also formalize the problem as a multi-objective optimization problem, and look at various pareto optimal solutions that trade off the various goals that each observer believes will be achieved.

The constraints provided to the problem, captures the exact evolution of the state and belief (for each agent) given the robot action, the observation generated by them and the fact that the robot needs to achieve its goal within the specified horizon limit. The exact IP encoding can be found in the paper Kulkarni et al. [2020b]. Also note that, the current objective function trades off goal obfuscation with goal legibility for the observers. However, the robot can ensure a predefined level of goal obfuscation (say obfuscate with *at least k* candidate goals) by adding an additional constraint that enforces a bound for goal obfuscation and goal legibility.

### 9.2.2.2 Search Algorithm

It is also possible to leverage search techniques that address goal obfuscation and goal legibility in isolation to solve MO-COPP. The search algorithm in Chapter 4 is adapted to address goal obfuscation and goal legibility simultaneously to two different observers. The bounds on the amount of goal obfuscation and goal legibility desired can be specified, similar to the ones seen in the IP: obfuscate with at least $k$ goals, make it legible with at most $j$ goals. These bounds, $\Phi = \langle \Phi_{\mathbf{X}}, \Phi_{\mathbf{C}} \rangle$, are given as input to the search algorithm.

---

**Algorithm 7** Heuristic-Guided Search

---

1: Initialize *open*, *closed* and *temp* lists; $\Delta = 1$
2: $\langle b_\Delta^X, b_\Delta^C \rangle \leftarrow \text{approx}(I, \mathcal{B}_0^X, \mathcal{B}_0^C)$
3: $open.\text{push}(I, \langle b_\Delta^X, b_\Delta^C \rangle, \langle \mathcal{B}_0^X, \mathcal{B}_0^C \rangle, priority = 0)$
4: **while** $\Delta \leqslant |\mathcal{S}|$ **do**
5:     **while** $open \neq \emptyset$ **do**
6:         $s, \langle b_\Delta^X, b_\Delta^C \rangle, \langle \mathcal{B}^X, \mathcal{B}^C \rangle, h_{node} \leftarrow open.\text{pop}()$
7:         **if** $|b_\Delta^X| > \Delta$ or $|b_\Delta^C| > \Delta$ **then**
8:             $temp.\text{push}(s, \langle b_\Delta^X, b_\Delta^C \rangle, \langle \mathcal{B}^X, \mathcal{B}^C \rangle, h_{node})$
9:             continue
10:         add $\langle b_\Delta^X, b_\Delta^C \rangle$ to *closed*
11:         **if** $s \models G_A$ and $\mathcal{B}^X \models \Phi_X$ and $\mathcal{B}^C \models \Phi_C$ **then**
12:             **return** $\pi, ObsSeq_X(\pi), ObsSeq_C(\pi)$
13:         **for** $s' \in \text{successors}(s)$ **do**
14:             $o^X \leftarrow \mathcal{O}_X(a, s'); o^C \leftarrow \mathcal{O}_C(a, s')$
15:             $\widehat{\mathcal{B}^X} = \text{Update}(\mathcal{B}^X, o^X); \widehat{\mathcal{B}^C} = \text{Update}(\mathcal{B}^C, o^C)$
16:             $\langle \widehat{b_\Delta^X}, \widehat{b_\Delta^C} \rangle \leftarrow \text{approx}(s', \widehat{\mathcal{B}^X}, \widehat{\mathcal{B}^C})$
17:             $h_{node} \leftarrow h_{G_A}(s') + h_{\mathcal{G}_{k-1}}(\widehat{\mathcal{B}^X}) - h_{\mathcal{G}_{\mathcal{G}-j}}(\widehat{\mathcal{B}^C})$
18:             add new node to *open* if $\langle \widehat{b_\Delta^X}, \widehat{b_\Delta^C} \rangle$ not in *closed*
19:     increment $\Delta$; copy items from *temp* to *open*; empty *temp*

---

Each search node maintains the associated beliefs for both observers. The *approx* function generates an approximate belief, $b_\Delta^i$, of size $\Delta$ (i.e. cardinality of $b_\Delta^i$ is $\Delta$). $b_\Delta^i$ is always inclusive of the true state of the robot, this is because the robot can only take actions that are consistent with its true state. If all such $\Delta$−sized beliefs (or smaller) are explored then $b_\Delta^i$ of $\Delta + 1$ size is computed, and this node gets put in the temporary list and is explored in the next outer iteration when $\Delta$ has been incremented. For each $\Delta$, all $\Delta$-sized unique combinations of belief (that include the actual state of the robot) are explored. This allows systematic and complete exploration of multiple paths to a given search node. The inner iteration performs heuristic guided forward search (we use greedy best first search) to find a plan while tracking at most $\Delta$ states in each $b_\Delta^i$. In the inner loop, the node expansion is guided by (1) customized heuristic function, which computes value of the node based on true goal and belief constraints given by $\Phi$ for the observers, and (2) goal test, which checks for satisfaction of true goal and satisfaction of the belief constraints given by $\Phi$. The algorithm stops either when a solution is found or when all the $\Delta$ iterations have been explored.

**Property** In both the solution approaches, we can assert a lower bound on the extent of goal obfuscation and goal legibility for a MO-COPP solution plan. In IP, we can specify the aforementioned goal constraints to assert this minimum value, while in the search, the goal tests allow us to assert it. By setting $k$, $j$ to desired values, we can eliminate solutions with low $GD$ score. This affords the following guarantee:

**Proposition 12.** *Let $\mathbf{X}$ and $\mathbf{C}$ be perfectly rational adversarial and cooperative observers respectively. Given a MO-COPP $= \langle \Lambda, \mathcal{M}^R, \mathcal{G}, \{\Omega_i\}_{i \in \Lambda}, \{\mathcal{O}_i\}_{i \in \Lambda}, \{\mathcal{B}_0^i\}_{i \in \{\mathbf{X}, \mathbf{C}\}} \rangle$ with equally probable and independent goals, $|\mathcal{G}| = n$, and goal constraints of at least $k$ goal obfuscation for $\mathbf{X}$ and at most $j$ goal legibility for $\mathbf{C}$, then a solution plan, $\pi$, gives the following guarantees:*

1. *$\mathbf{X}$ can infer $G_\mathbf{A}$ with probability $\leqslant 1/k$,*

2. *$\mathbf{C}$ can infer $G_\mathbf{A}$ with probability $\geqslant 1/j$, and*

3. *Goal difference metric, $GD(\pi) \geqslant \frac{k-j}{n-1}$*

Given $|\mathcal{G}_\mathbf{X}^\pi| \geqslant k$, $X$ can infer $G_\mathbf{A}$ with probability $\leqslant 1/k$. Similarly, given $|\mathcal{G}_\mathbf{C}^\pi| \leqslant j$, $C$ can infer $G_\mathbf{A}$ with probability $\geqslant 1/j$. Also, similarly, $GD(\pi) \geqslant \frac{k-j}{n-1}$. The above proposition states that, based on the observation equivalence there is no additional information revealed about the robot's true goal. Therefore, we can derive goal detection upper bound for $\mathbf{X}$ and lower bound for $\mathbf{C}$. Also this allows us to derive a lower bound on the plan quality.

## 9.3 Lies

Most accounts of explanations studied in this book consider cases where an automated agent shares information about its decision-making model, to accurately justify decisions it has made. It is of course possible that if a system would desire so, it could hijack such a model reconciliation communication process to shape the observer's expectation to be

divergent from the true model that is being used by the robot. In this section, we will refer to such manipulation of the model reconciliation process as *lies*. As with the implicit communication/behavioral manifestations of deception, lies, and explanations are inexorably connected. One could argue that any poct hoc reconstructive explanations discussed in Chapter 8 constitute a form of lies (insofar as the model being communicated was never explicitly used in its decision-making process), or by relying on minimal explanations, the agent is engaging in lies by omission (insofar that they may be relying on the human's false beliefs to simplify explanations). To simplify the discussion, when we refer to lies we are referring to scenarios wherein the agent is communicating information it knows to be false.

