# A Modal Characterization of
# Defeasible Deontic Conditionals and Conditional Goals

## Craig Boutilier
Department of Computer Science
University of British Columbia
Vancouver, British Columbia
CANADA, V6T 1Z2
**email:** cebly@cs.ubc.ca

## Abstract

We explore the notions of permission and obligation and their role in knowledge representation, especially as guides to action for planning systems. We first present a simple conditional deontic logic (or more accurately a *preference logic*) of the type common in the literature and demonstrate its equivalence to a number of modal and conditional systems for default reasoning. We show how the techniques of conditional default reasoning can be used to derive *factual preferences* from conditional preferences. We extend the system to account for the effect of beliefs on an agent's obligations, including beliefs held by default. This leads us to the notion of a *conditional goal*, goals toward which an agent should strive according to its belief state. We then extend the system (somewhat naively) to model the *ability* of an agent to perform actions. Even with this simple account, we are able to show that the deontic slogan "make the best of a bad situation" gives rise to several interpretations or *strategies* for determining goals (and actions). We show that an agent can improve its decisions and focus its goals by making observations, or increasing its knowledge of the world. Finally, we discuss how this model might be extended and used in the planning process, especially to represent planning under uncertainty in a qualitative manner.

## 1 Introduction

In the usual approaches to planning in AI, a planning agent is provided with a description of some state of affairs, a *goal state*, and charged with the task of discovering (or performing) some sequence of actions to achieve that goal. This notion of goal can be found in the earliest work on planning (see [24] for a survey) and persists in more recent work on intention and commitment [8]. In most realistic settings, however, an agent will frequently encounter goals that it cannot achieve. As pointed out by Doyle and Wellman [28] an agent possessing only simple goal descriptions has no guidance for choosing an alternative goal state toward which it should strive. Some progress has been made toward providing systems with the power to express sub-optimal goals. Recently, Haddawy and Hanks [12] have provided a very restricted model for partial fulfillment of deadline goals. Their goals are propositions that must be made true by a given deadline. If a goal cannot be met, *utility* is assigned to partial fulfillment, for instance, making "more" of the proposition true or getting close to the deadline. But we should not expect goals to be so well-behaved in general. There are many situations in which, say, missing a deadline by some small margin is worse than ignoring the task altogether; or in which an agent that fails to satisfy one conjunct of a conjunctive goal should not attempt to satisfy the rest of the goal. Indeed, goals can change drastically depending on the situation in which an agent finds itself. If it is raining, an agent's goal state is one it which it has its umbrella. If it's not raining, it should leave its umbrella home. Moreover, goals do not change simply because some ideal goal cannot be satisfied: a robot may be *able* to leave its umbrella home when it is raining. *Circumstances* play a crucial role in determining goals. We want the ability to express such goals as "Pick up a shipment at wholesaler X by 9AM," but state exceptions like "If it's a holiday, go to wholesaler Y." Permitting explicit exceptions to goal statements allows goals to be expressed naturally. Clearly, just as conditional plans allow circumstances to dictate the appropriate sequence of actions for satisfying a goal, *conditional goals* allow circumstances to dictate just what goal state an agent "desires."

It is profitable to think of an agent's goal as an *obligation*. It *ought* to do what it can to ensure the goals set for it are achieved. There have been a number of systems proposed to model the notions of permission and obligation, what ought to be the case [26, 7]. The earliest modal logics for obligation were unsuccessful, for they seem unable to represent *conditional obligations*. Conditional deontic logics have been developed to account for

the fact that obligations can vary in different circumstances [13, 17, 25]. This indicates that such systems are suitable for the representation of conditional goals.

In this paper we present a simple modal semantics for defeasible deontic conditionals and extend this in several ways to account for the influence of default knowledge, incomplete knowledge and ability on obligations and goals. Many systems have a rather complex semantics that attempts to capture the defeasibility of such conditionals [18, 15]. In Section 2, we present a family of simple modal logics that can be used to represent defeasible deontic conditionals. These logics deal only with preferences on deontic alternatives. This is a crucial abstraction for it allows us to concentrate solely on the structure of deontic preferences, ignoring important but separate concerns that must be considered in any true logic of obligation (see below). We show that our logic captures the semantic system of Hansson [13], and demonstrate that one can capture conditional obligations with a unary modal operator, contrary to "conventional wisdom" in the conditional logic literature.

Recently, the mechanisms of default reasoning have been applied to problems in deontic logic [14, 16]. Indeed, our conditional/modal logic was originally developed for applications in default reasoning [1, 2]; we show in Section 3 how techniques for conditional default reasoning can be used for deontic reasoning.

In Section 4, we turn our attention to the influence of beliefs on an agent's goals. Since goals are conditional on circumstances, the goals adopted by an agent will depend crucially (and defeasibly) on the agent's beliefs about the world. We (tentatively) adopt the deontic strategy that an agent should strive to bring about the best situations consistent with its beliefs. This provides a different conception of goals from that usually found in AI (e.g., the goals of [8]). Included in this analysis is a model of default beliefs. This takes the first steps toward defining a "qualitative decision theory."

Since an agent cannot be expected to do things beyond its ability, we introduce a naive model of ability and action in Section 5. We describe goals to be those propositions within an agent's control that ensure optimal outcomes. It turns out, however, that this model is sufficient to render useless the deontic slogan "bring about the best situations possible." Incomplete knowledge of the world gives rise to different *strategies* for deriving goals, in the tradition of game theory. We describe several of these and how *observations* can be used to improve decisions and focus goals.

Finally, we conclude with a brief summary and some future directions in which this work may be pursued, including applications to planning under uncertainty.

## 2 Conditional Preferences and Deontic Logic

Deontic logics have been proposed to model the concepts of obligation and permission [26, 27, 7, 13]. It is clear that a sentence "It ought to be the case that $A$" has non-extensional truth conditions, for what is *actually* the case need not have any influence on what (ideally) ought to be. Deontic logics usually have a modal connective $O$ where the sentence $O\alpha$ means that the proposition $\alpha$ is obligatory, or it ought to be the case that $\alpha$.[1] Semantically, the truth of such statements are usually determined with respect to a given set of *ideal* possible worlds, $\alpha$ being obligatory just when $\alpha$ is true in all ideal situations. While this set of ideal worlds (obligation) is often taken to be determined by some code of moral or ethical conduct, this need not be the case. "Obligations" may simply be the goals imposed on an agent by its designer, and the ideal worlds simply those in which an agent fulfills its specified duties.[2] We take the notion of ideality or preference in what follows to be determined by any suitable metric, and allow "obligation" to refer to the satisfaction of design goals, moral imperatives or anything similar.

