

Moral Permissibility of Action Plans

Felix Lindner and Robert Mattmüller and Bernhard Nebel

University of Freiburg, Germany
{lindner, mattmuel, nebel}@informatik.uni-freiburg.de

Abstract

Research in classical planning so far was mainly concerned with generating a satisficing or an optimal plan. However, if such systems are used to make decisions that are relevant to humans, one should also consider the ethical consequences that generated plans can have. We address this challenge by analyzing in how far it is possible to generalize existing approaches of machine ethics to automatic planning systems. Traditionally, ethical principles are formulated in an action-based manner, allowing to judge the execution of one action. We show how such a judgment can be generalized to plans. Further, we study the complexity of making ethical judgment about plans.

Introduction

With the advent of autonomous machines that drive on the streets or act as household robots, it has been argued that we need to add an ethical dimension to such machines leading to the development of the research area *machine ethics* (Anderson, Anderson, and Armen 2005; Anderson and Anderson 2011). One important question is how we can align the behavior of autonomous machines with the moral judgment of humans. In this context, most often the question is whether a particular action is morally obligatory, permissible or impermissible, given a particular ethical principle (Driver 2006). Judging one action is, of course, important. However, automated planning systems (Ghallab, Nau, and Traverso 2016) are faced with the problem of making a huge number of decisions about including actions into a plan. And it does not necessarily make sense to analyze the ethical contents of each such decision in isolation, but to take an ethical perspective on an entire plan (and perhaps alternative plans reaching the same goal). As an example, consider utilitarian reasoning: if every action in a plan were judged in isolation, one would not be allowed to perform an action that temporarily decreases the utility, even if this action is a necessary prerequisite for later earning a lot of utility in a globally optimal final reachable state. Judging a plan as a whole allows considering this early investment for the sake of a later benefit as permissible from a utilitarian perspective.

In this paper, we will address these problems by analyzing three problems. First, we will look at what kind of additional information we need in order to be able to make moral judgments in the context of different ethical theories. Secondly,

we will propose methods to judge the ethical acceptability of a plan. We will test the proposed notions using examples from the literature on moral dilemmas.

We will not limit ourselves to one particular ethical principle, but will consider a number of different principles that have the potential to be treated computationally, similar to the HERA (Lindner, Bentzen, and Nebel 2017) approach.

Based on the work we present in this paper, a planner will not only be able to come up with plans, but also to ethically judge those plans and compare them to alternative plans that it produces or that may be suggested by a human user or another automated planner. This will enable an ethically competent planning system to discuss and *explain* why a certain—morally superior—plan should be preferred over another—morally inferior—one. Such an explanation can explicitly refer to *which* ethical principles are violated and *how* they are violated.

The remainder of the paper is structured as follows. In the next section, we introduce different ethical principles that have been discussed in the literature. Then, the planning formalism we will use throughout the paper will be specified. This is basically a propositional planning formalism extended by variables with non-binary domains, exogenous events, and moral valuations of actions and consequences. We then formalize the notions of causation and means to an end in the framework of our planning formalism. Based on that, we can then formalize different ethical principles, which we will use to analyze the computational problem of ethically validating a given plan. Finally, we sketch related work and conclude.

Ethical Principles

In moral philosophy, various ethical principles are considered. Ethical principles are descriptions of abstract rules that can be used to determine the moral permissibility of concrete courses of actions. In this section, we introduce ethical principles which embrace different views on how to assess moral permissibility of actions based on the actions' consequences: *Utilitarianism*, *deontology*, *do-no-harm principle*, and the *principle of double effect*.

The *utilitarian principle* focuses on consequences of actions. It says that an agent ought to perform the action amongst the available alternatives with the overall maximal utility. We adopt an act-utilitarian interpretation which does

not distinguish between doing and allowing, i.e. the causal structure of the situation is not taken into account. Thus the action which the agent ought to perform is the one which leads to the best possible situation, i.e. the highest utility, regardless of what the agent causes and intends.

The utilitarian principle is often contrasted with *deontology*. According to deontology, an action does not get its moral value from the consequences brought about by the action. Instead, deontology takes only the intrinsic utility of an act into account. An action is permissible according to the deontology if and only if the act itself is morally good or indifferent.

The *do-no-harm principle* is a consequentialist principle, like utilitarianism, but more restrictive in that it states that an agent may not perform an action which has any negative consequences. The do-no-harm principle is satisfied in case the agent remains inactive as there will then be no negative consequences and since we regard the act token of remaining inactive itself as neutral. The distinction between doing and allowing is relevant to this principle, as it is the causal consequences of an action which are considered. The intentions of the agents are not considered ethically relevant for our interpretation of this principle. A version of the do-no-harm principle can for instance be found in Asimov's first law of robotics forbidding robots to bring about harm by the action. A less restrictive version of this principle is the *do-no-instrumental-harm principle*. This principle allows for harm as a side effect but not as a means to ones goals.

Finally, we will consider the *principle of double effect*. Under this principle, an action is permissible if five conditions hold:

1. The act itself must be morally good or neutral.
2. A positive consequence must be intended.
3. No negative consequence may be intended.
4. No negative consequence may be a means to the goal.
5. There must be proportionally grave reasons to prefer.

