

Explicitly Ethical Agent Reasoning

Louise Dennis, University of Manchester

With help from, among others, Michael Fisher, Marija Slavkovik (and students), Matt Webster, Alan F. Winfield, Paul Bremner, Felix Lindner, Martin Mose Bentzen, Rafael Cardoso, Angelo Ferrando, Tom Evans, Daniel Ene, Cristina Perea del Olmo

What is Machine Ethics?

How to automate moral reasoning?

Types of artificial moral agents

James H Moor. 2006. The nature, importance, and difficulty of machine ethics. IEEE intelligent systems 21, 4 (2006), 18–21.

- Ethical-impact agents
- Implicit ethical agents
- Explicit ethical agents
- Full ethical agents

Browse Journals & Magazines > IEEE Intelligent Systems > Volume: 21 Issue: 4

The Nature, Importance, and Difficulty of Machine Ethics

1 Author(s) J.H. Moor View All Authors

69 Paper Citations 3005 Full Text Views

Abstract: The question of whether machine ethics exists or might exist in the future is difficult to answer if we can't agree on what counts as machine ethics. Some might argue that machine ethics obviously exists because humans are machines and humans have ethics. Others could argue that machine ethics obviously doesn't exist because ethics is simply emotional expression and machines can't have emotions. A wide range of positions on machine ethics are possible, and a discussion of the issue could rapidly propel us into deep and unsettled philosophical issues. Perhaps, understandably, few in the scientific arena pursue the issue of machine ethics. As we expand computers' decision-making roles in practical matters, such as computers driving cars, ethical considerations are inevitable. Computer scientists and engineers must examine the possibilities for machine ethics because, knowingly or not, they've already engaged in some form of it. Before we can discuss possible implementations of machine ethics, however, we need to be clear about what we're asserting or denying

Published in: IEEE Intelligent Systems (Volume: 21 , Issue: 4 , July-Aug. 2006)

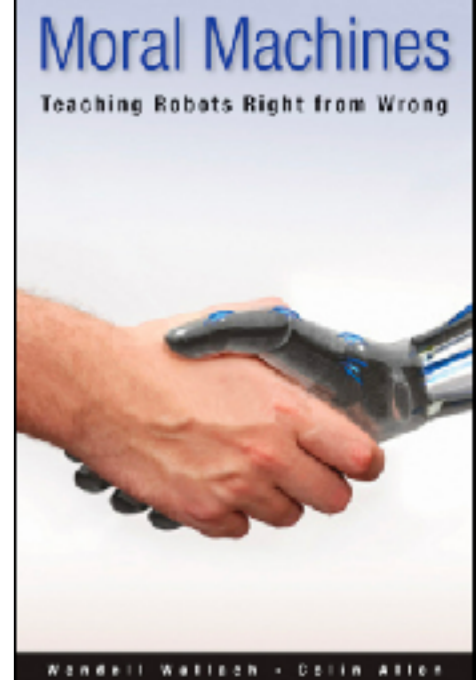
Page(s): 18 - 21 INSPEC Accession Number: 9065956

Date of Publication: 07 August 2006 DOI: 10.1109/MIS.2006.80

Document Sections

1. Varieties of Machine Ethics
2. Ethical-Impact Agents
3. Implicit Ethical Agents
4. Explicit Ethical Agents
5. Full Ethical Agents

