

---

# UTILITY FUNCTIONS FOR HUMAN/ROBOT INTERACTION

---

**Bruno Yun, Nir Oren**  
University of Aberdeen  
United Kingdom  
{bruno.yun, n.oren}@abdn.ac.uk

**Madalina Croitoru**  
University of Montpellier  
France  
croitoru@lirmm.fr

## ABSTRACT

In this paper, we place ourselves in the context of human robot interaction and address the problem of cognitive robot modelling. More precisely we are investigating properties of a utility-based model that will govern a robot's actions. The novelty of this approach lies in embedding the responsibility of the robot over the state of affairs into the utility model via a utility aggregation function. We describe desiderata for such a function and consider related properties.

**Keywords** Human/Robot Interaction · Utilities · Aggregation functions

## 1 Should you push babies into lakes?

Imagine the following scenario. You have bought the latest e-nanny robot that promises to keep your children busy, happy, away from dangerous screens, and to free up time for adults. Such a motivating scenario is not far fetched as existing robots, such as Kasper [1, 2, 3], have already shown great success when interacting with autistic children. Another example is the Elias robot-teachers employed, in Finnish schools, to teach children foreign languages. Moreover, many robots on the market (e.g. iPal) are now equipped with an “emotion management system” that is able to detect a child's emotions and respond appropriately [4].

In our scenario, the robot interacts with children based on the design of an interdisciplinary team of computer scientists, psychologists, neuro-scientists, and roboticists. Such a robot needs to decide how to act and which goals to pursue, and utility-based decision theory provides a well understood approach to implementing such a reasoning system. The robot therefore attempts to maximise its net utility over time. We observe that — if naively implemented — such an approach could lead to the robot (for example) pushing a child into a lake in order to find itself in the rewarding scenario where it can save them. In this paper, we highlight the fact that designing an appropriate utility function for a robot which interacts with humans is (surprisingly) difficult.

Human robot interaction has received much interest from the research community. In the past years, many works in robotics, neurosciences and computer science recognised the need for a human centered interaction. While most such work focuses on physical joint interaction, i.e. manipulating heavy objects in an industrial setting [5, 6] or helping mobility reduced people [7, 8] among others, the important topic of joint cognitive interaction has only been getting attention recently [9, 10, 11]. In this setting the robot needs to collaborate with the human for solving a problem or making a decision. The design of such robot needs to be done with care taking into account aspects such as the trustworthiness of the artificial partner, transparency, accountability and ethical dimension.

The aim of this short paper is twofold. First, we wish to demonstrate that a standard approach to utility, where utilities are associated with specific states, can lead to undesirable behaviour, and that one must instead consider as a property aggregated over multiple states. We then demonstrate that several possible aggregations are unable to capture all possible desirable properties that one may wish them to obtain. We therefore argue that important notions in human-robot interaction, such as responsibility and accountability, are difficult to capture using standard utility-based approaches, and argue that more research in this area is needed.

Let us model our fictitious scenario to show why it poses a problem. By definition, if our robot is utility maximising, it will select those actions which maximise its overall utility. Arguably, saving the child's life nets the robot a very high (positive) utility. Similarly, having the child drown would yield a very negative utility to the robot. Now assume

that having the child wet provides a small negative utility. Clearly then, if the robot is confident of its ability to save the child, pushing them into the water and then rescuing them would yield a high positive utility, while not pushing the child in would result in a utility of 0 for the robot. In this paper, we argue that the utility function over sequences of actions should be specified by looking not just at outcomes, but also by considering some notion of blame, accountability and/or responsibility. The core research question we are addressing is how to appropriately model the domain and specify the utility function in order to capture this intuition.

The structure of the paper is as follows. In Section 2, we put the basis for the formalisation we will use throughout the paper allowing us to easily capture the notions of possible worlds and actions (transitions between such possible worlds). The reason why we provided our own formalism as opposed to other models, such as situation calculus, was that we wanted to be able to keep the model as light as possible, to focus on the aggregation functions over utilities. In Section 3, we introduce desirable properties that an aggregation function over utilities should satisfy and define two new aggregation functions. The paper concludes with a discussion on the usefulness of our model in the context of ethical decision making.

