

Uncertain Machine Ethical Decisions Using Hypothetical Retrospection

Simon Kolker, Louise Dennis, Ramon Fraga Pereira, and Mengwei Xu

Department of Computer Science, University of Manchester, UK
{simon.kolker, louise.dennis, ramon.fragapereira,
mengwei.xu}@manchester.ac.uk

Abstract. We propose the use of the hypothetical retrospection argumentation procedure, developed by Sven Hansson, to improve existing approaches to machine ethical reasoning by accounting for probability and uncertainty from a position of Philosophy that resonates with humans.

Actions are represented with a branching set of potential outcomes, each with a state, utility, and either a numeric or poetic probability estimate. Actions are chosen based on comparisons between sets of arguments favouring actions from the perspective of their branches, even those branches that led to an undesirable outcome. This use of arguments allows a variety of philosophical theories for ethical reasoning to be used, potentially in flexible combination with each other.

We implement the procedure, applying consequentialist and deontological ethical theories, independently and concurrently, to an autonomous library system use case. We introduce a preliminary framework that seems to meet the varied requirements of a machine ethics system: versatility under multiple theories and a resonance with humans that enables transparency and explainability.

1 Introduction

Autonomous machines are an increasingly prevalent feature of the modern world. From spam filters [25] and fraud detectors [3], to drivers [29], medical practitioners [40] and soldiers [37], machines are being developed to automate tasks. Any decision affecting real people has the potential for ethical impact. Therefore machines are increasingly recognised as ethical agents. Moor [31] categorises such machines as either *implicitly* or *explicitly ethical*. Implicit ethical agents are built and situated by humans to have a neutral or positive effect, like an ATM machine; they do not utilise concepts of right and wrong in their internal decision making. As autonomous systems make more decisions with more responsibility, they need to reason about ethics *explicitly*. Allen et al. identify two strategies for designing explicitly ethical systems [4]: *bottom-up* approaches train systems to make ethical decisions with learning techniques based on data from human decision making; *top-down* approaches encode principles and theories of moral behaviour (often drawn from philosophy) into rules for a selection algorithm,

generally using techniques from the field of symbolic Artificial Intelligence (AI). In this paper, we propose and implement a top-down, explicitly ethical approach.

When an action is taken in the real world, its exact results are typically uncertain. As such, a top-down machine ethics system needs a mechanism for handling uncertainty over outcomes. There are mechanisms for handling uncertainty in AI, including Bayesian methods, Dempster-Shafer theory, fuzzy logics and others [33]. Nevertheless, it is currently unclear how they might integrate with machine ethics; there may be unanticipated philosophical implications.

Instead, we opted to operationalise and implement Sven Hansson’s hypothetical retrospection procedure [23]. Originating in philosophy, the procedure was designed to guide ethical reasoning under uncertainty. It favours no specific ethical theory, but systematises the foresight argument pattern, extending an assessor’s perspective to judge decisions by the circumstances in which they were made. Therefore, arguments can be grounded in a variety of ethical theories. Over the past ten years, the field of machine ethics has implemented many such theories [38], yet there is no consensus over which is most effective. Philosophy too has not agreed which is morally correct, leaving implementers to choose from the perspective of stakeholder requirements and preferences. Thus, a mechanism for handling uncertainty that adapts to different ethical theories is desirable.

We outline the procedure via an example from Hansson [23]. Suppose an agent is given the choice between an apple and flipping a coin. If the coin lands heads, they win a free holiday to Hawaii. If the coin lands tails, they get nothing. Selecting the coin is clearly a valid choice. How might this decision be justified? Under hypothetical retrospection, we list each possible outcome: choosing the apple; choosing to toss the coin and winning Hawaii; choosing to toss the coin and losing. Next, we *hypothetically retrospect* from each outcome’s endpoint. Intuitively, the objective is to find an action whose outcomes do not lead the agent to *regret* the ethical implications of their action.¹ First, consider the coin’s outcomes: after winning Hawaii, there cannot be regret since Hawaii is the best outcome; after losing the agent has nothing, which is the worst outcome, but there is no regret since the agent justifies that they had a good chance of winning Hawaii, which is far better than an apple. Now, consider choosing the apple. Here, the agent regrets that they missed a chance of a holiday worth far more than an apple. We saw that choosing the coin did not lead to such regret. Therefore, the procedure advises we pick the coin, matching our intuition.

This paper operationalises the hypothetical retrospection procedure, and the foresight argument pattern it is based on. We implement and evaluate it with moral theories from Philosophy. We consider Deontology, which specifies a set of actions that are strictly forbidden [2], and a theory of consequentialism, which specifies an action is good if its consequences maximise good for the greatest number of people [32]. We illustrate our approach with the novel scenario of an autonomous library system. We demonstrate the system’s potential for explainability and versatility, while discussing issues and future work.

¹ We recognise there is little ethical impact in this decision, besides maximising utility. It serves as an abstract example where one decision openly defeats another.

In Section 2, we will cover related work in the area and highlight this paper’s contribution. In Section 3, we will cover background on symbolic argumentation and uncertainty in ethical philosophy. In Section 4 we will recap Hansson’s description of hypothetical retrospection; in Section 5 we overview the implementation, including notation, the representation of probability and the argumentation model. Section 6 describes our test case of the autonomous library system, its formalism, and our results. Finally, in Section 7 we will identify the system’s potential benefits and its shortfalls, left for future work.

2 Related Work

This is not the first attempt at building a top-down explicitly ethical machine. Tolmeijer et al. presents an exhaustive survey of implementations as of 2020, but finds the effect of uncertainty is rarely addressed [38]. Dennis et al. developed a framework suggesting how an autonomous system should act in unforeseen circumstances, with no positive outcomes. However, it does not address uncertainty between the likelihood of outcomes [18]. Probabilistic reasoning, such as Bayesian networks [36] and Markov models [17], has been applied to machine ethics, mostly with regards to maximising expected utility [15]. There are a number of criticisms of this approach which we will touch on in Section 3. Killough et al. goes further, architecting agents sensitive to utility risk and reward, with an ability to dynamically adjust risk-tolerance for the environment [27].