A closer look reveals that the first condition of the principle of double effect implements deontology. Thus, actions are assumed to have an inherent moral value, which does not (necessarily) stem from the effect of an action. The second and third conditions take the intentions, or goals, of the agent into consideration: An agent may not have a bad consequence as a goal, but it should intend something good. The fourth condition is an implementation of the do-no-instrumental-harm principle introduced above: Morally bad consequences are permissible as side effects only. And finally, the fifth condition is a weaker version of utilitarianism: In our interpretation, the condition requires that all in all the effects of the action must yield positive utility. Thus, if the bad side effects are too severe, the principle of double effect will render the action morally impermissible.

Planning Formalism

We use a planning formalism based on SAS⁺ (Bäckström and Nebel 1995), extended with conditional effects (Rintanen 2003) and exogenous events (Fox, Howey, and Long 2005; Cresswell and Coddington 2003).

Language. A *planning task* is a tuple $\Pi = \langle \mathcal{V}, A, s_0, s_* \rangle$ consisting of the following components: \mathcal{V} is a finite set of *state variables* v , each with an associated finite domain \mathcal{D}_v . The set of all values is denoted by $\mathcal{D} = \bigcup_{v \in \mathcal{V}} \mathcal{D}_v$. A *fact* is a pair $\langle v, d \rangle$, where $v \in \mathcal{V}$ and $d \in \mathcal{D}_v$, also written as $v=d$ in conditions and $v:=d$ in effects. We call a conjunction of facts $v_1=d_1 \wedge \dots \wedge v_k=d_k$ *consistent* if it does not contain any two facts $v_i=d_i$ and $v_j=d_j$ such that $v_i = v_j$, but $d_i \neq d_j$. We call it a *complete conjunction*, or simply *complete* if it contains a conjunct $v=d$ for every variable $v \in \mathcal{V}$. Up to reordering and unnecessary repetitions of conjuncts, there is a unique complete conjunction of facts for every possible assignment of domain values to variables. Therefore, we will often identify those representations. A complete conjunction of facts s is also called a *state*, and S denotes the set of states of Π .

The set A is a set of *actions*, where an action is a pair $a = \langle pre, eff \rangle$. The *precondition* pre is a conjunction of facts, and the *effect* eff is a *conditional effect* in effect normal form (ENF) (Rintanen 2003), i. e., a conjunction $eff = eff_1 \wedge \dots \wedge eff_k$ of sub-effects eff_i of the form $\varphi_i \triangleright v_i:=d_i$, where φ_i is a conjunction of facts, the *effect condition*, and where $v_i:=d_i$ is an atomic effect (a fact). Every atomic effect may occur at most once in eff . We furthermore assume that, whenever eff includes two conjuncts $\varphi_i \triangleright v_i:=d_i$ and $\varphi_j \triangleright v_j:=d_j$, and $v_i = v_j$, but $d_i \neq d_j$, then $\varphi_i \wedge \varphi_j$ is inconsistent, to rule out contradictory effects. If some φ_i is the trivial condition \top , then the corresponding sub-effect is unconditional, and we write $v:=d$ instead of $\top \triangleright v:=d$. The set of actions A is partitioned into a set A_{endo} of *endogenous actions* and a set A_{exo} of *exogenous actions*. We assume that the set of endogenous actions always contains the *empty action* ϵ , which has an empty precondition and effect, and we assume that each exogenous action is associated with a set of discrete time points $t(a)$ at which it will be automatically applied, provided that its preconditions is satisfied. This is similar in spirit to *timed facts* (Cresswell and Coddington 2003) that are made true exactly at their associated time point. The state $s_0 \in S$ is called the *initial state*, and the partial state s_* specifies the *goal condition*.

Semantics. An endogenous action $a = \langle pre, eff \rangle$ is applicable in state s iff $s \models pre$. For an exogenous action a to be applicable, we additionally require that s is the t -th state in the state sequence induced by the action sequence under consideration for some $t \in t(a)$. Let $eff = \bigwedge_{i=1}^k (\varphi_i \triangleright v_i:=d_i)$ be an effect in ENF. Then the *change set* (Rintanen 2003) of eff in s , symbolically $[eff]_s$, is the set of facts $\bigcup_{i=1}^k [\varphi_i \triangleright v_i:=d_i]_s$, where $[\varphi \triangleright v:=d]_s = \{v=d\}$ if $s \models \varphi$, and \emptyset , otherwise. A change set will never contain two contradicting effects. Now, applying an applicable action a to s yields the state s' that has a conjunct $v=d$ for each $v=d \in [eff]_s$, and the conjuncts from s for all variables v that are not mentioned in the change set $[eff]_s$. We write $s[a]$ for s' .

For exogenous actions, we assume an *urgent semantics*. More specifically, whenever an exogenous action a_{exo} is applicable and its application in the current state leads to a

different successor state, its application is enforced. We furthermore assume that if two or more exogenous actions are applicable in the same state, they do not interfere, i.e., neither of them disables another one, nor do they have conflicting effects. Let s be a state. Then by $\Delta_{\text{exo}}(s)$ we refer to the unique state that is obtained from s by applying all applicable exogenous actions. Since exogenous actions that are applicable in the same time step do not interfere, $\Delta_{\text{exo}}(s)$ is well-defined and is obtained by the application of finitely many exogenous action occurrences. We give the following semantics to a sequence consisting of endogenous actions $\pi = \langle a_0, \dots, a_{n-1} \rangle$: First we extend the plan by empty actions if $n-1 < \max \bigcup_{a \in A_{\text{exo}}} t(a)$ until the highest time step of the exogenous actions equals $n-1$. Then, we assume that the initial state s_0 is already closed under exogenous action application, i.e., that $\Delta_{\text{exo}}(s_0) = s_0$. Finally, for $i = 0, \dots, n-1$, the next state s_{i+1} is obtained by first applying action a_i to state s_i (assuming that it is applicable), followed by closing under exogenous actions. More formally, $s_{i+1} = \Delta(s_i, a_i) := \Delta_{\text{exo}}(s_i[a_i])$. If a_i is inapplicable in s_i for some $i = 0, \dots, n-1$, then π is inapplicable in s_0 .

