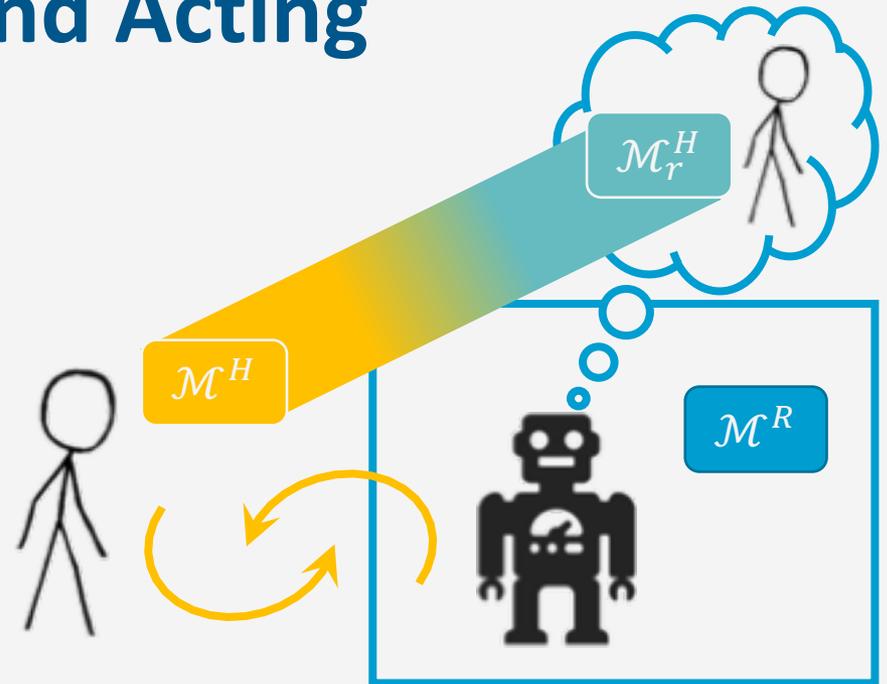


Automated Planning and Acting

Human-aware Planning

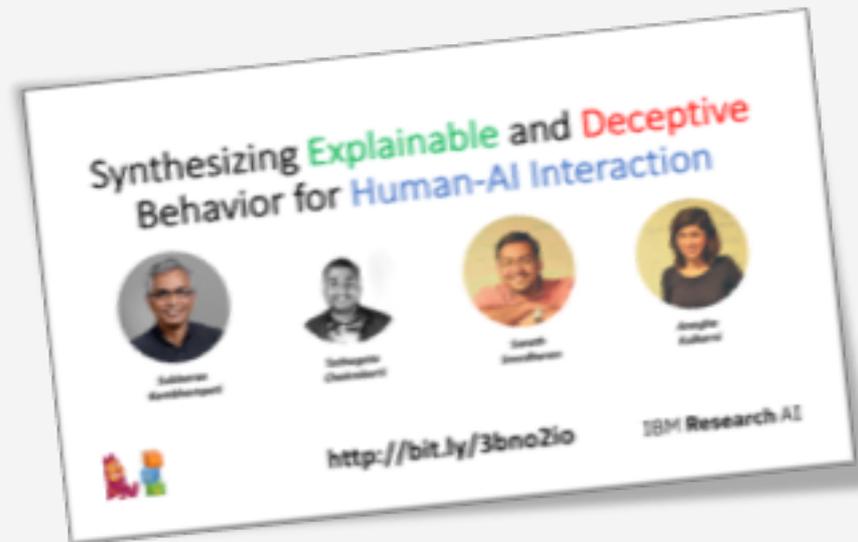


Content

1. Planning and Acting with **Deterministic** Models
2. Planning and Acting with **Refinement** Methods
3. Planning and Acting with **Temporal** Models
4. Planning and Acting with **Nondeterministic** Models
5. **Standard** Decision Making
6. Planning and Acting with **Probabilistic** Models
7. **Advanced** Decision Making
8. **Human-aware** Planning
 - a. Mental Models
 - b. Interpretable Behaviour
 - c. Explanations

Acknowledgements

- Slides based on material provided by Subbarao (Rao) Kambhampati and his colleagues (for more material on human-aware planning by Rao: <http://rakaposhi.eas.asu.edu>)



Outline

Mental Models

- Human-aware agent

Interpretable Behaviour

- Explicability
- Legibility
- Predictability

Explanations

- Model reconciliation

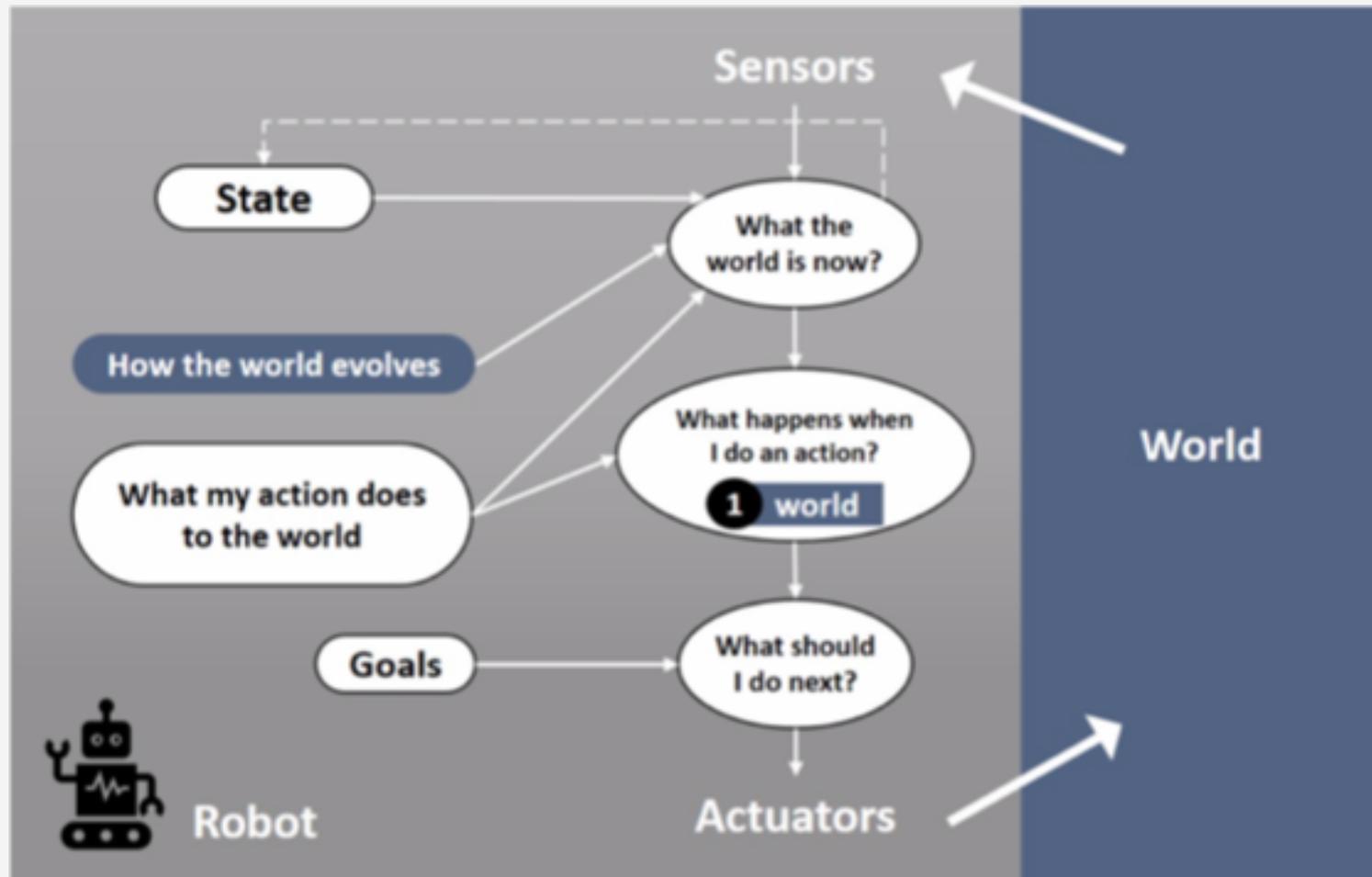
Motivation

- **Collaborations** between people and AI systems
 - I.e., systems with **humans in the loop**
 - Augment perception, cognition, problem-solving abilities of people
 - Examples
 - Help physicians make more timely and accurate diagnoses
 - Assistance provided to drivers of cars to help them avoid dangerous situations and crashes
- **Objective:** Systems that can interact intuitively with users and enable seamless machine-human collaborations
 - **Explainable** behaviour
 - Explainable AI = XAI

