



# Approval-directed agency and the decision theory of Newcomb-like problems

Caspar Oesterheld<sup>1,2</sup> 

Received: 17 January 2018 / Accepted: 15 February 2019  
© The Author(s) 2019

## Abstract

Decision theorists disagree about how instrumentally rational agents, i.e., agents trying to achieve some goal, should behave in so-called Newcomb-like problems, with the main contenders being causal and evidential decision theory. Since the main goal of artificial intelligence research is to create machines that make instrumentally rational decisions, the disagreement pertains to this field. In addition to the more philosophical question of what the right decision theory is, the goal of AI poses the question of how to implement any given decision theory in an AI. For example, how would one go about building an AI whose behavior matches evidential decision theory's recommendations? Conversely, we can ask which decision theories (if any) describe the behavior of any existing AI design. In this paper, we study what decision theory an approval-directed agent, i.e., an agent whose goal it is to maximize the score it receives from an overseer, implements. If we assume that the overseer rewards the agent based on the expected value of some von Neumann–Morgenstern utility function, then such an approval-directed agent is guided by two decision theories: the one used by the agent to decide which action to choose in order to maximize the reward and the one used by the overseer to compute the expected utility of a chosen action. We show which of these two decision theories describes the agent's behavior in which situations.

**Keywords** Reinforcement learning · Causal decision theory · Evidential decision theory · Newcomb's problem · AI safety · Philosophical foundations of AI

---

✉ Caspar Oesterheld  
caspar.oesterheld@foundational-research.org; caspar.oesterheld@duke.edu

<sup>1</sup> Foundational Research Institute, Berlin, Germany

<sup>2</sup> Present Address: Duke University, Durham, USA

## 1 Introduction

In decision theory, there is a large debate about how an instrumentally rational agent, i.e., an agent trying to achieve some goal or maximize some utility function, should decide in Newcomb's problem (introduced by Nozick 1969) and variations thereof (a list is given by Ledwig 2000, pp. 80–87). Consequently, different normative theories of instrumental rationality have been developed. The best known ones are evidential (sometimes also called Bayesian) decision theory (EDT) (Ahmed 2014; Almond 2010; Price 1986; Horgan 1981) and causal decision theory (CDT) (Gibbard and Harper 1981; Joyce 1999; Lewis 1981; Skyrms 1982; Weirich 2016), but many have attempted to remediate what they view as failures of the two theories by proposing further alternatives (Spohn 2003, 2012; Poellinger 2013; Arntzenius 2008; Gustafsson 2011; Wedgwood 2013; Dohrn 2015; Price 2012; Soares and Levinstein 2017).

Because the main goal of artificial intelligence is to build machines that make instrumentally rational decisions (Russell and Norvig 2010, Sects. 1.1.4, 2.2; Legg and Hutter 2007; Doyle 1992), this normative disagreement has some bearing on how to build these machines (cf. Soares and Fallenstein 2014a, Sect. 2.2; Soares and Fallenstein 2014b, Sect. 1; Bostrom 2014b, Chap. 13, Sect. "Decision theory"). The differences between these decision theories are probably inconsequential in most situations (Ahmed 2014, Sect. 0.5, Chap. 4; Briggs 2017),<sup>1</sup> but still matter in some (Ahmed 2014, Chap. 4–6; Soares 2014a; Bostrom 2014a). In fact, AI may expose the differences more often. For example, Newcomb's problem and the prisoner's dilemma with a replica (Kuhn 2017, Sect. 7) are easy to implement for agents with copyable source code (cf. Yudkowsky 2010 pp. 85ff. Soares and Fallenstein 2014b, Sect. 2; Soares 2014b; Cavalcanti 2010; Sect. 5). Indeed, the existence of many copies is the norm for (successful) software, including AI-based software. While copies of present-day software systems may only interact with each other in rigid, explicitly pre-programmed ways, future AI-based systems will make decisions in a more autonomous, flexible and goal-driven way. Overall, the decision theory of Newcomb-like scenarios is a central foundational issue which will plausibly become practically important in the longer term.

The problem for AI research posed by the disagreement among decision theorists can be divided into two questions:

1. What decision theory do we want an AI to follow?
2. How could we implement such a decision theory in an AI? Or: How do decision theories and AI frameworks or architectures map onto each other?

Although it certainly requires further discussion, there already is a large literature related to the first question.<sup>2</sup> In this paper, I would thus like to draw attention to the second question.

<sup>1</sup> In fact, Eells (1981) has argued that EDT and CDT always behave in the same way, but I disagree with this assessment based on the reasons given by Ahmed (2014, Sect. 4.3–4.6) and Price (1986).

<sup>2</sup> Of course, the existing literature asks about the right decision theory proper. The answer to that question might differ from the answer to the AI-specific question (cf. Kumar 2017; Treutlein 2018). After all, even if we have identified the right decision theory for ourselves, we may want to implement a different decision theory in an AI. One reason could be that the main contenders are not self-recommending—it has been pointed out that EDT and CDT both recommend to self-modify into slightly different decision theories

Specifically, I would like to investigate how approval-directed agents behave in Newcomb-like problems. By an approval-directed agent, I mean an agent that is coupled with an overseer. After the agent has chosen an action, the overseer scores the agent for that action. Rather than, say, trying to bring about particular states in the environment, the agent chooses actions so as to maximize the score it receives from the overseer (cf. Christiano 2014). A model of approval-directed agency that allows us to describe Newcomb-like situations is described and discussed in Sect. 2.