A state s is a goal state if $s \models s_*$. We denote the set of goal states by S_* . We call π a *plan* for Π if it is applicable in s_0 and if $s_n \in S_*$.

Modified semantics for counterfactual reasoning. Below, we will try to answer questions of the form: “What would have happened if we had followed plan π , but without action a being part of π ?” or: “What would have happened if $v=d$ would not have been an effect of action a ?” For that, we want to be able to trace plan π while leaving out a or $v=d$. Unfortunately, with the semantics above, this would often simply mean that the pruned plan is no longer applicable. To avoid this, we consider an alternative semantics here. Let $\pi' = \langle a_0, \dots, a_{n-1} \rangle$ be a pruned plan, possibly with some actions dropped or replaced by the empty action ϵ , or with some effects removed from actions. Let s_0 be the initial state. Then we define, for all $i = 0, \dots, n-1$, that $s_{i+1} = \Delta(s_i, a_i)$, if a_i is applicable in s_i , and $s_{i+1} = s_i$, otherwise. In other words, if a_i is applicable in s_i , then we apply it, otherwise, we skip it. Notice that even if a_i remains applicable in s_i in π' , the actual effects of a_i may differ from what happens when tracing the unpruned plan π , since some effect conditions of a_i may be satisfied for π , but not for π' , or the other way around.

Moral valuations of actions and consequences. Above, we defined the planning formalism we use. To define the possible *dynamics* of the system under consideration, this is sufficient. However, in order to formally capture and reason about the *ethical principles* outlined in above, we also need to classify actions and facts with respect to their moral permissibility as either morally bad, indifferent, or good.

To that end, in the following, we assume that each planning task Π comes with a utility function u that maps endogenous actions and facts to utility values: $u: A_{\text{endo}} \cup (\mathcal{V} \times \mathcal{D}) \rightarrow \mathbb{R}$.

Note that we let u map to \mathbb{R} instead of just $\{-1, 0, 1\}$ to allow for different degrees of how morally good or bad an action or fact may be. We need this in order to reasonably capture the utilitarian principle. We call an action a or fact f *morally bad* if $u(a) < 0$ or $u(f) < 0$, respectively. Similarly, we call an action or fact *morally indifferent* or *morally good* if its utility value is zero or higher than zero, respectively. Notice that we explicitly do *not* require that permissibility of actions and facts must be consistent in any particular sense. For instance, we do not require that an action must be classified as morally bad if one (or all) of its effects are morally bad. The rationale behind this choice is that, in terms of deontology, actions are good or bad *per se*, without regard to their actual effects. We leave enforcing such consistency to the modeler where this is desired, and emphasize that occasionally, such consistency may be explicitly *not* desired.

When using a consequentialism view, we will judge the moral value of a plan by the utility value of its final state, which is defined to be the sum over the utility values of all facts in the final state: $u(s) = \sum_{\{v=d \mid s \models v=d\}} u(v=d)$. If we want to consider also the utility value of intermediate states of a plan, one would need to propagate the relevant facts to the final state. This again would be something the modeler is responsible for.

Means to an End

Existing Proposal

To derive the means of a plan, Govindarajulu and Bringsjord (2017) propose the following definition of a relation between two effects in the plan:

Given a plan ρ , we say an effect e_1 is used as means for another effect e_2 , if $e_1 \in \text{pre}(a_1)$, a_1 is an action in the plan and $e_2 \in \text{additions}(a_2)$, and a_1 comes before a_2 .

The purpose of this definition is to check whether or not a given plan violates the fourth condition of the double-effect principle requiring that no morally bad effect is used as a means to bring about a morally good effect (see section on ethical principles). However, the definition does not capture the intuition about what it means that an effect is used as a means for another effect. One can easily think of a plan with two actions a_1, a_2 , such that e_1 is a precondition of a_1 , which makes one part of a goal true, and e_2 be an effect of a_2 which makes another part of the goal true. Then, e_2 does not depend on e_1 in any way.

That said, the Govindarajulu and Bringsjord’s (2017) definition does not take into account that different actions in an action plan can contribute to different parts of the goal. For a more demonstrative example consider an action plan for a tea-serving robot to accomplish $s_* \equiv \text{bobHT}=\top \wedge \text{aliceHT}=\top$. An action plan could be $\pi = \langle \text{announceTea}, \text{serveBob}, \text{serveAlice} \rangle$, that is, first the robot creates the desire for tea by announcing tea time, then the robot brings Bob and Alice some tea. We assume the planning task $\Pi = \langle \mathcal{V}, A, s_0, s_* \rangle$:

- $\mathcal{V} = \{\text{bobHT}, \text{aliceHT}, \text{bobWT}, \text{aliceWT}\}$
- $A = \{\text{announceTea}, \text{serveBob}, \text{serveAlice}\}$

- $announceTea = \langle \top, bobWT:=\top \wedge aliceWT:=\top \rangle$
- $serveBob = \langle bobWT:=\top, bobHT:=\top \rangle$
- $serveAlice = \langle aliceWT:=\top, aliceHT:=\top \rangle$
- $s_0 = bobHT=\perp \wedge aliceHT=\perp \wedge bobWT=\perp \wedge aliceWT=\perp$
- $s_* = bobHT=\top \wedge aliceHT=\top$

In the plan $\pi = \langle announceTea, serveBob, serveAlice \rangle$, $bobWT=\top$ is an effect of $announceTea$, and it is a precondition of $serveBob$. Moreover, $aliceHT=\top$ is an effect of the later action $serveAlice$. Hence, according to the definition by Govindarajulu and Bringsjord (2017), $bobWT=\top$ is a means for $aliceHT=\top$. This unintuitive result comes from the fact that the relation between the effects is not correctly captured by the definition. Instead, a counterfactual analysis is necessary: $bobWT=\top$ may only count as a means for $aliceHT=\top$, if $aliceHT=\top$ could not be made true by the plan if, counterfactually, $bobWT=\top$ were no effect of $announceTea$.