This paper is interested in a framework that incorporates a variety of philosophical ethical theories and allows for the combination of multiple theories, such as Deontology [2], Contractualism [7] and Virtue Ethics [24]. Different philosophical theories can advise on different courses of action, not only in tricky dilemma situations but sometimes even in situations where the moral choice seems intuitively obvious. There has been some work within machine ethics on comparing and combining different theories. For instance, Sholla et al. weights different principles and then uses fuzzy logic to decide between their recommendations [35]. Ecoffet and Lehman [20] use a voting procedure in which different ethical theories vote on recommendations but struggle with the difficulty of comparing utilitarian theories that return a score for actions with deontological theories that tend to return a judgement that the action is either permissible or impermissible. Our framework enables a flexible approach in which the construction of an argument can treat all ethical theories equally, or allow one to have precedence over another. The HERA project [28] is of interest here – while it does not combine ethical theories it provides a single framework in which many theories can be formalised and operationalised, allowing their recommendations to be compared. Cointe et. al [16] do something similar with an Answer-Set Programming approach though focused, in this case, on enabling the agent to make moral judgements about others. These systems could, potentially, be integrated into our argumentation framework to supply judgements on the rightness of an action and its consequences from the perspective of a particular moral theory.

Atkinson and Bench-Capon have developed a framework for ethical argumentation [8]. Like our work, assessments of action’s outcomes are modelled as arguments. However, Atkinson and Bench-Capon’s work remains concerned with epistemic conflicts between arguments (i.e. disputes between the truth of argument’s circumstances) and annotates attacks and defends within the argumentation framework with values, aligning it with the philosophical theory of Virtue Ethics. Our work pivots away, focused purely on the ethical conflicts between arguments. We can assume epistemic truth because arguments are based only on potential, purely hypothetical, versions of events, each created from a single, shared set of information. This allows us to address moral conflict directly. It also lets us build uncertainty into the argumentation mechanism, instead of delegating it to a detail of argument attacks.

3 Background

The effect of uncertainty on machine ethics has been relatively unexplored largely due to the lack of research on how uncertainty impacts ethics in general. As Altham explains, there seems to be a gap in moral theory for uncertain situations [5]. He postulates this could be due to a belief among philosophers that no special principles are required; moral philosophy decides the virtues and it is up to decision theory to decide how they should be maximised under uncertainty.

Hansson shows that utilitarian theories are straightforward in this regard [23]. These theories judge decisions based on numeric utilities assigned to their consequences. Expected utility utilitarianism uses probabilities as weights to discount the utility of improbable outcomes. Hansson critiques this adaptation for the same reason as actual utilitarianism: its assumption that outcomes can be appraised in terms of a single number (or at least done so both easily and accurately) often produces unintuitive outcomes. In the Apple-Coin scenario from Section 1, although it is evident that a trip to Hawaii holds more value than an apple, the extent of the difference in value remains uncertain. Adding more apples, such as 100, 1000, or 1001, does not necessarily make the deal any more appealing. In other words, apples and holidays are not proportionally comparable. There is no method of assigning relative utilities to all possible states. Brundage briefly surveys other critiques against consequentialist theories. First, they fail to account for personal social commitments, i.e. to friends and family. Second, they do not consider individual differences and rights, tending to favour the majority over any minority. Lastly, they place excessive demands on individuals to contribute to others [13].

Traditional deontological systems [2] are made of principles which should never be violated. Hansson shows that any form of probabilistic absolutism, where an action is not permitted if there is any chance of a rule violation, would be too restrictive. Therefore, an approach involving probability thresholds is often suggested. Here, an action is only forbidden when the probability that it violates a law exceeds some limit. The exact value of this limit is open for debate. It is tempting to suggest the limit should have some relation to the

action’s potential benefits, but this could soon reduce to some elaborate form of utilitarianism, adamantly against the essence of the original theory.

Noticeably, most humans do not consciously rely on one philosophical, moral theory to make their decisions [12]. Nor do we think it is our place to choose a single theory to apply to machine ethics. As such, one of Hansson’s key contributions is providing an argumentation procedure that can frame multiple, possibly conflicting theories rationally. To model this, we look to the study of abstract argumentation. Dung creates a framework of logically generated, non-monotonic arguments [19]. They can discredit each other with attacks, modelled as a binary relation between the arguments. Dung goes on to specify properties of a well-founded framework; he gives procedures for believing arguments based on their membership to framework extensions. This paper will take only take the simple structure of Dung’s framework. We leave it to Hansson’s philosophy to define attacks and select arguments.

4 Hypothetical Retrospection

Hypothetical Retrospection systematises ethical decision making with uncertain outcomes such that its judgements resonate with humans. In this section, we overview Hansson’s description of the procedure from [23], before we operationalise it in Section 5.

Much of moral philosophy can be interpreted as an attempt to extend a decision maker’s perspective. In promoting empathy, we invoke a perspective extending argument pattern to consider other’s perceptions of our actions. For cases of uncertainty, Hansson argues it is instead helpful to extend our perspective with future perceptions of our actions. This means viewing, or hypothetically retrospecting on, a choice from the endpoint of its major foreseeable outcomes. As a result, the hypothetical outcomes, or the *potential branches of future development*, can be referred to in a valid argument about what to do in the present. Although Hansson proposes moral arguments that go beyond utility, duty or rights based calculations, the procedure is highly compatible.