### 9.3.1 When Should an Agent Lie?

To start with, the first question we might want to ask is, why would we ever want an agent that lies? Outside purely adversarial scenarios, are there cases where we might want to employ such systems? Rather than look at scenarios where the agent is taking an explicitly adversarial stance, we will consider cases where the lies are used in service of improving the team utility in some capacity. There is enough evidence from social sciences that lies do play a (potentially positive) role in many real world human-human interaction scenarios. For example consider the interaction between a doctor and a patient. There are few scenarios where the role of lies have been thoroughly debated than in doctor-patient interactions. Whether it relates to the prescription of placebo or withholding certain information from the patient, there are many scenarios within a doctor-patient relationship wherein the people have argued for potential usefulness of deception or lies. In the case of human-machine teaming scenario, some potential use cases for such lies include

1. Belief Shaping: Here the system could choose to shape the beliefs of their potential teammate, through explicit communication, to persuade your teammate to engage behavior the autonomous agent believes would result in higher utility.

2. White Lies: This case could involve cases where the exact explanation may be too expensive to communicate, and the robot could choose to communicate a simpler 'explanation', which while technically contains untrue information is used to justify the optimal behavior. Such white lies could be particularly helpful when it may be expensive computationally or time consuming for the human teammate to process the original explanations.

### 9.3.2 How can an Agent Lie?

Now the question would be how one could generate lies. In general one could leverage the same model space search, and instead of restricting the model updates to only those constrained by the true robot model, in this case the search is free to make any update on the model. If the agent is not allowed to introduce new fluents or new actions, this is still a finite search space for classical planning domains. Though this would constitute a much larger search space than the one considered in the explanation, which brings up an

interesting point that *at least in computational terms telling the truth might be an easier strategy to following.*

One way to constrain the search space would be to limit the kind of changes that are allowed under the lies. One natural choice may be to limit lies to those that remove factors from the propositional representation of the model (i.e $\Gamma(\mathcal{M}_h^R)$) (Section 5.1). Such lies have been sometimes referred to as *lies by omission* in the literature. Another possibility may be to leverage theories of model evolution or drift when such information may be available to consider believable changes to models. Or if the lies need to include new possibilities (including new actions or fluents), consider leveraging large language models (like Brown et al. [2020]) or knowledge bases like WordNet [Fellbaum, 2010] to introduce them in a meaningful way. That is consider existing actions and fluents and use the knowledge base to look for related concepts and try to introduce them into the model.

### 9.3.3  Implications of Lies

User studies (as reported in Chakraborti and Kambhampati [2019b]) have shown that at the very least lay users seem to be open to the idea of an agent engaging in such deceptive behavior, when it is guaranteed to lead to higher team utility. In this book, we will not investigate the moral and ethical quandaries raised by designing an agent capable of lying, but rather look at a much simpler question – if we are allowing for deceptive behavior in the hope of higher utility, how does one even guarantee that such behavior would in fact lead to higher team utility? In some cases, the agent could very well be aware of the fact that it has access to more information than the other humans in the environment (owing to more sophisticated sensors or its location in the environment) and it may be confident in its ability to process and reason with the information it has access to. Though in many scenarios there is a small probability that the information it has access to is incorrect (owing to a faulty sensor or just change in information) or it may have overlooked some factor in its computation. As such the lie could lead to unanticipated negative outcomes. In such cases, it is quite possible that the human teammate could be a lot more critical of its automated teammate that chose to lie as compared to a case where the robot made an unintentional mistake. To us, this speaks for the need for significant further work to be done to not only understand under what conditions an automated system could confidently make such calls, but also better tools to model trust of teammates.

## 9.4  Bibliographical Remarks

The different obfuscatory behaviors discussed in this chapter have been formulated within the controlled observability planning framework introduced by Kulkarni et al. [2019b]. The work on deceptive planning was introduced by Masters and Sardina [2017]. The generalized extension of controlled observability framework (referred to as MO-COPP) used for balancing goal obfuscation for adversaries with the goal legibility for cooperative observers was introduced by Kulkarni et al. [2020b]. In terms of lies, the discussion provided is based on the papers, Chakraborti and Kambhampati [2019b] and Chakraborti and Kambhampati [2019a], where the works respectively considered when and why the agents

should consider generating lies and the computational mechanisms that could be employed towards generating such explanations. Outside of these specific papers, there is a lot of work investigating the utility of lies in various teaming scenarios. There exists a particular extensive literature on the role of lies in decision-making in medical literature Palmieri and Stern [2009]. In particular, many works have argued for the importance of a doctor withholding information from the patient (cf. Korsch and Harding [1998] and Holmes [1895]).

CHAPTER 10

# Applications

In this section, we will look at four different applications that leverage the ideas discussed in this book. In particular, all the systems discussed in this chapter will explicitly model the human's mental model of the task and among other things use it to generate explanations. In particular, we will look at two broad application domains. One where the systems are designed for collaborative decision-making, i.e systems designed to help user come up with decisions for a specific task and another system designed for helping users specify a declarative model of task (specifically in the context of dialogue planning for an enterprise chat agent).

## 10.1 Collaborative Decision-Making

Our first set of applications will be centered around systems that are designed to help end-users make decisions.

## 10.2 Humans as Actors

All the theoretical and formal discussions in the book until now have focused on scenarios, where the robot is the one acting in the world and the human is just an observer. The specific decision-support scenarios, we will look at in this section, requires us to consider a different scenario, one where the human is the primary actor and the robot is either an observer or an assistant to the robot. The setting is illustrated in figure 10.1. In this scenario, we have a human with a model $\mathcal{M}^H$ who is acting in the world, and a robot observing the actor. The robot has access to an approximation of the human model $\mathcal{M}_r^H$ and may also have access to a model $\mathcal{M}^*$ which they believe is the true model of the task that the human is pursuing. In addition to the decision-support settings, these settings are also present in settings where the AI system may be trying to teach the human (as in ITS systems VanLehn [2006]) and even cases where a robot may be trying to assist the human achieve their goal Chakraborti et al. [2015].

Model-reconciliation explanation in this setting would consist of the robot communicating information present in $\mathcal{M}^*$, but may be absent in $\mathcal{M}^H$ (or the robot believes it to be absent based on its estimate $\mathcal{M}_r^H$). One could also formalize a notion of explicable plan in this setting, particularly for cases where the robot may be suggesting plans to the human. In this case, the explicable plan consists of solving for the following objective

$$\text{Find: } \pi$$

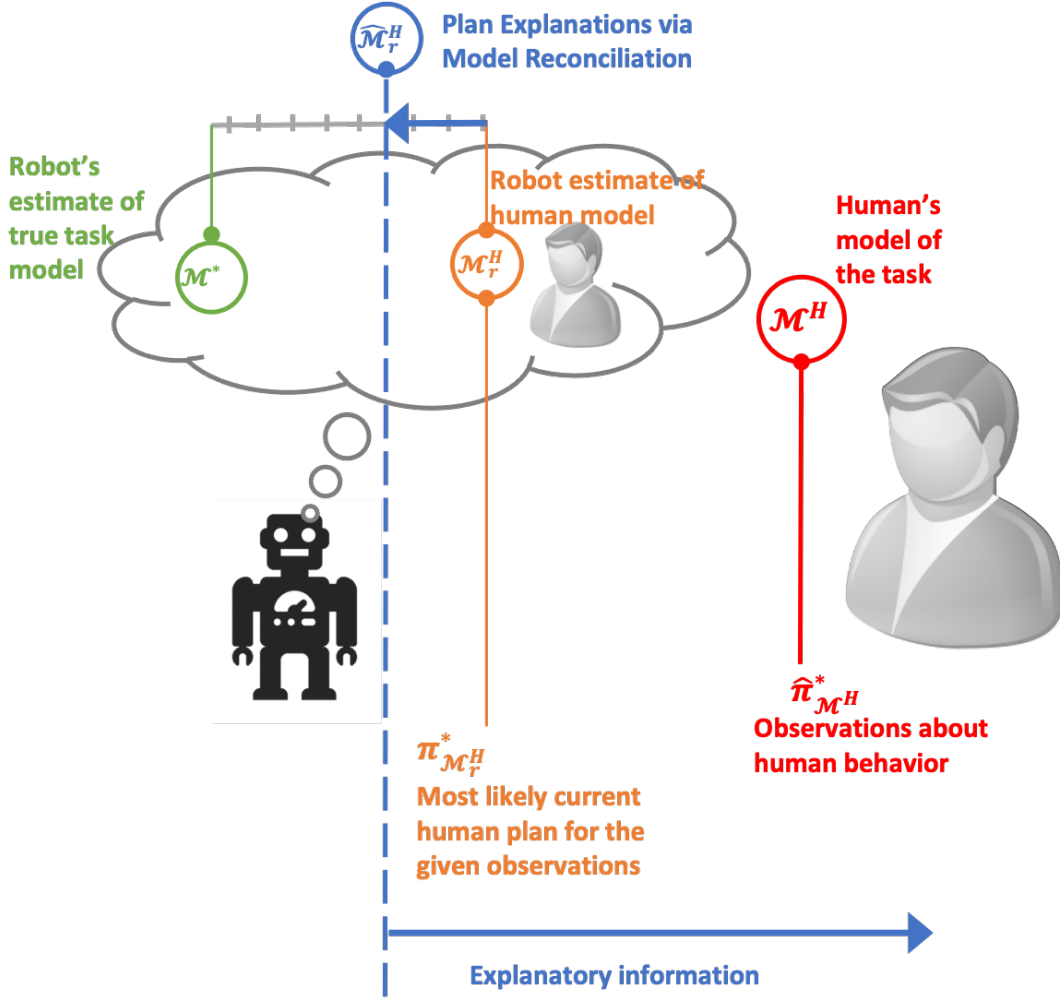$$\max_{\pi \in \Pi^{\mathcal{M}_r^H}} E(\pi, \mathcal{M}_r^H)$$