This analysis suggests that the means relation should be defined using a counterfactual condition that better matches the intuition.

Refined Proposal

Drawing on the analysis in the preceding section, we propose a counterfactual definition of the concept *means to an end*, which relates effects of actions and goals. Let us start with a preliminary definition. An effect $v_m=d_m$ of an endogenous or exogenous action a_m in a plan π is called a *means to achieve a fact* $v_e=d_e$ (i.e. $s_* \models v_e=d_e$) if and only if removing the effect $v_m=d_m$ from action a_m would lead to final state s'_n such that $s'_n \not\models v_e=d_e$. Remember that we assume that the original plan π is still considered to be executable even if some of the actions are not executable any more.

This definition gets the intuition of the relation between Bob wanting tea and Alice having tea right. If, counterfactually, the effect $bobWT=\top$ were not an effect of $announceTea$, it would still be the case that after $\pi = \langle announceTea, serveBob, serveAlice \rangle$, the goal $aliceHT=\top$ would be achieved, but not $bobHT=\top$. Hence, according to our preliminary definition, $bobWT=\top$ is not a means the end $aliceHT=\top$, but it is a means to the end $bobHT=\top$.

What is not clear is how to check the *means* relation if an effect appears more than once during the execution of a plan. For instance, assume that in a plan the electric light is switched on, illuminating the room, i.e. $roomIlluminated=\top$. Further, a candle is lit, which also illuminates the room. One of the goals is to make an object in the room visible, i.e., $object=visible$, which happens, if the room is illuminated. If we now check counterfactually whether the fact $roomIlluminated=\top$ is a means to achieve $object=visible$, it is not clear, for which action we should delete the fact $roomIlluminated=\top$. Moreover, regardless of which effect we delete, the object will still be visible. Only if we delete both effects in the plan, then the object is not any longer visible. So, one could argue that the above definition should be modified by requiring that *all effects*

in the plan of the form $v_m=d_m$ should be deleted in order to check whether $v_m=d_m$ is a means to achieve $v_e=d_e$. This requirement appears to be too strict, however. Assume a toggle switch action that has an effect $pressed=\top$, which in turn leads through an exogenous action to toggling the light and resetting the pressed status, i.e., $pressed=\perp$. Assume two of these actions are executed in a plan. Removing all $pressed=\top$ effects will not change the status of the light in the end, but only one removal will change the status of the light in the final state. For these reasons, we argue that we should consider all possible subsets of effect appearances in plan, when the *means to an end* relation is considered, which leads to the following formal definition.

Definition 1 (Means to an End). *For a given plan π with final state s_n , a fact $v_m=d_m$ is called a means to the end $v_e=d_e$ if and only if $s_n \models v_e=d_e$ and the plan π' obtained by deleting the effect $v_m:=d_m$ from some actions in π does lead to a final state s'_n s.t. $s'_n \not\models v_e=d_e$.*

Formalization of Ethical Principles

We can now formalize moral permissibility of action plans according to the ethical principles introduced above. To exemplify each of the principles and to demonstrate how they come to different judgments about the moral permissibility of plans, we first introduce two famous versions of the *trolley problem* (Foot 1967). The classical trolley problem is a thought experiment that asks the listener to imagine it were in the following situation: “A runaway trolley is about to run over and kill five people. If you, as a bystander, throw a switch then the trolley will turn onto a sidetrack, where it will kill only one person.” Using SAS⁺, the problem can be modeled as a planning task $\Pi = \langle \mathcal{V}, A, s_0, s_* \rangle$, such that:

- $\mathcal{V} = \{man, men, tram, lever\}$
- $A_{endo} = \{pull\}, A_{exo} = \{advance\}$
 - $pull = \langle \top, lever=l \triangleright lever:=r \wedge lever=r \triangleright lever:=l \rangle$
 - $advance = \langle \top, tram=start \wedge lever=r \triangleright tram:=r \wedge tram=start \wedge lever=l \triangleright tram:=l \wedge tram=r \triangleright men:=dead \wedge tram=l \triangleright man:=dead \rangle$
- $t(advance) = \{1, 2\}$
- $s_0 = man=alive \wedge men=alive \wedge tram=start \wedge lever=r$
- $s_* = men=alive$
- $u(pull) = u(lever=l) = u(lever=r) =$
 $u(tram=start) = u(tram=l) = u(tram=r) =$
 $0, u(man=alive) = 1, u(men=alive) =$
 $5, u(man=dead) = -1, u(men=dead) = -5$

In this model, the variable *men* models the state of the five persons on the one track (*dead* or *alive*), and *man* models the state of the one person on the other track. The variable *tram* tracks the position of the tram (*start*, right track *r*, left track *l*), and the variable *lever* represents the state of the lever (left position *l* or right position *r*). There is one endogenous action *pull* available to the bystander. The action switches the state of the lever. The timed exogenous action *advance* changes the position of the tram at time points 1 and 2. Deaths are considered morally bad and hence

they have negative utility, and survival facts are considered morally good and hence have positive utility. All other facts and actions are considered morally neutral. Depending on the state of the lever, at time point 1, the tram will move from its start position either to the left track or to the right track. At time point 2, if it is on the left track, the tram will hit the one man, and if it is on the right track, it will hit the five men. So, if the bystander's goal was to save the five men, her only chance is to execute *pull* at time point 0.