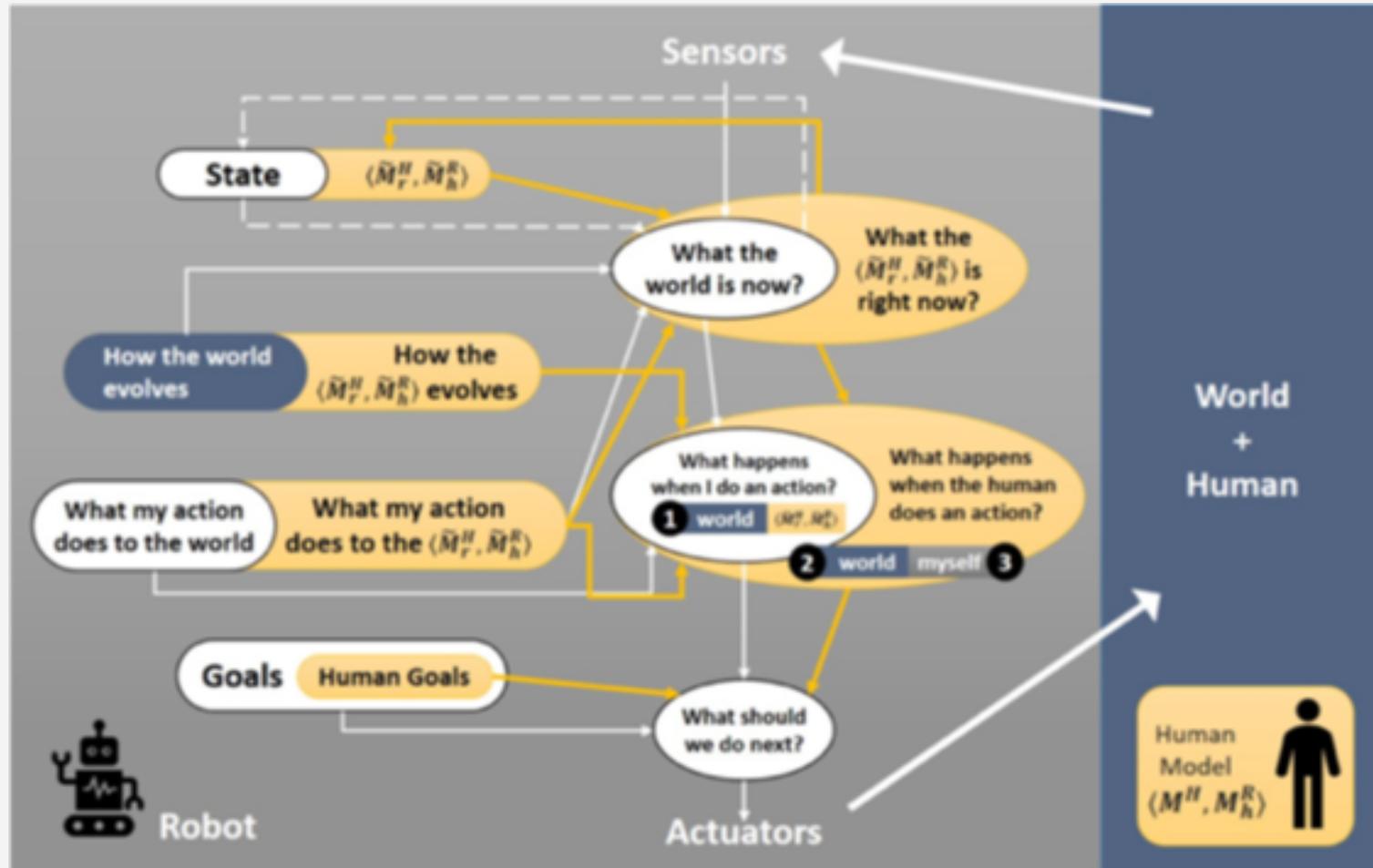
Proposed Solution

- Goal: Synthesise explainable behaviour
- Take into account the **mental model** of the human in the loop
 - Mental model:
 - Goals + capabilities of the humans in the loop
 - Human's model of AI agent's goals + capabilities

Classical Intelligent Agent

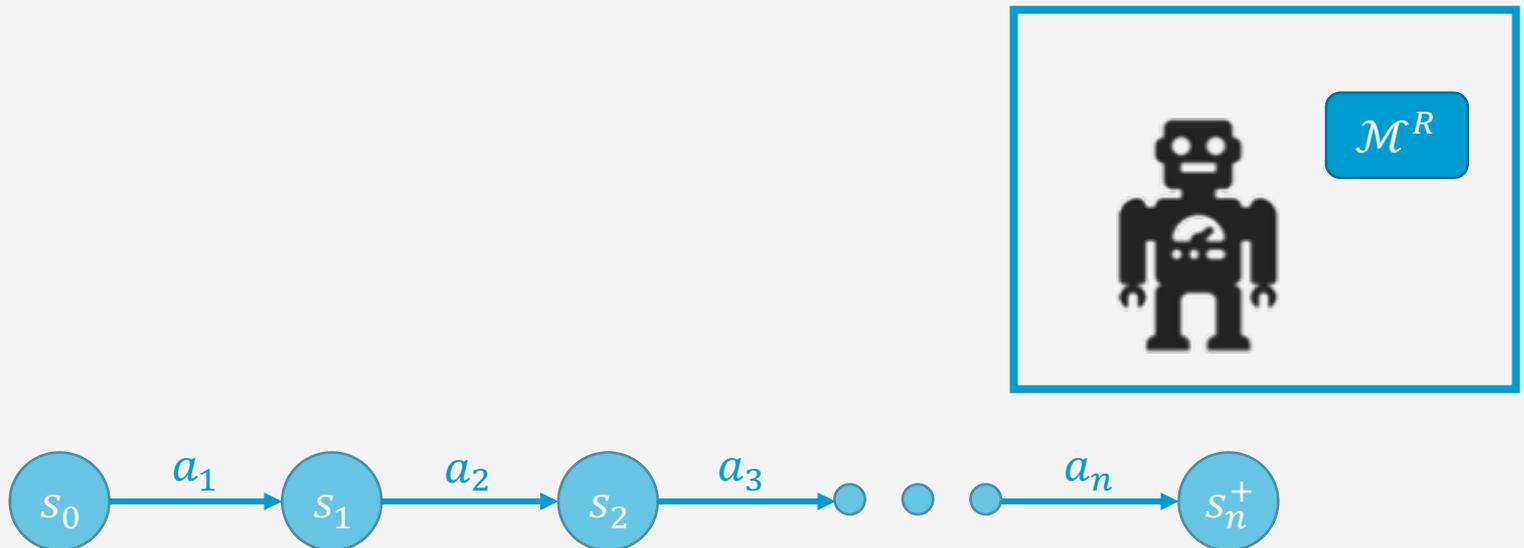


Human-aware Intelligent Agent



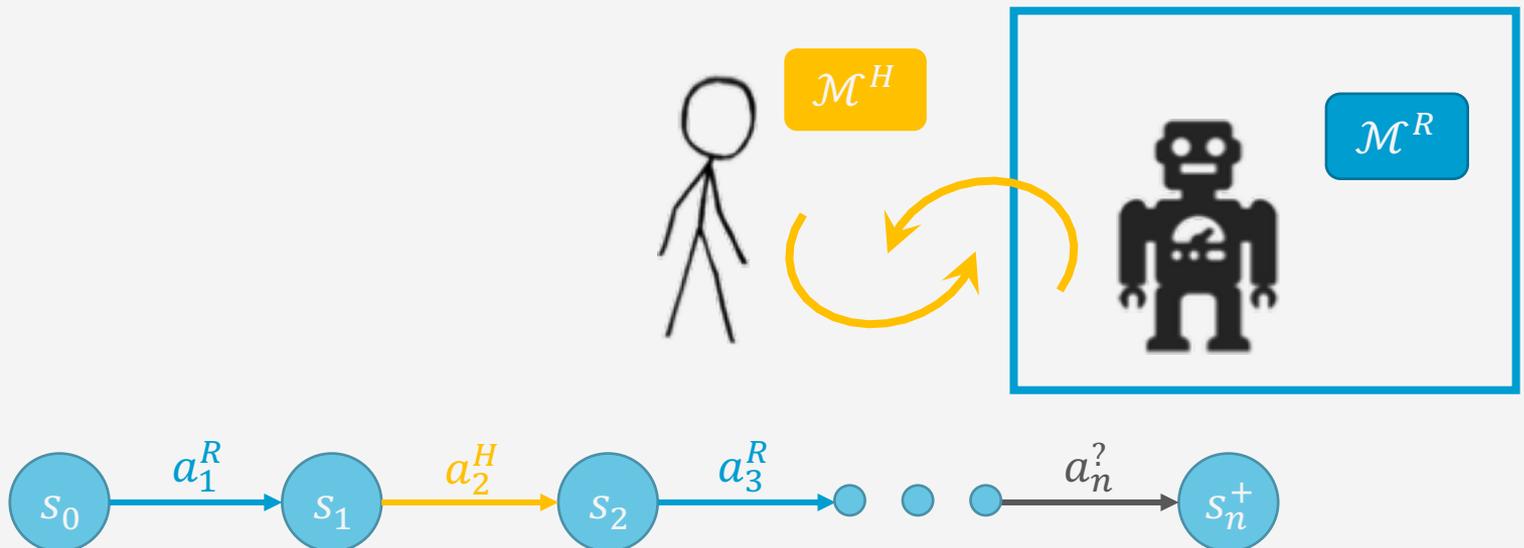
Classical Planning

- Given (Σ, s_0, S_g) , i.e., the agent's model \mathcal{M}^R
- Find a plan $\pi = \langle a_1, a_2, \dots, a_n \rangle$ that transforms s_0 to a state $s_n \in S_g$



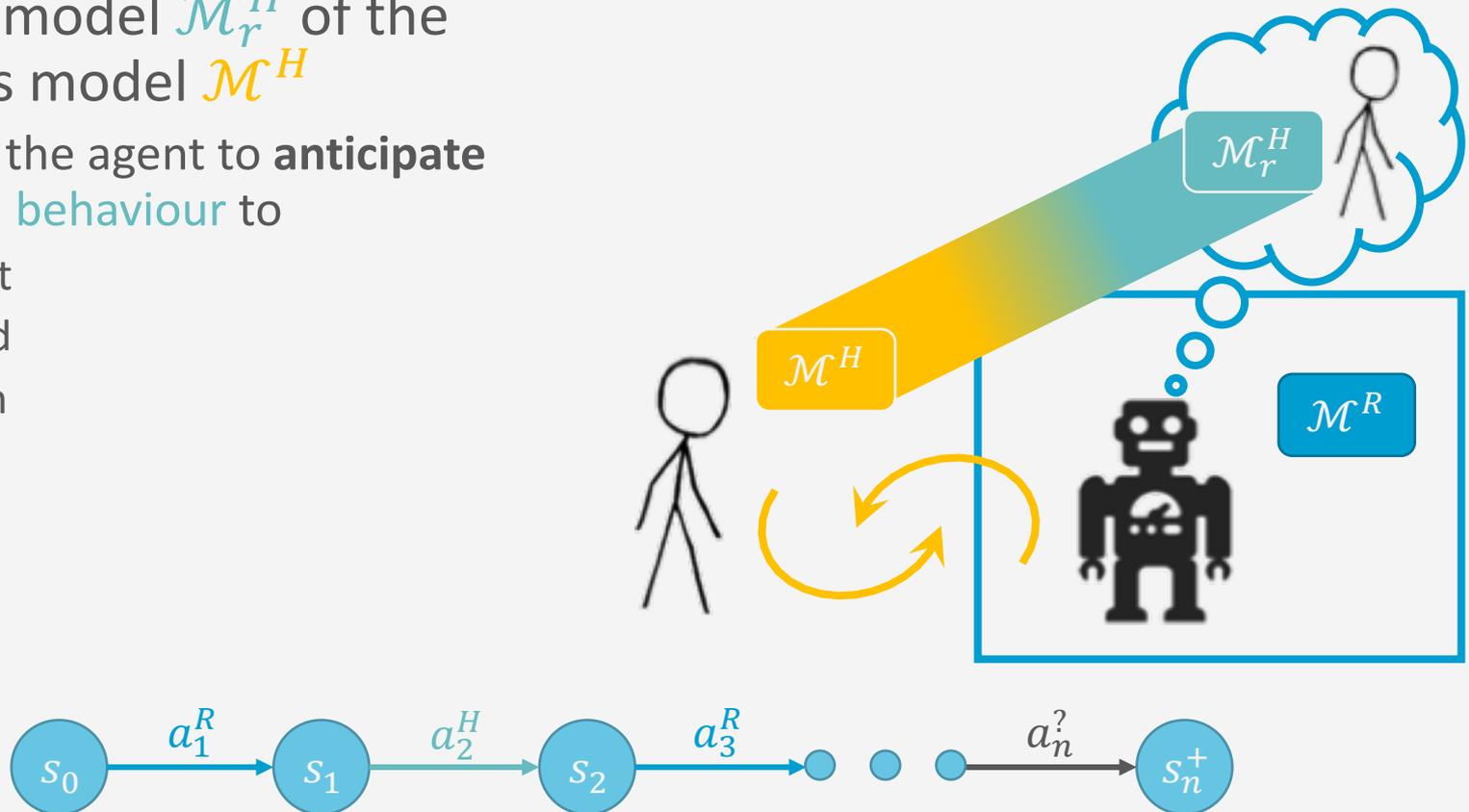
Collaborative Planning

- Given (Σ, s_0, S_g) , i.e., the agent's model \mathcal{M}^R
- Find a **joint plan** $\pi = \langle a_1^R, a_2^H, \dots, a_n^? \rangle$ that transforms s_0 to a state $s_n^+ \in S_g$