Approval-directed agency is intended as a model of reinforcement learning agents (see Sutton and Barto 1998; Russell and Norvig 2010; Chaps. 17, 21, for introductions to reinforcement learning), for whom the reward function is analogous to the approval-directed agent's overseer. Since reinforcement learning is such a general and commonly studied problem in artificial intelligence (Hutter 2005, e.g. Chap. 4.1.3; Russell and Norvig 2010, p. 831; Sutton and Barto 1998, Chap. 1), it is an especially attractive target for modeling.<sup>3</sup> However, because decision theories are usually defined only for single decisions, we will only discuss single decisions whereas reinforcement learning is usually concerned with sequential interactions of agent and environment. However, this decision can also be a policy choice to model sequential decision problems.<sup>4</sup> In addition to limiting our analysis to single decisions, we will not discuss the learning process and simply assume that the agent has already formed some model of the world.

If we assume that, after an action has been taken, the overseer rewards the agent based on the expected value of some von Neumann–Morgenstern utility function, the agent is implicitly driven by two decision theories: The overseer can use the regular conditional expectation or the causal expectation to estimate the value of its utility function; and the agent itself can follow CDT or EDT when maximizing the score it receives from the overseer (Sect. 3).

We then show how the overall decision theory depends on these two potentially conflicting decision theories. If the overseer bases its expected value calculations on looking only at the world, then the agent's decision theory is decisive. If the overseer

---

Footnote 2 continued

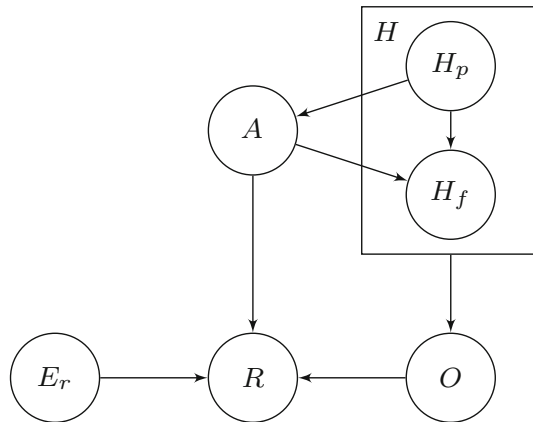
(Meacham 2010; Soares and Fallenstein 2014b, Sect. 3; Yudkowsky 2010, Sect. 2; Greene 2018). The same arguments imply that even if one is convinced of CDT or EDT one would not want the AI to use CDT and EDT. That said, one could also leave the self-modification to the AI.

<sup>3</sup> Reinforcement learning and approval-directed agency are also common outside of artificial intelligence. For example, Achen and Bartels (2016, Chap. 4) review evidence which shows that electorates often vote retrospectively to punish or reward incumbents.

<sup>4</sup> This is consistent with what reinforcement learning algorithms usually do—they choose policies rather than individual actions. This is because the utility of a single action usually cannot be evaluated without knowing how the agent will deal with situations that might arise as a result of taking that action.

When individual actions *can* be evaluated in isolation, the *ex ante* policy choice sometimes differs from the choice of individual actions (see the absent-minded driver, introduced by Piccione and Rubinstein 1997; cf. Aumann et al. 1997; the Newcomb-like scenarios discussed by, e.g., Hintze 2014; Soares and Levinstein 2017, Sect. 2; and the problems in anthropics discussed by Armstrong 2011). While it is rarely discussed in the debate between evidential and causal decision theorists, a few authors regard this discrepancy as crucial and have argued that a proper decision theory should be about optimal policy choices (e.g. Hintze 2014; Soares and Fallenstein 2014b, Sect. 2.1; Soares and Levinstein 2017, Sect. 2). However, this issue is beyond the scope of the present paper.

Further issues in sequential Newcomb-like problems are discussed by Everitt et al. (2015).



**Fig. 1** A causal model of an approval-directed agent in a Newcomb-like decision problem.  $A$  denotes the agent's action,  $H$  the environment history,  $O$  the observation on which the overseer bases the reward,  $R$  is that reward, and  $E_r$  is information about the way the reward is computed that is only available to the overseer. The box is used to indicate that  $H$  includes the two random variables  $H_p$  and  $H_f$ . All of  $H$  may have a causal influence on  $O$

bases its estimates only on the agent's action, then the overseer's decision (or perhaps rather action evaluation) theory is decisive.

## 2 Approval-directed agency

We first describe a model of approval-directed agency. To be able to apply both CDT and EDT, we will use causal models in Pearl's (2009) sense. Consequently, we use Pearl's *do*-calculus-based version of CDT (Pearl 2009, Chap. 4). We will, throughout this paper, assume that the agent has already formed a (potentially implicit) model of the world<sup>5</sup>—e.g., based on past interactions with the environment. Also, we will only consider single decisions rather than sequential problems of iterative interaction between agent and environment.

A causal model of such a one-shot Newcomb problem from the perspective of the approval-directed agent is given in Fig. 1. In this model, the agent decides to take some action  $A$ , which may causally affect some part of the environment history, i.e., the history of states,  $H$ . We will call that part of the history the agent's causal future  $H_f$ . Furthermore, the agent may be causally influenced by some other part of the environment history, which we will call the agent's causal past  $H_p$ .  $H$  may contain information other than  $H_f$  and  $H_p$ , which we will assume to be independent of  $A$ .<sup>6</sup> The *overseer*, physically realized by, e.g., some module physically attached to the agent or

<sup>5</sup> There is a broad philosophical literature on whether causal relationships exist and whether they can be inferred in cases where the agent is part of the environment. See, e.g., the edited volume by Price and Corry (2007).