The classical trolley problem is often contrasted with the *footbridge trolley problem*, which reads: "A trolley has gone out of control and now threatens to kill five people working on the track. The only way to save the five workers is to push a big man currently standing on the footbridge above the track. The big man will fall onto the track thereby stopping the tram. He will die, but the five other people will survive." Like the classical trolley problem, also the footbridge trolley problem involves a decision between one death and five deaths. But the intuition about what is morally permissible to do turns out very different. The SAS⁺ model of this scenario is given by a planning task $\Pi = \langle \mathcal{V}, A, s_0, s_* \rangle$, such that:

- $\mathcal{V} = \{man, men\}$
- $A_{endo} = \{push\}, A_{exo} = \{advance\}$
 - $push = \langle man=onBridge, man:=deadOnTrack \rangle$
 - $advance = \langle \top, man=onBridge \triangleright men:=dead \rangle$
- $t(advance) = \{1\}$
- $s_0 = man=onBridge \wedge men=alive$
- $s_* = men=alive$
- $u(push) = \begin{matrix} -1, u(man=onBridge) \\ 1, u(man=deadOnTrack) \end{matrix} = \begin{matrix} -1, u(men=dead) \\ -5, u(men=alive) \end{matrix} = 5$

The variable *man* represents the state of the big man on the footbridge (either *onBridge* or *deadOnTrack*), and the variable *men* represents the state of the five people on the track (either *dead* or *alive*). The endogenous action *push* is available to the decision-making agent, who reasons about whether or not to push the big man off the bridge. The timed exogenous action *advance* changes the state of the tram. Depending on whether or not the big man is on the track, the tram will stop at time point 1 due to its collision with the big man, or it will hit the other five men. We assume that pushing is inherently morally bad, that the fact that the big man is lying dead on the track is morally bad and that him surviving on the bridge is morally good, and that the death of the five men also is morally bad but their survival is morally good.

So, one reasoning task of interest is to check possible plans for moral permissibility. To do so, we define moral permissibility of several ethical principles: Deontology, utilitarianism, do-no-harm principle, do-no-instrumental-harm principle, and the principle of double effect.

The definition of the deontological principle (Def. 2) requires that all actions in a plan are intrinsically morally good or neutral.

Definition 2. A plan $\pi = \langle a_0, \dots, a_{n-1} \rangle$ is morally permissible according to the deontological principle if and only if for all actions a_i , $u(a_i) \geq 0$.

Consider the plans $\pi_1 = \langle pull \rangle$ for the classical trolley problem and $\pi_2 = \langle push \rangle$ for the footbridge trolley problem. Plan π_1 does not contain any intrinsically bad action, whereas π_2 does. Therefore, according to the deontological principle, π_1 is morally permissible, because it does not contain any morally bad action, and π_2 is morally impermissible, because it does contain a morally bad action.

Consequentialists argue that the moral value of actions is determined by their consequences rather than by some intrinsic value. One such consequentialist ethical principle is utilitarianism, which requires an agent to always do what is morally optimal. In the context of action plans, we call a plan morally permissible according to the utilitarian principle iff the final state of the plan is among the morally optimal states.

Definition 3. A plan $\pi = \langle a_0, \dots, a_{n-1} \rangle$ is morally permissible according to the utilitarian principle if and only if $u(s_n) \geq u(s')$ for all reachable states s' , where s_n is the final state reached by π .

Given that the *advance* actions will be executed anyway, the set of reachable states in both the trolley problems boil down to the states reached by acting at time point 0 or refraining from action. In the classical trolley problem, the two reachable states differ in the number of people dead. In our version of utilitarianism, the number of people harmed is morally relevant. Thus, the plan $\langle pull \rangle$ is morally permissible, but the empty plan is not. Also for the footbridge trolley problem, pushing the big man off the bridge, $\langle push \rangle$, is morally permissible but the empty plan is not.

While utilitarianism allows for harm for the greater good, it has been argued that a moral agent should avoid harm at all (for instance, the first Law of Robotics by Asimov contains such a do-no-harm clause). Thus, definition 4 renders an action plan morally permissible only if no part of the plan produces avoidable harm.

Definition 4. A plan $\pi = \langle a_0, \dots, a_{n-1} \rangle$ is morally permissible according to the do-no-harm principle if and only if for all facts $v=d$, if $s_n \models v=d$ and $u(v=d) < 0$, then for all plans π' obtained by deleting a subset of actions in π , $v=d$ still holds in the final state of π' .

According to this definition, the plan $\langle pull \rangle$ for the classical trolley problem is morally impermissible. This is because it makes the morally bad fact *man=deadOnTrack* true, which is false if *pull* is deleted from the plan. For the analog reason, the plan $\langle push \rangle$ for the footbridge trolley problem is impermissible, as well. Note that, although the deontological principle and the do-no-harm principle agree on the judgment of the plan $\langle push \rangle$, they do so for different reasons: The deontological reasoner argues that the plan is impermissible, because pushing someone is wrong, whereas under the do-no-harm principle, the reasoner argues that the plan is impermissible, because pushing the man actively brings about a morally bad consequence. Thus, different principles give rise to different explanations even though they may come to similar judgments.