It has been widely recognized that deontic logic, defined using a unary modal connective $O$, has certain limitations. In particular, it is difficult to represent *conditional obligations* [7]. For this reason, conditional deontic logic (CDL) has been introduced to represent the dependence of obligations on context [27]. Obligation is then represented by a two-place conditional connective. The sentence $O(B|A)$ is interpreted as "It ought to be that $B$ given $A$" or "If $A$ then it is obligatory that $B$," and indicates a conditional obligation to do $B$ in circumstances $A$. These logics can be interpreted semantically using an ordering on worlds that ranks them according to some notion of preference or ideality [13, 17]. Such a ranking satisfies $O(B|A)$ just in case $B$ is true at all most preferred of those worlds satisfying $A$. Thus, we can think of $B$ as a *conditional preference* given $A$. Once $A$ is true, the best an agent can do is $B$.

In this section, we present a simple modal logic and semantics for the representation of conditional obligations. The logic CO (and related systems) are presented below for this purpose. The presentation is brief and we refer to [1, 4] for further technical details and an axiomatization.

### 2.1 The Bimodal Logic CO

We assume a propositional bimodal language $\mathbf{L}_B$ over a set of atomic propositional variables $\mathbf{P}$, with the usual

---

[1]When discussing the obligations of a particular agent, we will often say that the agent has an obligation to "do $\alpha$." Though $\alpha$ is a proposition, we take "do $\alpha$" to be some (unspecified) action that brings about $\alpha$.

[2]Cohen and Levesque [8] analyze goals similarly. An agent has a goal $\alpha$ just in case $\alpha$ is true at all worlds that are "goal accessible."

Figure 1: A CO-model

classical connectives and two modal operators $\Box$ and $\overset{\leftarrow}{\Box}$. Our Kripkean possible worlds semantics for deontic preference will be based on the class of *CO-models*, triples of the form $M = \langle W, \leq, \varphi \rangle$ where $W$ is a set of possible worlds, or deontic alternatives, $\varphi$ is a valuation function mapping $\mathbf{P}$ into $2^W$ ($\varphi(A)$ is the set of worlds where $A$ is true), and $\leq$ is a reflexive, transitive connected binary relation on $W$.[3] Thus $\leq$ imposes a total preorder on $W$: $W$ consists of a set of $\leq$-equivalence classes, these being totally ordered by $\leq$. We take $\leq$ to represent an ordering of deontic preference: $w \leq v$ just in case $v$ is at least as preferable as $w$. This ordering is taken to reflect the preferences of an agent about complete situations, however these are to be interpreted (e.g., an ordering of moral acceptability, personal utility, etc.). We will sometimes speak of preferred situations as being more ideal or more acceptable than others. Each equivalence class, or *cluster* of worlds, consists of a set of equally preferred situations. Figure 1 illustrates a typical CO-model. The truth conditions for the modal connectives are

1. $M \models_w \Box\alpha$ iff for each $v$ such that $w \leq v$, $M \models_v \alpha$.

2. $M \models_w \overset{\leftarrow}{\Box}\alpha$ iff for each $v$ such that $w \not\leq v$, $M \models_v \alpha$.

$\Box\alpha$ is true at a world $w$ just in case $\alpha$ is true at all worlds at least as preferred as $w$, while $\overset{\leftarrow}{\Box}\alpha$ holds just when $\alpha$ holds at all less preferred worlds. The dual connectives are defined as usual: $\Diamond\alpha \equiv_{\mathrm{df}} \neg\Box\neg\alpha$ means $\alpha$ is true at some equally or more preferred world; and $\overset{\leftarrow}{\Diamond}\alpha \equiv_{\mathrm{df}} \neg\overset{\leftarrow}{\Box}\neg\alpha$ means $\alpha$ is true at some less preferred world. $\overset{\leftrightarrow}{\Box}\alpha \equiv_{\mathrm{df}} \Box\alpha \wedge \overset{\leftarrow}{\Box}\alpha$ and $\overset{\leftrightarrow}{\Diamond}\alpha \equiv_{\mathrm{df}} \Diamond\alpha \vee \overset{\leftarrow}{\Diamond}\alpha$ mean $\alpha$ is true at all worlds and at some world, respectively.

## 2.2 Deontic Conditionals

We now define a conditional connective $I(-|-)$ to express conditional preferences. Intuitively, $I(B|A)$ should hold just when $B$ holds at the most ideal worlds satisfying $A$. Of course, nothing in our models forces the existence of such minimal $A$-worlds (see Lewis [17] on the

---

[3]$\leq$ is *connected* iff $w \leq v$ or $v \leq w$ for each $v, w \in W$.

*Limit Assumption*). However, we may simply say that there should be some world satisfying $A \wedge B$ such that $A \supset B$ holds at all "accessible" (equally or more ideal) worlds. We also let the conditional hold vacuously when the antecedent $A$ is false in all situations. These truth conditions can be expressed in $\mathbf{L}_B$ as follows:

$$I(B|A) \equiv_{\mathrm{df}} \overset{\leftrightarrow}{\Box}\neg A \vee \overset{\leftrightarrow}{\Diamond}(A \wedge \Box(A \supset B)). \qquad (1)$$

$I(B|A)$ can be read as "In the most preferred situations where $A$ holds, $B$ holds as well," or "If $A$ then ideally $B$." This can be thought of, as a first approximation, as expressing "If $A$ then an agent ought to ensure that $B$," for making $B$ true (apparently) ensures an agent ends up in the best possible $A$-situation. There are problems with such a reading, as we discuss shortly; hence, we usually adopt the former reading. However, we will occasionally lapse and read $I(B|A)$ as "If $A$ then it ought to be that $B$." We note that an absolute preference $A$ can be expressed as $I(A|\top)$, or equivalently, $\overset{\leftrightarrow}{\Diamond}\Box A$. We abbreviate this as $I(A)$ ("ideally $A$"). We can also express the (analog of) the notion of conditional permission. If $\neg I(\neg A|B)$ holds, then in the most preferred $B$-situations it is not required that $\neg A$. This means there are ideal $B$-worlds where $A$ holds, or that $A$ is "tolerable" given $B$. We abbreviate this sentence $T(A|B)$. Loosely, we can think of this as asserting that an agent is *permitted* to do $A$ if $B$. Unconditional toleration is denoted $T(A)$ and stands for $\neg I(\neg A)$, or equivalently, $\overset{\leftrightarrow}{\Box}\Diamond A$.

Using CO we can express the conditional preferences involved in a number of classic deontic puzzles. Chisholm's paradox of contrary-to-duty imperatives [7] is one such puzzle. The preferences involved in the following account [18] cannot be adequately captured in with a unary modal deontic connective:

(a) It ought to be that Arabella buys a train ticket to visit her grandmother.

(b) It ought to be that if Arabella buys the ticket she calls to tell her she is coming.

(c) If Arabella does not buy the ticket, it ought to be that she does not call.

(d) Arabella does not buy the ticket.