**Fig. 10.1:** Illustration of the scenarios where human is the primary actor and robot the observer.

Such that $\pi$ is executable in both $\mathcal{M}_r^H$ and $\mathcal{M}^*$

That is the goal here becomes to suggest a plan that is explicable with respect to the model $\mathcal{M}_r^H$ in that it is close to the plan expected under the model $\mathcal{M}_r^H$, but is executable in both $\mathcal{M}_r^H$ and $\mathcal{M}^*$. Additionally, we may also want to choose plans whose cost in the model $\mathcal{M}^*$ is low (could be captured by adding a term $-1 \times C^*(\pi)$ to the objective).

### 10.2.1 RADAR

The first example, we will consider is a Proactive Decision Support system called RADAR [Sengupta et al., 2017a]. Proactive Decision Support (PDS) aims at improving the decision making experience of human decision makers by enhancing both the quality of the decisions and the ease of making them. RADAR leverages techniques from automated planning community that aid the human decision maker in constructing plans. Specifically, the system focuses on expert humans in the loop who share a detailed, if not complete, model of the domain with the assistant, but may still be unable to compute plans due to cognitive
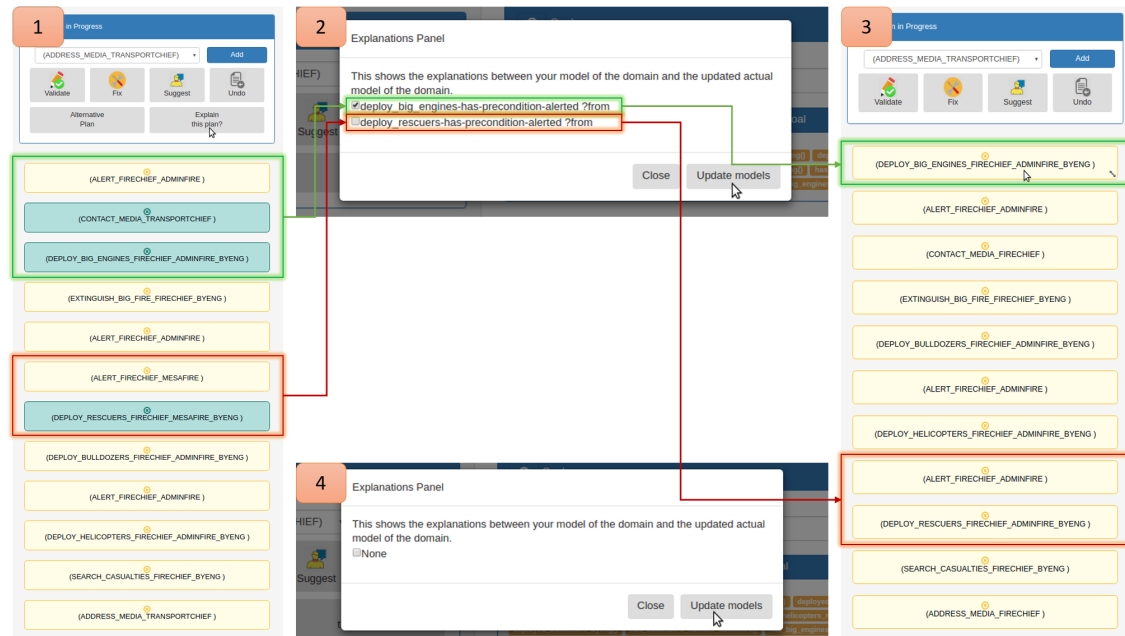
**Fig. 10.2:** RADAR system being applied to a firefighting domain, where the goal is to extinguish a number of fires. (1) RADAR knows that in the environment, the commander needs to inform the fire station's fire chief before deploying big engines and rescuers. In green, Adminfire's fire chief is alerted to deploy big engines from Admin Fire Station. In red, Mesa fire stations' fire chief is alerted to deploy rescuers from Mesa fire station. (2) The human's model believes that there is no need to inform fire chiefs and questions RADAR to explain his plan. RADAR finds these differences in the domain model and reports it to the human. The human acknowledges that before deploying rescuers one might need to alert the fire chief and rejects the update the fire chief needs to be alerted before deploying big engines. (3) In the alternative plan suggested by RADAR, it takes into account the humans knowledge and plans with the updated model. (4) Clicking on 'Explain This Plan' generates no explanations as there are none (with respect to the current plan) after the models were updated.

overload.

The system provides a number of useful features like plan suggestion, completion, validation and summarization to the users of the system. The entire system is built to allow users to perform naturalistic decision-making, whereby the system ensures the human is in control of the decision-making process. This means the proactive decision support system focuses on aiding and alerting the human in the loop with his/her decisions rather than generating a static plan that may not work in the dynamic worlds that the plan has to execute in. Figure 10.2, presents a screenshot of the RADAR system.

The component that is of particular interest to discussions in this book is the ability to support explanations. The system adapts model reconciliation explanations to address possible differences in the planner's model of the domain and the human expectation of it. Such differences could occur if the system may be automatically collecting information from external sources and thus the current estimate of the model diverges from the original specification provided/approved by the user. Here the system performs a model-space search to come up with Minimally Complete Explanation (Chapter 5 Section 5.2.1) to explain the plan being suggested. An important distinction in the use of explanations from previous chapter here is the fact that the human has the power to veto the model update if she believes that the planner's model is the one which is faulty, by choosing to approve or not approve individual parts of the explanation.

For example, consider the scenario highlighted in Figure 10.2, which presents a firefighting scenario where a commander is trying to come up with a plan to extinguish a fire in the city of Tempe. In this scenario, the RADAR system presents a plan (which itself is a completion of some action suggestions made by the commander), which contains unexpected steps to notify the Fire chief at multiple points. When the commander asks for an explanation, the system responds by pointing out that actions for both deploying big fire engines and the deploying rescuers has a precondition that the fire chief should be alerted in its model (which according to the system's model of the commander is missing from the commander's model). The commander responds by agreeing with the system on the fact that the fire chief needs to be informed before deploying rescuers, but also informs the system that fire chief doesn't need to be informed before deploying big fire engines. The system uses this new information to update its own models about the task and generates a new plan.