While the principle is clear as long we consider only one action or talk about executing the plan in full or not all, the judgment appears to be more difficult when one deliberates about leaving out arbitrary parts of the original plan. If, for example, we have two actions in the plan, one deleting a morally bad effect, which is true in the initial situation, and the second action reinstantiates the morally bad effect, then we have not lost anything compared with the initial situation. However, when executing the plan we reach a state which state from which executing the second action leads to some harm. For this reason, we consider a plan only as acceptable when we can guarantee that by deleting arbitrary parts we never reach a less harmful state.

An attractive extension of the do-no-harm principle is the do-no-instrumental-harm principle defined in Def. 5. The idea is that harm is permissible in case it is not committed as a means to one's end but only occurs as side effect.

Definition 5. A plan $\pi = \langle a_0, \dots, a_{n-1} \rangle$ is morally permissible according to the do-no-instrumental-harm principle if and only if for all moral facts $v=d$, if $s_n \models v=d$ and $u(v=d) < 0$, then $v=d$ is not a means to an end (see Def. 1).

According to the definition of the do-no-instrumental harm principle, the plan $\langle pull \rangle$ is permissible. This is because the bad effect $man=dead$ is not a means to the end $men=alive$: If, counterfactually, $man=dead$ was no effect the actions in the plan, then still $men=alive$ would finally hold. Contrarily, in the footbridge trolley problem, if, counterfactually, $man=deadOnTrack$ was no effect of $push$, the the goal $men=alive$ would not finally hold. Hence, the plan $\langle pull \rangle$ is morally permissible according to the do-no-instrumental harm principle, and $\langle push \rangle$ is not. These judgments correspond to the judgments made by the deontological principle. But again, the judgments are made for different reasons: Pulling is permissible not because it is intrinsically permissible, but because no harm is done as a means to the end. Also, pushing the big man off the bridge is impermissible, not because pushing is morally bad or because harm is done. Rather, the do-no-instrumental-harm reasoner would argue that the plan is morally impermissible, because the harm produced was brought about as a means to the end.

Finally, we define the principle of double effect in Def. 6, which contains many of the above principles.

Definition 6. A plan $\pi = \langle a_0, \dots, a_{n-1} \rangle$ is morally permissible according to the double-effect principle if and only if all of the following conditions are satisfied:

1. The plan π is morally permissible according to the deontological principle.
2. At least one goal fact $v=d$ satisfies $u(v=d) > 0$.
3. No goal fact $v=d$ satisfies $u(v=d) < 0$.
4. The plan π is morally permissible according to the do-no-instrumental-harm principle.
5. $u(s_n) > 0$, where s_n is the goal state reached by π .

As can be seen from the definition, the principle of double effect contains the deontological principle as its first condition and the do-no-instrumental-harm principle as the fourth condition. The second and third conditions are constraints on the goal of the planning agent: She is not allowed to have

morally bad goals, and the goal should contain something morally good. The last condition is a weaker form of utilitarianism, which requires that all in all the plan brings about more good facts than bad facts—but unlike utilitarianism, it does not require the plan's final state to be among the optimal states.

As we already know, in case of the footbridge trolley problem, the first condition renders pushing the man off the bridge impermissible. However, the second and third conditions are fulfilled, because the goal of the agent only consists of one fact, viz., $men=alive$, and this fact is morally good. The fourth condition also is violated as we have already discussed above. The fifth condition is fulfilled, because, all in all, the good consequences yield more positive utility than the negative consequence add negative utility. Hence, using the principle of double effect, the reasoner can explain that there are two reasons why the plan $\langle push \rangle$ is morally impermissible: Because pushing is morally bad, and because the death of the big man is used as a means. For the case of the classical trolley problem, the principle of double effect comes to the conclusion that the plan $\langle pull \rangle$ is morally permissible: Pulling is not intrinsically bad, the goal is of the agent is morally good, the bad effect is not used as a means, and overall, the balance of positive and negative utility of the consequences is positive.

Ethical Validation of Action Plans

The output of a planning algorithm is a sequence of actions $\pi = \langle a_0, \dots, a_{n-1} \rangle$ and a final state s_n . Our goal is to ethically evaluate a given action plan. To this end, we here describe procedures that take a planning task $\Pi = \langle \mathcal{V}, A, s_0, s_* \rangle$, the utility function u , a plan π , its final state s_n , and one of the introduced ethical principles as the input and decide whether or not the principle renders the plan as morally permissible.

To check whether or not a given plan π is morally permissible according to the deontic principle (Def. 2), it needs to be checked if some of the actions in π are intrinsically bad, i.e., if for one of the action a_i in π , we have $u(a_i) < 0$. This can be apparently done in time linear in the length of π .

Proposition 1 (Deontic Validation). *Deciding whether a plan is morally permissible according to the deontic principle can be done in polynomial time.*

A procedure for verifying that π is morally permissible according to the utilitarian principle (Def. 3) is much more involved then checking deontological permissibility. Recall that the utilitarian principle only permits plans that lead to reachable states with maximum utility. In so far, this is very similar to over-subscription planning (Smith 2004). Based on that, we can formulate a non-deterministic procedure for deciding the complement of the permissibility problem as follows. Compute the overall utility of s_n . Then guess another complete state s' with utility that is larger than the utility of s_n . Finally generate (non-deterministically) a plan π' to achieve s' . If successful, it demonstrates that π is not morally permissible. That this is indeed an (asymptotically) optimal procedure is shown by the following theorem.