We can, however, represent these sentences conditionally as $I(V)$, $I(C|V)$, $I(\neg C|\neg V)$ and $\neg V$. These give rise to no inconsistency in CO, and induce a natural ordering on worlds where only $V \wedge C$-worlds are most acceptable. Less preferred are certain $\neg V \wedge \neg C$-worlds, and still less preferred is any $\neg V \wedge C$-world:

$$VC < \overline{VC} < \overline{V}C$$

(The relative preference of $V \wedge \neg C$-worlds is left unspecified by this account, though we are assured that $VC < V\overline{C}$.) Notice that from this set we can derive

$I(C)$. More generally, CO satisfies the principle of *deontic detachment* [18]:

$$I(B|A) \wedge I(A) \supset I(B)$$

Another principle sometimes advocated in the deontic literature is that of *factual detachment*:

$$I(B|A) \wedge A \supset I(B)$$

This expresses the idea that if there is a conditional obligation to do $B$ given $A$, and $A$ is *actually* the case, then there is an actual obligation to do $B$. No extension of standard deontic logic can contain both principles as theorems [18]. Given our reading of $I(B|A)$, factual detachment should not be (and is not) valid in CO. However, clearly we require some mechanism for inferring "actual preferences" from factual statements and conditional preferences. In the next section we describe how techniques from conditional default reasoning can be used for just this purpose.

Another "puzzling" theorem of CO (as well as most conditional and standard deontic logics) is the following:

$$I(A) \supset I(A \vee B)$$

If we read $I(A)$ as "an agent is obligated to do $A$," this theorem seems somewhat paradoxical. From "It is obligatory that you help $X$" one can infer that "It is obligatory that you help $X$ or kill $X$." This suggests that one may fulfill this obligation by killing $X$ [13]. Of course, if we adhere strictly to the reading of $I$ as "In the most ideal situations you help $X$ or kill $X$" this is less problematic, for it does not suggest any means for fulfilling these obligations. Furthermore, seeing to it that $X$ is killed (presumably) removes one from the realm of ideal situations and did nothing to fulfill the original obligation (help $X$). In Sections 4 and 5 we suggest a means of capturing this distinction. This puzzle is closely related to the notion of *free choice permission*.

### 2.3 Representation Results

The idea of using an ordering on possible situations to reflect deontic preference and capture conditional obligations is not new. Hansson [13] provided a semantics for CDL that is much like ours.[4] While Hansson does not present an axiomatization of his system DSDL3, we can show that our modal semantics extends his and that a fragment of CO provides a sound and complete proof theory for his system. Let CO– denote the set of theorems in CO restricted to those sentences whose only non-classical connective is $I$.[5]

---

[4] One key difference is the fact that Hansson invokes the Limit Assumption: there must be a *minimal* (or most preferred) set of $A$-worlds for each satisfiable proposition $A$. This has no impact on the results below [4].

[5] Occurrences of $\square$ and $\overleftarrow{\square}$ must conform to the pattern in the definition of $I$.

**Theorem 1** $\vdash_{CO-} \alpha$  *iff*  $\models_{DSDL3} \alpha$.

Furthermore, the power of the second modal operator is not required. The results of [1, 3] show that the purely conditional fragment of CO can be captured using only the operator $\square$, that is, using the classical modal logic S4.3, simply by defining

$$I(B|A) \equiv_{\mathrm{df}} \square\neg A \vee \diamond(A \wedge \square(A \supset B)). \qquad (2)$$

Let S4.3– denote the conditional fragment of S4.3.

**Theorem 2** $\vdash_{S4.3-} \alpha$  *iff*  $\models_{DSDL3} \alpha$.

We note also that our system is equivalent to Lewis's conditional deontic logic VTA [17]. This shows that CDL, as conceived by Hansson and Lewis, can be captured using only a unary modal operator. The "trick", of course, lies in the fact that $\square$ is not treated as unconditional obligation.

Using a modal semantics of this sort suggests a number of generalizations of Lewis's logics. For instance, by using the modal logic S4, we can represent partially ordered or preordered preferences. Furthermore, we can define a version of the conditional that allows us to explicitly (and consistently) represent conflicting preferences of the form $I(B|A)$ and $I(\neg B|A)$ (see [1, 6]). The logic CO*, an extension of CO presented in [2, 4], is based on the class of CO-models in which every propositional valuation (every logically possible world) is represented. Using CO* we can ensure that each world is ranked according to our preference relation.

## 3 Defeasible Reasoning

While a standard modal logic (S4.3) suffices to represent Hansson's and Lewis's notion of conditional obligation, the added expressiveness of the logic CO can be used to great advantage. In particular, it allows us to express various assumptions about our premises and use the techniques of conditional default reasoning to infer actual preferences.

Loewer and Belzer [18] have criticized Lewis's semantics "since it does not contain the resources to express actual obligations and no way of inferring actual obligations from conditional ones." Clearly, in our example above, we should somehow be able to conclude that Arabella ought not call her grandmother; but we cannot infer logically in CO that $I(\neg C)$. This is to be expected, for the *actual fact* $\neg V$ does not affect the form taken by *ideal* situations. However, once $\neg V$ is realized, one ought to attempt to make the best of a bad situation. In other words, actual preferences (or obligations) should be simply those propositions true in the most preferred worlds that satisfy the actual facts.

Let $KB$ be a knowledge base containing statements of conditional preference and actual facts. Given that such facts actually obtain, the ideal situations are those most preferred worlds satisfying $KB$. This suggests a straightforward mechanism for determining actual preferences.

We simply ask for those $\alpha$ such that

$$\vdash_{CO} I(\alpha|KB)$$

If $KB$ is the representation of Chisholm's paradox above, we have that $I(\neg C|KB)$ is valid. Thus, we can quite reasonably model a type of factual detachment simply by using nested conditionals of this sort.

We notice that this is precisely the preliminary scheme for conditional default reasoning suggested in [9, 2]. This mechanism unfortunately has a serious drawback: seemingly *irrelevant* factual information can paralyze the "default" reasoning process. For instance, let $KB' = KB \cup \{R\}$, where $R$ is some distinct propositional atom (e.g., "it will rain"). We can no longer conclude that Arabella ought not call, for $I(\neg C|KB')$ is not valid. Intuitively, $R$ has nothing to do with Arabella's obligation to call, yet, from a logical perspective, there is nothing in the premises that guarantees that raining cannot affect Arabella's obligations. The following (incomplete) ordering is satisfies with the original premises:

$$VCR < \overline{VCR} < \overline{V}CR$$

Hence, it *could* be that Arabella should call if it's raining (whether she visits or not).