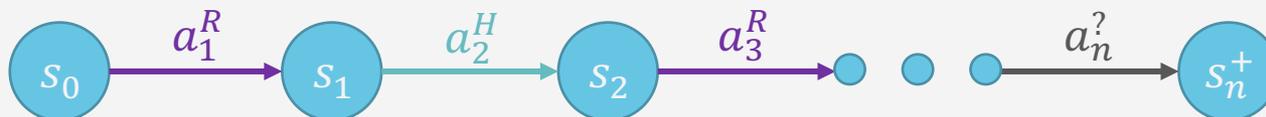
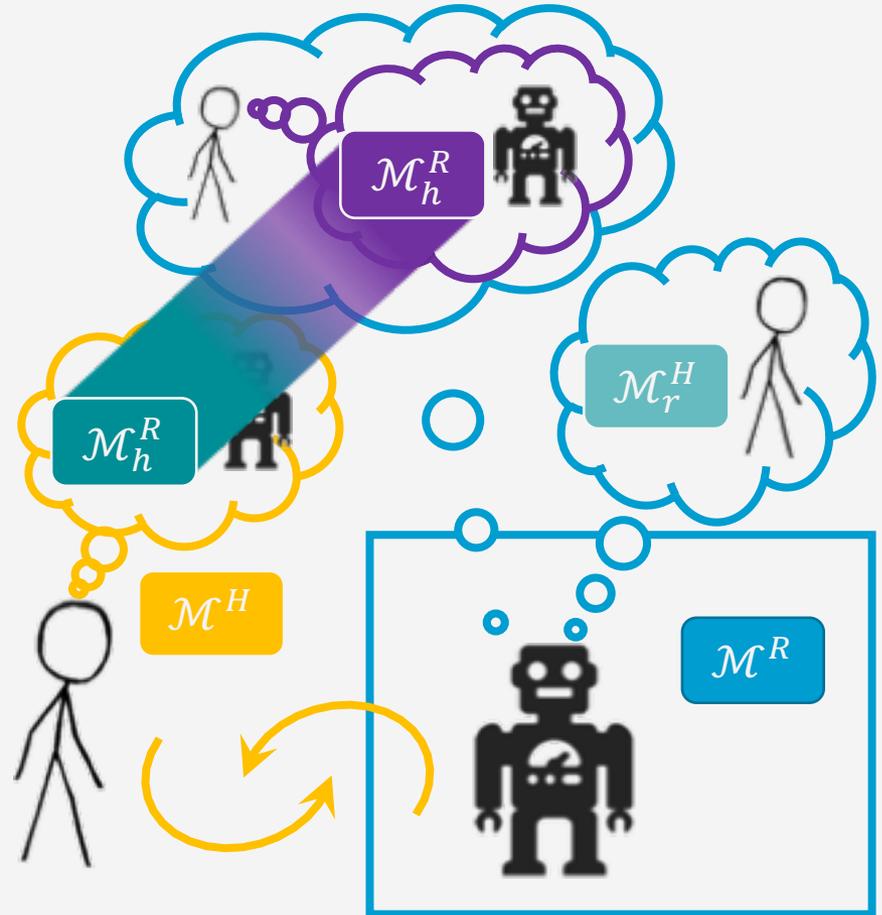
Human-aware Planning

- Next to \mathcal{M}^R
- Agent's model \mathcal{M}_r^H of the human's model \mathcal{M}^H
 - Allows the agent to **anticipate human behaviour** to
 - assist
 - avoid
 - team



Human-aware Planning

- Next to \mathcal{M}^R and \mathcal{M}_r^H
- Agent's model $\tilde{\mathcal{M}}_h^R$ that the agent expects the human to have of \mathcal{M}^R
 - Allows the agent to **anticipate human expectations** to
 - conform to those expectations
 - explain its own behaviour in terms of those expectations

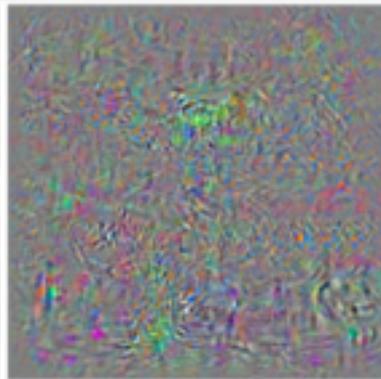


Generating Mental Models

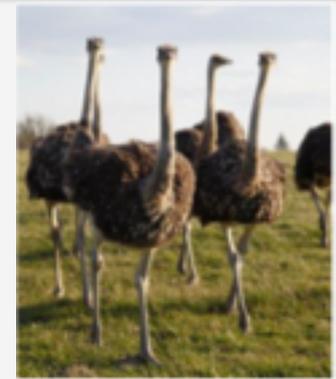
- Known beforehand (handcrafted/researched)
 - Urban Search and Rescue
 - Teaching
- Learning simple models for generating explanations/explicability
- Learning full models (transition functions, rewards)
 - Through interaction with users

XAI & Explanations

- Standard XAI: view of explanations too simple
 - Debugging tool for “inscrutable” representations
 - “**Pointing**” explanations (primitive)



Please point to
the “ostrich” part



Prediction:
School bus

Difference between left
and right magnified by 10

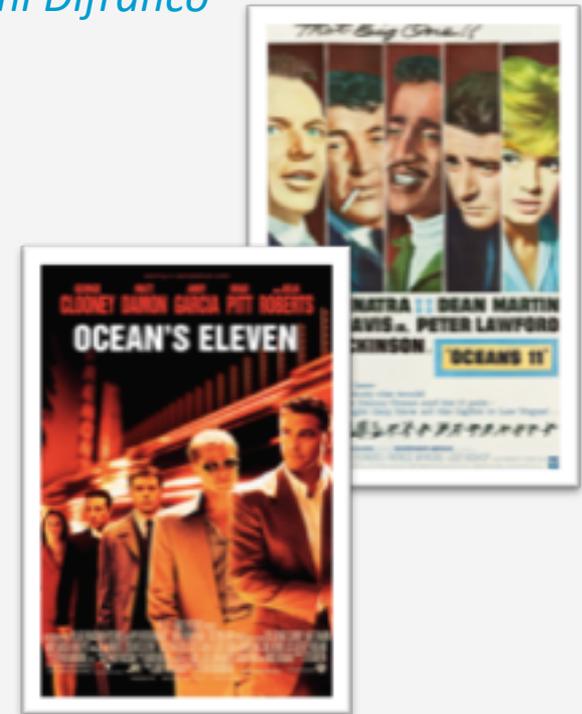
Prediction:
Ostrich

- Explaining decisions will involve pointing over space-time tubes
- Explanations critical for collaboration
 - But not as a monologue from the agent → **interaction**

Ethical Quandaries of Interaction

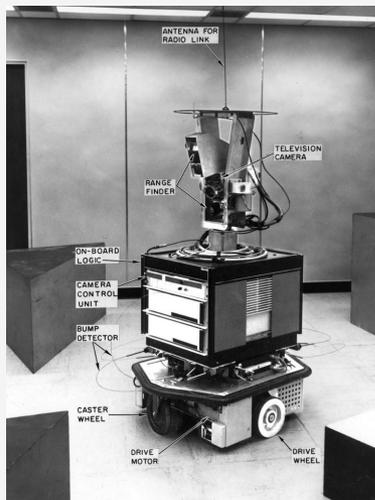
- Evolutionary, mental modelling allowed us to both cooperate or compete/sabotage each other
 - Lying is only possible because we can model others' mental states
- Human-aware AI systems with mental modelling capabilities bring additional ethical quandaries
 - E.g., automated negotiating agents that misrepresent their intentions to gain material advantage
 - Your personal assistant that tells you white lies to get you to eat healthy (or not...)

*Every tool is a
weapon, if you
hold it right.
--Ani Difranco*



Ethical Quandaries of Interaction

- Humans' example closure tendencies are more pronounced for emotional/social intelligence aspects
 - No one who saw Shakey the first time thought it could shoot hoops, yet the first people interacting with Eliza assumed it was a real doctor
 - Concerns about human-aware AI "toys" such as Cozmo (e.g., Sherry Turkle)



```

Welcome to
EEEEEE LL      IIII ZZZZZZZ AAAAA
EE      LL      II      ZZ      AA      AA
EEEEEE LL      II      ZZZ      AAAAAAA
EE      LL      II      ZZ      AA      AA
EEEEEE LLLLLL IIII ZZZZZZZ AA      AA

Eliza is a mock Rogerian psychotherapist.
The original program was described by Joseph Weizenbaum in 1966.
This implementation by Norbert Landsteiner 2005.

ELIZA: Is something troubling you ?
YOU:  Men are all alike.
ELIZA: What is the connection, do you suppose ?
YOU:  They're always bugging us about something or other.
ELIZA: Can you think of a specific example ?
YOU:  Well, my boyfriend made me come here.
ELIZA: Is it important to you that your boyfriend made you come here ?
YOU:  He says I'm depressed much of the time.
ELIZA: I am sorry to hear that you are depressed.
YOU:  It's true. I am unhappy.
ELIZA: Can you explain what made you unhappy ?
YOU:
  
```

Differences in Mental Models

- Expectations on **capabilities**
 - Human may have misconceptions about robot's actions
 - Certain actions in human's mental model may not be feasible for robot
- Expected state of the **world**
 - Human may assume certain facts are true (when they are not true)
- Expected **goals**
 - Human may have misconceptions about robot's objectives/intentions
- **Sensor** model differences
 - Human may have partial observability of robot's activities
 - Human may have incorrect beliefs about robot's observational capabilities
- Different **representations**
 - Robot's innate representation scheme might be too complex for human
 - Human may be thinking in terms of a different vocabulary