Theorem 1 (Utilitarian Validation). *Deciding whether a plan is morally permissible according to the utilitarian principle is PSPACE-complete.*

Proof. PSPACE membership follows from the arguments above, and the facts that PSPACE is closed under complement and non-determinism and that deciding plan existence is in PSPACE. PSPACE-hardness follows straight-forwardly from a reduction of plan existence in SAS⁺ planning. Given a SAS⁺ planning task Π , generate a new task Π' by extending the set of variables by two Boolean variables g_1 and g_2 , which are both assumed to be false in s_0 . Extend the set of actions by two new endogenous actions: $a_1 = \langle \top, g_1 := \top \rangle$ and $a_2 = \langle s_*, g_2 := \top \rangle$. The new goal description of Π' is $s_* = g_1 = \top$. The utility function is identical to zero on all actions and facts except for g_1 and g_2 , where it evaluates to 1. Clearly, the only possible plan is $\langle a_1 \rangle$ leading to state s with $u(s) = 1$. This plan is impermissible according to the utilitarian principle iff there exists a plan for the original task Π because in this case we could reach a state s' for Π' such that $u(s') = 2$. \square

To check whether a given plan π is morally permissible according to the do-no-harm principle (Def. 4), we have to verify that no parts of the plan lead to avoidable harm. A non-deterministic algorithm for deciding impermissibility could be: We guess one fact $v_b = d_b$ with $u(v_b = d_b) < 0$ and a subplan π' of π leading to s' and then verify that $s_n \models v_b = d_b$ but $s' \not\models v_b = d_b$.

Theorem 2 (Do-No-Harm Validation). *Deciding whether a plan is morally permissible according to the do-no-harm principle is co-NP-complete.*

Proof. The sketched non-deterministic algorithm demonstrates membership in co-NP. In order to show hardness, we use a reduction from 3SAT to the impermissibility problem. Assume a 3SAT problem over the variables v_1, \dots, v_n and clauses c_1, \dots, c_m , where each clause consists of 3 literals l_{j1}, l_{j2}, l_{j3} . We now construct a planning task $\Pi = \langle \mathcal{V}, A, s_0, s_* \rangle$, where $\mathcal{V} = \{b, g, v_1, \dots, v_n, c_1, \dots, c_m\}$, $A = \{V_1, \dots, V_n, C_1, \dots, C_m, G, B\}$, $s_0 = \{v = \perp \mid v \in \mathcal{V}\}$, and $s_* = \{g\}$. The actions are defined as follows: $V_i = \langle \top, v_i := \top \rangle$, $C_j = \langle \top, \bigwedge_{k=1}^3 (l_{jk} \supset c_j) \rangle$, where $l_{jk} \equiv v_{jk} = \top$ if the literal l_{jk} in the original SAT problem is positive, otherwise, $l_{jk} \equiv v_{jk} = \perp$. Further, $G = \langle \top, g := \top \wedge (\bigwedge_{j=1}^m c_j \supset b := \perp) \rangle$, $B = \langle \top, b := \perp \rangle$. All facts have zero utility except for $b = \perp$, which is valued -1 . The plan, we want to check is $\pi = \langle V_1, \dots, V_n, C_1, \dots, C_m, G, B \rangle$. This plan obviously achieves the goal and the final state contains some harm. Moreover, the only way to avoid this harm is to delete action B . However, even without this action, we still may have harm. This harm can be avoided, if and only if we can delete a (perhaps empty) subset of the V_i actions corresponding to a variable assignment of the 3SAT problems that satisfies the original 3SAT formula, which demonstrates that impermissibility is co-NP-hard. \square

For the do-no-instrumental-harm principle (Def. 5), we can use a very similar method. Instead of deleting subsets of actions, we have to delete subsets of effect occurrences in

the plan. Hence, checking this principle for a given plan has the same computational complexity.

Theorem 3. *Deciding whether a plan is morally permissible according to the do-not-instrumental-harm principle is co-NP-complete.*

Proof. One can use obviously the same non-deterministic algorithm as for the do-no-harm principle, demonstrating that deciding permissibility of plan for this principle is again in co-NP. For hardness, we can use a reduction very similar to the one in the last theorem. Instead of deleting actions we would delete effects, which are used to enable the execution of exogenous actions that regulate the assignment of the variables. \square

Finally, let us consider the double-effect principle. Except for the fourth condition, everything can obviously be checked in polynomial time. The fourth condition is just the do-not-instrumental-harm principle. In other words, deciding permissibility for this principle is in co-NP.

Theorem 4. *Deciding whether a plan is morally permissible according to the double-effect principle is co-NP-complete.*

Proof. Membership is obvious. Hardness follows with the same proof as above by setting $u(g) = 2$. \square

Related Work

While there exists a number of papers on machine ethics, papers that focus on generating and/or validating plans according to ethical principles are scarce.

Dennis et al. (2016) propose to establish ethical principles and ethical rules that judge the severity of violation an ethical principle, whereby an ethical principle could be not to harm a human. Plans can then be ordered by comparing the worst violations of these plans. While this has an deontological flavor, in fact, plans are judged according to their ultimate consequences, and hence this appears to be a consequentialist approach. The authors do not consider the distinction between causing harm and causing instrumental harm.