Several solutions have been proposed for the problem of irrelevance in default systems. We briefly describe one, Pearl's [22] System Z. Roughly, we want to assume that worlds are as preferred or as ideal as possible, subject to the constraints imposed by our theory of preferences. For instance, a world where it is raining, Arabella fails to buy a ticket and doesn't call her grandmother violates no more obligations (in $KB$) than a world where it is not raining and the other conditions obtain. It is consistent to assume that raining situations are no less ideal than non-raining situations. System Z provides a mechanism for determining the consequences of such assumptions, and it will allow us to conclude from $KB'$ that Arabella ought not call (when it is raining). Roughly, System Z chooses a preferred model from those satisfying $KB$. This model is the most "compact" model of $KB$, a model where worlds are "pushed down" in the preference ordering as far as possible (consistent with the constraints imposed by $KB$). For simple conditional theories[6] there is a unique preferred model, the *Z-model*. The Z-model for $KB$ (or $KB'$) above is

$$\{VCR, VC\overline{R}\} < \{\overline{VCR}, \overline{V}CR, V\overline{C}R, V\overline{C}R\} <$$
$$\{\overline{V}C\overline{R}, \overline{V}CR\}$$

Notice that $I(C|V \wedge R)$ is satisfied in this model. Also notice that visiting without calling is no "worse" than not visiting: since our premises do not specify exactly how bad failing to call is, it is assumed to be as "good"

as possible (though it cannot be ideal, since it violates the imperative $I(C|V)$).

Goldszmidt and Pearl [10, 11] have developed algorithms that compute (often efficiently) the conclusions that can be derived in this way. In [2] we describe how the expressive power of CO can be used to axiomatize the assumptions of System Z and describe an explicit preference relation on CO-models that captures the criterion of "compactness". Of course, any problems with a default reasoning scheme must be expected to carry over to our deontic approach. System Z has the drawback of failing to count situations that violate more conditional preferences as less ideal than those that violate fewer (at least, within a a fixed *priority level*). A simple solution has been proposed in [5] for the logic CO. Other solutions to the problem of irrelevance for conditional reasoning have also been proposed. We do not catalogue these here, but simply point out that deriving factual obligations from a conditional $KB$ has exactly the same structure as deriving default conclusions from a conditional $KB$ (see the next section); any problems and solutions for one will be identical for the other.

We note that recently other techniques for default reasoning have been proposed for deontic systems. Jones and Pörn [16] have also proposed an extension of their earlier deontic logics that uses some ideas from Delgrande's work in default reasoning. However, their system is rather complex and quite distinct from ours. Makinson [19] has made some preliminary suggestions for incorporating aspects of normality and agency in the formalization of obligations that relate to work in default reasoning. Horty [14] proposes representing imperatives and deriving obligations using Reiter's default logic. Situations are "scored" according to the imperatives they violate. This too is different from our approach, for we take preferences to be primitive and use these to derive (conditional) obligations. Naturally, constraints on a preference ordering can be derived from imperatives as well should we chose to view a statement $I(B|A)$ as an imperative. In fact, System Z can be viewed as ranking situations according to the *priority* of the imperatives violated by a situation. This is one advantage of the conditional approach over Horty's system: priorities on imperatives are induced by a conditional theory. In contrast, should one situation violate two rules of equal priority while a second violates just one, System Z views both situations as equally preferred (this is the drawback cited above), while Horty's system clearly "prefers" that fewer imperatives be violated. Again, the proposal in [5] for default reasoning can be seen as combining these two approaches.[7]

The view of conditionals (in the sense of CO) as im-

---

[6] Such theories may have propositions and conditionals of the form $I(B|A)$ where $B, A$ are propositional.

[7] While based on default logic, Horty's system shares some remarkable similarities because of the simple theories he uses and their connection to Poole's Theorist system, which in turn is related to our conditional logic [5].

peratives is tenable. However, imperatives reflect much more information than the preferences of an agent. Issues of action and ability must come into play, so it is hard to see just how a preference ordering can be derived from imperatives without taking into account such considerations. In the following sections, we explore these issues, showing how one might determine actions appropriate for an agent in various circumstances (the agent's "imperatives").

# 4 Toward a Qualitative Decision Theory

## 4.1 A Logic of Goals

While deontic logics concentrate on preferences related to notions of moral acceptability or the like, it is clear that preferences can come from anywhere. In particular, a system designer ought to be able to convey to an artificial agent the preferences according to which that agent ought to act. From such preferences (and other information) and agent ought to be able to derive *goals* and plan its actions accordingly.

Doyle and Wellman [28] have investigated a preferential semantics of goals motivated by a qualitative notion of utility. Roughly, $P$ is a goal just when any "fixed" $P$-situation is preferred to the corresponding $\neg P$-situation. A loose translation into our framework would correspond to a sentence schema $I(P|\alpha)$ for all propositions $\alpha$ where $\alpha \not\vdash \neg P$. This type of goal is unconditional in the sense that it ensures that $P$ is ideally true: $I(P)$. We can express conditional goals in their framework simply by fixing certain other propositions, but these can certainly not be defeasible.

Asserting $P$ as a goal in this sense is a very strong statement, in fact, so strong that very few goals will meet this criterion. For $P$ to be a goal, it must be that no matter what else is true the agent will be better off if $P$. We call such a goal *absolute* (whether conditional or unconditional). However, even such strong moral imperatives as "Thou shalt not kill" typically have exceptions (for instance, self-defense). Certainly, it is the case that we can postulate an absolute goal by considering all exceptional situations for a goal and making the goal conditional on the absence of exceptions. Unfortunately, just as with the qualification problem for default reasoning, this leads to goals that must be expressed in a very unnatural fashion. We should note that absolute goals have an advantage over the defeasible goals expressible in CO. If $Q$ is an absolute goal conditional on $P$, then once $P$ is known by an agent it is guaranteed that ensuring $Q$ is the best course of action. No matter what contingencies occur, $Q$ is better than $\neg Q$. However, as we will see below, such goals might provide very little guidance for appropriate behavior, especially in the presence of incomplete knowledge. Once again, very few realistic goals match this specification. To use an example developed below, an agent might have to decide whether or not to take its umbrella to work: if it rains, taking the umbrella

is best; if not, leaving the umbrella is best. Taking or leaving the umbrella cannot be an absolute goal for an agent, yet clearly there may be other considerations that make one or the other action a "real" goal. Such considerations include action, ability, expected outcomes and game-theoretic strategies.

## 4.2 Default Knowledge

As we have seen, the belief set of an agent provides the context in which its goals are determined. However, we should not require that goals be based only on "certain" beliefs, but on any reasonable default conclusions as well. For example, consider the following preference ordering with atoms $R$ (it will rain), $U$ (have umbrella) and $C$ (it's cloudy). Assuming $\overline{C} \wedge R$ is impossible, we might have the following preferences:

$$\{\overline{CRU}, C\overline{RU}\} < CRU < \{\overline{CR}U, C\overline{R}U\} < CR\overline{U}$$

Suppose, furthermore, that it usually rains when its cloudy. If $KB = \{C\}$, according to our notion of obligation in the last section, the agent's goals are $\overline{R}$ and $\overline{U}$. Ideally, the agent ought to ensure that it doesn't rain and that it doesn't bring its umbrella. Ignoring the absurdity of the goal $\overline{R}$ (we return to this in the next section), even the goal $\overline{U}$ seems to be wrong. Given $C$, the agent should *expect* $R$ and act accordingly.

