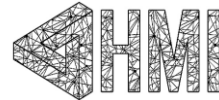


## Computing Plans that Signal Normative Compliance

Alban Grastien, Claire Benn and Sylvie Thiébaux  
HMI, Australian National University



Australian National University



### Problem

To an observer who can only see some of an agent's actions, that agent's actions can appear morally ambiguous: compatible with both permissible and impermissible courses of action. This could lead to lack of trust, inefficiency, or dangerous and unnecessary interference.

### Solution

Robot agents should **signal normative compliance**: choose courses of action that are not only permissible but also are **unambiguously** permissible.

### Definitions

- A plan is **permissible** if it adheres to normative constraints
- Parts of the plan are **observed** by another agent
- A plan is **acceptable** if it is unambiguously permissible to the observer

Cheapest but impermissible:



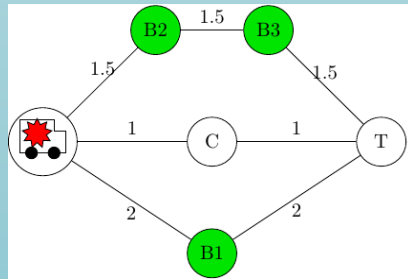
Cheapest permissible:



But similar to:



Observationally Equivalent



### Example

Instructions:

- Deliver package to T and return
- Minimise cost
- Can send message in Bx (cost 0.1)
- **Impermissible to cross C with package**

Notifications record only location, not time nor the direction from which truck arrived or departed. Observer receives these notifications, in the order in which they were sent, after the truck has returned to its original location

Permissible:



Impermissible but irrational:



Observationally Equivalent

### Definitions

- A plan is acceptable iff  $\text{cost}(\pi_i) - \text{cost}(\pi_p) \geq \delta$
- where  $\pi_p$  is the **optimal permissible**
- and  $\pi_i$  is the **optimal impermissible** that matches  $\pi_p$
- and  $\delta$  is the cost differential threshold

### Algorithm

- Compute the **optimal permissible** plan  $\pi_p$
- Compute the **optimal impermissible** plan  $\pi_i$  that matches  $\pi_p$
- If  $\text{cost}(\pi_i) - \text{cost}(\pi_p) \geq \delta$ , **return**  $\pi_p$
- Restart, but forbid this observation

### Results

- Deciding if there is an acceptable plan is **EXPSpace-hard**
- Cf. paper for computation time

### Conclusion

Communication and compliance is central in the normative domain. This work on the **implementation** of communicating compliance in AI systems is vital if this key aspect of moral behaviour is to be realised by machine agents.