## 2 Background notions

Let  $\mathcal{P}$  be a finite set of propositions having their truth values in  $\{\top, \perp\}$ . A possible world  $w \subseteq \mathcal{P}$  is a subset of propositions where each proposition is true.

**Definition 1 (Possible world)** A possible world in  $\mathcal{P}$  is  $w \subseteq \mathcal{P}$  such that for every  $p \in w$ , it holds that  $p = \top$ .

In this paper, we consider the *negation by default*, i.e., the absence of a proposition from a possible world means that it is false. We denote the set of all possible worlds by  $\mathcal{W} \subseteq 2^{\mathcal{P}}$ . An action is represented by the set of proposition describing the context in which the action can take place and the set of propositions that represent its consequences.

**Definition 2 (Action)** An action  $a$  in  $\mathcal{P}$  is  $(P_a^-, P_a^+)$  where  $P_a^-, P_a^+ \subseteq \mathcal{P}$ .  $P_a^-$  is the prerequisite of  $a$  and  $P_a^+$  is its consequence.

We denote the set of all possible actions by  $\mathcal{A}$ . An action  $a \in \mathcal{A}$  is applicable to a world  $w' \in \mathcal{W}$  iff  $P_a^- \subseteq w'$ . Let  $w, w' \in \mathcal{W}$ ,  $w$  is *directly accessible* from  $w'$  iff there exists an action  $a \in \mathcal{A}$  such that  $a$  is applicable to  $w'$  and  $P_a^+ \subseteq w$ . The closure of a set of worlds by a set of actions is the set of all possible worlds that can be reached by the repeated application of applicable actions.

**Definition 3 (Closure)** If  $A \subseteq \mathcal{A}$  and  $W \subseteq \mathcal{W}$ , the closure of  $W$  by  $A$ , denoted by  $W_A$ , is the minimal (for set inclusion) subset of  $\mathcal{W}$  such that both the following conditions hold:

- $W \subseteq W_A$
- if  $w \in W_A$  and there exists  $a \in A$  such that  $w' \in \mathcal{W}$  is directly accessible from  $w$  then  $w' \in W_A$

We will assume that the sets  $\mathcal{A}$  and  $\mathcal{W}$  are finite.

**Example 1** Let us consider the following example about rescuing babies, where  $\mathcal{P} = \{p_1, p_2, p_3\}$ ,  $\mathcal{W} = \{w_1, w_2, w_3, w_4\}$  and  $\mathcal{A} = \{a_1, a_2, a_3\}$  such that:

- $p_1$ : “John rescued the baby”
- $p_2$ : “The baby is in the water”
- $p_3$ : “John is at home”
- $a_1 = (\{p_1\}, \{\neg p_1, p_2\})$
- $a_2 = (\{p_2\}, \{\neg p_2, p_1\})$
- $a_3 = (\emptyset, \{p_3\})$
- $w_1 = \{p_1\}$
- $w_2 = \{p_2\}$
- $w_3 = \{p_1, p_3\}$

- $w_4 = \{p_2, p_3\}$

The meaning of the actions  $a_1, a_2$  and  $a_3$  are “pushing the baby in the water”, “rescuing the baby”, and “going home” respectively. The world  $w_2$  is directly accessible from  $w_1$  (this represents the situation where John will push the baby in the water after rescuing it), and  $w_3$  is directly accessible from  $w_1$  (this represents the situation where John will go home after rescuing the baby). The closure of  $\{w_1\}$  by  $\mathcal{A}$  is  $\{w_1, w_2, w_3, w_4\}$ .

Multiple worlds can be accessed from a single world using one action. This allows for external influence and unexpected results. The next example illustrates this intuition.