Much like in decision theory, actions should be based not just on the utilities (preferences) over outcomes, but also on the likelihood (or typicality or normality) of outcomes. In order to capture this intuition in a qualitative setting, we propose a logic that has two orderings, one representing preferences over worlds and one representing the degree of *normality* or *expectation* associated with a world. The presentation is again brief and we refer to [6] for further motivation and technical details.

The logic QDT, an attempt at a qualitative decision theory, is characterized by models of the form $M = \langle W, \leq_P, \leq_N, \varphi \rangle$, where $W$ is a set of worlds (with valuation function $\varphi$), $\leq_P$ is a transitive, connected *preference ordering* on $W$, and $\leq_N$ is a transitive, connected *normality ordering* on $W$. We interpret $w \leq_P v$ as above, and take $w \leq_N v$ to mean $w$ is at least as *normal* a situation as $v$ (or is at least as *expected*). The submodels formed by restricting attention to either relation are clearly CO-models. The language of QDT contains four modal operators: $\Box_P$, $\overleftarrow{\Box}_P$ are given the usual truth conditions over $\leq_P$ and $\Box_N$, $\overleftarrow{\Box}_N$ are interpreted using $\leq_N$. The conditional $I(B|A)$ is defined as previously, using $\Box_P$, $\overleftarrow{\Box}_P$. A new *normative conditional* connective $\Rightarrow$ is defined in exactly the same fashion using $\Box_N$, $\overleftarrow{\Box}_N$:

$$A \Rightarrow B \equiv_{\mathrm{df}} \overleftrightarrow{\Box}_N \neg A \vee \overleftrightarrow{\Diamond}_N (A \wedge \Box_N(A \supset B)). \quad (3)$$

The sentence $A \Rightarrow B$ means $B$ is true at the most normal $A$-worlds, and can be viewed as a default rule. This

conditional is exactly that defined in [2, 4], and the associated logic is equivalent to a number of other systems (e.g., the qualitative probabilistic logic of [21, 11]).[8]

Given a QDT-model and a (finite) set of facts $KB$, we define the *default closure* of $KB$ to be (where $\mathbf{L}_{CPL}$ is our propositional sublanguage)

$$Cl(KB) = \{\alpha \in \mathbf{L}_{CPL} : KB \Rightarrow \alpha\}$$

That is, those propositions $\alpha$ that are normally true given $KB$ form the agent's set of default conclusions. It seems natural to ask how the ordering $\leq_N$ is determined. Typically, we will have a set of conditional premises of the form $A \Rightarrow B$, plus other modal sentences that constrain the ordering. We note that conditional deontic sentences may also be contained in $KB$; but these impose no constraints on the normality ordering. Unless these premises form a "complete" theory, there will be a space of permissible normality orderings. Many default reasoning schemes will provide a "preferred" such ordering and reason using that ordering. System Z, described above, is one such mechanism, forming the most compact (normality) ordering consistent with the $KB$. We make no commitment to the default reasoning scheme we use to determine the ordering, simply that the closure $Cl(KB)$ is semantically well-defined. We assume for simplicity (though this is relaxed in [6]), that the default closure $Cl(KB)$ is finitely specifiable and take it to be a single propositional sentence.[9]

We remark that similar considerations apply to the preference ordering $\leq_P$. One can take the ordering to be the most compact ordering satisfying the premises $KB$ or use some other strategy to determine allowable models. However, we do not *require* the use of a single ordering — the definitions presented below can be reinterpreted to capture truth in all permissible orderings or all QDT-models of a given theory [6]. It is sufficient, though, to define goals relative to a single model, and then use simple logical consequence (truth in all QDT-models of $KB$) to derive goals, if desired.

An agent ought to act not as if only $KB$ were true, but also these default beliefs $Cl(KB)$. Assume a particular QDT-model $M$. As a first approximation, we define an *ideal goal* (w.r.t. $KB$) to be any $\alpha \in \mathbf{L}_{CPL}$ such that

$$M \models I(\alpha | Cl(KB))$$

The *ideal goal set* is the set of all such $\alpha$. In our previous example, where $KB = \{C\}$, we have that $Cl(KB) \equiv$

---

[8]QDT can be axiomatized using the axioms of CO for each pair of connectives $\Box_P$, $\overleftarrow{\Box}_P$ and $\Box_N$, $\overleftarrow{\Box}_N$, plus the single interaction axiom $\overleftarrow{\Box}_P \alpha \equiv \overleftarrow{\Box}_N \alpha$. Thus, $\leq_P$ and $\leq_N$ are completely unrelated except that they must be defined over the same set of worlds.

[9]A sufficient condition for this property is that each "cluster" of equally normal worlds in $\leq_N$ corresponds to a finitely specifiable theory. This is the case in, e.g., System Z [2].

---

$C \wedge R$ and the agent's goals are those sentences entailed by $C \wedge R \wedge U$. It should be clear that goals are *conditional* and *defeasible*; for instance, given $C \wedge \overline{R}$, the agent now has as a goal $\overline{U}$.[10]

# 5 Ability and Incomplete Knowledge

The definition of ideal goal in the previous section is somewhat unrealistic as it fails to adequately account for the ability of an agent. Given $C$ in our example above, the derived goal $U$ seems reasonable while the goal $R$ seems less so: we should not expect an agent to make it rain! More generally, there will be certain propositions $\alpha$ over which the agent has no control: it is not within its power to make $\alpha$ true or false, regardless of the desirability of $\alpha$. We should not require an agent to have as a goal a proposition of this type. Ideal goals are best thought of as the "wishes" of an agent that finds itself in situation $KB$ (but cannot change the fact that it is in a $KB$ situation).

## 5.1 Controllable Propositions

To capture distinctions of this sort, we introduce a naive model of action and ability and demonstrate its influence on conditional goals. While this model is certainly not very general (see the concluding section), it is sufficient to illustrate that conditional goals will require more structure than we have suggested above using ideal goals.

We partition our atomic propositions into two classes: $\mathbf{P} = \mathcal{C} \cup \overline{\mathcal{C}}$. Those atoms $A \in \mathcal{C}$ are *controllable*, atoms over which the agent has direct influence. We assume that the only actions available are $do(A)$ and $do(\overline{A})$, which make $A$ true or false, for each controllable $A$. We assume also that these actions have no effects other than to change the truth value of $A$. The atom $U$ (ensure you have your umbrella) is an example of a controllable atom. Atoms in $\overline{\mathcal{C}}$ are *uncontrollable*. $R$ (it will rain) is an example of an uncontrollable atom.