### 10.2.2 MA-RADAR

Next we will consider yet another extension of RADAR, but one that is focused on supporting multiple user types. In particular, MA-RADAR [Sengupta et al., 2018] considers a scenario where there are multiple users of possibly different backgrounds working together to construct a single solution. Such decision-making scenarios are further complicated when the different users may have differing levels of access and may have privacy concerns limiting the sharing some of the information with other decision-makers in the loop. As such, one of the focuses of MA-RADAR was to use augmented reality (AR) to allow the different users to view differing views of the task while working together on the same interface. The fact that they are using AR techniques means that in addition to the public information (i.e. the information that all the users can access), each user can view their specific private information which may not be accessible to others. For example, revisiting the firefighting domain, let us consider the case where there are two commanders working together to come up with a plan for controlling the fires. Here in addition to each commanders personal understanding of the task which may be inconsistent with each other and even what the system knows (for example status of some resources etc.), but may have knowledge about certain fluents and actions that are private to each commander and may not necessarily want the other commander to know. In terms of the explanations, each user could have different background knowledge and the system should strive to establish a common ground whenever possible. As such, the system uses the multi-model explanation method discussed in Chapter 6 (Section 6.2). Specifically, MA-RADAR combines the individual models into a single annotated model and explanations are generated with respect to this annotated model. The generated explanation is then filtered with respect to each user, so they only view information relevant to them (i.e the information is not redundant as per their model and is not private to any of the other users). Each user can then use their AR interface to view their specific explanations. Figure 10.3, presents the augmented reality interface that is shown to the user of the system.
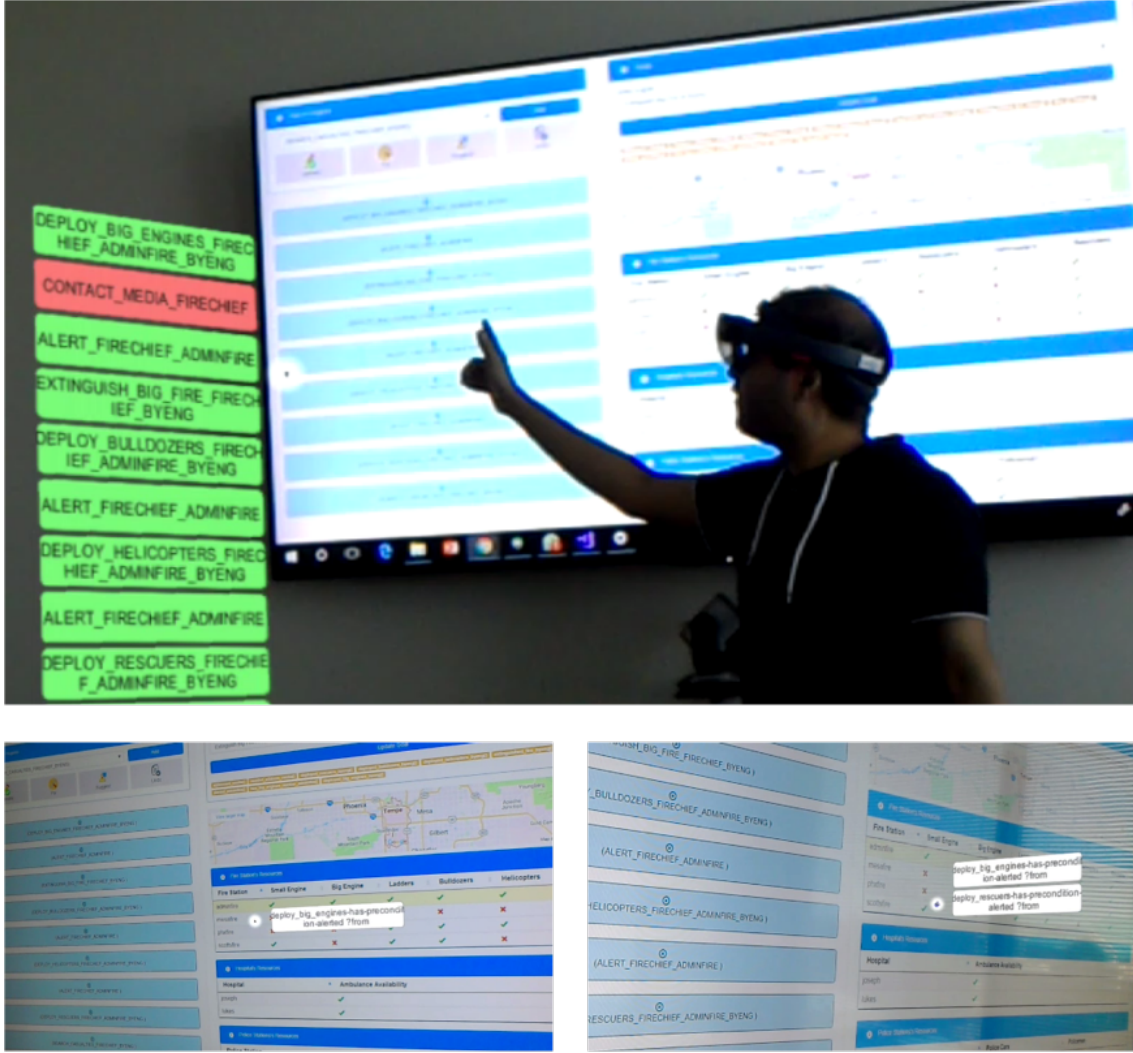
**Fig. 10.3:** A user wearing an augmented reality (AR) device to use the MA-RADAR interface. The figures show, how in addition to the common information shown on the screen the AR device allow users to view information private to each user.

### 10.2.3   RADAR-X

Next we will consider a variant of the basic RADAR system was extended to support explanatory dialogue and iterative planning. The system, named RADAR-X [Valmeekam et al., 2020], assumes that the user may have latent preferences that may not be specified to the system upfront or need not be completely realizable. The explanatory dialogue, particularly contrastive questions (i.e the user asks *Why this plan P as opposed to plan Q?*) thus becomes a way for the user to expose their underlying preferences. This system assumes that as a response to a specific plan, the users responds with a partial specification of the plans they prefer. In particular the system expects users to provide partial plans that can be defined by a tuple of the form $\hat{\pi} = \langle \hat{A}, \prec \rangle$, where $\hat{A}$ specifies a multi-set of actions the user expects to see in the plan and $\prec$ specifies the set of ordering constraints defined over $\hat{A}$ that the user expects to be satisfied. A sequential plan is said to satisfy a given partial plan $\hat{\pi}$ if the plan contains each action specified in $\hat{A}$ and they satisfy all the ordering constraints specified in $\prec$.

If the partially specified plan is feasible, the system suggests one of the possible instantiations of the partial plan (i.e. a plan that satisfies the given specification) to the decision-maker. The user could possibly further refine their choice by adding more information into the partial specification, until she is happy with the choice made by the system. If the partial specification is not feasible, then the system responds by first providing an explanation as to why the foil is not feasible. This is similar to the explanation generation method specified in Section 5.3.1, but with the end condition being the case of identifying the updated model where the foils are infeasible.

Once the explanation is provided, the system tries to generate plans that are closer to the specified foil. The system currently tries to find a subset of the original partial plan that the user specified that can be realized by the system, where for a given partial plan $\hat{\pi} = \langle \hat{A}, \prec \rangle$ a partial plan $\hat{\pi}' = \langle \hat{A}, \prec \rangle$ is considered a subset if $\hat{A}' \subseteq \hat{A}$, and for every $a_1, a_2 \in \hat{A}'$, $a_1 \prec' a_2$ if and only in $a_1 \prec a_2$. The system currently considers three strategies for coming up with such subsets.

1. Present the closest plan: Among the possible maximal subsets the system chooses one of them and presents a plan that corresponds to this subset.

2. Present plausible subsets: The user is presented with all possible maximal subsets and they can select the one that most closely represents their preferences.

3. Present conflict sets: The user with minimal conflict sets. That is minimal subset of actions and their corresponding ordering constraints that cannot be achieved together. So the user is asked to make a choice to remove one of the conflicting action from the set.

The system also looks at possible approximations that could be used to speed up the calculations.

## 10.3 Model Transcription Assistants

In this section, we will consider the problem of systems designed to allow domain experts to transcribe the models in some declarative form. Here the difference in the mental models actually comes from possible mistakes that the user may make while writing the declarative model. Clearly in this case, the system doesn't have access to the human's mental model and the direction of reconciliation is to get closer to the user's mental model. So here the process involves generating explanations for specific user queries by assuming that the human's model is an abstract version of the current model. Thus exposing relevant fragments of models that correspond to the relevant behavior and thus letting user directly fix any inconsistencies in the exposed fragment.