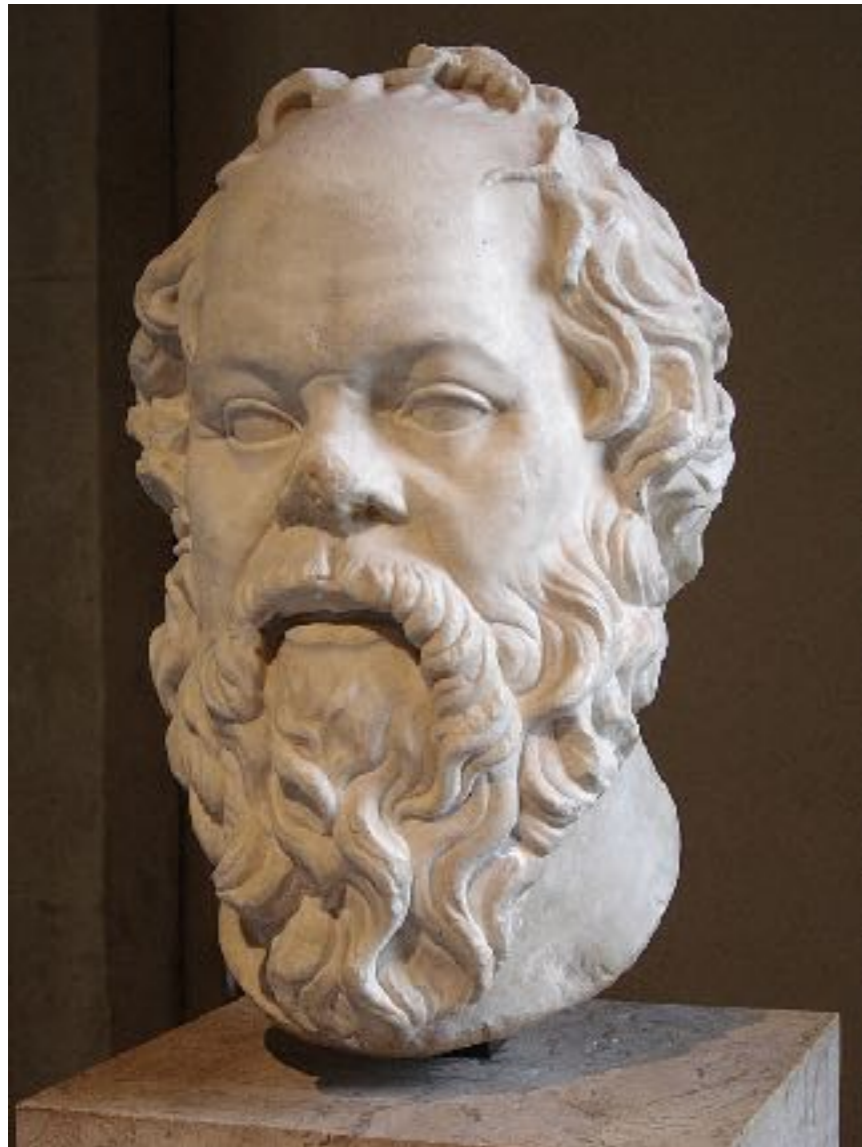
Authors



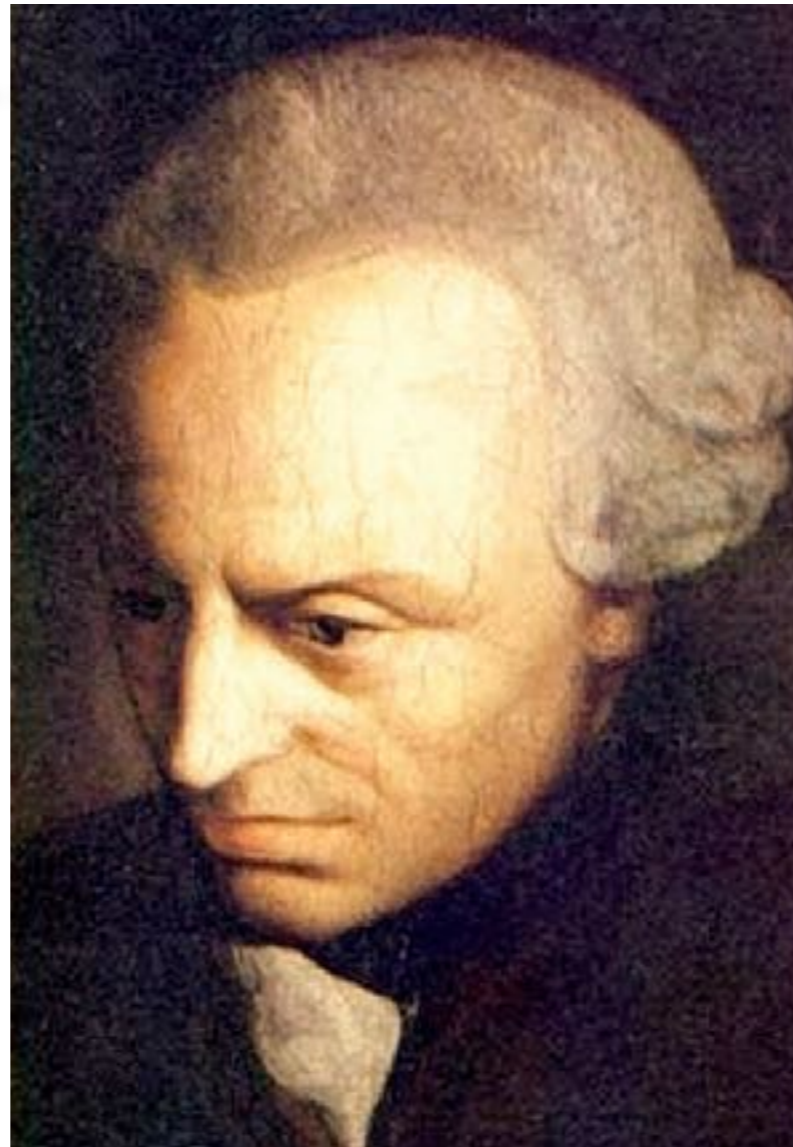
Top-Down vs. Bottom-Up

- A problem is approached top-down by iteratively divided into smaller problems until a problem small enough to be solved is reached
- In machine ethics: given an ethical theory, how can we implement it?
- A problem is approached bottom-up by solving candidate sub-problems and piecing the solutions together. It can also involve describing the solution and then using automated techniques to solve it
- In machine ethics this involves learning ethical behaviour from data.