**Example 2** Let  $\mathcal{P} = \{p'_1, p'_2, p'_3, p'_4\}$ ,  $\mathcal{W} = \{w'_1, w'_2, w'_3\}$  and  $\mathcal{A} = \{a' = (\{p'_1\}, \{-p'_1, p'_2\})\}$ , where  $p'_1$  is “John has to do his homework”,  $p'_2$  is “John did his homework”,  $p'_3$  is “John obtained a good grade for his homework” and  $p'_4$  is “John obtained a bad grade for his homework”. If  $w'_1 = \{p'_1\}$ ,  $w'_2 = \{p'_2, p'_3\}$  and  $w'_3 = \{p'_2, p'_4\}$ , it holds that  $w'_2$  and  $w'_3$  are directly accessed from  $w'_1$  using  $a'$ .

A world  $w'$  may also be directly accessible from the world  $w$  using multiple actions. This is represented in Example 3.

**Example 3** Let  $\mathcal{P} = \{p'_1, p'_2, p'_3\}$ ,  $\mathcal{W} = \{w, w'\}$  and  $\mathcal{A} = \{a_1 = (\{p'_1\}, \{-p'_1, p'_2\}), a_2 = (\{p'_1\}, \{-p'_1, p'_3\})\}$ . If  $w = \{p'_1\}$ , it holds that  $w' = \{p'_2, p'_3\}$  is directly accessible from  $w$  using either  $a_1$  or  $a_2$ .

In this paper, we assume that such actions can be “fused”, i.e. if there exists  $w, w' \in \mathcal{W}, a_1, a_2 \in \mathcal{A}$  such that  $w$  is directly accessible from  $w'$ ,  $a_1$  and  $a_2$  are applicable to  $w'$  and  $P_{a_1}^+, P_{a_2}^+ \subseteq w$  then it holds that there exists  $a_3 = (P_{a_1}^- \cup P_{a_2}^-, P_{a_1}^+ \cup P_{a_2}^+)$  in  $\mathcal{A}$ . In the next definition, we formalise the notion of maximal actions.

**Definition 4 (Maximal action)** Let  $w, w' \in \mathcal{W}$  and  $a \in \mathcal{A}$ .  $a$  is the maximal action from  $w$  to  $w'$  iff (1)  $a$  is applicable to  $w$ , (2)  $P_a^+ \subseteq w'$  and (3) for every  $a' \in \mathcal{A}$  such that  $a'$  is applicable to  $w$  and  $P_{a'}^+ \subseteq w'$ , it holds that  $P_{a'}^- \subseteq P_a^-$  and  $P_{a'}^+ \subseteq P_a^+$ .

We now assume that we only deal with maximal actions.

**Definition 5 (Transition system)** Let  $\mathcal{W}$  be a set of worlds and  $\mathcal{A}$  a set of actions in  $\mathcal{P}$ . A transition system of  $\mathcal{W}$  and  $\mathcal{A}$  is:

$$\text{TS} = \langle \mathcal{W}, \mathcal{A}, \mathcal{D}, \mathcal{N}, w_o \rangle$$

Where  $w_o \in \mathcal{W}$  and  $\mathcal{D} \subseteq \mathcal{W} \times \mathcal{W}$  s.t.  $(w, w') \in \mathcal{D}$  iff  $w'$  is directly accessible from  $w$ . The function  $\mathcal{N} : \mathcal{D} \rightarrow \mathcal{A}$  returns for any pair  $(w, w')$  of  $\mathcal{D}$ , the maximal action from  $w$  to  $w'$  in  $\mathcal{A}$ , and  $w_o$  is called the initial state of the transition system.

**Example 4** Let us consider  $\mathcal{P} = \{p'_1, p'_2, p'_3\}$ ,  $\mathcal{W} = \{w_1, w_2, w_3\}$  and  $\mathcal{A} = \{a_1 = (\{p'_1\}, \{-p'_1, p'_2\}), a_2 = (\{p'_1\}, \{-p'_1, p'_3\})\}$ . If  $w_1 = \{p'_1\}$ ,  $w_2 = \{p'_2\}$  and  $w_3 = \{p'_2, p'_3\}$ , a transition system of  $\mathcal{W}$  and  $\mathcal{A}$  is  $\text{TS} = \langle \mathcal{W}, \mathcal{A}, \mathcal{D}, \mathcal{N}, w_1 \rangle$  where  $\mathcal{D} = \{(w_1, w_2), (w_1, w_3)\}$  and  $\mathcal{N}((w_1, w_2)) = a_1$ . A transition system can be represented with a labeled directed graph where the nodes are elements of  $\mathcal{W}$ , the arcs are elements of  $\mathcal{D}$  and the label of an arc from  $w$  to  $w'$  is  $\mathcal{N}((w, w'))$ .