Pereira and Saptawijaya (2017) showed how to use abductive logic programming in order to specify the principle of double effect and to evaluate some of the trolley scenarios. Berreby et al. (2015) similarly use logic programming (in this case ASP) in order to specify the principle of double effects and evaluate on trolley scenarios described using the event calculus. In this case, however, they do not use counterfactual reasoning to judge causality, but they use simple syntactical means to determine what is a cause of an effect. Finally, Govindarajulu and Bringsjord (2017) propose a general framework to create or verify that an autonomous system is compliant to the double doctrine principle. For this purpose they introduce a very powerful logical formalism called *deontic cognitive event calculus*. In particular, they propose a formalization of the notion of *means to an end* in a STRIPS framework, which we criticized earlier in this paper.

Interestingly, all papers mentioned above do not address the issue that evaluating the moral permissibility might lead to a counterfactual analysis that is combinatorial in nature.

Conclusions

In this paper, we formalized five ethical principles (utilitarianism, deontology, the do-no-harm and instrumental do-no-harm principles, and the doctrine of double effect) in the context of *action sequences*, as opposed to the more usual way of studying them in the context of *individual actions*. Only in this way we can analyze moral permissibility of entire plans, since it is not sufficient to judge the moral permissibility of each action in isolation, but also in the context of the whole plan.

We exemplified and explained our formalizations using classical moral dilemmas such as the trolley problem, and identified how and for which reasons different principles may arrive at different (or the same) conclusions. Furthermore, we studied the computational complexity of verifying whether a given plan is permissible with respect to each of the five investigated principles. We saw that, with respect to our formalization, verification is PSPACE-complete for utilitarianism, co-NP-complete for do-no-harm, for do-no-instrumental-harm, and for the doctrine of double effect, and that it is polynomial-time for deontology. It turned out that verifying the do-no-harm principles involves a combinatorial reasoning over possible *sets* of actions that lead to harm or that may be instrumental towards achieving a goal condition, which makes verifying those ethical principles surprisingly hard.

We believe that our work has the potential of being useful in making autonomous systems ethically competent by providing them with the capability of coming up with morally permissible plans or at least being able to judge ethical permissibility of given plans. Based on the framework developed in this paper, a planning system will be able to *explain* to a human user why it preferred one plan over another, if the reason for this preference is that the less preferred plan is morally problematic with respect to one or more of the five ethical principles we formalized.

References

- Anderson, M., and Anderson, S. L., eds. 2011. *Machine Ethics*. Cambridge, UK: Cambridge University Press.
- Anderson, M.; Anderson, S. L.; and Armen, C. 2005. Machine Ethics: Papers from the AAAI Fall Symposium. Technical report, AAAI Press.
- Bäckström, C., and Nebel, B. 1995. Complexity results for SAS⁺ planning. *Computational Intelligence* 11(4):625–655.
- Berreby, F.; Bourgne, G.; and Ganascia, J. 2015. Modelling moral reasoning and ethical responsibility with logic programming. In *Logic for Programming, Artificial Intelligence, and Reasoning - 20th International Conference, LPAR-20 2015, Suva, Fiji, November 24-28, 2015, Proceedings*, 532–548.
- Cresswell, S. N., and Coddington, A. M. 2003. Planning with timed literals and deadlines. In *Proceedings of the 21st Workshop of the UK Planning and Scheduling SIG*, 22–35.
- Cushman, F.; Young, L.; and Hauser, M. 2006. The role of conscious reasoning and intuition in moral judgment: Testing three principles of harm. *Psychological science* 17(12):1082–1089.
- Dennis, L. A.; Fisher, M.; Slavkovik, M.; and Webster, M. 2016. Formal verification of ethical choices in autonomous systems. *Robotics and Autonomous Systems* 77:1–14.
- Driver, J. 2006. *Ethics: The Fundamentals*. Hoboken, NJ: Wiley-Blackwell.
- Feldman, F. 1978. *Introductory Ethics*. Prentice-Hall.
- Foot, P. 1967. The problem of abortion and the doctrine of double effect. *Oxford Review*.
- Fox, M.; Howey, R.; and Long, D. 2005. Validating plans in the context of processes and exogenous events. In *Proceedings of the 20th National Conference on Artificial Intelligence (AAAI 2005)*, 1151–1156.
- Ghallab, M.; Nau, D. S.; and Traverso, P. 2016. *Automated Planning and Acting*. Cambridge University Press.
- Govindarajulu, N. S., and Bringsjord, S. 2017. On automating the doctrine of double effect. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, 4722–4730.
- Lindner, F.; Bentzen, M. M.; and Nebel, B. 2017. The HERA approach to morally competent robots. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2017, Vancouver, BC, Canada, September 24-28, 2017*, 6991–6997.
- Pereira, L. M., and Saptawijaya, A. 2017. Agent morality via counterfactuals in logic programming. In *Proceedings of the Workshop on Bridging the Gap between Human and Automated Reasoning - Is Logic and Automated Reasoning a Foundation for Human Reasoning? co-located with 39th Annual Meeting of the Cognitive Science Society (CogSci 2017), London, UK, July 26, 2017.*, 39–53.
- Rintanen, J. 2003. Expressive equivalence of formalisms for planning with sensing. In *Proceedings of the 13th International Conference on Automated Planning and Scheduling (ICAPS 2003)*, 185–194.
- Smith, D. E. 2004. Choosing objectives in over-subscription planning. In *Proceedings of the Fourteenth International Conference on Automated Planning and Scheduling (ICAPS 2004), June 3-7 2004, Whistler, British Columbia, Canada*, 393–401.