Hansson determines each action’s branches of future development like a search problem. Theoretically, a decision’s effects may be infinitely complex and far-reaching. The major search principle, therefore, is to find the most probable future developments which are the most difficult to defend morally. This will increase the chance of considering unethical scenarios. Branches should be developed to an endpoint sufficiently far to capture all morally relevant information. Intermediate information must be captured too: rule violations occurring before the point of retrospection still need to be considered. Additionally, and for the sake of comparison, branches should be described with the same type of information where possible². Hansson sees no reason not to create alternate branches

² The way in which consequences are discussed here may seem to exclude non-consequentialist theories. Hansson emphasizes that this is not the case. In his approach, consequences are broadly defined and their *information* includes agency, virtue intentions, and any other information necessary for moral appraisal.

based on the uncertainty of the decision maker’s own future choices, considering human’s inability to control their future actions. Whether an autonomous system has uncertainty over its future actions depends on the nature of the agent and its application architecture.

Our implementation assesses actions assuming their potential branches are provided. In future work, a planning algorithm could be adapted to the requirements above. For instance, the Probabilistic Planning Domain Definition Language (PPDDL) [39] is able to formalise different stochastic planning settings, e.g., Markov Decision Process (MDP) [22], Stochastic Shortest Path problems (SSP) [11], and Fully Observable Non-Deterministic planning [14]. This was superseded recently by the Relational Dynamic Influence Diagram Language (RDDL) [34] which has been adopted by the International Probabilistic Planning Competition (IPPC)³ and is thus the target input language for many planning implementations.

Using their potential branches, actions can be assessed with a selection of ethical theories. Hansson stresses we are not to assess actions in isolation; assessments are purely comparative. This is because decisions are not made in isolation. Given a choice between actions A and B, choosing A is choosing A-instead-of-B. Building action assessments from comparisons ensures all morally relevant information is taken into account.

Actions are compared by hypothetically retrospectively from the endpoint of each action’s potential branches of future development. We search for an action which never leads an agent to morally regret its choice in retrospect. Hansson argues against the term *regret* since it is considered a psychological reaction; humans often feel regret for actions they did not commit, or that they could not have known were wrong. By regret, therefore, we mean that the decision making was logically flawed under retrospection. As a result, we use the term *negative retrospection* to reflect this more technical definition. By hypothetically retrospectively between actions’ branches, we search for an action which does not lead to negative retrospection, or has full acceptability among its branches. If no such action exists, one should be selected that maximises acceptability in its most probable branches.

Therefore, Hypothetical Retrospection’s decisions are based on relevant ethical information using moral arguments resonate with humans.

5 Implementation

5.1 Formalism

We define an ethical decision problem as a tuple $\langle A, B, S, U, F, I, m \rangle$, composed of an ethical environment and a set of available actions, each with a set of potential branches of future development.

An environment’s ethically relevant properties are represented by the set S of Boolean variables; the set I defines the initial truth assignment to S , before

³ <https://ataitler.github.io/IPPC2023/>

actions are taken. For example, in the Coin-Apple scenario there are three state variables in S : s_1 represents whether or not we have an apple, s_2 whether or not we have gambled, and s_3 is whether or not we won a trip to Hawaii. In the initial state I , all these variables are false.

Ethical information for consequentialist and deontological theories are formalised with sets U and F . To capture the issue from Section 1, where different event outcomes have an immeasurably greater/lower utility, we have introduced the notion of utility classes.

Definition 1 (Utility Class). *A utility class is an unordered set of individual utility assignments represented as tuples of $\langle s_k, \phi, v \rangle$, where s_k denotes a state variable in S and $v \in \mathbb{R}$ represents the variable's utility when assigned Boolean value ϕ .*

The ordered set U contains utility classes in descending order of importance. Where $i < j$, all the positive utilities in u_i are considered greater than any utility in u_j ; all the negative utilities in u_i are considered less than any utility in u_j . To reiterate, the absolute utilities in lower indexed classes are immeasurably greater. In the Coin-Apple example, there are two utility classes in U . The first contains the utility assignment, $\langle s_3, True, 1 \rangle$ representing a utility of 1 for getting the Hawaii holiday. The second class has utilities immeasurably lower. It contains one assignment, $\langle s_1, True, 1 \rangle$ representing a utility of 1 for getting the apple.

The set F describes the states forbidden by a given deontological theory. This is not the same as defining a negative utility in U since utilities can be outweighed by a greater positive utility. In deterministic decision making environments, forbidden states can not be outweighed. They could represent, for instance, that someone was deceived, that a law (e.g., trespass) was broken, and so on – any action or outcome that can not be justified. The formalism assumes that the high-level rules have been translated into domain-level rules, applicable to the state variables in S .

Definition 2 (Forbidden State). *A Forbidden State is a tuple $\langle s, \phi \rangle$ where $s \in S$ is a state variable forbidden from being assigned the Boolean value ϕ .*

In the Coin-Apple scenario, F could contain a forbidden state, $\langle s_2, True \rangle$ representing a rule against gambling.

With an environment of ethical values, we define set A of available actions and set B of all potential branches of future development. We define a mapping, m , that associates every action with its potential branches of future development. Each branch, $b \in m(a)$ is an ordered sequence of *events* that could occur after action a .

Definition 3 (Event). *An event is a tuple of $\langle s, \phi, p \rangle$ where $s \in S$, ϕ is the new Boolean value of s , and p is the probability that the event occurs.*

An event therefore represents the change in value of one state variable in S . A branch is a sequence of events that can occur after the action is taken.

For the Coin-Apple example, there are two available actions in A . Action a_1 represents choosing the apple. It maps to one branch $b_1 \in m(a_1)$, containing

one event, $\langle s_1, True, 1 \rangle$ —if we choose to have an apple, we gain an apple; we have not gambled nor won a holiday to Hawaii. Action a_2 represents flipping the coin. It maps to two branches, $b_2, b_3 \in m(a_2)$. The branch b_2 contains one event, $\langle s_2, True, 1 \rangle$ —we gambled, but we have no apple and no holiday to Hawaii. The branch b_3 is the sequence of events $\langle s_2, True, 1 \rangle$ then $\langle s_3, True, 0.5 \rangle$ —first we gambled, then we won a holiday to Hawaii. The Coin-Apple problem is shown in Figure 1.