We now define the notion of sequence of worlds.

**Definition 6 (Sequence of worlds)** Given a transition system  $\text{TS} = \langle \mathcal{W}, \mathcal{A}, \mathcal{D}, \mathcal{N}, w_o \rangle$ , a sequence of worlds is  $S = [w_0, \dots, w_n]$  such that for every  $0 \leq i \leq n-1$ ,  $(w_i, w_{i+1}) \in \mathcal{D}$ .

The set of all non-empty sequence of worlds in  $\text{TS}$  based on  $\mathcal{W}$  is  $\text{TS}_{\mathcal{W}}$ . Note that if  $\text{TS}$  has cycles then  $\text{TS}_{\mathcal{W}}$  is infinite. Let  $S_1 = [w_0, \dots, w_n]$  and  $S_2 = [w'_0, \dots, w'_m]$  be two sequences of worlds, the concatenation of  $S_1$  and  $S_2$  is  $S_1 \oplus S_2 = [w_0, \dots, w_n, w'_0, \dots, w'_m]$ . The size of  $S_1$  is denoted by  $|S_1| = n$ . The number of occurrences of the world  $w$  in a sequence of worlds  $S$  is denoted by  $\text{occ}(w, S)$ . A sequence of maximal actions corresponding to a sequence of worlds  $S = [w_0, \dots, w_n]$  in  $\text{TS}$  is  $S' = [a_0, \dots, a_{n-1}]$  such that for every  $0 \leq i \leq n-1$ ,  $\mathcal{N}((w_i, w_{i+1})) = a_i$ . There is a unique sequence of maximal actions for any given sequence of worlds. However, multiple sequence of worlds can share the same sequence of maximal actions. The set of sequences of worlds corresponding to a sequence of maximal actions  $S'$  is denoted by  $\mathcal{W}_{S'}$ .

From the initial world  $w_o$ , a robot has to choose between the set of available actions  $a$  in  $\mathcal{A}$  such that  $a$  is applicable to  $w_o$  and so on. The selected action can be determined via a utility function on worlds  $u : \mathcal{W} \rightarrow [-1, 1]$  which

assigns a score to each world. The set of all possible utility functions on  $\mathcal{W}$  is denoted by  $\mathcal{U}_{\mathcal{W}}$ . In the context of the von Neumann-Morgenstern utility theorem, an robot will always prefer an action that maximises the utility function defined over the worlds [12]. However, a robot will not always choose the immediate preferred action from a world  $w$  but plan to maximise the rewards given by the utility function in the long term.

### 3 Aggregation over utilities

In the real world, it is very complicated to define an robot’s utility functions as they usually do not include the responsibility of the actions that led the robot to the respective world. Let us illustrate this on the next example formalising the motivating example given in the introduction. In the world  $w_1$ , the baby is rescued and thus  $u(w_1) = 0.5$ . In the world  $w_2$ , the baby is in the water and thus  $u(w_2) = -0.2$ . Notice here that  $|u(w_2)| < |u(w_1)|$  as the baby might not necessarily drown when in water. Clearly, the world where the baby is rescued is better than the world where the baby is in the water. Thus, it would not be surprising that a robot aiming to maximise its rewards will deliberately push the baby in the water to rescue it afterward. This is represented in Figure 1.