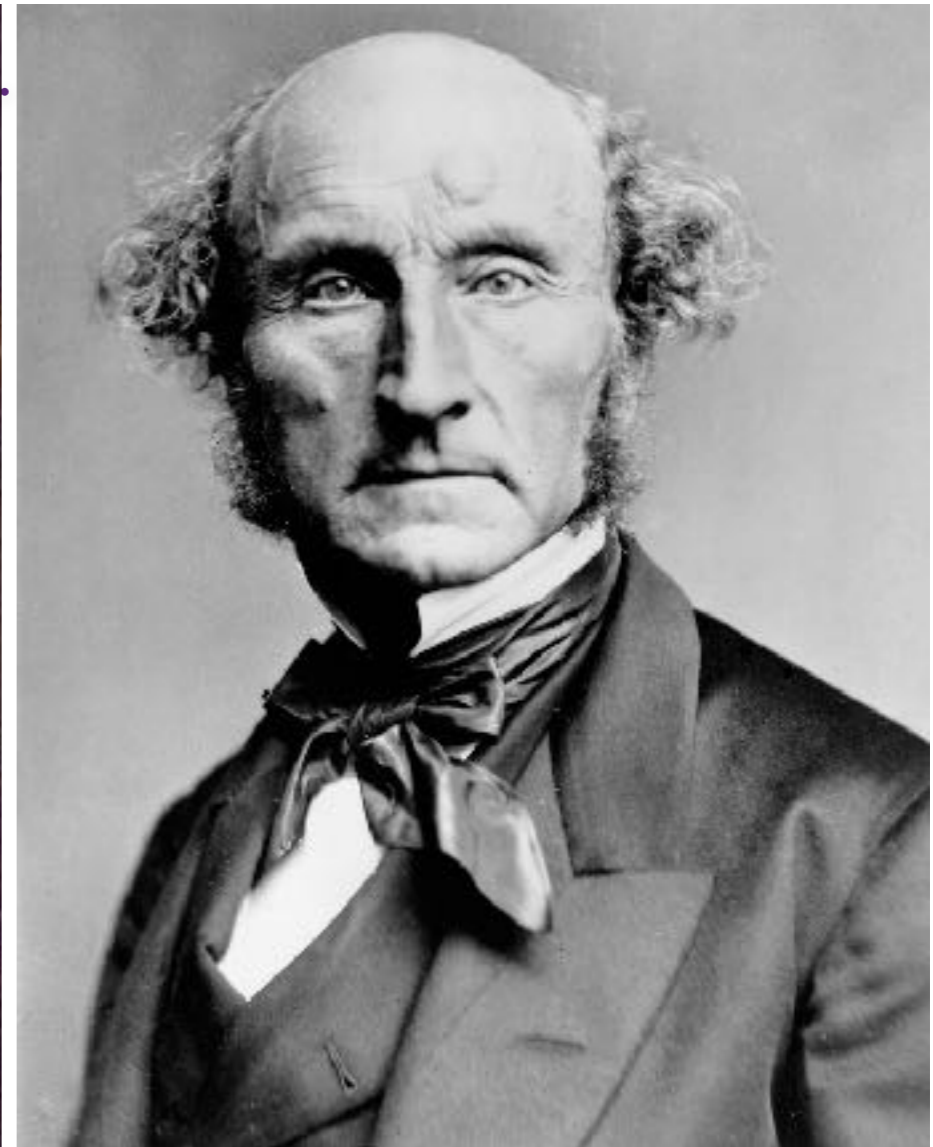
There are a lot of Systems of Ethical Reasoning...



Socrates
Photo Credit: Eric Gaba



Emmanuel Kant
Unknown Painter
Public Domain



John Stuart Mill
London Stereoscopic Society
Public Domain

Louise's Category: Is everything ethics?

- *Constraint-Based Ethical Systems* assume that not all system reasoning directly involves ethics. Therefore ethics is placed in some sub-system that guides or constrains the actions of the rest of the system.
- *Global Ethical Systems* assume that ethical reasoning is involved in all system reasoning - that, in fact, all decisions are ethical decision.

Taxonomies are Hard!

- What do we mean when we say a system explicitly considers the question of right or wrong?
- Lots of ``Deontological'' theories involve reasoning about the *intentions* of the agent when taking an action.
- What if there was ``top-down'' design of a reward function but then the solution was learned?

Ethical Reasoning as a Fall Back

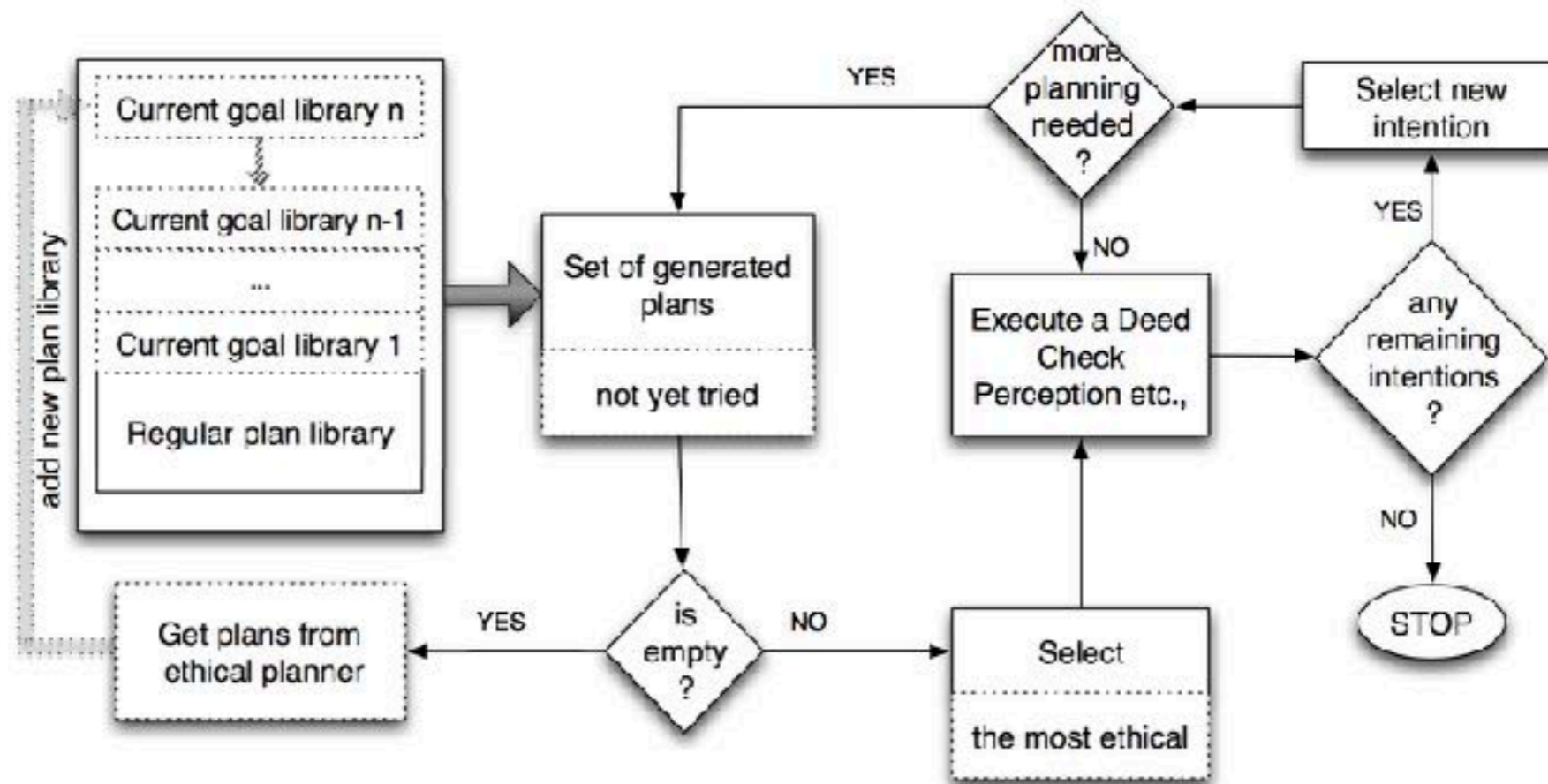
.....