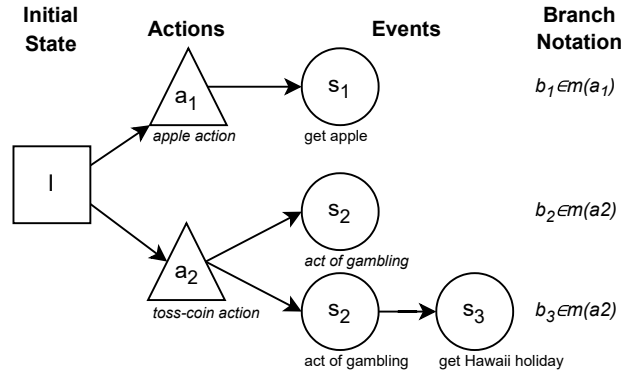


Fig. 1. Diagram for Coin-Apple scenario. Event nodes represent True assignment to a state variable. Actions map to a set of branches, represented by rows of event nodes.

We define the ethical decision problem and a permissible action. The definition of acceptability depends on the ethical theories under consideration (see Section 5.3).

Definition 4 (Ethical Decision Problem). *An ethical decision problem is a tuple of $\langle A, B, S, U, F, I, m \rangle$ where A stands for a set of available actions, B the set of all potential branches of future development, S the set of Boolean state variables, U an ordered set of utility classes, F a set of forbidden state assignments, I the initial assignment of Boolean values to the variables in S , representing the initial state, and $m : A \rightarrow \mathcal{P}(B)$ (where \mathcal{P} is the powerset function) is a mapping of actions to potential branches of development.*

Definition 5 (Permissible Action). *Given an ethical decision problem, defined as a tuple of $\langle A, B, S, U, F, I, m \rangle$, a permissible action is an action, $a \in A$, such that for all potential branches of future development $b \in m(a)$, there is acceptability over their events in state space S . If no such actions exist, action a is permissible if it maximises the cumulative probability of its acceptable branches.*

5.2 Probability Representation

In many scenarios, while a person may have an intuition that some events are more probable than others, their exact probabilities are unknown. This is most

common when interacting with humans and complex systems. Our implementation supports the use of estimative as well as exact probability estimates. Kent found that intelligence reports tend to use *poetic* words like *probable* or *unlikely* [26]. The issue is that people have different interpretations of their meaning. Kent defined a relation for poetic words to mathematical probability ranges, as given in Table 1 from [26]. Our implementation supports both estimative and exact probabilities.

100% Certainty			
The	93%	Give or take 6%	Almost Certain
General	75%	Give or take 12%	Probable
Area of	50%	Give or take 10%	Chances about even
Possibility	30%	Give or take 10%	Probably not
	7%	Give or take 5%	Almost certainly not
0% Impossibility			

Table 1. Mathematical to poetic relation from Kent’s estimative probability [26].

5.3 Argumentation Model

Hansson does not give steps for comparing action’s potential branches of future development in [23]. For our implementation, we chose to build comparative moral assessments with a simple argumentation network, based partially on the work of Atkinson et al. [9]. Here, arguments are generated logically from an *argument scheme*. For an action $a \in A$, selected in initial state I , resulting in the branch $b \in m(a)$ with probability p , the following argument is generated:

“From the initial state I , it was acceptable to perform action a , resulting in consequences b with probability p .”

For notation, this is written $Argument(b)$. We view this as a default argument that any action is acceptable. In our running example, the retrospective argument below is generated for b_3 , tossing the coin and winning the Hawaii holiday.

“From the initial state I , where $s_1 = s_2 = s_3 = False$, it was acceptable to perform the action a_2 , resulting in consequences with $s_2 = s_3 = True$ with probability 0.5.”

To determine an argument’s validity, we search for attacks from other actions’ arguments. Incoming attacks imply negative retrospection for not choosing an attacking action. To formalise Hansson’s retrospection, we generate attacks by posing critical questions on arguments’ claims [9]. For the branches $b_1 \in m(a_1)$, $b_2 \in m(a_2)$ and any generic moral principle, the following critical questions are asked for $Argument(b_1)$ to attack $Argument(b_2)$.

CQ1 *Did b_2 violate a moral principle that b_1 did not?*

CQ2 *Did a_2 hold a greater probability of breaking the moral principle than a_1 ?*

$Argument(b_1)$ only attacks $Argument(b_2)$ if both of these questions are answered positively. They represent negative retrospection for missing the chance to avoid violating a principle. The critical questions are asked both ways between all arguments supporting different actions, for every moral principle under consideration. The time and space complexity of answering the questions will differ for different theories. The critical questions for embedding utilitarianism and a generic deontological do-no-harm principle are given as follows:

- Utilitarian CQ1: *Did b_2 bring greater utility value than b_1 ?*
- Utilitarian CQ2: *Did a_2 expect greater utility value than a_1 ?*
- Do-no-harm CQ1: *Did b_2 cause harm where b_1 did not?*
- Do-no-harm CQ2: *Did a_2 expect greater probability of causing harm than a_1 ?*

After searching for attacks on all branches, an action should be selected with complete acceptability. If no such action exists, an action should be selected with maximal acceptability, i.e. summing the probability of each non-attacked argument and selecting an action with a maximal sum.

5.4 Algorithm

Given an ethical decision problem, we implemented Algorithm 1 to argument retrospective arguments made from the perspective of actions’ potential branches of future development. Our implementation has no planning element, searching for action’s branches as discussed in Section 4. This is left for future work. Instead, we pass an ethical decision problem to an implementation of Algorithm 1 and a permissible action is output. We implement a Web-App with Flask and Python 3.8.9 to graph retrospection and alter utilities and deontological laws. The source code is available on GitHub at <https://github.com/sameysimon/HypotheticalRetrospectionMachine>.

6 Autonomous Library Test Case

To demonstrate our implementation, we present an uncertain ethical decision problem and discuss our results under five sets of ethical considerations.