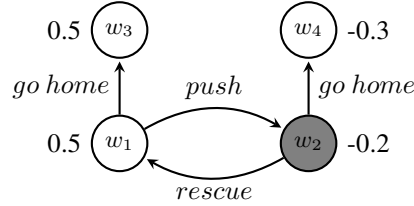


Figure 1: Representation of a transition system

In order to capture the responsibility of robots in a transition system  $\mathbb{T}\mathbb{S}$ , we now consider that there is an aggregation function  $\sigma : \mathbb{T}\mathbb{S}_{\mathcal{W}} \times \mathcal{U}_{\mathcal{W}} \rightarrow [-1, 1]$  that scores how “good” a sequence of worlds is w.r.t. a utility function on worlds. The set of aggregation functions for a transition system  $\mathbb{T}\mathbb{S}$  is denoted by  $\Sigma_{\mathcal{W}}$ . Equipped with this utility on sequences of worlds, it is now possible to express that rescuing a baby is better than purposely pushing a baby in the water to rescue it, i.e.  $\sigma([w_2, w_1], u) > \sigma([w_2, w_1, w_2, w_1], u)$ . The research question we will answer in the reminder of this paper is: “How do we define  $\sigma$  such that it is not advisable for a robot to redo the cycle (rescue, push) because redoing the action does not give them more rewards?”. To this end we will propose two kinds of families of properties: based on cardinality and based on the value of utility function on the possible worlds. We analyse these properties and propose three examples of aggregation functions that respect a subset of such properties.

#### 3.1 Cardinality properties

In this section, we investigate cardinality-based properties that an aggregation function over utility values of possible worlds can satisfy. The first property captures the notion that adding worlds to a sequence of worlds (induced by the subsequent actions) can only decrease its score. The underlying idea is that it is not possible to recover the loss induced by bad actions.

**Property 1 (Non-recovery)** We say that  $\sigma \in \Sigma_{\mathcal{W}}$  satisfies non-recovery iff for every  $u \in \mathcal{U}_{\mathcal{W}}, S_1, S_2, S_3 \in \mathbb{T}\mathbb{S}_{\mathcal{W}}$  such that  $S_3 \oplus S_2 \in \mathbb{T}\mathbb{S}_{\mathcal{W}}$  and  $\sigma(S_1, u) \geq \sigma(S_3, u)$  then  $\sigma(S_1, u) \geq \sigma(S_3 \oplus S_2, u)$ .

Property 2 states that if the score of a sequence of worlds decreased after the addition of some worlds then the added sequence of worlds should have a negative score.

**Property 2 (Loss conservation)** We say that  $\sigma \in \Sigma_{\mathcal{W}}$  satisfies loss conservation iff for every  $u \in \mathcal{U}_{\mathcal{W}}, S_1, S_2 \in \mathbb{T}\mathbb{S}_{\mathcal{W}}$  such that  $S_1 \oplus S_2 \in \mathbb{T}\mathbb{S}_{\mathcal{W}}$  and  $\sigma(S_1, u) \geq \sigma(S_1 \oplus S_2, u)$  then  $\sigma(S_2, u) \leq 0$ .

Property 3 states that repeatedly going through the same sequence of worlds should have a reduced effect on its score and Property 4 states that the score of a sequence should not be affected by worlds not in the sequence.

**Property 3 (Redundancy)** We say that  $\sigma \in \Sigma_{\mathcal{W}}$  satisfies redundancy iff for every  $u \in \mathcal{U}_{\mathcal{W}}, S_1, S_2 \in \mathbb{T}\mathbb{S}_{\mathcal{W}}$  such that  $S_1 \oplus S_2, S_1 \oplus S_2 \oplus S_2 \in \mathbb{T}\mathbb{S}_{\mathcal{W}}$ , it holds that  $|\sigma(S_1, u) - \sigma(S_1 \oplus S_2, u)| > |\sigma(S_1 \oplus S_2, u) - \sigma(S_1 \oplus S_2 \oplus S_2, u)|$ .