**Definition** For any set of atomic variables $\mathcal{P}$, let $V(\mathcal{P})$ be the set of truth assignments to this set. If $v \in V(\mathcal{P})$ and $w \in V(\mathcal{Q})$ for distinct sets $\mathcal{P}$, $\mathcal{Q}$, then $v; w \in V(\mathcal{P} \cup \mathcal{Q})$ denotes the obvious extended assignment.

We can now distinguish three types of propositions:

**Definition** A proposition $\alpha$ is *controllable* iff, for every $u \in V(\overline{\mathcal{C}})$, there is some $v \in V(\mathcal{C})$ and $w \in V(\mathcal{C})$ such that $v; u \models \alpha$ and $w; u \models \neg \alpha$.

A proposition $\alpha$ is *influenceable* iff, for some $u \in V(\overline{\mathcal{C}})$, there is some $v \in V(\mathcal{C})$ and $w \in V(\mathcal{C})$ such that $v; u \models \alpha$ and $w; u \models \neg \alpha$.

---

[10]The "priority" given to defaults can be thought of as assuming arbitrarily high conditional probabilities. We are currently investigating the ability to give priority to certain preferences (e.g., infinitely low or high utility).

Finally, $\alpha$ is *uninfluenceable* iff it is not influenceable.

Intuitively, since atoms in $\mathcal{C}$ are within complete control of the agent, it can ensure the truth or the falsity of any controllable proposition $\alpha$, according to its desirability, simply by bringing about an appropriate truth assignment. If $A, B \in \mathcal{C}$ then $A \vee B$ and $A \wedge B$ are controllable. If $\alpha$ is influenceable, we call the assignment $u$ to $\overline{\mathcal{C}}$ a *context* for $\alpha$; intuitively, should such a context hold, $\alpha$ can be controlled by the agent. If $A \in \mathcal{C}$, $X \in \overline{\mathcal{C}}$ then $A \vee X$ is influenceable but not controllable: in context $X$ the agent cannot do anything about the truth of $A \vee X$, but in context $\overline{X}$ the agent can make $A \vee X$ true or false through $do(A)$ or $do(\overline{A})$. Note that all controllables are influenceable (using a tautologous context). In this example, $X$ is uninfluenceable. It is easy to see that these three types of propositions are easily characterized according to properties of their prime implicates [6].

## 5.2 Complete Knowledge

Given the distinction between controllable and uncontrollable propositions, we would like to modify our definition of a goal so that an agent is obligated to do only those things within its control. A first attempt might simply be to restrict the goal set as defined in the last section to controllable propositions. In other words, we determine the those propositions that are ideally true given $Cl(KB)$, and then denote as a goal any such proposition within the agent's control (i.e., any proposition that is *influenceable* in context $KB$). The following example shows this to be inadequate.

Consider three atoms: $L$ (my office thermostat is set low); $W$ (I want it set low); and $H$ (I am home this morning). My robot has control only over the atom $L$ (it can set the thermostat), and possesses the following default information: $\top \Rightarrow W$, $\overline{L} \Rightarrow \overline{W}$ and $H \wedge \overline{L} \Rightarrow W$. (I normally want the thermostat set low; but if it's high, I probably set it myself and want it high — unless I stayed at home and the caretaker turned it up.) The robot has the "factual" knowledge $KB = \{\overline{L}, H\}$, namely, that the thermostat setting is high and I'm at home. The default closure of its knowledge is $Cl(KB) = \{\overline{L}, H, W\}$: most likely I want the thermostat set low, even though it is currently high. Finally, the robot's preference ordering is designed to respect my wishes:

$$\{WL, \overline{WL}\} < \overline{W}L < W\overline{L}$$

(we assume $H$ does not affect preference).

It should be clear that the robot should not determine its goals by considering the ideal situations satisfying $Cl(KB)$. In such situations (since $\overline{L}$ is known), $\overline{L}$ is true and the robot concludes that $\overline{L}$ *should be true*.[11] This is clearly mistaken, for considering only the best situations in which one's knowledge of controllables is

---

[11] This is a simple theorem of our conditional logic: $I(\alpha|\alpha)$.

---

true prevents one from determining whether changing those controllables could lead to a better situation. Since any controllable proposition can be changed if required, we should only insist that the best situations satisfying *uninfluenceable* known propositions be considered. We shouldn't allow the fact that $\overline{L}$ is known unduly influence what we consider to be the best alternatives — we can make $L$ true if that is what's best. But notice that we should not ignore the truth of controllables when making default predictions. The prior truth value of a controllable might provide some indication of the truth of an uncontrollable; and we *must* take into account these uncontrollables when deciding which alternatives are *possible*, before deciding which are best. In this example, the fact $\overline{L}$ might provide an indication of my (uncontrollable) wish $\overline{W}$ (though in this case defeated by $H$).[12]

This leads to the following formulation of goals that account for ability. We again assume a QDT-model $M$ and sets $\mathcal{C}$, $\overline{\mathcal{C}}$. The closure of $KB$ is defined as usual. The *uninfluenceable belief set* of an agent is

$$UI(KB) = \{\alpha \in Cl(KB) : \alpha \text{ is uninfluenceable}\}$$

This set of beliefs is used to determine an agent's goals. We say $\alpha$ is a *complete knowledge goal* (CK-goal) iff

$$M \models I(\alpha|UI(KB)) \qquad \text{and} \qquad \alpha \text{ is controllable}$$

In our example above, the only atomic goal the robot has is $L$ (low thermostat). In the earlier example, given $C$ (cloudy) the goal will be $U$ (take umbrella).

As with ideal goals, the set of CK-goals is deductively closed. We can think of CK-goals as *necessary conditions* for an agent achieving some ideal state. Usually we are interested in *sufficient conditions*, some sentence that, if known, guarantees that the agent is in an ideal state (given $UI(KB)$).[13] This can be captured in modal terms. We say proposition $G$ is *CK-sufficient* with respect to $KB$ (or "guarantees ideality") just when $G \wedge UI(KB)$ is satisfiable and

$$M \models \overset{\leftrightarrow}{\boxminus}_P(UI(KB) \supset \overset{\leftarrow}{\boxminus}_P(UI(KB) \supset \neg G))$$

This simply states that any world satisfying $G \wedge UI(KB)$ is at least as preferred as any other $UI(KB)$-world. Thus

---

[12] If a controllable provides some indication of the truth of an uncontrollable or another controllable, (e.g., $\overline{L} \Rightarrow \overline{W}$) we should think of this as an *evidential rule* rather than a *causal rule*. Given our assumption about the independence of atoms in $\mathcal{C}$, we must take all such rules to be evidential (e.g., changing the thermostat will not change my wishes). We discuss this further in the concluding section. Note the implicit temporal aspect here; propositions should be thought of as *fluents*. The theorem $I(\alpha|\alpha)$ should not be viewed as paradoxical for precisely this reason. It does not suggest that $\alpha$ *ought* to be true if it is true, merely that it must be true in the best situations where it *is* true. As we see below, $\neg\alpha$ can be a *goal* or obligation even if $\alpha$ is known.