<sup>6</sup> For simplicity, we will ignore dependences not resulting from causation (Arntzenius 2010). For example, if you play against a copy, there is a logical dependence between your and your copy's decision. Even if you know a set of nodes in the causal graph that *d*-separates your and your copy's decision (e.g., if you

a human supervisor, observes the agent's action and partially, via some percept  $O$ , the state of the world<sup>7</sup>. The overseer then calculates the reward  $R$ . To set proper incentives to the agent, we will assume the overseer to know not only the action and observation, but also everything that the agent knows (cf. Christiano 2016). The overseer may also have access to some additional piece of information  $E_r$  about the way the reward is to be calculated.<sup>8</sup> Lastly, we assume that the sets of possible values of  $A$ ,  $O$  and  $E_r$  are finite.

In principle, the overseer could reward the agent in all kinds of ways. E.g., it could reward the agent "deontologically" (Alexander and Moore 2016) for taking a particular action independently of the consequences of taking that action. In this paper, we will assume that the reward estimates the value of some von Neumann–Morgenstern utility function  $U$  that only depends on states of the world. I use the capital  $U$  to indicate that the utility function, too, is a random variable (in the Bayesian sense). For simplicity's sake, we will, again, assume that the set of possible values of  $U$  is finite.

We will view  $U$  as representing the system designer's preferences over world states.<sup>9</sup> While other ways of assigning the reward are possible, this is certainly an attractive way of getting an approval-directed agent to achieve goals that we want it to achieve. After all, in real-world applications, we will usually care about the outcomes of the agent's decisions, such as whether a car has reached its destination in time or whether a human has been hurt.

The standard way of estimating  $U(H)$  (or any quantity for that matter) is the familiar conditional expectation. Thus, the overseer may compute the reward as

$$r = \mathbb{E}[U(H) \mid e_r, a, o], \quad (1)$$

---

Footnote 6 continued

know your common source code), the dependence persists. We exclude these dependences because such situations cannot be modeled by standard causal graphs.

However, we could adapt causal graphs to accommodate for these kinds of dependences. First, we could modify our definition of causality in such a way that dependence does imply causation, as has been proposed by Spohn (2003, 2012), Yudkowsky (2010) and others. For instance, we could model the dependence between the outputs of two instances of an algorithm by introducing a logical node as a common cause of the two. This logical node would then represent the output of the abstract algorithm that the two copies implement. While changes to the concept of causation may affect CDT's implied behavior, the results from this paper can be directly transferred to such modifications.

Alternatively, we could extend causal graphs to also include non-causal dependences (cf. Poellinger 2013). Such extension necessitates a new CDT formalism, so the proofs from this paper do not directly transfer to this case. That said, I expect our results to generalize given that both EDT and CDT would probably treat non-causal dependences on the action just like they treat causal arrows directed toward the action.

<sup>7</sup> Christiano (2014) does not define approval-directed agency formally, but judging from a comment he made at <https://medium.com/paulfchristiano/i-agree-that-the-key-feature-of-approval-directed-agents-is-that-the-causal-picture-is-736b4474910e>, he considers it crucial to his conception that the overseer only looks at the agent's action and does not observe the action's consequences (cf. the distinction introduced in Sect. 3).

<sup>8</sup> One reason for the overseer to have access to such additional information is that some of the human supervisor's values may not be expressible in a way that the approval-directed agent's algorithm can utilize (cf. Muehlhauser and Helm 2012, Sects. 3, 4, 5.3).

<sup>9</sup> Some have tried to modify the reward relative to the designer's preferences to make the reinforcement learning problem easier to solve (Sorg 2011), although Sutton and Barto (1998, Sect. 3.2) explicitly discourage such tricks in their reinforcement learning textbook.

where  $r$ ,  $a$ ,  $e_r$ , and  $o$  are values of  $R$ ,  $A$ ,  $E_r$ , and  $O$ , respectively.<sup>10</sup>

A causal decision theorist overseer agrees that after an action  $a$  is taken the right-hand side of Eq. 1 most accurately estimates how much utility is achieved. She merely thinks that this term should not be used to decide which action  $a$  to take in the first place.<sup>11</sup> However, this puts a causal decision theorist overseer in a peculiar situation. Whatever formula she uses to compute the reward will also be used by the reward-maximizing agent to decide which action to take. A causal decision theorist overseer might therefore worry (rightfully, as we will see) that providing rewards according to Eq. 1 will make the agent EDT-ish. Hence, she either has to incorrectly estimate how much utility was achieved; or live with the agent using an—in her mind—incorrect way of weighing her options. If she prefers the latter, she would reward according to Eq. 1. But arguably getting the agent to choose correctly is the overseer’s primary goal. Thus, she might prefer to compute the reward according to

$$r = \mathbb{E}[U(H) \mid e_r, do(a), o]. \quad (2)$$

Here,  $do(a)$  refers to Pearl’s do-calculus, where conditioning on  $do(a)$  roughly means intervening from outside the causal model to set  $A$  to  $a$ . For an introduction to the do-calculus, see Pearl (2009). Although a causal decision theorist overseer may prefer computing rewards according to Eq. 1, we will from now on say “the overseer uses CDT” if rewards are computed according to Eq. 2 and “the overseer uses EDT” if rewards are calculated according to Eq. 1.