### 10.3.1 D3WA+

The tool D3WA+ Sreedharan et al. [2020b] implemented this idea in the context of transcribing declarative models for dialogue planning for an automated chat agent. It was an

extension of a previous system called D3WA Muise et al. [2019] developed by IBM. D3WA allowed dialogue editors to encode their knowledge about the dialogue tree in the form of a non-deterministic planning domain. D3WA+ took this base system and extended by providing debug tools to the domain designers that allowed them to query the system to better understand why the system was generating the current dialogue tree. Figure 10.4 presents a screenshot of the D3WA+ interface.

In particular, the system focused on providing the domain writer with the ability to raise two types of queries:

1. Why are there no solutions?

2. Why does the generated dialogue tree not conform to their expectation?

The first question is expected to be raised when there exists no possible trace from initial state to goal under the current model specification, and the latter when the domain writer was expecting the dialogue tree to take a particular form that is not satisfied by the one generated by the system. Thus the system would need to explain why the current problem is unsolvable. In the second case, the domain writer has to specify their expected dialogue flow on the storyboard panel. A dialogue flow in this case would consist of possible questions the chat agent could raise and possible outcomes. Note the specified flow doesn't need to contain a complete dialogue sequence (which starts at the beginning, ends with the end-user, i.e. the one who is expected to use the chat bot, getting the desired outcome and contains every possible intermediate steps) but could very well be a partial specification of the dialogue flow that highlights some part of the dialogue. Assuming that the expected flow can not be supported by the current model specification, the system would need to explain why this particular flow isn't possible. If the user had specified a complete dialogue sequence, then the system can merely test the sequence in the model and provide the failure point. Though this will no longer be possible if the user only specified a partial foil. In such cases, we would need to respond to why any possible sequence that satisfies the partial specification provided by the domain writer will be impossible. This can now be mapped into explaining the unsolvability of a modified planning problem, one that constrains the original problem to only allow solutions that align with the specified flow.

Thus the answer to both these questions maps into explanations of problem unsolvability. The system here leverages approaches discussed in Sreedharan et al. [2019b], to find the minimal subset of fluents for which the problem is unsolvable. A minimal abstraction over this set of fluents (with others projected out) is then presented to the domain writer as an unsolvable core of the problem that they can then try to fix. In addition to the model abstraction, two additional debugging information is provided to the user. An unreachable landmark and the failure information for an example trace. The unsolvable landmark is extracted by considering the delete relaxation of the unsolvable abstract model. The example trace is generated from the most concrete solvable model abstraction for the original model, so that more detailed plans are provided to the user. Such concrete models are generated by searching for the minimum number of fluents to be projected out to make the problem solvable.
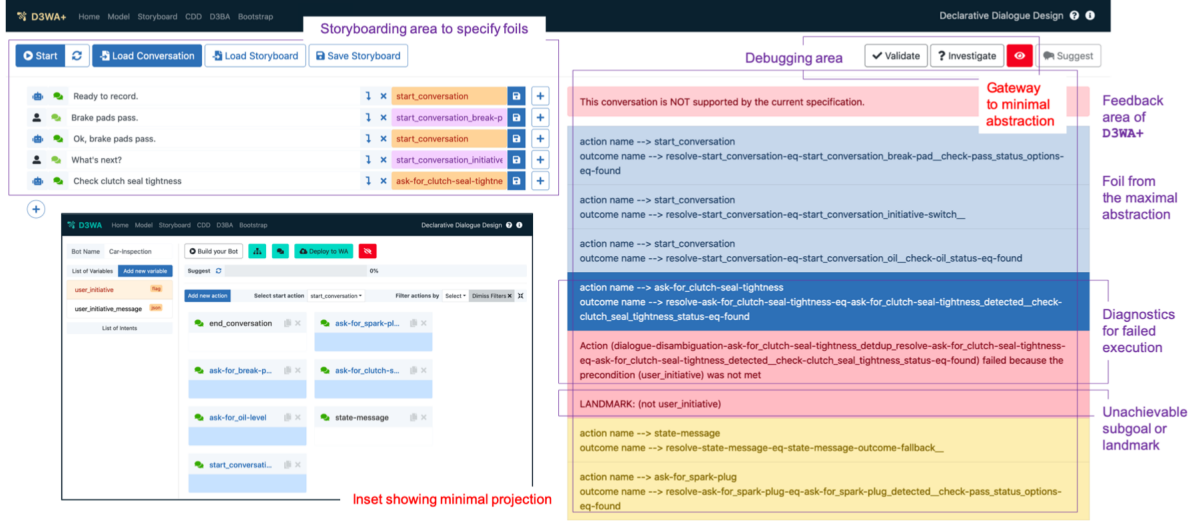
**Fig. 10.4:** An overview of the interface provided by D3WA+ to its end users. Which includes an option to specify foils or expected behavior, a panel that shows an abstract version of the domain and also additional debugging information provided to the user that includes information like: example plan failure, unachievable landmarks etc.

## 10.4  Bibliographic Remarks

Many of the applications we consider in this chapter have been demonstrated at various conferences including ICAPS and AAAI. RADAR system was first described in Sengupta et al. [2017a]. A version of the system that focused on plan of study generation was described in Grover et al. [2020]. The paper also presented a user study which validated various components that were part of the system. The multi agent version MA-RADAR was presented in Sengupta et al. [2018] and was demonstrated in ICAPS 2018 demo track. The contrastive version, i.e., RADAR-X was presented in Valmeekam et al. [2020] and was also presented as a demo at AAAI-21. The D3WA+ system was introduced in Sreedharan et al. [2020b] , the system was presented as a demo in ICAPS-20 and was awarded the best demo award. Apart from the ones listed in this chapter the ideas presented in the book have also been used in other applications. One prominent one is the FRESCO system [Chakraborti et al., 2017b] that was part of an IBM project for a smart room. The system made use of a variant of the model reconciliation for selective plan visualization. They try to model plan visualization as a process of model reconciliation against an empty user model. The process involved identifying the minimum number of components of the actions in the current plan (i.e., the preconditions and effects) that needed to be communicated to the user for the plan to make sense (i.e the plan is valid and optimal). As mentioned earlier, many of the decision-support applications look at settings where the human is the actor and the system has access to a model $\mathcal{M}_r^H$ of the human. All the specific applications discussed assume that $\mathcal{M}_r^H$ is accurate in that it is close to $\mathcal{M}^H$, but this may not be true in general. While this is quite similar to the asymmetry between $\mathcal{M}^R$ and $\mathcal{M}_h^R$. There is a major difference in that unlike the earlier case, the agent with access to the true model, in this case, the human, may not be invested in performing the reconciliation even if she is aware of the asymmetry. In this scenario, the robot may have to initiate the process of reconciliation, one possible way may be to make use of questions

to identify information about the human model Grover et al. [2020].

CHAPTER 11

# Conclusion

This book presents a concise introduction to recent research on human-aware decision-making, particularly ones focused on the generation of behavior that a human would find explainable or deceptive. Human-aware AI or HAAI techniques are characterized by the acknowledgment that for automated agents to successfully interact with humans, they need to explicitly take into account the human's expectations about the agent. In particular, we look at how for the robot to successfully work with humans, it needs to not only take into account its model $\mathcal{M}^R$, which encodes the robot's beliefs about the task and their capabilities but also take into account the human's expectation of the robot model $\mathcal{M}^R_h$, which captures what the human believes the task to be and what the robot is capable of. It is the model $\mathcal{M}^R_h$ that determines what the human would expect the robot to do and as such if the robot expects to adhere to or influence the human expectations, it needs to take into account this model. Additionally, this book also introduces three classes of interpretability measures, which capture certain desirable properties of robot behavior. Specifically, we introduce the measures *Explicability*, *Legibility*, and *Predictability*.

In the book, we mostly focused on developing and discussing methods to address the first two of these interpretability measures. We looked at specific algorithms that allow us to generate behaviors that boost explicability and legibility. We also looked at the problem of generating explanations for a given plan, and how it could be viewed as the use of communication to boost the explicability score of a selected plan by updating human expectations. Additionally, we looked at variations of this basic explanation framework under differing settings, including cases where the human models of the robot may be unknown or where the robot's decision-making model may be expressed in terms that the human doesn't understand. We also saw a plan generation method that is able to incorporate reasoning about the overhead of explanation into the plan selection process, thereby allowing for a method that is able to combine the benefits of purely explicable plan generation methods and those that identify explanations after the plan has been selected.