[13] Hector Levesque (personal communication) has suggested that this is the crucial "operator."

ensuring $G$ is true (assuming as usual that $UI(KB)$ cannot be affected) guarantees the agent is among the best possible $UI(KB)$-worlds.

Of course, changes in the world can only be effected through atomic actions, so we are most interested in sufficient conditions described using atomic actions. We say an (atomic) *action set* is any set of controllable literals (drawn from $\mathcal{C}$). If $\mathcal{A}$ is such a set we use it also to denote the conjunction of its elements. An *atomic goal set* is any action set $\mathcal{A}$ that guarantees each CK-goal (see [6] for further details). We can show that any atomic goal set determines a reasonable course of action.

**Theorem 3** *Let $\mathcal{A}$ be an atomic goal set for KB. Then $\mathcal{A}$ is CK-sufficient for KB.*

We note that usually we will be interested in *minimal* atomic goal sets, since these require the fewest actions to achieve ideality. We may wish to impose other metrics and preferences on such goals sets as well (e.g., associating costs with various actions).

### 5.3   Incomplete Knowledge

The goals described above seem reasonable, in accord with the general deontic principle "do the best thing possible consistent with your knowledge." We dubbed such goals "CK-goals" because they seem correct when an agent has complete knowledge of the world (or at least uncontrollables). But CK-goals do not always determine the best course of action if an agent's knowledge is *incomplete*. Consider our preference ordering above for the umbrella example and an agent with an empty knowledge base. For all the agent knows it could rain or not (it has no indication either way). According to our definition of a CK-goal, the agent ought to $do(\overline{U})$, for the best situation consistent with its $KB$ is $\overline{RU}$. Leaving its umbrella at home will be the best choice should it turn out not to rain; but should it rain, the agent has brought about the *worst* possible outcome. It is not clear that $\overline{U}$ should be a goal. Indeed, one might expect $U$ to be a goal, for no matter how $R$ turns out, the agent has avoided the worst outcome.

It is clear, in the presence of incomplete knowledge, that there are various *strategies* for determining goals. CK-goals (the "deontic strategy") form merely one alternative. Such a strategy is opportunistic, optimistic or adventurous. Clearly, it *maximizes potential gain*, for it allows the possibility of the agent ending up in the best possible outcome. In certain domains this might be a prudent choice (for example, where a cooperative agent determines the outcome of uncontrollables). Of course, another strategy might be the cautious strategy that *minimizes potential loss*. This corresponds precisely to the *minimax* procedure from game theory, described formally here.

*Complete action sets* (complete truth assignments to the atoms in $\mathcal{C}$) are all of the alternative courses of action available to an agent. To minimize (potential) loss, we must consider the worst possible outcome for each of these alternatives, and pick those with the "best" worst outcomes. If $\mathcal{A}_1$, $\mathcal{A}_2$ are complete action sets, we say $\mathcal{A}_1$ is *as good as* $\mathcal{A}_2$ ($\mathcal{A}_1 \leq \mathcal{A}_2$) iff

$$M \models \overset{\leftrightarrow}{\Diamond}_P(\mathcal{A}_2 \wedge UI(KB) \wedge \neg\overset{\leftarrow}{\Diamond}_P(\mathcal{A}_1 \wedge UI(KB)))$$

Intuitively, if $\mathcal{A}_1 \leq \mathcal{A}_2$ then the worst worlds satisfying $\mathcal{A}_1$ are at least as preferred in $\leq_P$ as those satisfying $\mathcal{A}_2$ (considered, of course, in the context $UI(KB)$). It is not hard to see that $\leq$ forms a transitive, connected preference relation on complete action sets. The *best* actions sets are those minimal in this ordering $\leq$. To determine the best action sets, however, we do not need to compare all action sets in a pairwise fashion:

**Theorem 4** $\mathcal{A}_i$ *is a best action set iff* $M \models \mathcal{A}_i \leq \neg\mathcal{A}_i$.

This holds because the negation of a complete action set is consistent with any other action set. We say $\alpha$ is a *cautious goal* iff

$$\vee\{\mathcal{A}_i : \mathcal{A}_i \text{ is a best action set }\} \models \alpha$$

In this way, if (say) $A \wedge B$ and $A \wedge \neg B$ are best action sets, then $A$ is a goal but $B$ is not. Simply doing $A$ (and letting $B$ run its natural course) is sufficient. This notion of goal has controllability built in (ignoring tautologies). In our example above, $U$ is a cautious goal. Of course, it is the action sets that are most important to an agent. We cannot expect best action sets, in general, to be sufficient in the same sense that CK-goal sets are. The potential for desirable and undesirable outcomes makes this impossible. However, we can show that if there does exist some action set that is sufficient for $KB$ that it will be a best action set.

**Theorem 5** *If some complete action set $\mathcal{A}$ is CK-sufficient for KB, then every best action set is CK-sufficient.*

Note that the concept of CK-sufficiency can be applied even in the case of incomplete knowledge. When it is meaningful, it must be that possible outcomes of unknown uncontrollable have no influence on preference (given best action sets): all *relevant* factors are known.

In [6] we describe various properties of these two strategies. In particular, we show that the adventurous and cautious strategies do indeed maximize potential gain and minimize potential loss, respectively. The cautious strategy seems applicable in a situation where one expects the worst possible outcome, for example, in a game against an adversary. Once the agent has performed its action, it expects the worst possible outcome, so there is no advantage to discriminating among the candidate (best) action sets: all have equally good worst outcomes. However, it's not clear that this is the best strategy if the outcome of uncontrollables is essentially "random." If outcomes are simply determined by the natural progression of events, then one should be more

selective. We think of nature as neither benevolent (a co-operative agent) nor malevolent (an adversary). Therefore, even if we decide to be cautious (choosing from *best* action sets), we should account for the fact that the worst outcome might not occur: we should choose the action sets that take advantage of this fact. We are currently investigating such strategies in relation to, for instance, the absolute goals of Doyle and Wellman [28]. Indeed, it's not hard to see that being a cautious goal is a necessary condition for being an absolute goal, but not sufficient. Other strategies provide useful (closer) approximations to absolute goals. Such considerations also apply to games where an opponent might not be able to consistently determine her best moves and an agent wants to exploit this fact. It should be clear that the combination of ability and incomplete knowledge also has a profound impact on how obligations (in the traditional deontic sense) must be defined.

## 5.4 Observations

It should be clear that if an agent can *observe* the truth values of certain unknown propositions before it acts, it can improve its decisions. Eliminating situations cannot make the worst outcomes of action sets any worse (all contingencies are accounted for in cautious goals); but in many cases, it will make the worst outcomes better and change the actions chosen. To continue our previous example, suppose $R$ and $C$ are unknown. The agent's cautious goal is then $U$. If it were in the agent's power to determine $C$ or $\overline{C}$ before acting, its actions could change. Observing $\overline{C}$ indicates the impossibility of $R$, and the agent could then decide to $do(\overline{U})$. In [6] we formalize this notion by distinguishing two types of uncontrollable atoms: *observables* (like "cloudy") and *unobservables* (like "it will rain soon"). We describe the value of observations in terms of their effect on decisions and possible outcomes (much like "value of information" in decision theory [21]).