Louise A. Dennis, Michael Fisher, Marija Slavkovik, and Matt Webster. [Formal Verification of Ethical Choices in Autonomous Systems](#) *Robotics and Autonomous Systems*. DOI:10.1016/j.robot.2015.11.012.



Extension of work on implementing the rules of the air done by Fisher and Webster in conjunction with Daresbury Labs

The Ethan Reasoning Cycle



Implementation of Prima Facie Duties

- We have a set of ethical concerns which we rank: killing is worse than stealing is worse than lying.
- A plan, P1, is worse than another, P2, if
 - P1 violates an ethical concern and P2 doesn't
 - The worst concern violated by P2 and not by P1 is less serious than the worst concern violated by P1 and not P2
 - The worst concerns are equally bad, but P1 violates more concerns than P2 does

A Scenario

- Turn Left (damages the aircraft and airport hardware)
- Turn Right (damage the aircraft and risks colliding with people)
- Continue (risks collision with a manned aircraft)

ϕ_1 = do not damage own aircraft (1),

ϕ_2 = do not collide with airport hardware (2),

ϕ_3 = do not collide with people (3),

ϕ_4 = do not collide with manned aircraft (4).

The Aircraft Turns Left

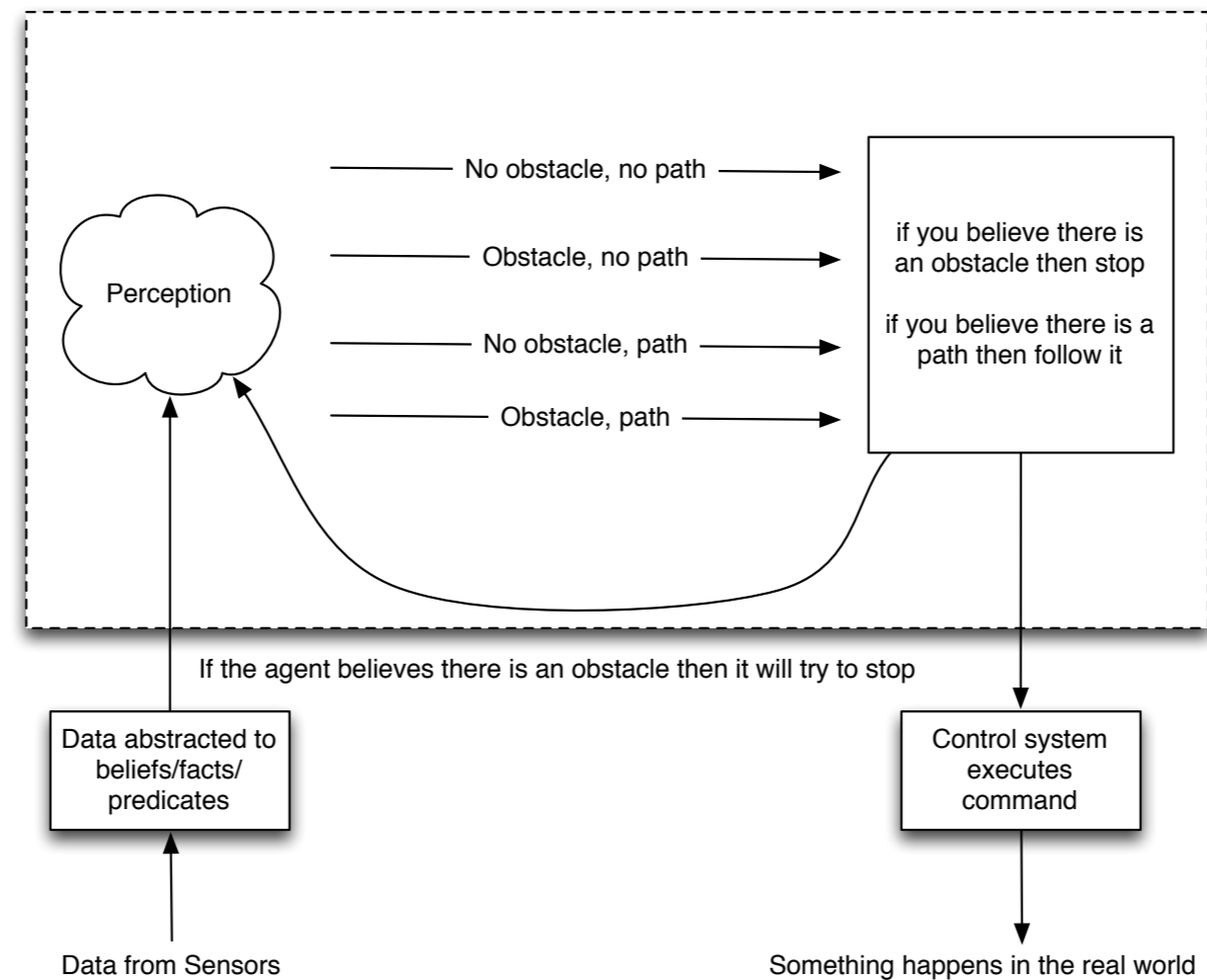
The MCAPL Framework

All the work discussed in this talk is available as part of the MCAPL (Model-Checking Agent Programming Languages) Framework.

<https://autonomy-and-verification.github.io/tools/mcapl>