Additionally, we also saw that modeling human expectations not only allows us to create interpretable behaviors but also provides us with the tools needed to generate deceptive and adversarial behaviors. In particular, we saw how one could leverage the modeling of the other agent to create behaviors that obfuscate certain agent model information from an adversarial observer or even deceive them about the model component. We also saw how model reconciliation methods could be molded to create lies, which even in non-adversarial scenarios could help the agent to achieve higher team utility at the cost of providing some white lies to the human.

While the problem of generating explanations for AI decisions or developing deceptive AI agents has been studied for a while, the framing of these problems within the context of human-aware AI settings is a relatively recent effort. As such, there exists a number of exciting future directions to be explored and technical challenges to overcome, before we can have truly human-aware systems. Some of these challenges are relatively straightforward (at least in their conception), as in the case of scaling up the methods and applying them within more complex scenarios that are more faithful to real-world scenarios. Then

there are challenges that may require us to rethink our current strategies and whose formalization itself presents a significant challenge. Here we would like to take a quick look at a few of these challenges.

**Creating a Symbolic Middle Layer**   The necessity of symbols in intelligent decision-making has been a topic that has been widely debated within the field of AI for a very long time. Regardless of whether symbols are necessary for intelligence, it remains a fact that people tend to communicate in terms of symbols and concepts. As such, if we want to create systems that can interact with people effectively they should be capable of communicating using symbols and concepts people understand. In chapter 8, we have already seen an example, where an agent translates information about its model into terms that are easier for a human to understand, but if we want to create successful systems that are truly able to collaborate with humans then we need to go beyond just creating explanation generation systems. They need to be able to take input from the human in symbolic terms even when the model may not be represented in those terms. Such advice could include information like instructions the agent should follow, possible domain, and preference information. Interestingly the advice provided by the human would be colored by what they believe the agent model to be, and as such correct interpretation of the specified information may require analyzing the human input in the light of their expectation about the agent model. Some preliminary works in this direction include Guan et al. [2020] and for a discussion on the overall direction, the readers can refer to Kambhampati et al. [2021].

**Trust and Longitudinal Interaction**   Most of the interactions discussed in this book are single-step interactions, in so far that they focus on the agent proposing a plan for a single task and potentially handling any interaction requirements related to that plan. But we are not focused on creating single-use robots. Rather we want to create automated agents that we can cohabit with and work with on a day-to-day basis. This means the robot's choice of actions can no longer be made in isolation, rather it should consider the implications of choosing a certain course of action on future interaction with humans. We saw some flavors of such consideration in methods like Minimally Monotonic Explanations (MME) (Chapter 5, Section 5.2.1) and the discounting of explicability over a time horizon (Chapter 3, Section 3.4.2), though these are still limited cases. As we move forward, a pressing requirement is for the robots to be capable of modeling the level of trust the human holds for the robot and how the robot's actions may influence the human trust. Such trust-level modeling is of particular importance in longitudinal settings, as it would be the human trust on the robot that would determine whether or not the human would choose to work with it on future tasks. Thus the impact of the robot action on human trust should be a metric it should consider while coming up with its actions. At the same time, reasoning about trust levels also brings up the question of trust manipulation and how to design agents that are guaranteed to not induce undeserved trust in its capabilities. Some preliminary works in this direction are presented in Zahedi et al. [2021].

**Learning Human Models**   The defining feature of many of the techniques discussed in this book is the inclusion of the human's expectations, sensory capabilities, and in general their models into the reasoning process. In many scenarios, such models may

not be directly available to the robot and it may be required to learn such models either by observing the human or by receiving feedback from the human. We have already seen examples of learning specific model proxies that are sufficient to generate specific classes of behaviors. But these are specialized models meant for specific applications. More general problems may require the use of additional information and even the complete model. While the model $\mathcal{M}^H$ could be learned by observing the human behavior, the model $\mathcal{M}_h^R$ is usually more expensive to learn as it requires the human to either provide feedback on the robot behavior or provide the full plan they are expecting from the robot. As such learning a fully personalized model for a person may require too much information to be provided by a single person. A more scalable strategy may be to learn approximate models for different user types from previously collected data. As and when the robot comes into contact with a new human, they can use the set of learned models to identify the closest possible model. This model can act as a starting point for the interaction between the system and the user and can be refined over time as the robot interacts with the human. Some initial work in this direction can be found in Soni et al. [2021].

**Super Human AI and Safety**   Currently, most successful AI systems are generally less competent than humans overall but may have an edge over people on specific narrow tasks. There is no particular reason to believe that this condition should always persist. In fact, one could easily imagine a future where AI systems outpace humans in almost all tasks. It may be worth considering how the nature of human-robot interaction may change in such a world. For one thing, the nature and goal of explanation may no longer be about helping humans understand the exact reason for selecting a decision but rather about giving a general sense of why the decisions make sense. For example, it may be enough to establish why the decision is better than any alternative the human could come up with. Going back to the problem of vocabulary mismatch introduced in Chapter 8, it may very well be the case that there may be concepts that the system makes use of that have no equivalent term in human vocabulary and as such explanation might require the system teaching new concepts to the human. The introduction of the ability to reason and model the human mental model also raises additional safety and ethical questions in the context of such superhuman human-aware AI. For one, the questions of white lies and deception for improving team utility (Chapter 9, Section 9.3) takes on a whole new dimension and could potentially head into the realm of wireheading. It is very much an open question as to how one can build robust methods and safeguards to control for and avoid such potential safety concerns that may arise from the deployment of such systems. As we see more AI systems embrace human-aware AI principles, it becomes even more important to study and try ameliorate unique safety concerns that arise in systems capable of modeling and influencing human's mental models and beliefs.

# Bibliography

James F Allen. Mixed initiative planning: Position paper. In *ARPA/Rome Labs Planning Initiative Workshop*, 1994.

Saleema Amershi, Daniel S. Weld, Mihaela Vorvoreanu, Adam Fourney, Besmira Nushi, Penny Collisson, Jina Suh, Shamsi T. Iqbal, Paul N. Bennett, Kori Inkpen, Jaime Teevan, Ruth Kikin-Gil, and Eric Horvitz. Guidelines for human-ai interaction. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems, CHI 2019, Glasgow, Scotland, UK, May 04-09, 2019*, page 3, 2019.

Dan Amir and Ofra Amir. Highlights: Summarizing agent behavior to people. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, pages 1168–1176, 2018.

Chris L Baker and Joshua B Tenenbaum. Modeling human plan recognition using bayesian theory of mind. *Plan, activity, and intent recognition: Theory and practice*, 7:177–204, 2014.

Pascal Bercher, Susanne Biundo, Thomas Geier, Thilo Hoernle, Florian Nothdurft, Felix Richter, and Bernd Schattenberg. Plan, Repair, Execute, Explain – How Planning Helps to Assemble Your Home Theater. In *ICAPS*, 2014.

Blai Bonet and Hector Geffner. Belief tracking for planning with sensing: Width, complexity and approximations. *Journal of Artificial Intelligence Research*, 50:923–970, 2014.

Cynthia Breazeal. Toward sociable robots. *Robotics and autonomous systems*, 42(3-4): 167–175, 2003.

Cynthia L Breazeal. *Designing sociable robots*. MIT press, 2004.

Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.

Dan Bryce, J Benton, and Michael W Boldt. Maintaining evolving domain models. In *Proceedings of the twenty-fifth international joint conference on artificial intelligence*, pages 3053–3059, 2016.

Tathagata Chakraborti and Subbarao Kambhampati. (How) Can AI Bots Lie? In *XAIP Workshop*, 2019a.

Tathagata Chakraborti and Subbarao Kambhampati. (when) can ai bots lie? In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, pages 53–59, 2019b.

Tathagata Chakraborti, Gordon Briggs, Kartik Talamadupula, Yu Zhang, Matthias Scheutz, David Smith, and Subbarao Kambhampati. Planning for serendipity. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5300–5306. IEEE, 2015.

Tathagata Chakraborti, Sarath Sreedharan, Yu Zhang, and Subbarao Kambhampati. Plan Explanations as Model Reconciliation: Moving Beyond Explanation as Soliloquy. In *IJCAI*, 2017a.