An approval-directed agent is characterized by maximizing the reward it receives from the overseer.<sup>12</sup> However, decision theory offers us, again, (at least) two different expected values, the regular expected value of EDT

$$\mathbb{E}[R \mid a], \quad (3)$$

and CDT’s causal expected value

$$\mathbb{E}[R \mid do(a)]. \quad (4)$$

<sup>10</sup> At first sight this may be confusing to some readers, because in reinforcement learning, utility sometimes refers to expected cumulative reward (Russell and Norvig 2010, Chap. 17, 21), although others use the term *value function* instead (Sutton and Barto 1998, Sect. 3.7). Here,  $U$  does not refer to utility in that sense but in the decision-theoretical sense of representing intrinsic values. So, in the present case, we have two “layers” of goals: first, the agent maximizes the reward  $r$ . Second, the agent as incentivized by the overseer’s way of calculating rewards maximizes utility  $U(H)$ .

One cause of confusion is that in model applications of reinforcement learning, the reward function possesses full knowledge of the world state and thus does not require the use of the expectation operator.

<sup>11</sup> If the disagreement in Newcomb’s problem is to be about different theories of rational choice (EDT, CDT and so forth) rather than the predictive abilities of “the being”, Omega or the psychologist, then after requesting both boxes a proponent of two-boxing has to believe that she will probably receive only \$1000. Causal and evidential decision theorists agree that regular conditional expectation is the correct way of updating one’s beliefs about the state of the world after an action has been taken (cf. the distinction between “acts” and “actions” in Pearl 2009 Sect. 4.1.1).

<sup>12</sup> In reinforcement learning, some have proposed alternative optimization targets that incorporate, e.g., risk aversion (García and Fernández 2015, Sect. 3).

We leave the interesting question of which (if any) decision theory describes the behavior of current reinforcement learning algorithms to future research<sup>13</sup> and in the following assume that the agent is known to implement either CDT or EDT.

### 3 The conflict of the decision theories of agent and overseer

When viewed together with the overseer, our agent may now be seen as containing two decision theories, one for computing the reward and one in the algorithm that tries to find the action to maximize that reward. These decision theories may not always be the same. Given this potential discrepancy, the question is which of the two decision theories prevails, i.e., for which configurations of the two decision theories the overall agent acts like a CDT agent and for which it acts like an EDT agent w.r.t.  $U$ .

As it turns out, the answer to this question depends on the decision problem in question. In particular, it depends on whether the overseer updates its estimate of  $U(H)$  primarily based on the action taken by the agent or on its observation of the environment.

For illustration, consider two versions of Newcomb's problem. In both versions, the predictor is equally reliable—e.g., correct with 90% probability—and the potential box contents are the same—e.g., the standard \$1K and \$1M. As usual, the content of the opaque box cannot be causally influenced by one's decision. In the first version, the overseer eventually sees the payoff, i.e., how much money the agent has made. In this case, as soon as the money is observed, the overseer's estimate of  $U(H)$  becomes independent of the agent's action. More generally,  $O$  may tell the overseer so much about  $U(H)$  that it becomes independent of  $A$  even if  $U(H)$  is not yet fully observed. That is,

$$\mathbb{E}[U(H) \mid e_r, a, o] = \mathbb{E}[U(H) \mid e_r, o] \quad (5)$$

and

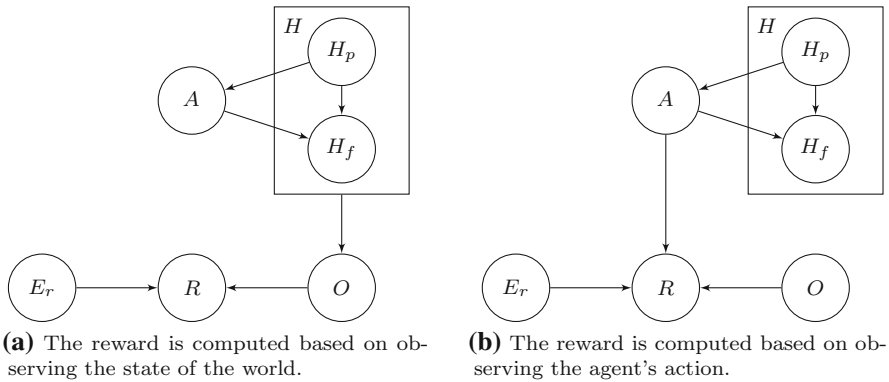
$$\mathbb{E}[U(H) \mid e_r, do(a), o] = \mathbb{E}[U(H) \mid e_r, o] \quad (6)$$

for all  $e_r$ ,  $a$  and  $o$ . Note that neither of these two implies the other.<sup>14</sup> Intuitively speaking, these two mean that the reward is ultimately determined by  $U(H)$ .

In the second version of Newcomb's problem, the monetary payoff is not observed but covertly invested into increasing the agent's utility function. Only the agent's choice can then inform the overseer about  $U(H)$ . Formally, it is both

<sup>13</sup> For preliminary work on this question, see Mayer et al. (2016), Oesterheld (2018a) and perhaps Albert and Heiner (2001).

<sup>14</sup> We give a brief justification of this claim. If all of  $a$ 's causal influence on  $H$  can be discerned from  $O$ , then, of course,  $a$  could still be diagnostically relevant for one's estimate of  $U(H)$ . The other direction is more complicated. The idea is that Eq. 5 can be true if the causal and non-causal implications of  $a$  exactly cancel each other out. An example is a version of Newcomb's problem in which one-boxing ensures with certainty that both boxes contain the same amount of money. Then if  $O$  and  $E_r$  do not contain any information, the expected value of two-boxing and one-boxing is the same and so learning of the action is irrelevant for estimating  $U(H)$ . However, two-boxing is causally better than one-boxing, so Eq. 6 is violated.