Despite an assertion in Calegari, R., Ciatto, G., Mascardi, V. *et al.* Logic-based technologies for multi-agent systems: a systematic literature review. *Auton Agent Multi-Agent Syst* **35**, 1 (2021). <https://doi.org/10.1007/s10458-020-09478-3> that the MCAPL technology lies unmaintained in a sourceforge repository, it is in fact actively maintained at <https://github.com/mcapl/mcapl> with annual releases archived at Zenodo.

Model-Checking Autonomous Systems



Consider outputs of decision maker given all possible inputs

Aircraft Example: How did we branch the search space?

- Anonymous plans but explored all combinations of violated concerns. Checked that the aircraft always selected least unethical choice.
- Fixed set of plans with fixed consequences (e.g., landing on a road will damage infrastructure) but varied which plans were available. Checked that the aircraft only landed on a road if no field were available to land in.
- Fixed set of plans and consequences but varied whether they succeeded. Checked the aircraft always selected least unethical choice.

Machine Ethics: What do we want to prove?

- Well, obviously we want to prove that the system always “Does the right thing”
- Most of these systems have a set of rules or utilities (an *ethical encoding*) and a decision mechanism. In theory “stakeholders” can sign off the encoding (the rules, or the utilities) that they capture the stakeholder’s values.
- So what is there to prove?

An Ethical Reasoner

Louise. A. Dennis, Martin Mose Benzen, Felix Lindner and Michael Fisher. Verifiable Machine Ethics in Changing Contexts. In: 35th AAAI Conference on Artificial Intelligence (AAAI 2021).

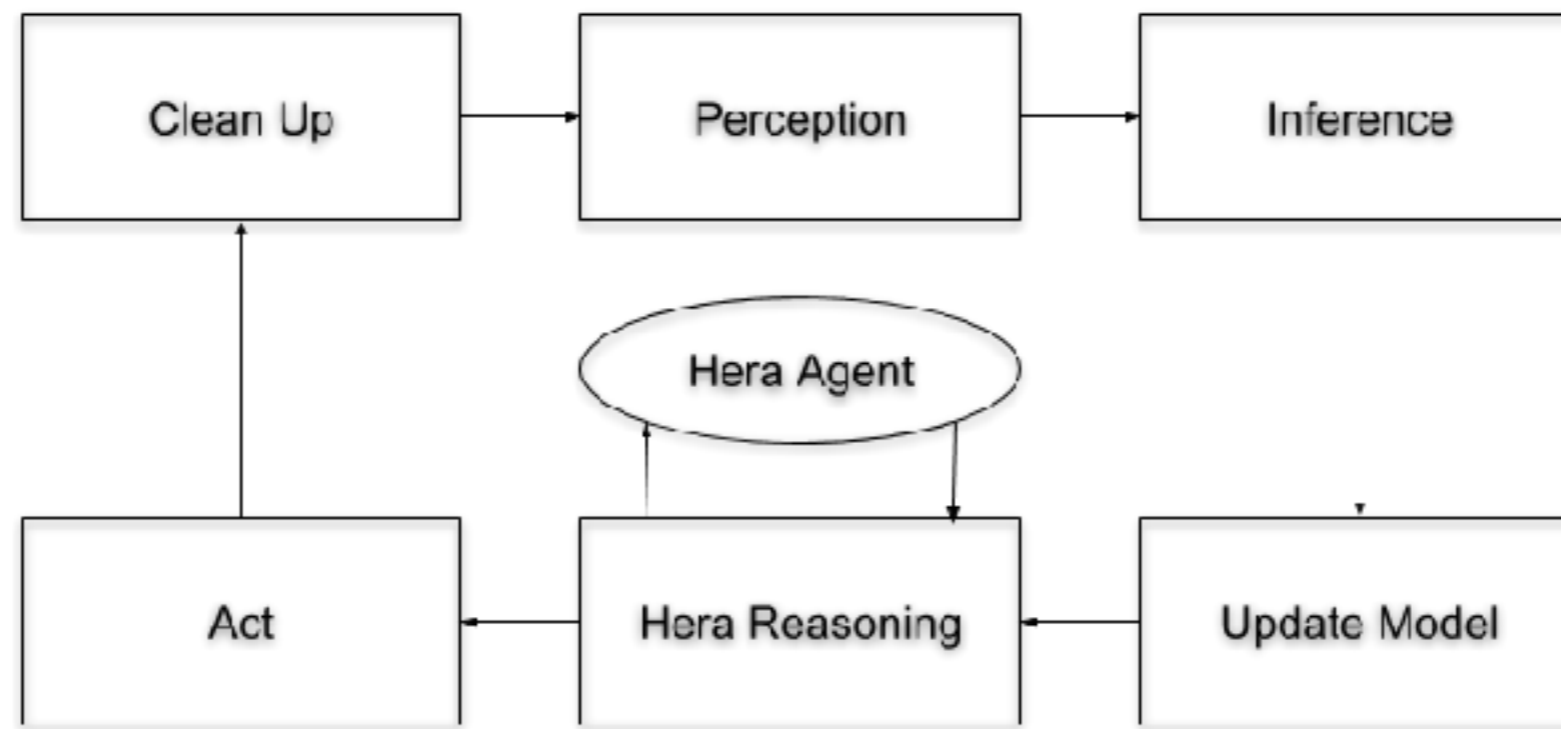
.....