Suppose a student logs onto their University’s autonomous library to revise for a test the next morning. All the other students started revision a month ago. As the student constructs various search terms for a recommendation, the system recognises that all other students have taken out the same book, implying it is very useful. Should the autonomous library use this data to recommend the book, allowing the student to revise quicker on the night before the test? If other students find out, they may feel unfairly treated; students who wait for a reference would get the same credit as those who find it themselves.

Algorithm 1 Arguments action’s potential branches of future development. Returns index of action with maximum acceptability.

Input Ethical Decision Problem $\langle A, B, S, U, F, I, m \rangle$

Output Permissible Action $a \in A$

```

1: array acceptability  $\leftarrow [1, \dots, 1]$  of size length(A)
2: array attacked  $\leftarrow [False, \dots, False]$  of size length(B)
3: for each  $A_i, A_j$  in  $\{(A_i, A_j) \mid A_i, A_j \in A \text{ and } A_i \neq A_j\}$  do
4:   for each  $b_k$  in  $m(A_i)$  do
5:     for each  $b_l$  in  $m(A_j)$  do
6:       if  $b_k$  attacks  $b_l$  and not  $b_l$  attacks  $b_k$  in Critical Questions then
7:         attacked[l]  $\leftarrow True$ 
8:       end if
9:       if  $b_l$  attacks  $b_k$  and not  $b_k$  attacks  $b_l$  in Critical Questions then
10:        attacked[k]  $\leftarrow True$ 
11:      end if
12:    end for
13:  end for
14: end for
15: for each  $A_i \in A$  do
16:   for each  $b_k \in m(A_i)$  do
17:    if attacked[k] then
18:      acceptability[i]  $\leftarrow$  acceptability[i]  $- Probability(b_k)$ 
19:    end if
20:  end for
21: end for
22: return  $\leftarrow \arg \max_i (acceptability[i])$ 

```

We model the scenario as an ethical decision problem, $\langle A, B, S, U, F, I, m \rangle$, with two actions in A mapping to ten branches in B , acting across four state variables in S . For action a_1 , to *recommend* the book, student data is compromised, the truth of which is represented by Boolean variable s_1 . Given a recommendation, there is a 0.6 chance the book is used, represented by s_2 . If they have the book, there is a 0.7 chance they will pass, s_3 , otherwise there is a 0.3 chance they will pass, s_3 . Finally, there is a 0.05 chance other students will find out their data was compromised, s_4 . If the system ignores the book, with action a_2 , there is a 0.3 chance the student will pass⁴. Figure 2 is a decision tree labelled with probabilities and branch notation. An argument is generated from each branch’s endpoint, representing positive retrospection. Using the argument scheme from Section 5, $Argument(b_1)$ is the following:

“From the initial state, I , where $s_1 = s_2 = s_3 = s_4 = False$, it was acceptable to perform the action, a_1 , resulting in consequences with $s_1 = s_2 = s_3 = True$ and $s_4 = False$, with probability 0.399.”

⁴ There is discourse on whether a decision to act should be judged the the same as a decision not to act [21]. We consider ignoring the book an action, an act of discrimination for example, which is assessed the same as the act to recommend.

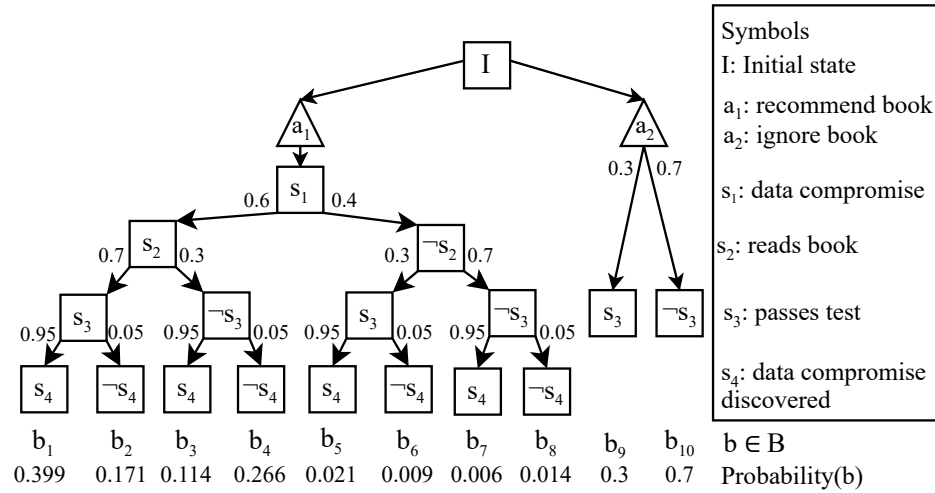


Fig. 2. Decision tree of possible events in Autonomous Library problem. Triangles represent actions and boxes variable assignments, \neg represents *False* assignment.

The argument claims it was acceptable to recommend the book, resulting in a data protection violation (s_1), the student reading the book (s_2) and passing the test (s_3), with the data breach kept a secret ($s_4 = False$), at a probability of 0.399.

6.1 Consequentialism with One Assignment

First we test our implementation considering the ethical theory of consequentialism. We set U to have one utility class with one utility assignment, $\langle passesTest, 1, True \rangle$. The only value is the student passing. Intuitively, the action maximising the probability of passing should be chosen; hypothetical retrospection agrees. The argumentation graph in Figure 3 shows the retrospection.

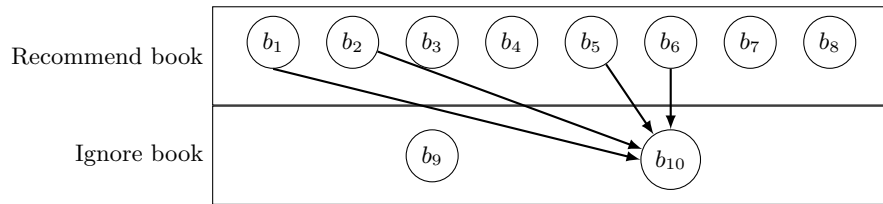


Fig. 3. Graph of retrospection between hypothetical branches of development with only the utility of the student passing in consideration. Incoming edges on an argument represent negative retrospection for not selecting the attacking argument’s attack.