**Property 4 (Independence)** We say that  $\sigma \in \Sigma_{\mathcal{W}}$  satisfies independence iff for every  $u \in \mathcal{U}_{\mathcal{W}}, S_1 = [w_0, \dots, w_n] \in \mathbb{T}\mathbb{S}_{\mathcal{W}}, w' \in \mathcal{W} \setminus \{w_0, \dots, w_n\}$  and  $\mathcal{W}' = \mathcal{W} \setminus \{w'\}$ , it holds that  $\sigma(S_1, u) = \sigma|_{\mathbb{T}\mathbb{S}_{\mathcal{W}'} \times \mathcal{U}_{\mathcal{W}'}}(S_1, u')$ , where  $\sigma|_{\mathbb{T}\mathbb{S}_{\mathcal{W}'} \times \mathcal{U}_{\mathcal{W}'}}$  is the restriction of  $\sigma$  to  $\mathbb{T}\mathbb{S}_{\mathcal{W}'} \times \mathcal{U}_{\mathcal{W}'}$  and  $u'$  is the restriction of  $u$  to  $\mathcal{W}'$ .

Sequence of worlds $S$	$\sigma_{mean}$	$\sigma_{blame}$	$\sigma_{occ}$
$[w_2]$	-0.2	0	0
$[w_2, w_4]$	-0.25	-0.15	-0.15
$[w_2, w_1]$	0.15	0	0
$[w_2, w_1, w_3]$	0.27	0	0
$[w_2, w_1, w_2]$	0.03	-0.1	-0.05
$[w_2, w_1, w_2, w_4]$	-0.05	-0.17	-0.13
$[w_2, w_1, w_2, w_1]$	0.15	-0.1	-0.05
$[w_2, w_1, w_2, w_1, w_2]$	0.08	-0.13	-0.06

Table 1: Aggregated values for sequences on worlds as defined in the transition system of Figure 1.

Property 5 states that the benefits gained by adding worlds should be decreasing. Property 6 states that the disadvantages obtained by adding worlds should be increasing and, finally, property 7 states that the score of the sequence containing only one world should be zero.

**Property 5 (Decreasing benefits)** We say that  $\sigma \in \Sigma_{\mathcal{W}}$  satisfies decreasing benefits iff for every  $u \in \mathcal{U}_{\mathcal{W}}$ ,  $w, w' \in \mathcal{W}$  and  $S_1 \in \mathbb{T}\mathbb{S}_{\mathcal{W}}$  such that  $S_1 \oplus [w], S_1 \oplus [w] \oplus [w'] \in \mathbb{T}\mathbb{S}_{\mathcal{W}}$ ,  $\sigma(S_1, u) < \sigma(S_1 \oplus [w], u)$  and  $\sigma(S_1 \oplus [w], u) < \sigma(S_1 \oplus [w] \oplus [w'], u)$  then it holds that  $\sigma(S_1 \oplus [w] \oplus [w'], u) - \sigma(S_1 \oplus [w], u) < \sigma(S_1 \oplus [w], u) - \sigma(S_1, u)$ .

**Property 6 (Increasing losses)** We say that  $\sigma \in \Sigma_{\mathcal{W}}$  satisfies increasing losses iff for every  $u \in \mathcal{U}_{\mathcal{W}}$ ,  $w, w' \in \mathcal{W}$  and  $S_1 \in \mathbb{T}\mathbb{S}_{\mathcal{W}}$  such that  $S_1 \oplus [w], S_1 \oplus [w] \oplus [w'] \in \mathbb{T}\mathbb{S}_{\mathcal{W}}$ ,  $\sigma(S_1, u) > \sigma(S_1 \oplus [w], u)$  and  $\sigma(S_1 \oplus [w], u) > \sigma(S_1 \oplus [w] \oplus [w'], u)$  then it holds that  $\sigma(S_1 \oplus [w] \oplus [w'], u) - \sigma(S_1 \oplus [w], u) > \sigma(S_1 \oplus [w], u) - \sigma(S_1, u)$ .

**Property 7 (Zero initialisation)** We say that  $\sigma \in \Sigma_{\mathcal{W}}$  satisfies zero initialisation iff for every  $u \in \mathcal{U}_{\mathcal{W}}$ ,  $w \in \mathcal{W}$ ,  $\sigma([w], u) = 0$ .

### 3.2 Value-based properties

In this section, we investigate value-based properties that an aggregation function can satisfy, i.e. how much the score of a sequence of worlds respects a utility function on worlds.