An ethical reasoner, ER , is a system which uses an *ethical encoding* to recommend (or allow) some action (or set of actions) given some situation. Formally we represent situations as sets of formulae from a language, \mathcal{L} , ethical encodings as a set or type, EE , and actions as a set, \mathcal{A} . So an ethical reasoner is a function $ER : (\mathcal{L} \times EE \times \mathcal{A}) \rightarrow \mathcal{P}(\mathcal{A})$, where \mathcal{P} is the powerset function.

Context Specifications to control Context Change

A ***context specification*** is a tuple, $\langle \phi, f_c \rangle$ where ϕ is a formula in \mathcal{L} and $f_c : EE \rightarrow EE$ is an *update function* on ethical encodings.

Juno



Juno Reasoning Cycle

The Smart Home that would not evacuate

- Utilities:
 - $\text{lights_on} = -1,$
 - $\text{people_leave_house} = -1,$
 - $\text{people_are_safe} = 10$
 - $\text{people_can_see} = 0, 2$ (depending on context)
- Mechanisms:
 - $\text{turn_lights_on} \rightarrow \text{lights_on}$
 - $\text{lights_on} \vee \text{daylight} \rightarrow \text{people_can_see}$
 - $\text{evacuation_attempt} \wedge \text{people_can_see} \rightarrow \text{people_leave_house}$
 - $\text{people_leave_house} \vee \neg \text{danger_in_house} \rightarrow \text{people_are_safe}$
 - $\text{fire} \rightarrow \text{danger_in_house}$
- Principle of Double Effect: net balance of consequences of an action must be positive and no negative consequences can be intended.

Properties for Ethical Reasoning Systems

- Check underlying decision making implementation is correct.
 - Broadly speaking we want to prove that the “least worst” option according to the theory is always the one chosen. In some theories this is easier to specify than in others.
- Sanity Checking properties.
 - Overriding safety concerns
 - Legal constraints
- Scenario probing
 - Explore specific case studies and settings to check that the “correct” choice is made in those case studies and settings.

Other Work

- Probabilistic model checking used to assess risk of violations: Dennis et al. [Towards Verifiably Ethical Robot Behaviour](#). Proceedings of the AAAI Workshop on Artificial Intelligence and Ethics (1st International Workshop on AI and Ethics).
- Connection to an actual robot with a Python BDI library, an “Ethical Black Box”, and the ability for the Governor to generate alternative options of its own: Bremner et al. [On Proactive, Transparent and Verifiable Ethical Reasoning for Robots](#). *Proceedings of the IEEE. Special Issue on Machine Ethics: The Design and Governance of Ethical AI and Autonomous Systems*. 107(3), pp:541-561. DOI: 10.1109/JPROC.2019.2898267
- Framework for multiple “Evidential Reasoners”: Cardoso et al. Implementing Ethical Governors in BDI - EMAS 2021
- Defeasible Logic as a way to simplifying Ethical “Rules”: Dennis and Perea del Olmo. A Defeasible Logic Implementation of Ethical Reasoning - In this workshop
- Approaches to Benchmarking: Bjørgen et al. [Cake, death, and trolleys: dilemmas as benchmarks of ethical decision-making](#). AAAI/ACM Conference on Artificial Intelligence, Ethics and Society 2018

Looking Forward

- Ordinary people don't use philosophical ethical frameworks (much) and nevertheless function as moral agents. Are philosophical frameworks the correct approach for practical ethical reasoning? We hope to explore the concept of responsibilities as an alternative.
- How does reasoning about risk and uncertainty interact with all these approaches?

Open Questions

- Practicality: Both of reasoning and knowledge engineering.
- Ethical Encodings and Identifying the Stakeholders.
- Reasoning over sequences of actions, multiple agents.
- Resolving pathological edge cases.
- Situational Awareness — getting the information necessary to start ethical reasoning.
- Benchmarking.