Model Differences

- Robot and human may have different models of same task
 - Divergence in models can lead to expectation mismatch
 - Consequence: Plans that are optimal to robot may not be so in model of human
 - **Inexplicable** plans
- Robot has two options
 - **Explicable planning** – sacrifice optimality in own model to be explicable to human
 - *interpretable behaviour*
 - **Plan Explanations** – resolve perceived suboptimality by revealing relevant model differences
 - *model reconciliation*

Intermediate Summary

- Different mental models
 - Mental model of the human
 - Mental model that the human has of the agent
 - Mental model that the agent assumes the human has of the agent
- Differences between mental models
 - May lead to inexplicable behaviour
- Ethical quandaries
 - Modelling mental state of humans requires ethical behaviour of agent

Outline

Mental Models

- Human-aware agent

Interpretable Behaviour

- Explicability
- Legibility
- Predictability

Explanations

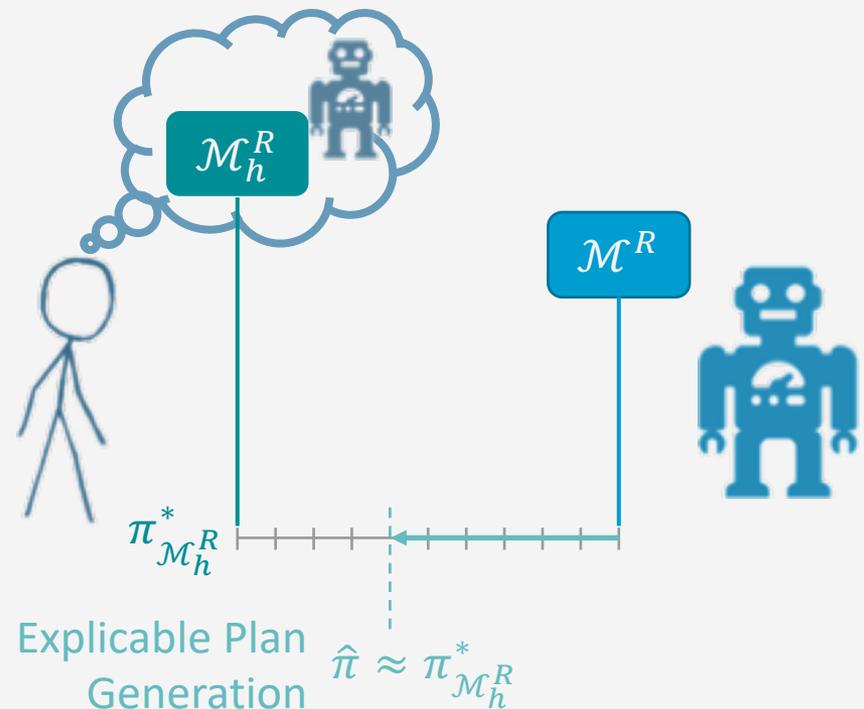
- Model reconciliation

Interpretable Behaviour

- **Explicable** behaviour
 - Acting in a way that make sense to the user
- **Legible** behaviour
 - Acting in a way that convey necessary information to the user
- **Predictable** behaviour
 - Acting in a way that allow users to accurately anticipate future behaviour

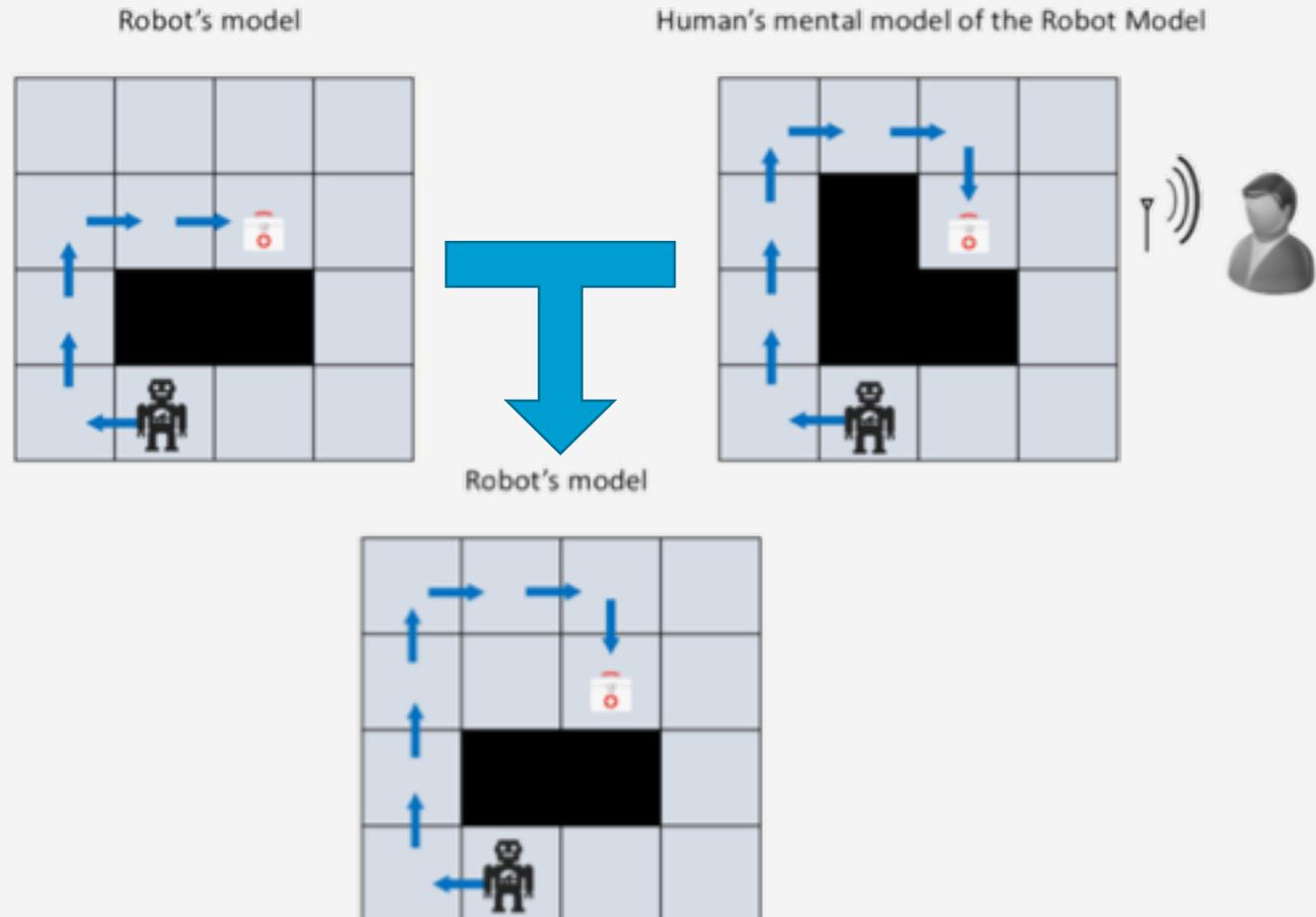
Why *Explicable* Behaviour?

- Robot's behaviour may diverge from human's expectations of it
- Human may get surprised by robot's inexplicable behaviour
- One way to avoid surprising a human involves generating **explicable behaviour by conforming to human's expectations**
 - Account for human's mental model



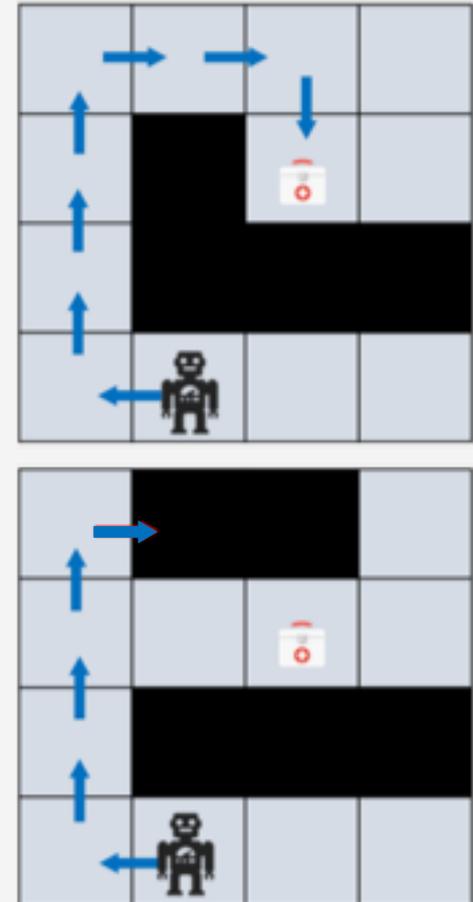
Explicable Behaviour

- Example:
Robot may have to sacrifice its optimality to improve explicability



Model-based Explicable Behaviour

- Human's mental model is available to the robot
- But robot may not be able to plan directly with human mental model
- Find a valid plan that is “closest” to the expected plan
- Involves minimising distance w.r.t. expected plans
 - Cost difference in human model
 - Action set difference



Model-free Explicable Planning

- Problem to solve:

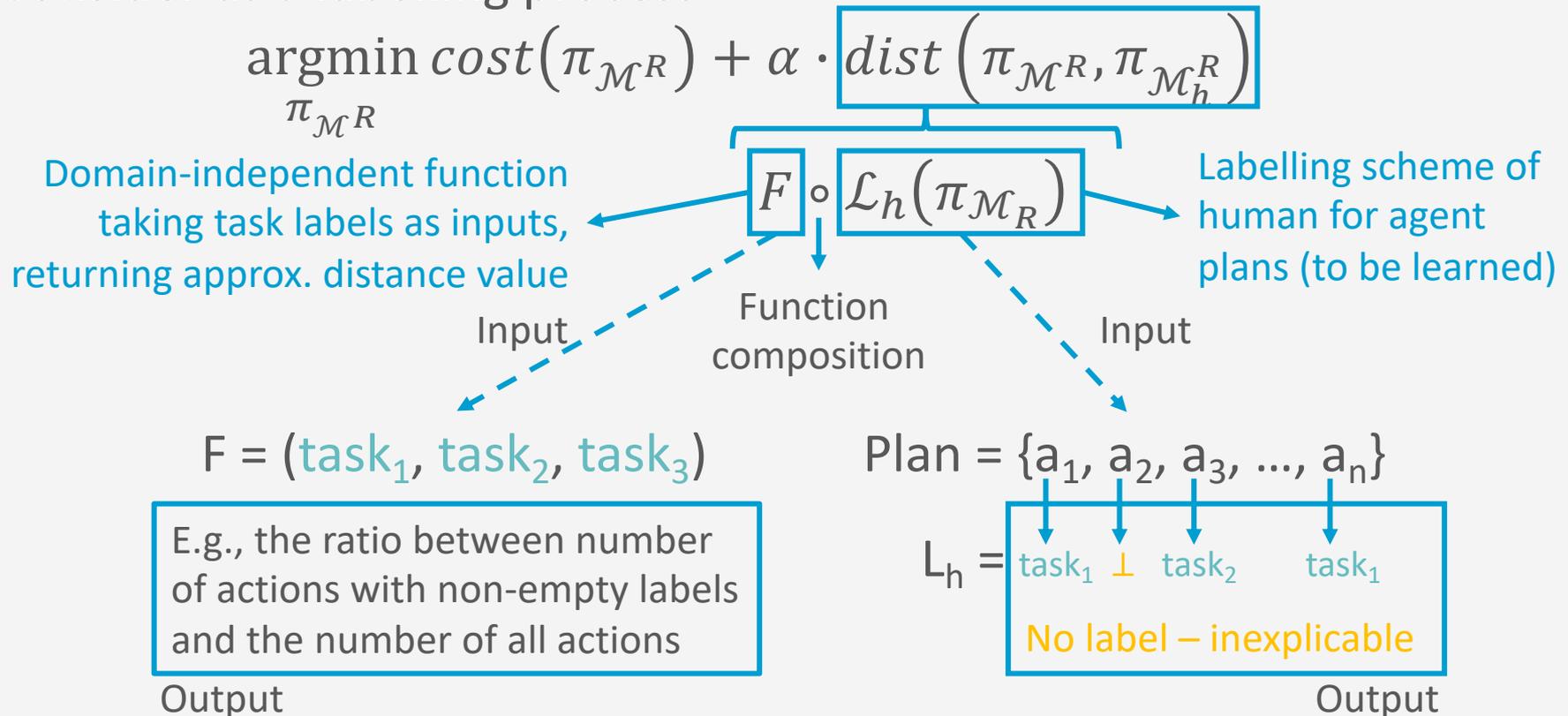
$$\operatorname{argmin}_{\pi_{\mathcal{M}^R}} \boxed{\text{cost}(\pi_{\mathcal{M}^R})} + \alpha \cdot \boxed{\text{dist}(\pi_{\mathcal{M}^R}, \pi_{\mathcal{M}_h^R})}$$

Cost of robot plan
Distance between robot plan and human's expectation of robot plan

- Robot may not have human's mental model \mathcal{M}_h^R upfront
 - But: We do not necessarily need to learn the full model

Model-free Explicable Planning

- Understand = Associate abstract tasks with actions
- Consider as a labelling process

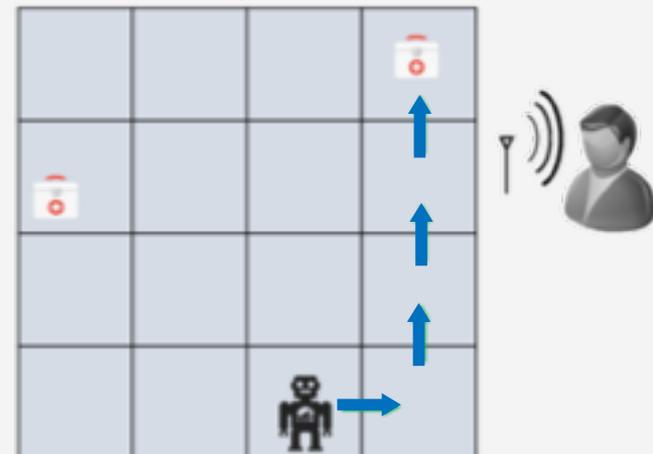
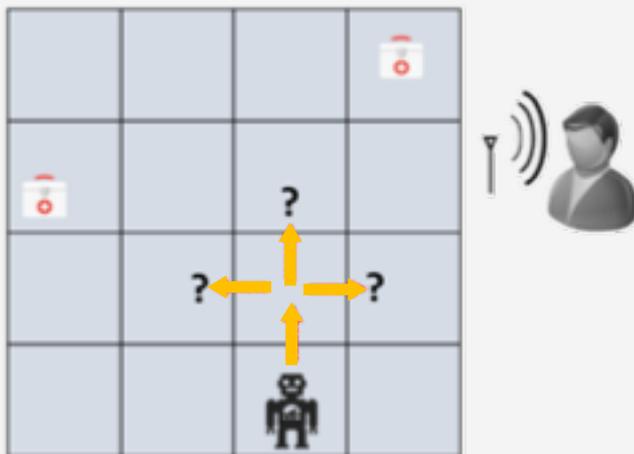


Why *Legible* Behaviour?

- In human-robot teams, essential for the robot to communicate its intentions and objectives to the human
 - Explicitly communicate its intentions to the human
 - Generating a behaviour which **implicitly** reveals robot's intentions to the human
 - Might be easier for the human teammate

Legible Behaviour

- In general, involves a setting where
 - Human has access to **candidate goals** but does not know true goal
- Robot's objective: Convey true goal implicitly through its behaviour
- Human updates its belief on set of candidate goals when it receives observations
- **By synthesising legible behaviour, robot reduces human's uncertainty over candidate goals**



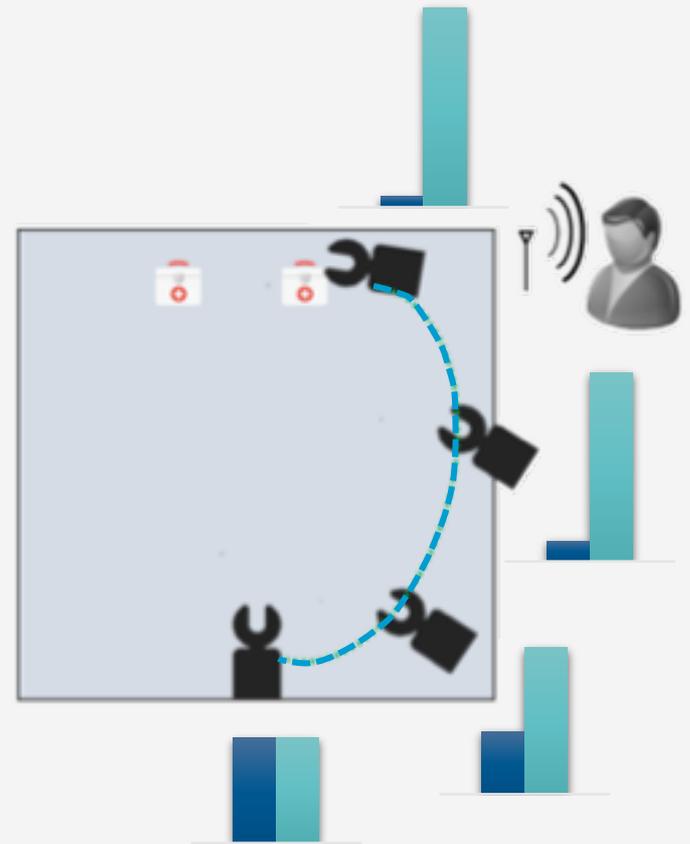
Online Legible Behaviour

- Enables human to **quickly** and **confidently** infer robot's true goal
- Human's belief update is captured using a probabilistic goal recognition system
- Actions that maximise the posterior probability of the true goal G are favoured

$$\operatorname{argmax}_{G \in \mathcal{G}} P(G | \text{Observations})$$

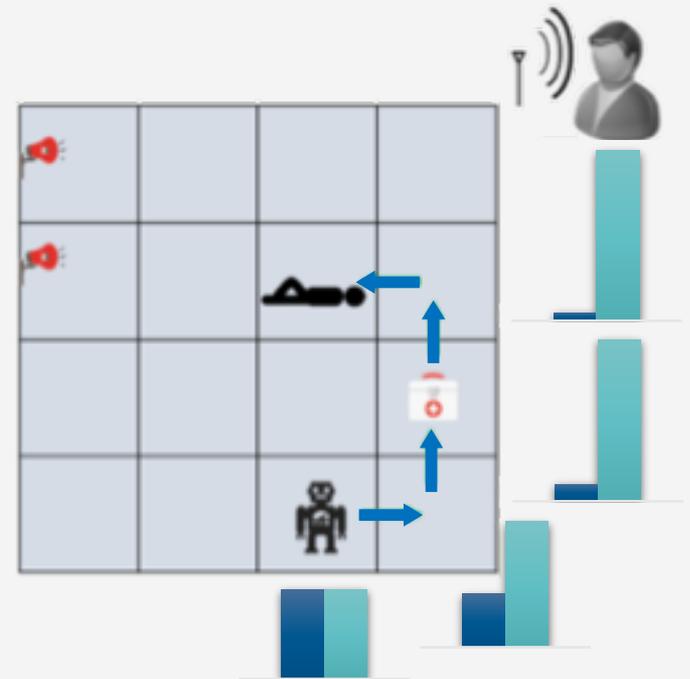
Legible Robot Motion

- Example: Which med kit will the robot pick up?
- While performing goal recognition, human considers shortest distances
- Approach involves finding a trajectory endpoint between start point and true goal such that posterior probability of true goal is maximised
 - The sooner the goal is recognised in the trajectory, the better is the trajectory's legibility



Transparent Planning

- Example: Is the robot surveying the rooms or performing triage?
- Whenever an action is performed, goal recognition system is used to update human's belief
- Objective: Reach a target belief where true goal is more probable than other goals
- Take the first applicable action associated with a belief of highest utility (closest to target belief)