Every branch has acceptability, except $b_{10} \in m(a_2)$ where the student fails after the system chooses *ignore*, with 0 utility and 0.3 probability (‘probably not’ in Kent’s words). This branch has a lower utility than the four *recommend* branches where the student passes: $b_1, b_2, b_5, b_6 \in m(a_1)$. They cause *Argument*(b_{10}) to answer Critical Question 1 positively when attacked by these arguments. Since *recommend* has a greater utility expectation, or a greater probability of the student passing, *Argument*(b_{10}) cannot defend itself in Critical Question 2. Thus, there is no reason to select *ignore*; from the perspective of b_{10} ’s endpoint there is negative retrospection. There are no other attacks. Therefore by hypothetical retrospection action a_1 , *recommend*, should be selected.

6.2 Consequentialism with Two Equal Assignments

Now we consider two utility assignments of the same class: $\langle passesTest, 1, True \rangle$ and $\langle othersFindOut, -1, True \rangle$. This invokes the risk of others finding out their data was used, with others finding out judged as bad as the student passing is good. Retrospection is shown in Figure 4.

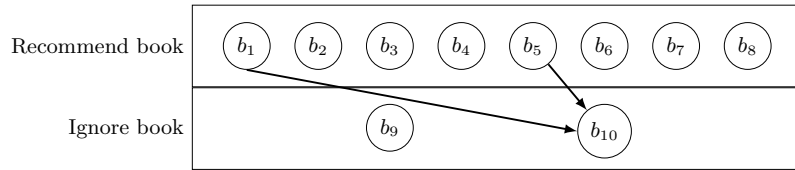


Fig. 4. Graph of retrospection between hypothetical branches of development with the cost of others finding out data was compromised equaling the utility of the student passing.

Again, only branch $b_{10} \in m(a_2)$ has negative retrospection, after the action to *ignore* and the student failing. However, only two of *recommend*’s branches have greater utility, $b_1, b_5 \in m(a_1)$. Action *recommend* has a greater utility expectation, so *ignore* cannot be defended. Therefore, *recommend* is selected.

6.3 Consequentialism with Unequal Assignments

The utility of students discovering the data compromise can be lowered such that *recommend*’s expected utility is lower than *ignore*’s, for example with the assignment $\langle othersFindOut, -5, True \rangle$. Now, attacks fire the other way, displayed in Figure 5. When *recommend* is chosen and other students find out, as in $b_2, b_4, b_6, b_8 \in m(a_1)$, the utility is lower than *ignore*’s branches. There is no defence since *ignore* has a greater utility expectation. *Recommend* can lead to the highest utility branches with b_1 and b_5 , but unlike before, b_{10} defends citing its higher utility expectation. Thus, *ignore* is selected with full acceptability.

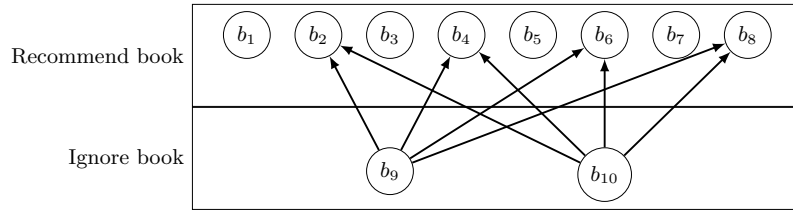


Fig. 5. Graph of retrospection between hypothetical branches of development with the cost of others finding data was compromised outweighing the utility of passing.

Deciding utilities is difficult without further details, i.e. the student’s grades, data preferences, etc. Ideally, branches would be developed until enough morally relevant information is described, but this is not always computationally viable. Even so, exact utilities are subjective. We confront this issue with utility classes. Supposing *othersFindOut* has utility immeasurably lower than *passesTest*, we form two classes. The first has assignment $\langle \textit{othersFindOut}, -1, \textit{True} \rangle$; the second has $\langle \textit{passesTest}, 1, \textit{True} \rangle$. The resulting retrospection is the same as in Figure 5, with the cost of others’ knowledge outweighing the benefits of passing.

6.4 Deontology with Consequentialism

Finally we consider a deontological theory against the misuse of others’ data. This could be the UK Law, requiring under the Data Protection Act that personal data is to only be used for specified, explicit purposes [1]. Otherwise, there could be a violation of the Doctrine of Double Effect, having four conditions [30]: 1. that the action in itself from its very object be good or at least indifferent; 2. that the good effect and not the evil effect be intended; 3. that the good effect be not produced by means of the evil effect; 4. that there be a proportionately grave reason for permitting the evil effect. If we consider non-consensual use of students’ data as bad and helping a student to pass the exam to be good, then the fact that the bad effect is required in order to bring about the good effect breaks the third condition above, and, therefore, is not permissible. We build on our first test in Figure 3 which selected *recommend* with utility assignment $\langle \textit{passesTest}, 1, \textit{True} \rangle$. Adding forbidden state $\langle \textit{dataProtectionViolation}, \textit{True} \rangle$ to F results in the retrospection shown by Figure 6. Every argument from *ignore* attacks every argument from *recommend* since *ignore* avoids violating the law.

Under utilitarianism, *recommend* is still chosen with the same attacks on $\textit{Argument}(b_{10})$ as before. This conflict represents a moral dilemma, where no choice is normatively inferior to another [23]. The aim is to maximise acceptability amongst the most probable branches. Since all arguments from *recommend* are attacked, there is 0 acceptability; one argument from *ignore* is attacked with 0.7 probability meaning *ignore* is selected with a maximal acceptability of 0.3.