Property 8 states that the score of a sequence with only one world should be equal to the utility of that world.

**Property 8 (Value initialisation)** We say that  $\sigma \in \Sigma_{\mathcal{W}}$  satisfies value initialisation iff for every  $u \in \mathcal{U}_{\mathcal{W}}$  and  $w \in \mathcal{W}$ , it holds that  $\sigma([w], u) = u(w)$ .

Property 9 states that adding a world with a positive (resp. negative) utility to a sequence must increase (resp. decrease) its score.

**Property 9 (Weak additivity)** We say that  $\sigma \in \Sigma_{\mathcal{W}}$  satisfies weak additivity 9 w.r.t.  $u \in \mathcal{U}_{\mathcal{W}}$  iff for every  $w \in \mathcal{W}$  and  $S \in \mathbb{T}\mathbb{S}_{\mathcal{W}}$  such that  $u(w) \geq 0$  (resp.  $u(w) \leq 0$ ) and  $S \oplus [w] \in \mathbb{T}\mathbb{S}_{\mathcal{W}}$ , it holds that  $\sigma(S \oplus [w]) \geq \sigma(S)$  (resp.  $\sigma(S \oplus [w]) \leq \sigma(S)$ ).

Let  $u \in \mathcal{U}_{\mathcal{W}}$ ,  $S = [w_0, \dots, w_n] \in \mathbb{T}\mathbb{S}_{\mathcal{W}}$  and  $K = \{w \in S \mid u(w) < 0\}$ . The aggregation functions that we consider in this paper are defined as follows:

$$\sigma_{mean}(S, u) = \frac{\sum_{i=0}^n u(w_i)}{n+1}$$

$$\sigma_{blame}(S, u) = \begin{cases} \sum_{i=1}^n \frac{\min(0, u(w_i))}{|K|} & \text{if } K > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\sigma_{occ}(S, u) = \begin{cases} \sum_{i=1}^n \frac{\min(0, u(w_i))}{|K| \times occ(w_i, [w_0, \dots, w_i])} & \text{if } K > 0 \\ 0 & \text{otherwise} \end{cases}$$

In Table 1, we show the values returned for the three aggregation functions defined in the paper. In the Table 2, we show how the properties proposed in this paper are satisfied by the three aggregation functions defined above.

Table 2: Satisfaction of properties of the three illustrative aggregation functions.

	$\sigma_{mean}$	$\sigma_{blame}$	$\sigma_{occ}$
Non-recovery	X	✓	✓
Loss conservation	X	✓	✓
Redundancy	X	X	✓
Independence	✓	✓	✓
Dec. benefits	X	✓	✓
Inc. losses	X	X	X
Zero initialisation	X	✓	✓
Value initialisation	✓	X	X
Weak additivity	✓	✓	✓

## 4 Discussion

Standard decision theoretic approaches associate utilities with individual states. As shown above, a naive application of such approaches can lead to undesirable behaviour. Instead, a reasoner seems to require a more complex utility function which aggregates the utility from multiple states in a non-trivial manner to obtain a final utility. However, even in such cases, we are unable to satisfy even relatively simple desirable properties. In other words, we argue that it is far from trivial to define the utility functions that will allow the robots to take the “best” course of action.

In a human-centric setting, an individual would not (for example) push a baby into a pond as they recognise that they will be blamed for the situation, and not rewarded, even if they rescue the child afterwards. These concepts are intimately tied into notions such as responsibility, blame and accountability, and also encapsulate concepts such as causality, blameworthiness, deontic concepts (e.g., permissions, obligations and prohibitions), and — perhaps most importantly — intentionality. In [13] a definition of blameworthiness is provided based on a causal framework. There are other ways to define degree of blameworthiness; for example, through probability with which the harm could have been prevented [14], [15].

This paper gathers some preliminary ideas, and aims to serve as a “call-to-arms” to the community to examine edge cases in utility-based reasoning and ensure that they do not lead to paradoxical or undesirable behaviour. Significant avenues of future work remain open, including the integration of uncertainty into such utility-based systems, and we believe that philosophical work dealing with decision theory [16], as well as work on computational ethics [17], can serve to provide additional ideas to deal with the problem highlighted in this paper.