**Fig. 2** Two different ways in which the overseer can calculate the reward

$$\mathbb{E}[U(H) | e_r, a, o] = \mathbb{E}[U(H) | e_r, a] \quad (7)$$

and

$$\mathbb{E}[U(H) | e_r, do(a), o] = \mathbb{E}[U(H) | e_r, do(a)]. \quad (8)$$

Intuitively speaking, these two equations mean that the reward is not determined by  $U(H)$  but by what the overseer believes  $U(H)$  will be given  $a$  or  $do(a)$ .

Again, we assume that this is known to the agent. An example class of cases is that in which the agent's decisions are correlated with those of agents in far-away parts of the environment (cf. Treutlein and Oesterheld 2017; Oesterheld 2018b). The two versions are depicted in Fig. 2.

Of course, these are only the two extremes from the set of all possible situations. In real-world Newcomb-like scenarios, the overseer may also draw some information from both sources. Nonetheless, it seems useful to understand the extreme cases, as this may also help us understand mixed ones.

In the following subsections, we will show that in the first type, the decision theory of the agent is decisive, whereas in the second type, the overseer's decision theory is<sup>15</sup>. Roughly, the reason for that is the following: As noted earlier, the reward in the first type depends directly on  $U(H)$ . Thus, the agent will try to maximize  $U(H)$  according to its own decision theory. In the second type, the overseer takes the agent's action  $a$  and then considers what either  $a$  or  $do(a)$  says about  $U(H)$ . Thus, the agent has to pay careful attention to whether the overseer uses EDT's or CDT's expected value.

We prove this formally by considering all possible configurations of the type of the problem, the overseer's decision theory and the agent's decision theory. While we will limit our analysis to EDT and CDT, the results can easily be generalized to variants of these that arise from modifying the causal model or conditional credence distribution (e.g. Yudkowsky 2010; "Disposition-based decision theory"; Spohn 2012; Dohrn 2015). The analysis is summarized in Table 1.

<sup>15</sup> The dominance of the overseer's decision theory in the second type of Newcomb's problem is mentioned (though not proven) by Christiano (2014, Sect. "Avoid lock-in").



**Table 1** An overview of the results of the calculations in Sect. 3

Type of Newcomb problem	Agent's DT	Overseer's DT	Resulting DT
First type	CDT	EDT	CDT
		CDT	CDT
	EDT	EDT	EDT
		CDT	EDT
Second type	CDT	EDT	EDT
		CDT	CDT
	EDT	EDT	EDT
		CDT	CDT

### 3.1 First type

#### 3.1.1 The EDT agent

The EDT agent judges its action by

$$\mathbb{E}[R | a]. \tag{9}$$

If the overseer calculates regular conditional expectation, then it is

$$\mathbb{E}[R | a] = \mathbb{E}[\mathbb{E}[U(H) | E_r, O, a] | a] \tag{10}$$

$$= \mathbb{E}[U(H) | a], \tag{11}$$

where the last line is due to what is sometimes called the law of total expectation (LTE) or the tower rule (see, e.g., Ross 2007, Sect. 3.4; Billingsley 1995, Theorem 34.4). Intuitively, you cannot expect that gaining more evidence (i.e.,  $E_r$  and  $O$  in addition to  $a$ ) moves your expectation of  $U(H)$  into any particular direction.

Because the overseer knows more than the agent, we will need this rule in all of the following derivations. Its application makes it hard to generalize these results to other decision theories, since LTE does not apply if the two decision theories do not both compute a form of expected utility.

Equations 10 and 11 show that if the overseer computes regular expected value and the agent maximizes the reward according to EDT, then the agent as a whole maximizes  $U$  according to EDT.

If the overseer computes CDT's expected value, it is

$$\mathbb{E}[R | a] = \mathbb{E}[\mathbb{E}[U(H) | E_r, do(a), O] | a] \tag{12}$$

$$= \sum_{e_r, o} P(e_r, o | a) \cdot \mathbb{E}[U(H) | e_r, do(a), o] \tag{13}$$

$$\stackrel{\text{eq. 5 and 6}}{=} \sum_{e_r, o} P(e_r, o | a) \cdot \mathbb{E}[U(H) | e_r, a, o] \tag{14}$$

$$= \mathbb{E} [\mathbb{E} [U(H) \mid E_r, a, O] \mid a] \tag{15}$$

$$\stackrel{\text{LTE}}{=} \mathbb{E} [U(H) \mid a] \tag{16}$$

### 3.1.2 The CDT agent

The CDT agent judges its action by

$$\mathbb{E} [R \mid do(a)]. \tag{17}$$

If the overseer uses regular expected value (EDT), then

$$\mathbb{E} [R \mid do(a)] = \mathbb{E} [\mathbb{E} [U(H) \mid a, O, E_r] \mid do(a)] \tag{18}$$

$$= \sum_{e_r, o} P(e_r, o \mid do(a)) \cdot \mathbb{E} [U(H) \mid a, o, e_r] \tag{19}$$

$$\stackrel{\text{eq. 5 and 6}}{=} \sum_{e_r, o} P(e_r, o \mid do(a)) \cdot \mathbb{E} [U(H) \mid do(a), o, e_r] \tag{20}$$

$$= \mathbb{E} [\mathbb{E} [U(H) \mid do(a), O, E_r] \mid do(a)] \tag{21}$$

$$\stackrel{\text{LTE}}{=} \mathbb{E} [U(H) \mid do(a)] \tag{22}$$

Learning about an intervention  $do(a)$  cannot always be treated in the same way as learning about other events. Hence, the application of the law of total expectation is not straightforward. However,  $P(\cdot \mid do(x))$  is always a probability distribution. Because the law of total expectation applies to all probability distributions, it also applies to ones resulting from the application of the do-calculus.