Tathagata Chakraborti, Kartik Talamadupula, Mishal Dholakia, Biplav Srivastava, Jeffrey O Kephart, and Rachel KE Bellamy. Mr. jones–towards a proactive smart room orchestrator. In *AAAI Fall Symposia*, 2017b.

Tathagata Chakraborti, Sarath Sreedharan, Anagha Kulkarni, and Subbarao Kambhampati. Projection-Aware Task Planning and Execution for Human-in-the-Loop Operation of Robots. In *IROS*, 2018.

Tathagata Chakraborti, Anagha Kulkarni, Sarath Sreedharan, David Smith, and Subbarao Kambhampati. Explicability? Legibility? Predictability? Transparency? Privacy? Security?: The Emerging Landscape of Interpretable Agent Behavior. In *ICAPS*, 2019a.

Tathagata Chakraborti, Sarath Sreedharan, Sachin Grover, and Subbarao Kambhampati. Plan explanations as model reconciliation–an empirical study. In *2019 14th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pages 258–266. IEEE, 2019b.

Tathagata Chakraborti, Sarath Sreedharan, and Subbarao Kambhampati. Balancing Explicability and Explanation in Human-Aware Planning. In *IJCAI*, 2019c.

Tathagata Chakraborti, Sarath Sreedharan, and Subbarao Kambhampati. The emerging landscape of explainable ai planning and decision making. In *IJCAI*, 2020.

Anca Dragan and Siddhartha Srinivasa. Generating Legible Motion. In *RSS*, 2013.

Anca D Dragan. Robot Planning with Mathematical Models of Human State and Action. *arXiv:1705.04226*, 2017.

Anca D Dragan, Kenton CT Lee, and Siddhartha S Srinivasa. Legibility and Predictability of Robot Motion. In *HRI*, 2013.

Anca D Dragan, Shira Bauman, Jodi Forlizzi, and Siddhartha S Srinivasa. Effects of Robot Motion on Human-Robot Collaboration. In *HRI*, 2015.

Rebecca Eifler, Michael Cashmore, Jörg Hoffmann, Daniele Magazzeni, and Marcel Steinmetz. A New Approach to Plan-Space Explanation: Analyzing Plan-Property Dependencies in Oversubscription Planning. In *AAAI*, 2020.

Christiane Fellbaum. Wordnet. In *Theory and applications of ontology: computer applications*, pages 231–243. Springer, 2010.

Jaime F Fisac, Chang Liu, Jessica B Hamrick, S Shankar Sastry, J Karl Hedrick, Thomas L Griffiths, and Anca D Dragan. Generating Plans that Predict Themselves. In *WAFR*, 2018.

Maria Fox, Derek Long, and Daniele Magazzeni. Explainable Planning. In *IJCAI XAI Workshop*, 2017.

Hector Geffner and Blai Bonet. A Concise Introduction to Models and Methods for Automated Planning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 2013.

Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.

Sachin Grover, David Smith, and Subbarao Kambhampati. Model elicitation through direct questioning. *arXiv preprint arXiv:2011.12262*, 2020.

Lin Guan, Mudit Verma, Sihang Guo, Ruohan Zhang, and Subbarao Kambhampati. Explanation augmented feedback in human-in-the-loop reinforcement learning. *arXiv preprint arXiv:2006.14804*, 2020.

Bradley Hayes and Julie A Shah. Improving robot controller transparency through autonomous policy explanation. In *2017 12th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pages 303–312. IEEE, 2017.

Malte Helmert. The Fast Downward Planning System. *JAIR*, 2006.

Jörg Hoffmann. Ff: The fast-forward planning system. *AI magazine*, 22(3):57–57, 2001.

Jörg Hoffmann. Where'ignoring delete lists' works: Local search topology in planning benchmarks. *Journal of Artificial Intelligence Research*, 24:685–758, 2005.

Oliver Wendell Holmes. *Medical Essays, 1842-1882*, volume 9. Houghton, Mifflin, 1895.

Subbarao Kambhampati, Sarath Sreedharan, Mudit Verma, Yantian Zha, and Lin Guan. Symbols as a lingua franca for bridging human-ai chasm for explainable and advisable ai systems. In *AAAI Senior Member Track*, 2021.

Sarah Keren, Avigdor Gal, and Erez Karpas. Goal Recognition Design. In *ICAPS*, 2014.

Sarah Keren, Avigdor Gal, and Erez Karpas. Privacy Preserving Plans in Partially Observable Environments. In *IJCAI*, 2016.

Sarah Keren, Avigdor Gal, and Erez Karpas. Strong Stubborn Sets for Efficient Goal Recognition Design. In *ICAPS*, 2018.

Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *International conference on machine learning*, pages 2668–2677. PMLR, 2018.

Ross A Knepper, Christoforos I Mavrogiannis, Julia Proft, and Claire Liang. Implicit communication in a joint action. In *Proceedings of the 2017 ACM/IEEE International Conference on Human-Robot Interaction*, pages 283–292. ACM, 2017.

George Konidaris, Leslie Pack Kaelbling, and Tomas Lozano-Perez. From skills to symbols: Learning symbolic representations for abstract high-level planning. *Journal of Artificial Intelligence Research*, 61:215–289, 2018.

Barbara M Korsch and Caroline Harding. *The intelligent patient's guide to the doctor-patient relationship: learning how to talk so your doctor will listen*. Oxford University Press, 1998.

Benjamin Krarup, Michael Cashmore, Daniele Magazzeni, and Tim Miller. Model-Based Contrastive Explanations for Explainable Planning. In *XAIP Workshop*, 2019.

Anagha Kulkarni, Tathagata Chakraborti, Yantian Zha, Satya Gautam Vadlamudi, Yu Zhang, and Subbarao Kambhampati. Explicable Robot Planning as Minimizing Distance from Expected Behavior. In *AAMAS Extended Abstract*, 2019a.

Anagha Kulkarni, Siddharth Srivastava, and Subbarao Kambhampati. A Unified Framework for Planning in Adversarial and Cooperative Environments. In *AAAI*, 2019b.

Anagha Kulkarni, Sarath Sreedharan, Sarah Keren, Tathagata Chakraborti, David Smith, and Subbarao Kambhampati. Designing environments conducive to interpretable robot behavior. In *IROS*, 2020a.

Anagha Kulkarni, Siddharth Srivastava, and Subbarao Kambhampati. Signaling friends and head-faking enemies simultaneously: Balancing goal obfuscation and goal legibility. In *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems*, AAMAS '20, page 1889–1891. International Foundation for Autonomous Agents and Multiagent Systems, 2020b.

John Lafferty, Andrew McCallum, and Fernando CN Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *ICML*, 2001.

Isaac Lage, Daphna Lifschitz, Finale Doshi-Velez, and Ofra Amir. Exploring Computational User Models for Agent Policy Summarization. In *IJCAI*, 2019.

Aleck M MacNally, Nir Lipovetzky, Miquel Ramirez, and Adrian R Pearce. Action Selection for Transparent Planning. In *AAMAS*, 2018.

Peta Masters and Sebastian Sardina. Deceptive Path Planning. In *IJCAI*, 2017.

Tim Miller. Explanation in Artificial Intelligence: Insights from the Social Sciences. *Artificial Intelligence*, 2019.

Christian Muise, Vaishak Belle, Paolo Felli, Sheila McIlraith, Tim Miller, Adrian Pearce, and Liz Sonenberg. Planning over multi-agent epistemic states: A classical planning approach. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2015.

Christian Muise, Tathagata Chakraborti, Shubham Agarwal, Ondrej Bajgar, Arunima Chaudhary, Luis A Lastras-Montano, Josef Ondrej, Miroslav Vodolan, and Charlie Wiecha. Planning for goal-oriented dialogue systems. *arXiv preprint arXiv:1910.08137*, 2019.

Tuan Nguyen, Sarath Sreedharan, and Subbarao Kambhampati. Robust planning with incomplete domain models. *Artificial Intelligence*, 245:134–161, 2017.

Tuan Anh Nguyen, Minh Do, Alfonso Emilio Gerevini, Ivan Serina, Biplav Srivastava, and Subbarao Kambhampati. Generating diverse plans to handle unknown and partially known user preferences. *Artificial Intelligence*, 190(0):1 – 31, 2012.