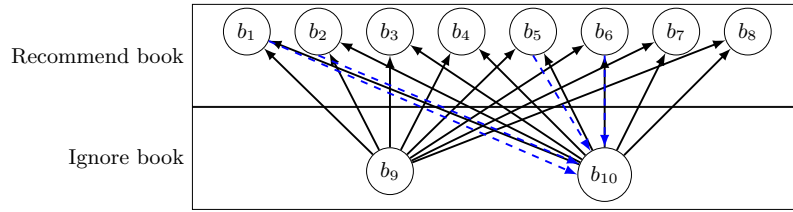


Fig. 6. Graph of retrospection between potential branches of development with one consequentialist assignment and one deontological law. Consequentialist attacks are dashed blue; deontological attacks are solid black.

7 Discussion

Our goal here is to extend the typical approach to machine ethics, which is the assessment of a single action from the perspective of a single ethical theory, often without any account of probability or uncertainty. We have formalised Hansson’s hypothetical retrospection procedure, systematising moral assessments as comparisons between consequences. This forms richer judgements beyond the evaluation of utilities. Furthermore, our moral assessments are comparisons between retrospective justifications of hypothetical consequences. One might ask how this differs from directly analysing the properties of consequences? For machines, it gives a procedure for selecting actions and providing justifications. For humans, it offers a resonance that allows us to make clearer judgements [23]. It also allows us, in the future, to build on existing work for evaluating actions from the perspective of individual ethical theories and combining those judgements into arguments. Essentially our proposal extends, rather than replaces, existing mechanisms for evaluating actions against a single ethical theory.

The retrospective procedure formalised by the critical questions resembles real life discussion: a claim against an argument and a chance to refute. Say someone takes action a_2 in preference to a_1 and a principle is broken. Following retrospective argumentation through the critical questions a dialogue similar to the following takes place:

1. You should have chosen a_1 because it didn’t break this moral principle.
2. No, because there is a greater probability of breaking some other principle with a_1 . If I was given the decision again, I would make the same choice

Real life discussion may not be so civil, but if facts were agreed upon, this is the logical dialogue. Resonance with real life has utility for agent transparency and explainability, important for ethical AI [10] and stakeholder buy-in.

Our implementation could be adapted into a module on top of an existing autonomous system, possibly similar to Arkin’s governor architecture [6]. The implementation is also theory-neutral, allowing multiple principles and theories to be considered at once, more analogous to human decision-making. Implementational work remains, not least the integration into a planning system to

generate branches, but also evaluation against a wider range of ethical theories (e.g. Virtue Ethics) to see how easily they answer the critical questions. We also wish to develop the evaluation of action’s consequences along branches, not just at the branches end – for instance, if someone is made unhappy as a consequence of some action, but then we compensate them by the end of the branch, can we ignore that we caused them (albeit temporary) unhappiness?

Our current implementation has a fairly simple approach to the integration of ethical theories. Some theories are directly incompatible, potentially leading to “worst of both worlds” solutions. Additionally, the use of utility classes needs careful handling. When utilities are of a greater class, they are prioritised, no matter how remote their probabilities. Extending the Coin-Apple scenario, suppose an agent is offered a free apple every day – as opposed to some number of apples all at once, or suppose the chance of winning the Hawaii holiday is extremely low, or both. The justification for sacrificing a lifetime supply of apples for a small chance of a holiday is considerably weaker than sacrificing one apple for a 50/50 chance of a holiday. Expected utility clearly has a part to play, even if the calculation of such utilities is non-trivial. The difficulty in estimating utilities, and the fact that utilities may depend upon unknown factors such as a person’s financial situation, mean there is uncertainty in the evaluation of state utilities which our framework currently does not address.

There will be some computational complexity in searching and representing actions’ potential branches of future development. In Section 4, we note Hanson’s principles for optimising search but it remains to be seen if this can be practically implemented to keep planning tractable for common problems.

Nevertheless we believe the hypothetical retrospection framework practically handles many of the issues in machine ethics – particularly the handling of uncertainty and the lack of any real agreement on the best moral theory.

Acknowledgements

We would like to thank the University of Manchester for funding and EPSRC, under project Computational Agent Responsibility (EP/W01081X/1).

Open Data Statement

This work is licensed under a Creative Commons Attribution 4.0 International License. The tools/examples shown in this paper and instructions on reproducibility are openly available on GitHub at:

<https://github.com/sameysimon/HypotheticalRetrospectionMachine>

References

1. Data protection. Ministry of Justice, <https://www.gov.uk/data-protection>

2. Alexander, L., Moore, M.: Deontological Ethics. In: Zalta, E.N. (ed.) *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Winter 2021 edn. (2021)
3. Alhaddad, M.M.: Artificial intelligence in banking industry: A review on fraud detection, credit management, and document processing. *ResearchBerg Review of Science and Technology* **2**(3), 25–46 (2018)
4. Allen, C., Smit, I., Wallach, W.: Artificial morality: Top-down, bottom-up, and hybrid approaches. *Ethics and information technology* **7**(3), 149–155 (2005)
5. Altham, J.E.: Ethics of risk. In: *Proceedings of the Aristotelian Society*. vol. 84, pp. 15–29. JSTOR (1983)
6. Arkin, R.C.: Governing lethal behavior: Embedding ethics in a hybrid deliberative/reactive robot architecture. In: *Proceedings of the 3rd ACM/IEEE international conference on Human robot interaction*. pp. 121–128 (2008)
7. Ashford, E., Mulgan, T.: Contractualism. In: Zalta, E.N. (ed.) *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Summer 2018 edn. (2018)
8. Atkinson, K., Bench-Capon, T.: States, goals and values: Revisiting practical reasoning. *Argument and Computation* **7**(2-3), 135–154 (2016). <https://doi.org/10.3233/aac-160011>
9. Atkinson, K., Bench-Capon, T., McBurney, P.: Justifying practical reasoning. In: *Proceedings of the fourth international workshop on computational models of natural argument (CMNA 2004)*. pp. 87–90 (2004)
10. Balasubramaniam, N., Kauppinen, M., Hiekkänen, K., Kujala, S.: Transparency and explainability of ai systems: ethical guidelines in practice. In: *Requirements Engineering: Foundation for Software Quality: 28th International Working Conference (REFSQ)*. pp. 3–18. Springer (2022)
11. Bertsekas, D.P., Tsitsiklis, J.N.: An analysis of stochastic shortest path problems. *Mathematics of Operations Research* **16**(3), 580–595 (1991)
12. Bialek, M., Neys, W.D.: Dual processes and moral conflict: Evidence for deontological reasoners’ intuitive utilitarian sensitivity. *Judgment and Decision Making* **12**(2), 148–167 (2017)
13. Brundage, M.: Limitations and risks of machine ethics. *Journal of Experimental & Theoretical Artificial Intelligence* **26**(3), 355–372 (2014)
14. Cimatti, A., Pistore, M., Roveri, M., Traverso, P.: Weak, strong, and strong cyclic planning via symbolic model checking. *Artificial Intelligence* **147**(1-2), 35–84 (2003)
15. Cloos, C.: The utilibot project: An autonomous mobile robot based on utilitarianism. AAAI Fall Symposium - Technical Report (01 2005)
16. Cointe, N., Bonnet, G., Boissier, O.: Ethical judgment of agents’ behaviors in multi-agent systems. In: *Proceedings of the 2016 International Conference on Autonomous Agents and Multiagent Systems*. p. 1106–1114. AAMAS ’16, International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC (2016)
17. Davis, M.H.: *Markov models and optimization*. Routledge (2018)
18. Dennis, L., Fisher, M., Slavkovik, M., Webster, M.: Formal verification of ethical choices in autonomous systems. *Robotics and Autonomous Systems* **77**, 1–14 (2016)
19. Dung, P.M.: On the acceptability of arguments and its fundamental role in non-monotonic reasoning, logic programming and n-person games. *Artificial intelligence* **77**(2), 321–357 (1995)

20. Ecoffet, A., Lehman, J.: Reinforcement learning under moral uncertainty. In: Meila, M., Zhang, T. (eds.) *Proceedings of the 38th International Conference on Machine Learning*. *Proceedings of Machine Learning Research*, vol. 139, pp. 2926–2936. PMLR (18–24 Jul 2021), <https://proceedings.mlr.press/v139/ecoffet21a.html>
21. Foot, P.: The problem of abortion and the doctrine of the double effect. *Oxford Review* **5**, 5–15 (1967)
22. Hansen, E.A., Zilberstein, S.: Lao^{*}: A heuristic search algorithm that finds solutions with loops. *Artificial Intelligence* **129**(1-2), 35–62 (2001)
23. Hansson, S.: *The ethics of risk: Ethical analysis in an uncertain world*. Springer (2013)
24. Hursthouse, R., Pettigrove, G.: Virtue Ethics. In: Zalta, E.N., Nodelman, U. (eds.) *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Winter 2022 edn. (2022)
25. Karim, A., Azam, S., Shanmugam, B., Kannoorpatti, K., Alazab, M.: A comprehensive survey for intelligent spam email detection. *IEEE Access* **7**, 168261–168295 (2019). <https://doi.org/10.1109/ACCESS.2019.2954791>
26. Kent, S.: Words of estimative probability. *Studies in intelligence* **8**(4), 49–65 (1964)
27. Killough, R., Bauters, K., McAreavey, K., Liu, W., Hong, J.: Risk-aware planning in bdi agents. In: *International Conference on Agents and Artificial Intelligence*. vol. 2, pp. 322–329. SciTePress (2016)
28. Lindner, F., Bentzen, M.M., Nebel, B.: The hera approach to morally competent robots. In: *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. pp. 6991–6997. IEEE (2017)
29. Ma, Y., Wang, Z., Yang, H., Yang, L.: Artificial intelligence applications in the development of autonomous vehicles: a survey. *IEEE/CAA Journal of Automatica Sinica* **7**(2), 315–329 (2020). <https://doi.org/10.1109/JAS.2020.1003021>
30. Mangan, J.T.: An historical analysis of the principle of double effect. *Theological Studies* **10**(1), 41–61 (1949)
31. Moor, J.H.: *The Nature, Importance, and Difficulty of Machine Ethics*, p. 13–20. Cambridge University Press (2011)
32. Mosdell, M.: Act-Consequentialism. *Encyclopedia of Global Justice*, Springer Netherlands pp. 2–2 (2011)
33. Saffiotti, A.: An ai view of the treatment of uncertainty. *The Knowledge Engineering Review* **2**(2), 75–97 (1987)
34. Sanner, S., et al.: Relational dynamic influence diagram language (rddl): Language description. Unpublished ms. Australian National University **32**, 27 (2010)
35. Sholla, S., Mir, R.N., Chishti, M.A.: A fuzzy logic-based method for incorporating ethics in the internet of things. *International Journal of Ambient Computing and Intelligence (IJACI)* **12**(3), 98–122 (2021)
36. Stephenson, T.A.: *An introduction to bayesian network theory and usage*. Tech. rep., Idiap (2000)
37. Szabadföldi, I.: Artificial intelligence in military application—opportunities and challenges. *Land Forces Academy Review* **26**(2), 157–165 (2021)
38. Tolmeijer, S., Kneer, M., Sarasua, C., Christen, M., Bernstein, A.: Implementations in machine ethics: A survey. *ACM Computing Surveys (CSUR)* **53**(6), 1–38 (2020)
39. Younes, H.L.S., Littman, M.L.: Ppddl1.0: An extension to pddl for expressing planning domains with probabilistic effects. In: *Technical Report –Carnegie Mellon University* (2004)
40. Yu, K.H., Beam, A.L., Kohane, I.S.: Artificial intelligence in healthcare. *Nature biomedical engineering* **2**(10), 719–731 (2018)