If the overseer uses CDT’s expected value, then

$$\mathbb{E} [R \mid do(a)] = \mathbb{E} [\mathbb{E} [U(H) \mid E_r, O, do(a)] \mid do(a)] \tag{23}$$

$$\stackrel{\text{LTE}}{=} \mathbb{E} [U(H) \mid do(a)]. \tag{24}$$

## 3.2 Second type

### 3.2.1 The EDT agent

The EDT agent judges its actions by

$$\mathbb{E} [R \mid a]. \tag{25}$$

If the overseer is based on regular conditional expectation (EDT), then it is again

$$\mathbb{E} [R \mid a] = \mathbb{E} [\mathbb{E} [U(H) \mid E_r, a] \mid a] \tag{26}$$

$$\stackrel{\text{LTE}}{=} \mathbb{E} [U(H) \mid a]. \tag{27}$$

If the overseer is based on CDT-type expectation, then

$$\mathbb{E}[R | a] = \mathbb{E}[\mathbb{E}[U(H) | E_r, do(a)] | a] \quad (28)$$

$$= \sum_{e_r} P(e_r | a) \cdot \mathbb{E}[U(H) | do(a), e_r] \quad (29)$$

$$= \sum_{e_r} P(e_r) \cdot \mathbb{E}[U(H) | do(a), e_r] \quad (30)$$

$$= \sum_{e_r} P(e_r | do(a)) \cdot \mathbb{E}[U(H) | do(a), e_r] \quad (31)$$

$$= \mathbb{E}[\mathbb{E}[U(H) | E_r, do(a)] | do(a)] \quad (32)$$

$$\stackrel{\text{LTE}}{=} \mathbb{E}[U(H) | do(a)]. \quad (33)$$

### 3.2.2 The CDT agent

The CDT agent judges actions by

$$\mathbb{E}[R | do(a)]. \quad (34)$$

Because of Rule 2 in Theorem 3.4.1 of Pearl (2009, Sect. 3.4.2) applied to the causal graph of Fig. 2b, it is

$$\mathbb{E}[R | do(a)] = \mathbb{E}[R | a]. \quad (35)$$

Thus, the analysis of the CDT agent is equivalent to that of the EDT agent.

## 4 Conclusion

In this paper, we have taken a step to map reinforcement learning architectures onto decision theories. We found that in Newcomb-like problems, if the overseer rewards the agent purely on the basis of the agent's action, then the overall system's behavior is determined by the decision theory implicit in the overseer's reward function. If the overseer judges the agent based on looking at the world, however, then the agent's decision theory is decisive.

This has implications for how we should design approval-directed agents. For instance, if we would like to leave decision-theoretical judgements to the overseer, we must ensure that the overseer assigns rewards before making new observations about the world state (cf. Christiano 2014, Sect. "Avoid lock-in"). Of course, this makes the reward less accurate and may thus slow down the agent's learning process. If we want the overseer to look at both the world and the agent's action, then we need to align both the overseer's and the agent's decision theory.

Much more research is left to be done at the intersection of decision theory and artificial intelligence. For instance, what (if any) decision theories describe the way modern reinforcement learning algorithms maximize reward? Do the results of this paper generalize to sequential decision problems? Moving away from the reinforcement learning

framework, what decision theories do other frameworks in AI implement? What about decision theories other than CDT and EDT?

**Acknowledgements** I am indebted to Max Daniel, Johannes Treutlein, Tom Everitt, Lukas Gloor, Sören Mindermann, Brian Tomasik and Tobias Baumann for valuable comments and discussions.