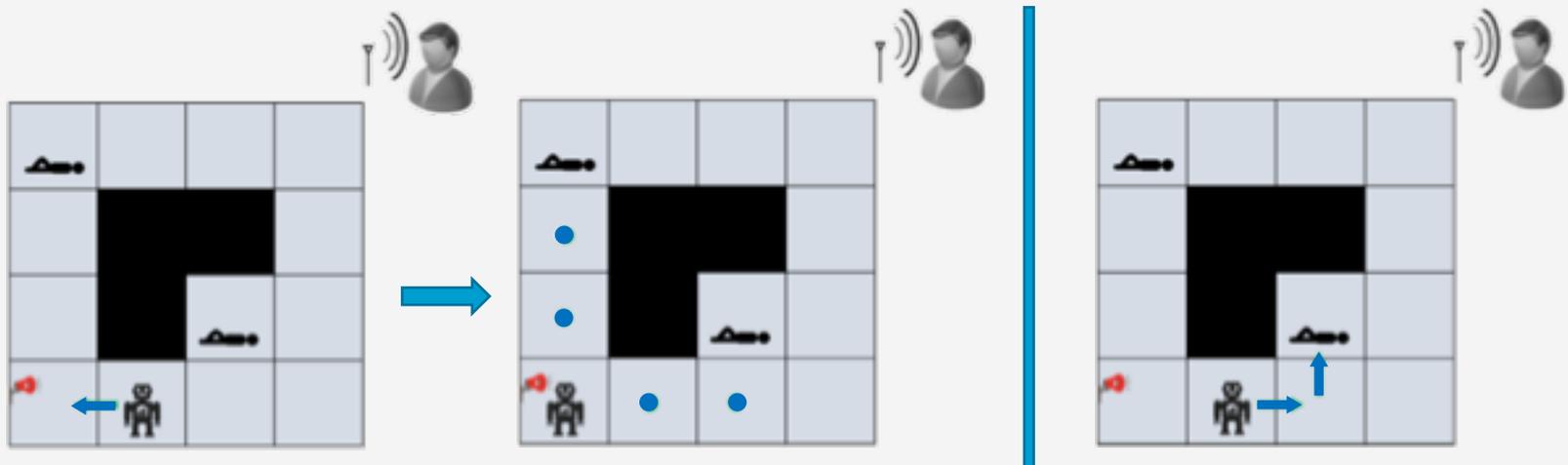


Offline Goal Legibility

- Generalises problem of goal legibility in terms of
 - Partial observability of the human
 - Amount of goal legibility achieved
- Partial observability:
 - Multiple action and state pairs may yield the same observation
 - Human's belief update consists of all possible states that emit given observation and are valid considering previous belief
 - $b_{i+1} = \text{update}(b_i, o_{i+1})$

Offline Goal Legibility

- Example: Robot has to survey and treat a victim
 - Has to convey which victim it is treating
- Key idea: Limit number of candidate goals (at most j goals) possible in observer's final belief
- Explores legible behaviour that satisfies predetermined amount of goal legibility, i.e., the plan is *j -legible*

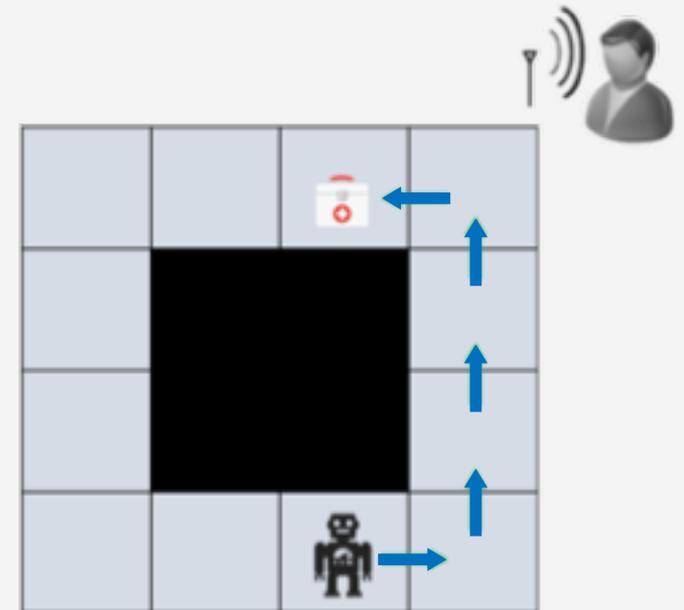


Why *Predictable* Behaviour?

- In human-robot teams, if robot's behaviour cannot be anticipated by human, it can hamper team performance
- Predictable robot behaviours are easy for the human to understand and help in engendering trust in the robot
- *Predictability and legibility are fundamentally different and often contradictory properties of motion*

Predictable Behaviour

- In general, involves a setting where
 - Human knows start state and goal but does not know which plan will be executed
- Robot's objective is to behave in a way that can be anticipated by the human
- Observer updates its belief on set of valid plans when it receives observations
- By synthesising predictable behaviour, robot reduces human's uncertainty over possible behaviours



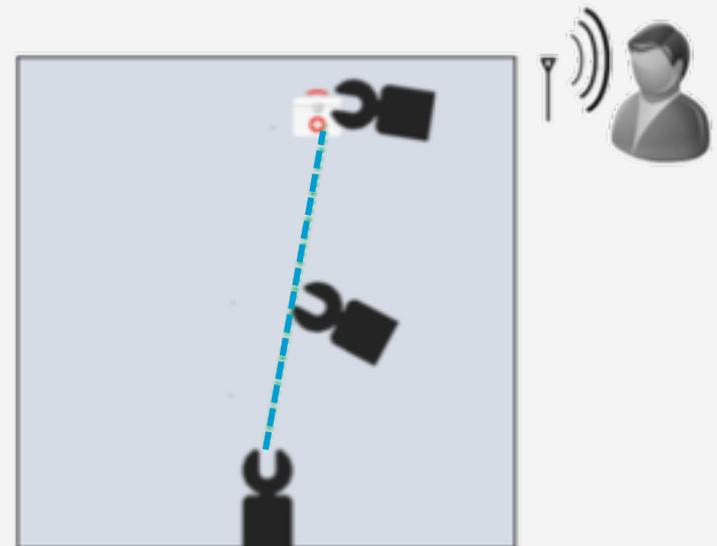
Predictable Robot Motion

- Example: What trajectory will robot take?
- Human assumes that robot is rational and that it prefers short length trajectory
- Most predictable trajectory optimises path towards the goal (C cost fct. modelling human's expectation)

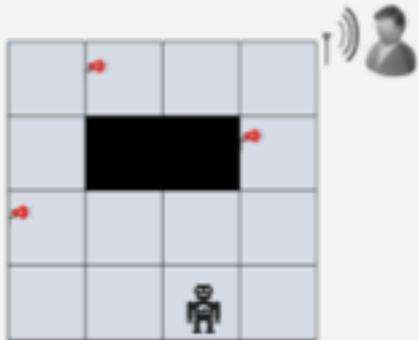
$$\underset{traj}{\operatorname{argmin}} C(traj)$$

- There are two aspects of generating predictable motion:

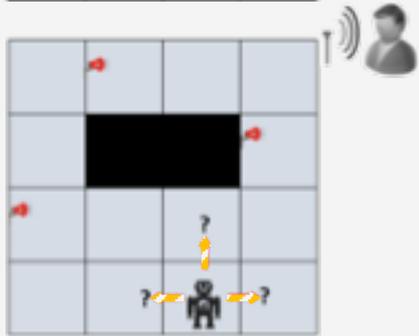
- Learning C
- Minimising C



t -Predictability



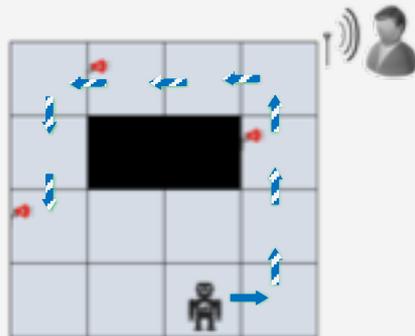
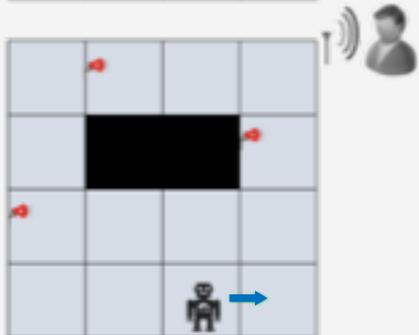
- Key idea: first t actions should foreshadow rest of actions
- Example: What route would the robot take to survey the rooms?



- t -predictability score $P_t =$ probability of sequence $a_{t+1} \dots a_T$, given start state, goal and $a_1 \dots a_t$

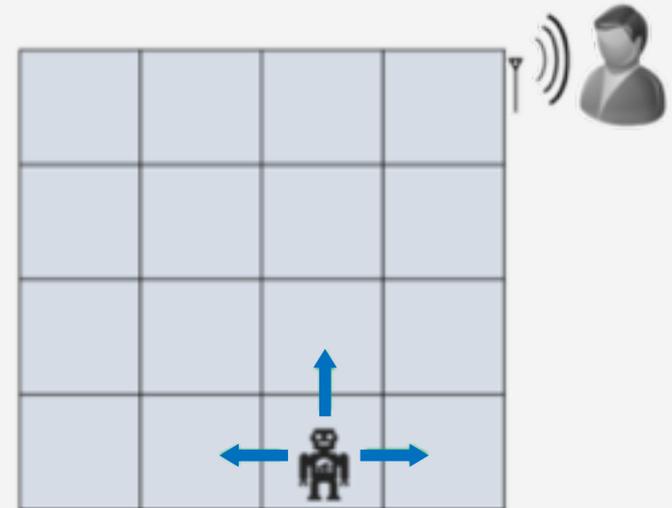
- t -predictable planner finds action sequence \mathbf{a}^* such that

$$\mathbf{a}^* = \operatorname{argmax}_{\mathbf{a} \in A} P_t(\mathbf{a})$$



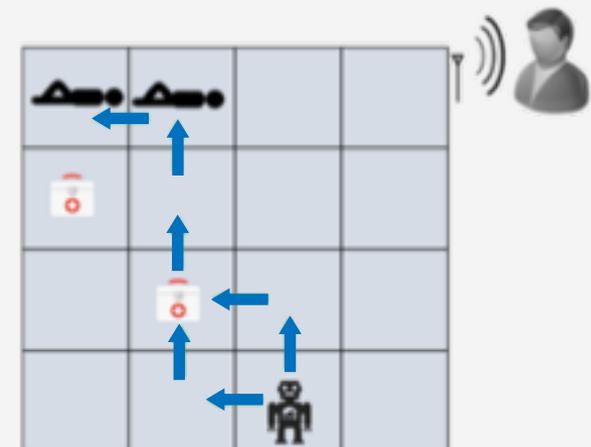
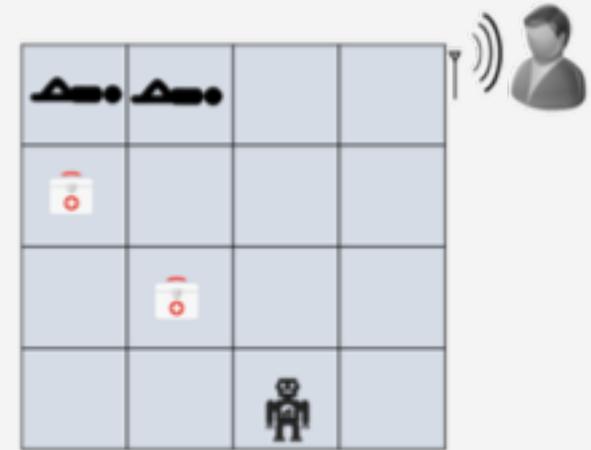
Offline Plan Predictability

- Assume offline setting
 - Human has partial observability
 - Belief update performed after receiving all observations
- Human guesses robot's actions based on plans that
 - Are consistent with observation sequence
 - Achieve goal
- Generalises the problem of conveying actions to observer