Nils J Nilsson. *Principles of artificial intelligence*. Morgan Kaufmann, 2014.

John J Palmieri and Theodore A Stern. Lies in the doctor-patient relationship. *Primary care companion to the Journal of clinical psychiatry*, 11(4):163, 2009.

Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. " why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016.

Stuart Russell and Peter Norvig. *Artificial intelligence: a modern approach*. Prentice Hall, 2002.

Brian Scassellati. Theory of mind for a humanoid robot. *Auton. Robots*, 12(1):13–24, 2002.

Bastian Seegebarth, Felix Müller, Bernd Schattenberg, and Susanne Biundo. Making Hybrid Plans More Clear to Human Users – A Formal Approach for Generating Sound Explanations. In *ICAPS*, 2012.

Sailik Sengupta, Tathagata Chakraborti, Sarath Sreedharan, and Subbarao Kambhampati. RADAR – A Proactive Decision Support System for Human-in-the-Loop Planning. In *AAAI Fall Symposium*, 2017a.

Sailik Sengupta, Tathagata Chakraborti, Sarath Sreedharan, Satya Gautam Vadlamudi, and Subbarao Kambhampati. Radar-a proactive decision support system for human-in-the-loop planning. In *AAAI Fall Symposia*, pages 269–276, 2017b.

Sailik Sengupta, Tathagata Chakraborti, and Subbarao Kambhampati. Ma-radar–a mixed-reality interface for collaborative decision making. *ICAPS UISP*, 2018.

Utkarsh Soni, Sarath Sreedharan, and Subbarao Kambhampati. Not all users are the same: Providing personalized explanations for sequential decision making problems. *arXiv preprint arXiv:2106.12207*, 2021.

Sarath Sreedharan, Subbarao Kambhampati, et al. Handling model uncertainty and multiplicity in explanations via model reconciliation. In *Proceedings of the International Conference on Automated Planning and Scheduling*, 2018a.

Sarath Sreedharan, Siddharth Srivastava, and Subbarao Kambhampati. Hierarchical Expertise Level Modeling for User Specific Contrastive Explanations. In *IJCAI*, 2018b.

Sarath Sreedharan, Alberto Olmo, Aditya Prasad Mishra, and Subbarao Kambhampati. Model-free model reconciliation. In *AAAI*, 2019a.

Sarath Sreedharan, Siddharth Srivastava, David Smith, and Subbarao Kambhampati. Why can't you do that hal? explaining unsolvability of planning tasks. In *International Joint Conference on Artificial Intelligence*, 2019b.

Sarath Sreedharan, Tathagata Chakraborti, Christian Muise, and Subbarao Kambhampati. Planning with Explanatory Actions: A Joint Approach to Plan Explicability and Explanations in Human-Aware Planning. In *AAAI*, 2020a.

Sarath Sreedharan, Tathagata Chakraborti, Christian Muise, Yasaman Khazaeni, and Subbarao Kambhampati. –d3wa+–a case study of xaip in a model acquisition task for dialogue planning. In *Proceedings of the International Conference on Automated Planning and Scheduling*, volume 30, pages 488–497, 2020b.

Sarath Sreedharan, Utkash Soni, Mudit Verma, Siddharth Srivastava, and Subbarao Kambhampati. Bridging the gap: Providing post-hoc symbolic explanations for sequential decision-making problems with black box simulators. *arXiv preprint arXiv:2002.01080*, 2020c.

Sarath Sreedharan, Siddharth Srivastava, and Subbarao Kambhampati. Tldr: Policy summarization for factored ssp problems using temporal abstractions. In *Proceedings of the International Conference on Automated Planning and Scheduling*, volume 30, pages 272–280, 2020d.

Sarath Sreedharan, Tathagata Chakraborti, and Subbarao Kambhampati. Foundations of explanations as model reconciliation. *Artificial Intelligence*, 301:103558, 2021a.

Sarath Sreedharan, Anagha Kulkarni, David E Smith, and Subbarao Kambhampati. A unifying bayesian formulation of measures of interpretability in human-ai interaction. In *International Joint Conference on Artificial Intelligence*, pages 4602–4610, 2021b.

Sarath Sreedharan, Siddharth Srivastava, and Subbarao Kambhampati. Using state abstractions to compute personalized contrastive explanations for ai agent behavior. *Artificial Intelligence*, 301:103570, 2021c.

Biplav Srivastava, Tuan Anh Nguyen, Alfonso Gerevini, Subbarao Kambhampati, Minh Binh Do, and Ivan Serina. Domain independent approaches for finding diverse plans. In *IJCAI*, pages 2016–2022, 2007.

Nicholay Topin and Manuela Veloso. Generation of Policy-Level Explanations for Reinforcement Learning. In *AAAI*, 2019.

Karthik Valmeekam, Sarath Sreedharan, Sailik Sengupta, and Subbarao Kambhampati. Radar-x: An interactive interface pairing contrastive explanations with revised plan suggestions. In *XAIP ICAPS*, 2020.

Kurt VanLehn. The behavior of tutoring systems. *I. J. Artificial Intelligence in Education*, 16(3):227–265, 2006.

Stylianos Loukas Vasileiou, Alessandro Previti, and William Yeoh. On exploiting hitting sets for model reconciliation. In *AAAI*, 2021.

Manuela M Veloso. Learning by Analogical Reasoning in General Problem Solving. *Doctoral Thesis*, 1992.

J Waa, J van Diggelen, K Bosch, and M Neerincx. Contrastive Explanations for Reinforcement Learning in Terms of Expected Consequences. In *IJCAI Workshop on explainable AI (XAI)*, 2018.

H Wimmer and J Perner. Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children's understanding of deception. *Cognition*, 1983.

Tom Zahavy, Nir Ben-Zrihem, and Shie Mannor. Graying the Black Box: Understanding DQNs. In *ICML*, 2016.

Zahra Zahedi, Mudit Verma, Sarath Sreedharan, and Subbarao Kambhampati. Trust-aware planning: Modeling trust evolution in longitudinal human-robot interaction. *arXiv preprint arXiv:2105.01220*, 2021.

Yantian Zha, Lin Guan, and Subbarao Kambhampati. Learning from ambiguous demonstrations with self-explanation guided reinforcement learning, 2021.

Haoqi Zhang, Yiling Chen, and David C Parkes. A General Approach to Environment Design with One Agent. In *IJCAI*, 2009.

Yu Zhang, Sarath Sreedharan, Anagha Kulkarni, Tathagata Chakraborti, Hankz Hankui Zhuo, and Subbarao Kambhampati. Plan Explicability and Predictability for Robot Task Planning. In *ICRA*, 2017.
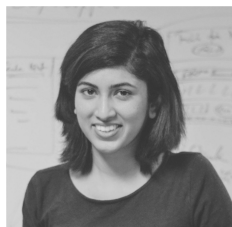
# Authors' Biographies

## Sarath Sreedharan

**Sarath Sreedharan** is a Ph.D. student at Arizona State University working with Prof. Subbarao Kambhampati. His primary research interests lie in the area of human-aware and explainable AI, with a focus on sequential-decision making problems. Sarath's research has been featured in various premier research conferences, including IJCAI, AAAI, AAMAS, ICAPS, ICRA, IROS, etc, and journals like AIJ. He was also the recipient of Outstanding Program Committee Member Award at AAAI-2020.

## Anagha Kulkarni

**Anagha Kulkarni** is an AI Research Scientist at Invitae. Before that, she received her Ph.D. in Computer Science from Arizona State University. Her Ph.D. thesis was in the area of human-aware AI and automated planning. Anagha's research has featured in various premier conferences like AAAI, IJCAI, ICAPS, AAMAS, ICRA and IROS.

## Subbarao Kambhampati

**Subbarao Kambhampati** is a professor in the School of Computing & AI at Arizona State University. Kambhampati studies fundamental problems in planning and decision making, motivated in particular by the challenges of human-aware AI systems. He is a fellow of Association for the Advancement of Artificial Intelligence, American Association for the Advancement of Science, and Association for Computing machinery, and was an NSF Young Investigator. He was the president of the Association for the Advancement of Artificial Intelligence, trustee of International Joint Conference on Artificial Intelligence, and a founding board member of Partnership on AI. Kambhampati's research as well as his views on the progress and societal impacts of AI have been featured in multiple national and international media outlets.