**Open Access** This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

## References

- Achen, C. H., & Bartels, L. M. (2016). *Democracy For realists. Why elections do not produce responsive government*. Princeton: Princeton University Press.
- Ahmed, A. (2014). *Evidence, decision and causality*. Cambridge: Cambridge University Press.
- Albert, M., & Heiner, R. A. (2001). *An indirect-evolution approach to Newcomb's problem*. CSLE discussion paper, no. 2001-01. [https://www.econstor.eu/bitstream/10419/23110/1/2001-01\\_newc.pdf](https://www.econstor.eu/bitstream/10419/23110/1/2001-01_newc.pdf). Accessed 22 Feb 2019.
- Alexander, L., & Moore, M. (2016). Deontological ethics. In E. N. Zalta (Ed.), *The Stanford encyclopedia of philosophy*. Winter 2016. Metaphysics Research Lab, Stanford University. <https://plato.stanford.edu/archives/win2016/entries/ethics-deontological/>. Accessed 22 Feb 2019.
- Almond, P. (2010). *On causation and correlation part 1: Evidential decision theory is correct*. [https://casparosterheld.files.wordpress.com/2016/12/almond\\_edt\\_1.pdf](https://casparosterheld.files.wordpress.com/2016/12/almond_edt_1.pdf). Accessed 22 Feb 2019.
- Armstrong, S. (2011). *Anthropic decision theory*. Future of Humanity Institute. [arXiv: 1110.6437](https://arxiv.org/abs/1110.6437).
- Arntzenius, F. (2008). No regrets, or: Edith Piaf revamps decision theory. *Erkenntnis*, 68(2), 277–297. <https://doi.org/10.1007/s10670-007-9084-8>.
- Arntzenius, F. (2010). Reichenbach's common cause principle. In E. N. Zalta (Ed.), *The Stanford encyclopedia of philosophy*. Fall 2010. Metaphysics Research Lab, Stanford University.
- Aumann, R. J., Hart, S., & Perry, M. (1997). The absent-minded driver. *Games and Economic Behavior*, 20, 102–116.
- Billingsley, P. (1995). *Probability and measure* (3rd ed.). Hoboken: Wiley.
- Bostrom, N. (2014a). *Hail mary, value porosity, and utility diversification*. <http://www.nickbostrom.com/papers/porosity.pdf>. Accessed 22 Feb 2019.
- Bostrom, N. (2014b). *Superintelligence. Paths, dangers, strategies* (1st ed.). Oxford: Oxford University Press.
- Briggs, R. (2017). Real-life Newcomb problems? In *Talk at the 1st workshop on decision theory & the future of artificial intelligence in Cambridge, UK*.
- Cavalcanti, E. G. (2010). Causation, decision theory, and Bell's theorem: A quantum analogue of the Newcomb problem. *The British Journal for the Philosophy of Science*, 61(3), 569–597. <https://doi.org/10.1093/bjps/axp050>.
- Christiano, P. (2014). *Model-free decisions*. <https://ai-alignment.com/model-free-decisions-6e6609f5d99e>. Accessed 22 Feb 2019.
- Christiano, P. (2016). *Adequate oversight*. <https://ai-alignment.com/adequate-oversight-25fadf1edce9>. Accessed 22 Feb 2019.
- Dohrn, D. (2015). Egan and agents: How evidential decision theory can deal with Egan's dilemma. *Synthese*, 192(6), 1883–1908.
- Doyle, J. (1992). Rationality and its roles in reasoning. *Computational Intelligence*, 8(2), 376–409.
- Eells, E. (1981). Causality, utility, and decision. *Synthese*, 48(2), 295–329.
- Everitt, T., Leike, J., & Hutter, M. (2015). Sequential extensions of causal and evidential decision theory. In T. Walsh (Ed.), *Algorithmic decision theory: 4th international conference, ADT 2015, Lexington, KY, USA, September 27–30, 2015, Proceedings* (pp. 205–221). Springer. [https://doi.org/10.1007/978-3-319-23114-3\\_13](https://doi.org/10.1007/978-3-319-23114-3_13).
- Fisher, J. C. *Disposition-based decision theory*. <https://casparosterheld.files.wordpress.com/2019/02/dbdt.pdf>. Accessed 22 Feb 2019.