Offline Plan Predictability

- Example:
 - Robot has to perform triage
 - Which med kit should the robot pick?
- Solution: Generate a plan whose observation sequence is associated with
 - At least m plans to the same goal,
 - And the plans have high similarity.
 - i.e., m plans that are at most d distance from each other – m -similar plans
- Using plan distance metrics
 - Action set distance gives the number of similar actions given two plans



Summary

- Aspects of interpretable behaviour
- Explicability
 - Act in a way that is comprehensible to the human agent
- Legibility
 - Act in a way such that a human agent can determine which goal is pursued by agent
- Predictability
 - Act in a way such that a human agent can predict the next steps given the previous steps

Outline

Mental Models

- Human-aware agent

Interpretable Behaviour

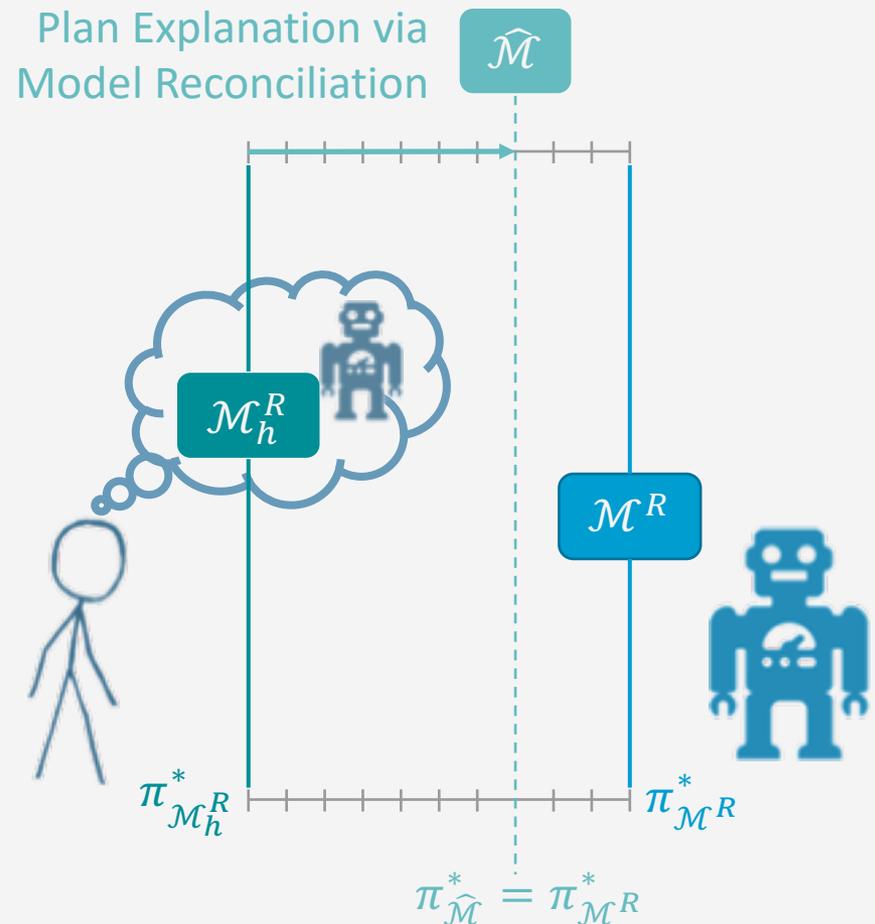
- Explicability
- Legibility
- Predictability

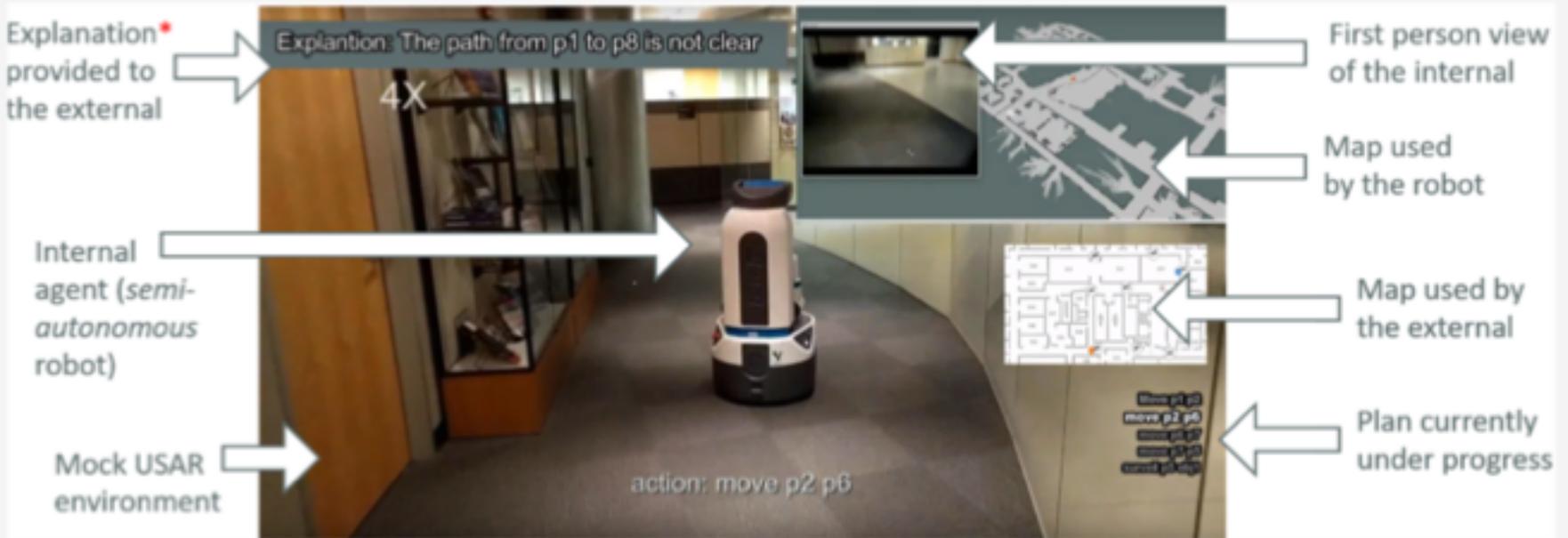
Explanations

- Model reconciliation

Plan Explanations

- Conforming to expectations of human
 - E.g., by explicable planning; may involve giving up optimal plan
 - But: May *not* be feasible
- **Model reconciliation:** Bring mental model closer by explanations
 - Planner is optimal in own but not in human's model
 - Given a plan, explanation is a model update
 - After explanation, plan is also optimal in the updated human model





Example

- Mock search and reconnaissance scenario with internal robot and external human



Aspects to Explanations

- **Completeness:** No better explanation exists, no aspect of plan remains inexplicable
 - Requires explanations of a plan to be comparable
- **Conciseness:** Explanations are easily understandable to the explainee
 - The larger an explanation, the harder for the human to incorporate information into deliberation process
- **Monotonicity:** Remaining model differences cannot change completeness of explanation, i.e., all aspects of model that yielded plan are reconciled
 - Subsumes completeness
- **Computability:** Ease of computing explanation from robot's point of view

Types of Explanations

- **Plan Patch Explanation (PPE)**
 - Provide model differences pertaining to only the actions present in the plan that needs to be explained
- **Model Patch Explanation (MPE)**
 - Provide all model differences to the human
- **Minimally Complete Explanation (MCE)**
 - Shortest complete explanation
 - Can be rendered invalid given further updates
- **Minimally Monotonic Explanation (MME)**
 - Shortest explanation preserving monotonicity
 - Not necessarily unique as there may be model differences supporting the same causal links in the plan; exposing one link is enough (to guarantee optimality in the updated model)

Aspects of Types of Explanations

- Plan Patch Explanation (PPE)
- Model Patch Explanation (MPE)
- Minimally Complete Explanation (MCE)
- Minimally Monotonic Explanation (MME)

Explanation Type	Completeness	Conciseness	Monotonicity	Computability
PPE	X	✓*	X	✓
MPE	✓	X	✓	✓
MCE	✓	✓	X	?
MME	✓	✓	✓	?

$$|\text{approx. } MCE| \leq |\text{exact } MCE| < |MME| \ll |MPE|$$

* In the sense that it focuses on the differences w.r.t. the plan but not necessarily a short explanation

Example – FetchWorld

- Fetch robot whose design requires it to tuck its arms and lower its torso or crouch before moving – not obvious to human navigating



Robot's Model

```
(: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)))

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

Human's Model

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

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

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

Example – FetchWorld

- Initial state and goal:
- Robot’s optimal plan:
- Human’s expected plan:

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

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

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

Robot’s Model