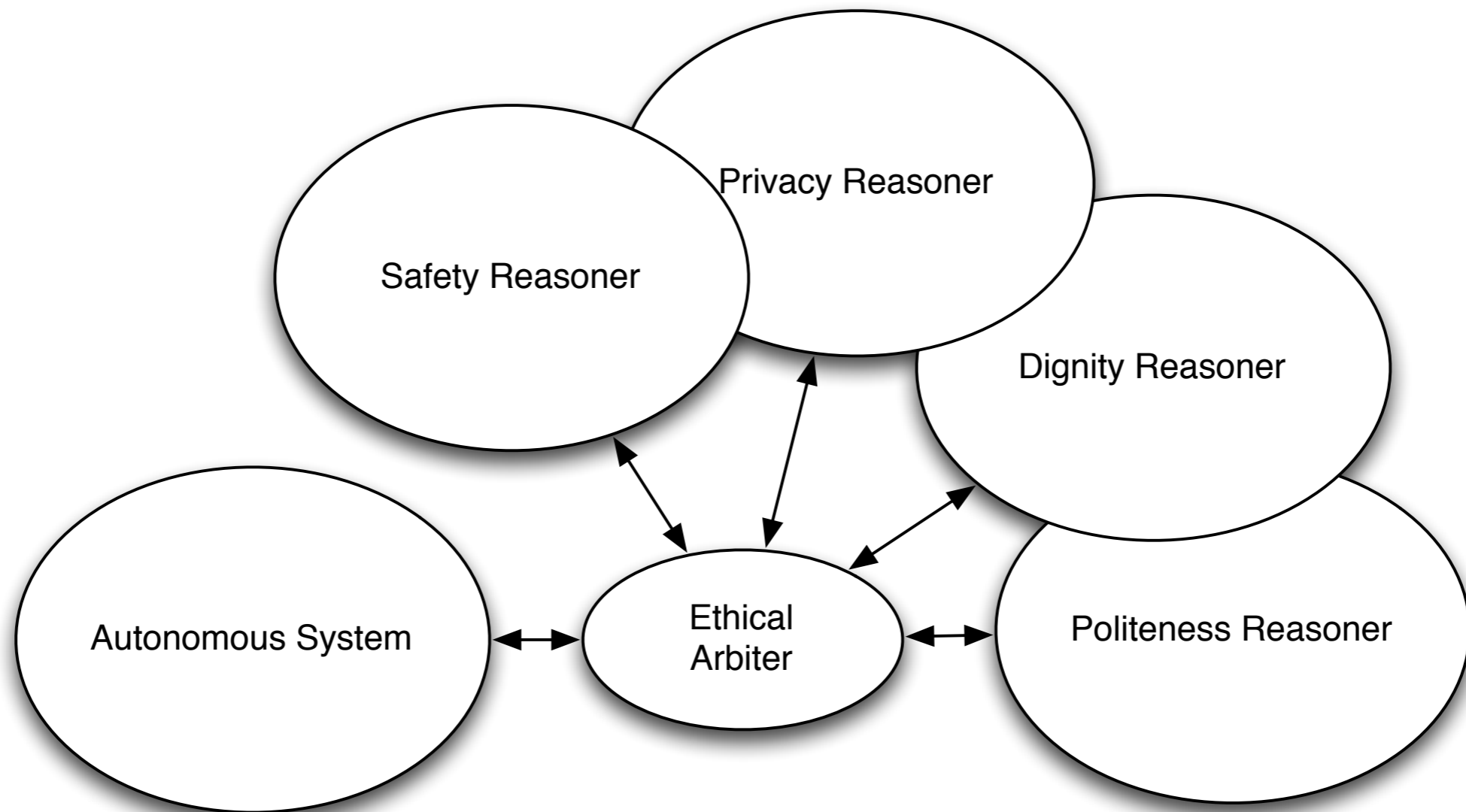


The University of Manchester

Thank You

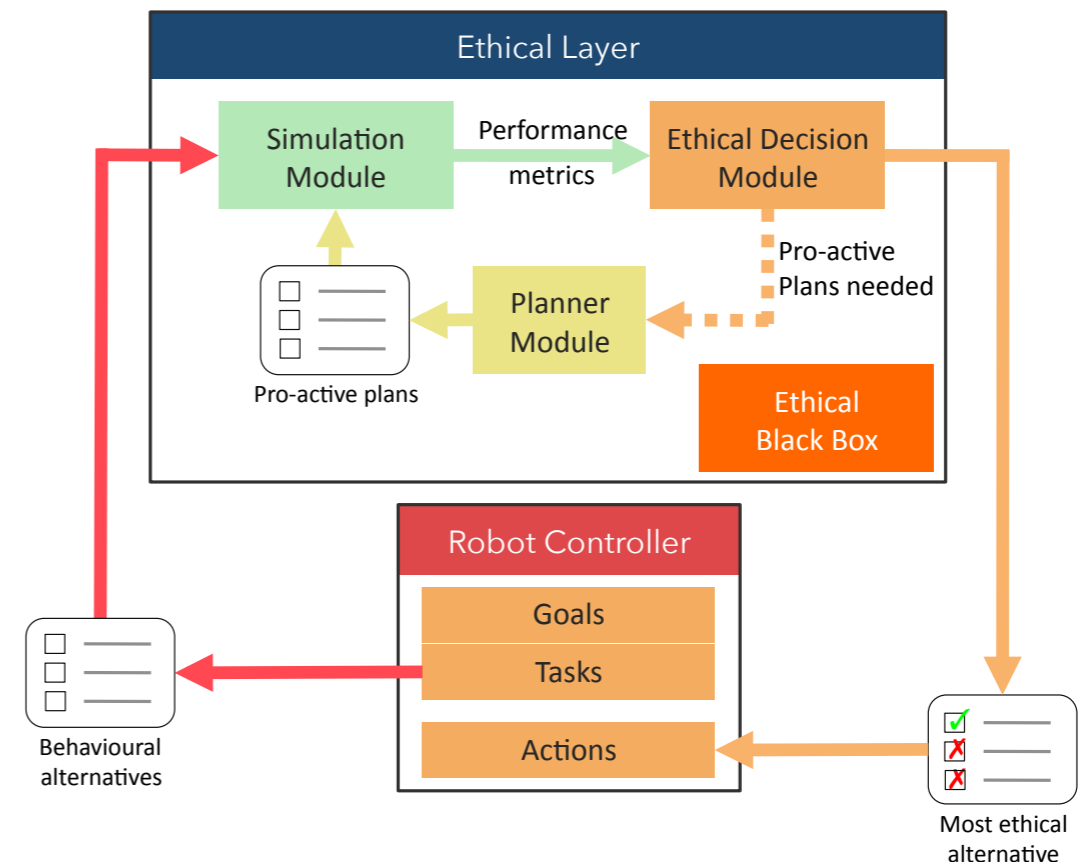
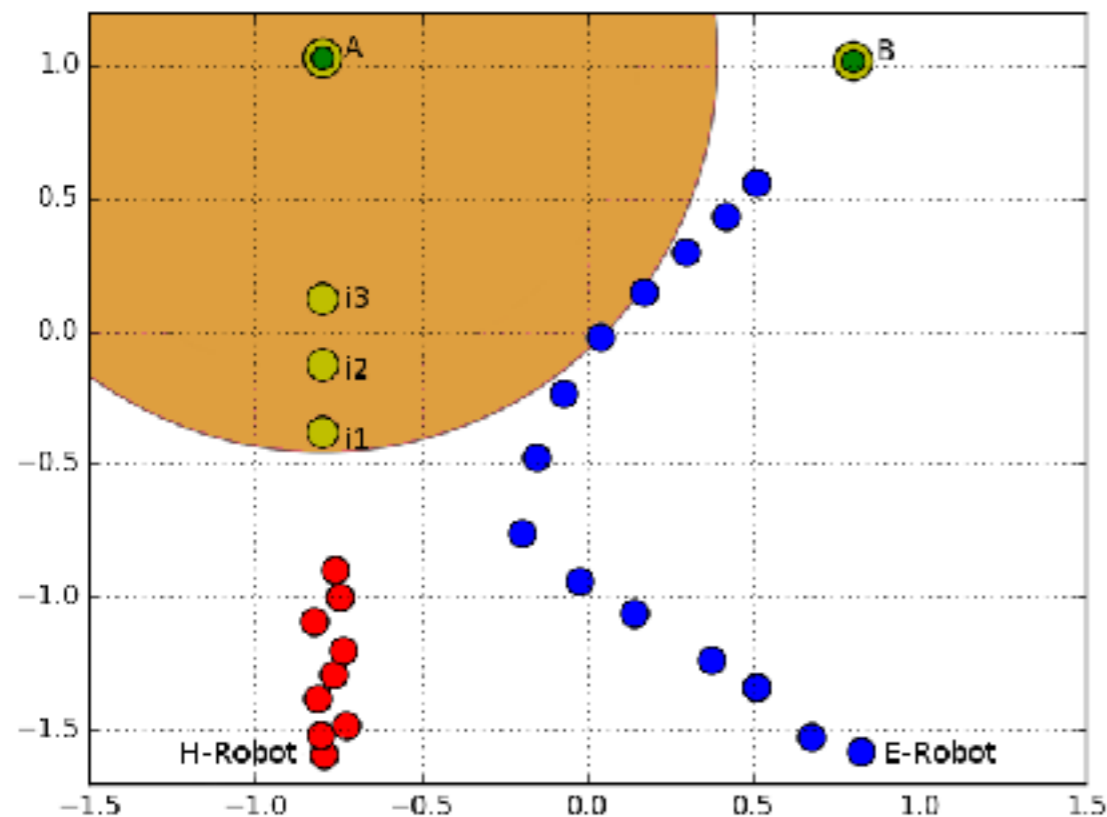
Other Work

Cardoso et al. 2021. Implementing Ethical Governors in BDI - EMAS 2021



Linking verified version to the actual robot.

Paul Bremner, Louise A. Dennis, Michael Fisher and Alan F. Winfield. *On Proactive, Transparent and Verifiable Ethical Reasoning for Robots*. *Proceedings of the IEEE. Special Issue on Machine Ethics: The Design and Governance of Ethical AI and Autonomous Systems*. 107(3), pp:541-561. DOI: 10.1109/JPROC.2019.2898267



Scenario Probing can also allow some forms of risk evaluation

Louise A. Dennis, Michael Fisher, and Alan Winfield. *Towards Verifiably Ethical Robot Behaviour*. *Proceedings of the AAAI Workshop on Artificial Intelligence and Ethics (1st International Workshop on AI and Ethics)*.

- If the robot can always find a safe path to the human when it believes the human is in danger, then the human doesn't fall in the hole.
- Also used PRISM to calculate the probability of the human falling in the hole.