- García, J., & Fernández, F. (2015). A comprehensive survey on safe reinforcement learning. *Journal of Machine Learning Research*, 16, 1437–1480.
- Gibbard, A., & Harper, W. L. (1981). Counterfactuals and two kinds of expected utility. In W. L. Harper, R. Stalnaker, & G. Pearce (Eds.), *Ifs. Conditionals, belief, decision, chance and time* (Vol. 15). The University of Western Ontario Series in Philosophy of Science. A series of books in philosophy of science, methodology, epistemology, logic, history of science, and related fields (pp. 153–190). Springer. [https://doi.org/10.1007/978-94-009-9117-0\\_8](https://doi.org/10.1007/978-94-009-9117-0_8).
- Greene, P. (2018). Success-first decision theories. In A. Ahmed (Ed.), *Newcomb's problem*. Classic Philosophical Arguments. Cambridge University Press. <https://doi.org/10.1017/9781316847893.007>.
- Gustafsson, J. E. (2011). A note in defence of ratificationism. *Erkenntnis*, 75(1), 147–150.
- Hintze, D. (2014). *Problem class dominance in predictive dilemmas*. <http://intelligence.org/files/ProblemClassDominance.pdf>. Accessed 22 Feb 2019.
- Horgan, T. (1981). Counterfactuals and Newcomb's problem. *The Journal of Philosophy*, 78(6), 331–356.
- Hutter, M. (2005). *Universal artificial intelligence. sequential decision based on algorithmic probability*. In W. Brauer, G. Rozen-berg, & A. Salomaa (Eds.), *Texts in theoretical computer science*. Springer.
- Joyce, J. M. (1999). *The foundations of causal decision theory. Cambridge studies in probability, induction, and decision theory*. Cambridge: Cambridge University Press.
- Kuhn, S. (2017). Prisoner's dilemma. In E. N. Zalta (Ed.), *The Stanford encyclopedia of philosophy*. Spring 2017. Metaphysics Research Lab, Stanford University. <https://plato.stanford.edu/archives/spr2017/entries/prisoner-dilemma/>. Accessed 22 Feb 2019.
- Kumar, R. (2017). New work for decision theorists. In *Talk at the 1st workshop on decision theory & the future of artificial intelligence in Cambridge, UK*.
- Ledwig, M. (2000). *Newcomb's problem*. Ph.D. thesis, University of Constance. <https://kops.uni-konstanz.de/bitstream/handle/123456789/3451/ledwig.pdf>. Accessed 22 Feb 2019.
- Legg, S., & Hutter, M. (2007). Universal intelligence: A definition of machine intelligence. *Minds and Machines*, 17(4), 391–444.
- Lewis, D. (1981). Causal decision theory. *Australasian Journal of Philosophy*, 59(1), 5–30.
- Mayer, D., Feldmaier, J., & Shen, H. (2016). Reinforcement learning in conflicting environments for autonomous vehicles. In *International workshop on robotics in the 21st century: Challenges and promises*. arXiv: 1610.07089.
- Meacham, C. J. G. (2010). Binding and its consequences. *Philosophical Studies*, 149(1), 49–71. <https://doi.org/10.1007/s11098-010-9539-7>.
- Muehlhauser, L., & Helm, L. (2012). *Intelligence explosion and machine ethics*. Machine Intelligence Research Institute. <https://intelligence.org/files/IE-ME.pdf>. Accessed 22 Feb 2019.
- Nozick, R. (1969). Newcomb's problem and two principles of choice. In N. Rescher, et al. (Eds.), *Essays in honor of Carl G. Hempel* (pp. 114–146). Berlin: Springer.
- Oesterheld, C. (2018a). *Doing what has worked well in the past leads to evidential decision theory*. <https://casparoesterheld.files.wordpress.com/2018/01/learning-dt.pdf>. Accessed 22 Feb 2019.
- Oesterheld, C. (2018b). Newcomb's problem, the Prisoner's dilemma and large universes: A consideration for consequentialists. In *Talk at the 15th conference of the international society for utilitarian studies*. Karlsruhe Institute of Technology (KIT), July 24–26, 2018.
- Pearl, J. (2009). *Causality. Models, reasoning, and inference*. Cambridge: Cambridge University Press.
- Piccione, M., & Rubinstein, A. (1997). On the interpretation of decision problems with imperfect recall. *Games and Economic Behavior*, 20, 3–24.
- Poellinger, R. (2013). *Unboxing the concepts in Newcomb's paradox: Causation, prediction, decision*. [http://philsci-archive.pitt.edu/9887/7/newcomb\\_in\\_ckps.pdf](http://philsci-archive.pitt.edu/9887/7/newcomb_in_ckps.pdf). Accessed 22 Feb 2019.
- Price, H. (1986). Against causal decision theory. *Synthese*, 67, 195–212.
- Price, H. (2012). Causation, chance, and the rational significance of supernatural evidence. *Philosophical Review*, 121(4), 483–538.
- Price, H., & Corry, R. (Eds.). (2007). *Causation, physics, and the constitution of reality: Russell's republic revisited*. Oxford: Oxford University Press.
- Ross, S. M. (2007). *Introduction to probability models* (9th ed.). Cambridge: Academic Press.
- Russell, S., & Norvig, P. (2010). *Artificial intelligence. A modern approach* (3rd ed.). London: Pearson Education, Inc.
- Skyrms, B. (1982). Causal decision theory. *The Journal of Philosophy*, 79(11), 695–711.
- Soares, N. (2014a). *Newcomblike problems are the norm*. <http://mindingourway.com/newcomblike-problems-are-the-norm/>. Accessed 22 Feb 2019.

- Soares, N. (2014b). *Why Ain't you rich?* <https://intelligence.org/2014/10/07/nate-soares-talk-aint-rich/>. Accessed 22 Feb 2019.
- Soares, N., & Fallenstein, B. (2014a). *Aligning superintelligence with human interests: A technical research agenda*. Technical report. 2014-8. Machine Intelligence Research Institute. <https://intelligence.org/files/TechnicalAgenda.pdf>. Accessed 22 Feb 2019.
- Soares, N., & Fallenstein, B. (2014b). *Toward idealized decision theory*. Technical report 2014-7. Machine Intelligence Research Institute. [arXiv: 1507.01986](https://arxiv.org/abs/1507.01986).
- Soares, N., & Levinstein, B. A. (2017). Cheating death in damascus. In *Formal epistemology workshop (FEW) 2017*. University of Washington, Seattle, USA. <https://intelligence.org/files/DeathInDamascus.pdf>. Accessed 22 Feb 2019.
- Sorg, J. D. (2011). *The optimal reward problem: Designing effective reward for bounded agents*. PhD thesis, University of Michigan. [https://deepblue.lib.umich.edu/bitstream/handle/2027.42/89705/jdsorg\\_1.pdf](https://deepblue.lib.umich.edu/bitstream/handle/2027.42/89705/jdsorg_1.pdf). Accessed 22 Feb 2019.
- Spohn, W. (2003). Dependency equilibria and the causal structure of decision and game situation. *Homo Oeconomicus*, 20, 195–255.
- Spohn, W. (2012). Reversing 30 years of discussion: Why causal decision theorists should one-box. *Synthese*, 187(1), 95–122.
- Sutton, R. S., & Barto, A. G. (1998). *Reinforcement learning: An introduction*. Cambridge: MIT Press.
- Treutlein, J. (2018). How the decision theory of Newcomb like problems differs between humans and machines. In *Talk at the 2nd workshop on decision theory & the future of artificial intelligence in Munich, Germany*.
- Treutlein, J., & Oesterheld, C. (2017). *A wager for evidential decision theory*. Unpublished manuscript.
- Wedgwood, R. (2013). Gandalph's solution to the Newcomb problem. *Synthese*, 190(14), 2643–2675. <https://doi.org/10.1007/s11229-011-9900-1>.
- Weirich, P. (2016). Causal decision theory. In *The Stanford encyclopedia of philosophy*. Spring 2016.
- Yudkowsky, E. (2010). *Timeless decision theory*. The Singularity Institute. <http://intelligence.org/files/TDT.pdf>. Accessed 22 Feb 2019.

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.