```
(: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)))
```

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

Human’s Model

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

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

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

Example – FetchWorld

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

- Robot's optimal plan:

Robot's Model

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

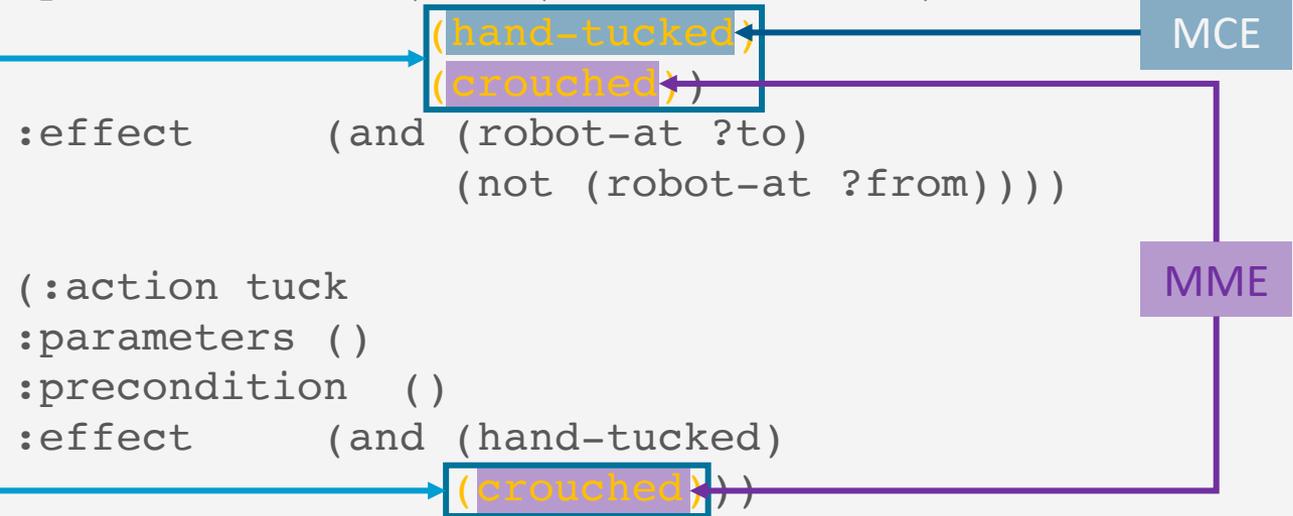
(:action tuck
:parameters ()
:precondition ()
:effect (and (hand-tucked)
(crouched)))

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

PPE = MPE

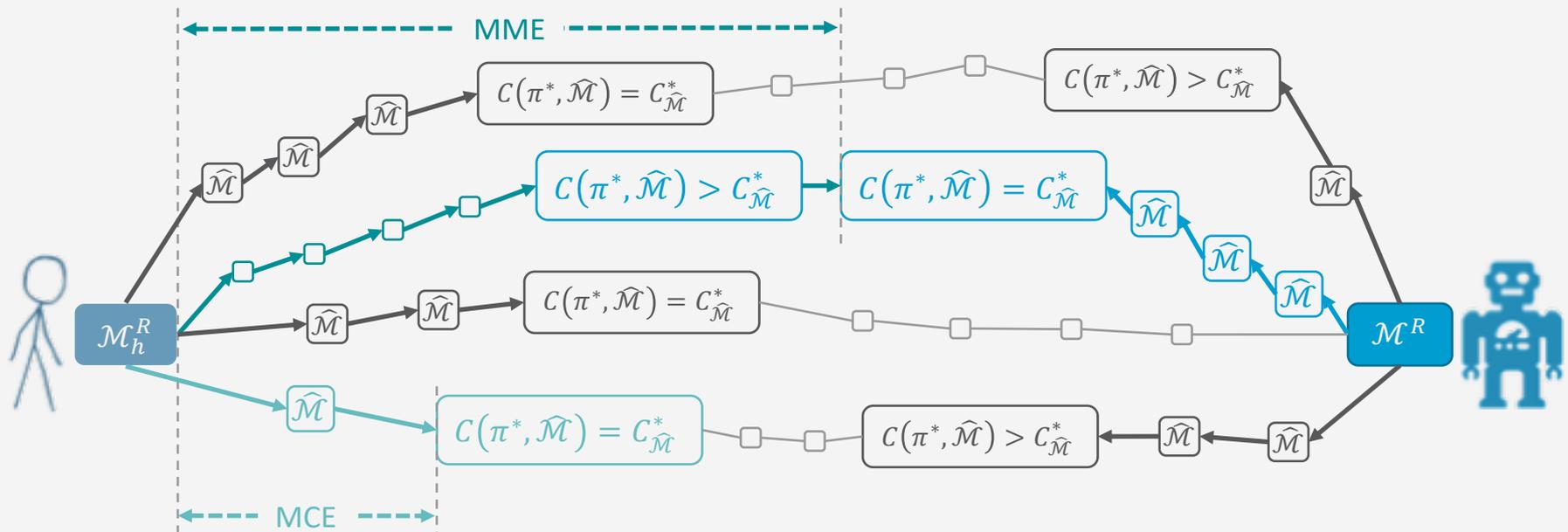
MCE

MME



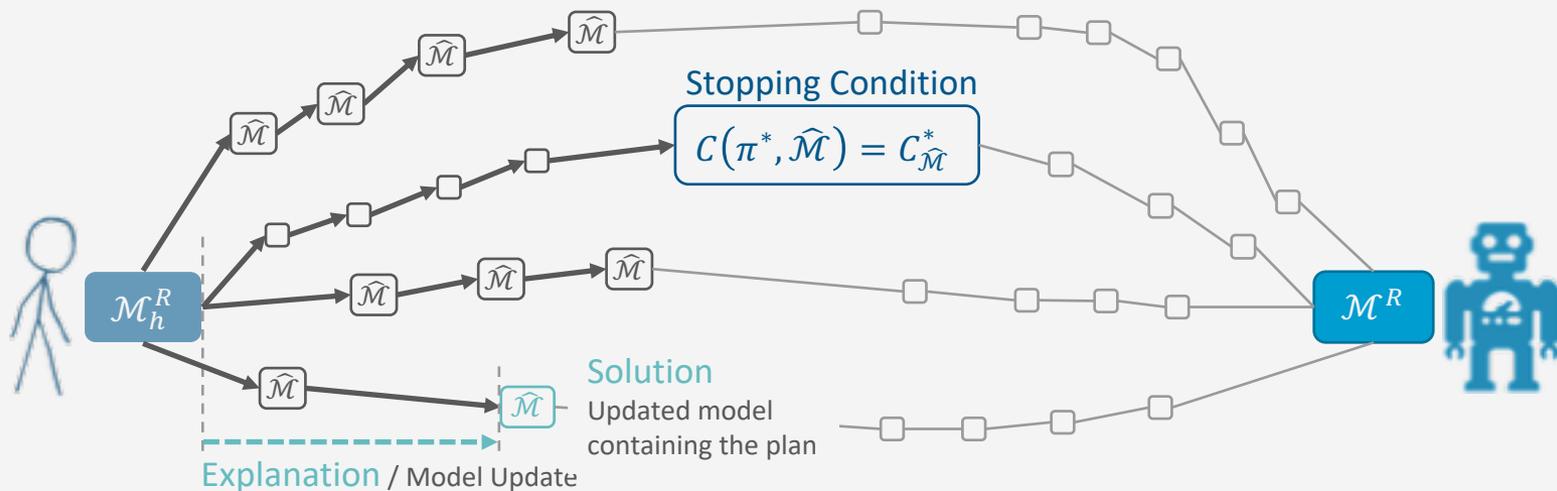
Model Space Search

- Search algorithms for finding MCEs and MMEs



Model Space Search

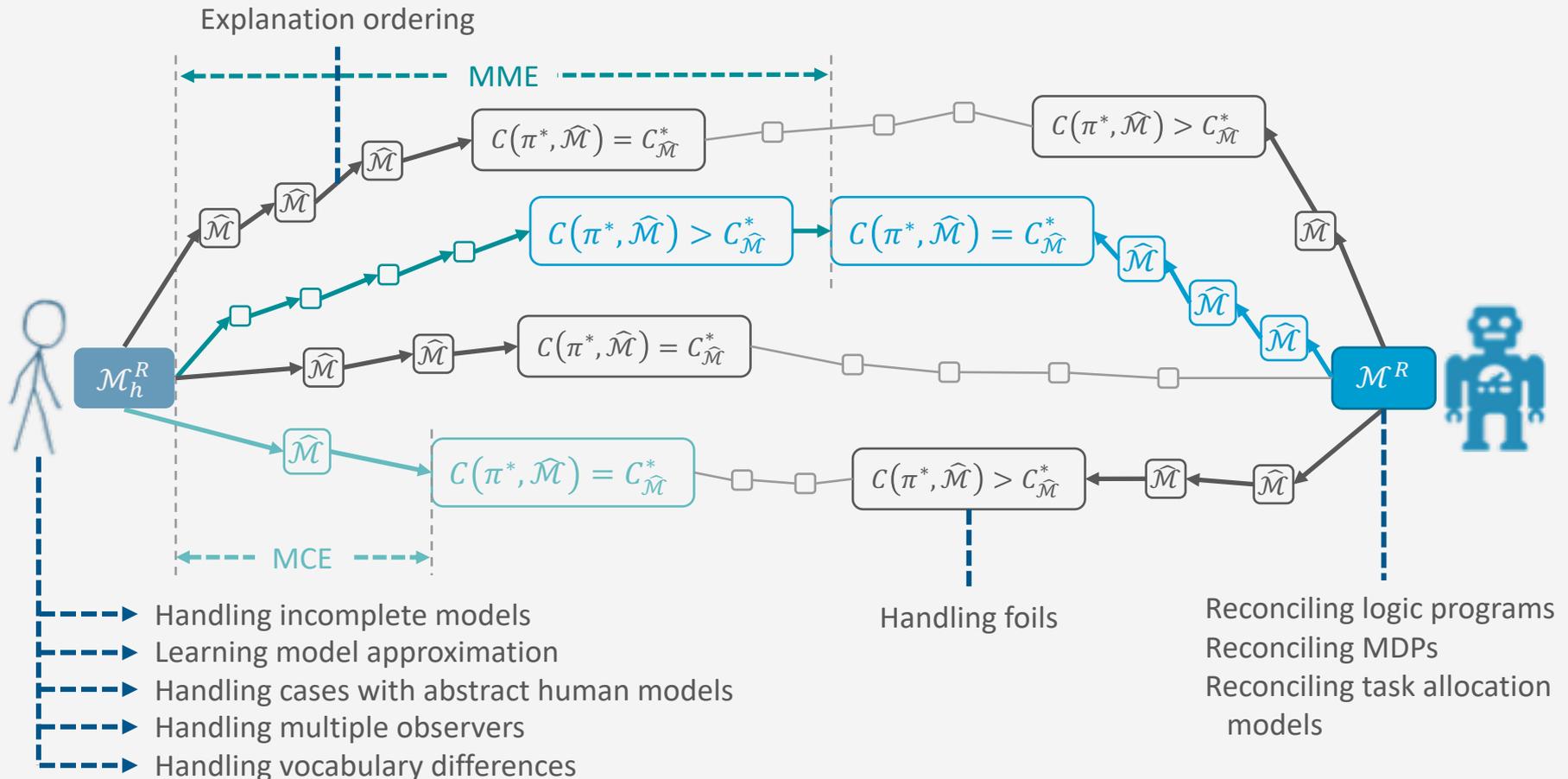
- Human-aware planning: Given model of planning problem and mental model of human, find right model to plan in
 - Trade-off explicability and explanation



- Minimise

cost/length of explanations + $\alpha \cdot$ departure from optimality

Extensions (Outlook)



Summary

- Model reconciliation
 - Explain differences in model
 - PPE, MPE, MCE, MME

Outline

Mental Models

- Human-aware agent

Interpretable Behaviour

- Explicability
- Legibility
- Predictability

Explanations

- Model reconciliation

The End

Content

1. Planning and Acting with **Deterministic** Models
 - Conventional AI planning
2. Planning and Acting with **Refinement** Methods
 - Abstract activities → collections of less-abstract activities
3. Planning and Acting with **Temporal** Models
 - Reasoning about time constraints
4. Planning and Acting with **Nondeterministic** Models
 - Actions with multiple possible outcomes
5. **Standard** Decision Making
 - Utility theory
 - Markov decision process (MDP)
6. Planning and Acting with **Probabilistic** Models
 - Actions with multiple possible outcomes, with probabilities
7. **Advanced** Decision Making
 - Hidden goals
 - Partially observable MDP (POMDP)
 - Decentralised POMDP
8. **Human-aware** Planning
 - Planning with a human in the loop