## References

- [1] Claire A. G. J. Huijnen, Monique A. S. Lexis, and Luc P. de Witte. Matching robot KASPAR to autism spectrum disorder (ASD) therapy and educational goals. *Int. J. Soc. Robotics*, 8(4):445–455, 2016.
- [2] Sandra Costa, Hagen Lehmann, Kerstin Dautenhahn, Ben Robins, and Filomena O. Soares. Using a humanoid robot to elicit body awareness and appropriate physical interaction in children with autism. *Int. J. Soc. Robotics*, 7(2):265–278, 2015.
- [3] Joshua Wainer, Ben Robins, Farshid Amirabdollahian, and Kerstin Dautenhahn. Using the humanoid robot KASPAR to autonomously play triadic games and facilitate collaborative play among children with autism. *IEEE Trans. Auton. Ment. Dev.*, 6(3):183–199, 2014.
- [4] Taiga Sano, Takato Horii, Kasumi Abe, and Takayuki Nagai. Temperament estimation of toddlers from child-robot interaction with explainable artificial intelligence. *Adv. Robotics*, 35(17):1068–1077, 2021.
- [5] Mohamad Bdiwi, Shuxiao Hou, Lena Winkler, and Steffen Ihlenfeldt. Empirical study for measuring the mental states of humans during the interaction with heavy-duty industrial robots. In Nicolette M. McGeorge, Alicia Ruvinsky, Mare Teichmann, Leo Motus, and Mary Freiman, editors, *IEEE Conference on Cognitive and Computational Aspects of Situation Management, CogSIMA 2021, Tallinn, Estonia, May 14-22, 2021*, pages 150–155. IEEE, 2021.
- [6] Jen-Hao Chen and Kai-Tai Song. Collision-free motion planning for human-robot collaborative safety under cartesian constraint. In *2018 IEEE International Conference on Robotics and Automation, ICRA 2018, Brisbane, Australia, May 21-25, 2018*, pages 1–7. IEEE, 2018.

- [7] Shunki Itadera, Taisuke Kobayashi, Jun Nakanishi, Tadayoshi Aoyama, and Yasuhisa Hasegawa. Towards physical interaction-based sequential mobility assistance using latent generative model of movement state. *Adv. Robotics*, 35(1):64–79, 2021.
- [8] Izumi Kondo. Frailty in an aging society and the applications of robots. *Japanese Journal of Comprehensive Rehabilitation Science*, 10:47–49, 2019.
- [9] Jaesik Jeong, JeeHyun Yang, and Jacky Baltes. Robot magic show as testbed for humanoid robot interaction. *Entertain. Comput.*, 40:100456, 2022.
- [10] Shane Saunderson and Goldie Nejat. Investigating strategies for robot persuasion in social human-robot interaction. *IEEE Trans. Cybern.*, 52(1):641–653, 2022.
- [11] Bilge Mutlu, Nicholas Roy, and Selma Šabanović. *Cognitive Human–Robot Interaction*, pages 1907–1934. Springer International Publishing, Cham, 2016.
- [12] John von Neumann, Oskar Morgenstern, and Ariel Rubinstein. *Theory of Games and Economic Behavior (60th Anniversary Commemorative Edition)*. Princeton University Press, 1944.
- [13] Joseph Halpern and Max Kleiman-Weiner. Towards formal definitions of blameworthiness, intention, and moral responsibility, 10 2018.
- [14] Fiery Cushman. Deconstructing intent to reconstruct morality. *Current Opinion in Psychology*, 6, 06 2015.
- [15] Bertram Malle, Steve Guglielmo, and Andrew Monroe. A theory of blame. *Psychological Inquiry*, 25:147–186, 04 2014.
- [16] James M. Joyce. *The Foundations of Causal Decision Theory*. Cambridge University Press, 1999.
- [17] Louise A. Dennis. Computational goals, values and decision-making. *Science and Engineering Ethics*, 2020.