## 6 Concluding Remarks

We have presented a modal logic of preferences that captures the notion of "obligation" defined in the usual deontic logics. We have extended this logic with the ability to represent normality and added to the system a naive account of action and ability. Within this qualitative framework we have proposed various methods for determining the obligations and goals of an agent that account for the beliefs of an agent (including its default beliefs), the agent's ability, and the interaction of the two. We have shown that goals and obligations cannot be uniquely defined in the presence of incomplete knowledge, rather that *strategies* for determining goals arise.

There are a number of ways in which this framework must be extended. Clearly, the account of actions is naive. True actions have preconditions, default or un-

certain effects, and so on.[14] We are currently extending the framework to include a representation of events and actions that have such properties. In this way, we can account for planning under uncertainty with a qualitative representation, and plan for conditional goals.

We would like to include observations as actions themselves [**?**] and use these to characterize (dynamic) conditional goals and conditional plans. We hope to characterize the changes in an agent's belief set, due to observations and its knowledge of the effects of actions, using belief revision and update semantics. One drawback of this system is the extralogical account of action and ability. We hope to embed the account of action and ability directly in the object language (e.g., using the methods of dynamic logic). Temporal notions also have a role to play (see, e.g., McCarty [20]).

Finally, we need to explore extensions in which certain preferences or "utilities" can be given precedence over expectations or "probabilities." The definition of goal in our system has the property that an agent should act as if every belief (including default beliefs) is true. This seems to be reasonable for the most part. But the consequences of being wrong for certain acts may outweigh the "probability" of being right (in the classic decision theoretic sense). For example, even if I believe that I can safely run across the busy freeway, the drastic consequences of being wrong greatly outweigh the benefit of being right. We would like to capture this type of trade-off in a qualitative way. Possible methods include the "stratification" of propositions or situations to give priority to preferences in some cases, expectations in others; or explicitly using the ranking information implicit in the ordering structure (as is done by Goldszmidt and Pearl [10, 11]) and comparing the qualitative degree of belief/expectation to the qualitative degree of preference.[15] This may lead to a concrete qualitative proposal for *decision-theoretic defaults* [23], where a de-

---

[14]The default "effects" of actions need not be causal. Consider the the preference relation induced by the payoff matrix for the classic *Prisoner's Dilemma* (where $A$ means "our agent" cooperates and $O$ means the other agent cooperates):

$$\overline{A}O < AO < \overline{A}\,\overline{O} < A\overline{O}$$

Our agent's best course of action *for any given choice by the other agent* is not to cooperate. However, on the assumption that the other agent reasons similarly, our agent ends up in a suboptimal situation $\overline{A}\overline{O}$. Hence the dilemma: mutual cooperation would have been better ($AO$). To incorporate the assumption that the other agent reasons similarly, our agent might hold two defaults: $A \Rightarrow O$ and $\neg A \Rightarrow \neg O$. For these defaults to fill the appropriate role, they must be "applied" *after* the agent has made its choice: though apparently "uninfluenceable," the truth value $O$ must not persist after action. Given this, (should both agent's use the same defaults) we have a model of the Prisoner's Dilemma that accounts for mutual cooperation.

[15]This final suggestion is due to Judea Pearl (personal communication).

fault rule $A \to B$ means that "acting as if $B$" has higher expected utility than not, given $A$.

# References

[1] Craig Boutilier. Conditional logics of normality as modal systems. In *Proc. of AAAI-90*, pages 594–599, Boston, 1990.

[2] Craig Boutilier. Inaccessible worlds and irrelevance: Preliminary report. In *Proc. of IJCAI-91*, pages 413–418, Sydney, 1991.

[3] Craig Boutilier. Conditional logics for default reasoning and belief revision. Technical Report KRR-TR-92-1, University of Toronto, Toronto, January 1992. Ph.D. thesis.

[4] Craig Boutilier. A logic for revision and subjunctive queries. In *Proc. of AAAI-92*, pages 609–615, San Jose, 1992.

[5] Craig Boutilier. What is a default priority? In *Proceedings of Canadian Society for Computational Studies of Intelligence Conference*, pages 140–147, Vancouver, 1992.

[6] Craig Boutilier. Beliefs, ability and obligations: A framework for conditional goals. Technical report, University of British Columbia, Vancouver, 1994. (Forthcoming).

[7] Roderick M. Chisholm. Contrary-to-duty imperatives and deontic logic. *Analysis*, 24:33–36, 1963.

[8] Philip R. Cohen and Hector J. Levesque. Intention is choice with commitment. *Artificial Intelligence*, 42:213–261, 1990.

[9] James P. Delgrande. An approach to default reasoning based on a first-order conditional logic: Revised report. *Artificial Intelligence*, 36:63–90, 1988.

[10] Moisés Goldszmidt and Judea Pearl. System Z+: A formalism for reasoning with variable strength defaults. In *Proc. of AAAI-91*, pages 399–404, Anaheim, 1991.

[11] Moisés Goldszmidt and Judea Pearl. Reasoning with qualitative probabilities can be tractable. In *Proceedings of the Eighth Conference on Uncertainty in AI*, pages 112–120, Stanford, 1992.

[12] Peter Haddawy and Steve Hanks. Representations for decision-theoretic planning: Utility functions for deadline goals. In *Proc. of KR-92*, pages 71–82, Cambridge, 1992.

[13] Bengt Hansson. An analysis of some deontic logics. *Noûs*, 3:373–398, 1969.

[14] John F. Horty. Moral dilemmas and nonmonotonic logic. *J. of Philosophical Logic*, 1993. To appear.

[15] Andrew J. I. Jones and Ingmar Pörn. Ideality, subideality and deontic logic. *Synthese*, 65:275–290, 1985.

[16] Andrew J. I. Jones and Ingmar Pörn. On the logic of deontic conditionals. In *Workshop on Deontic Logic in Computer Science*, Amsterdam, 1991.

[17] David Lewis. *Counterfactuals*. Blackwell, Oxford, 1973.

[18] Barry Loewer and Marvin Belzer. Dyadic deontic detachment. *Synthese*, 54:295–318, 1983.

[19] David Makinson. Five faces of minimality. *Studia Logica*, 1992. To appear.

[20] L. Thorne McCarty. Defeasible deontic reasoning. In *Fourth International Workshop on Nonmonotonic Reasoning*, pages 139–147, Plymouth, VT, 1992.

[21] Judea Pearl. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann, San Mateo, 1988.

[22] Judea Pearl. System Z: A natural ordering of defaults with tractable applications to default reasoning. In M. Vardi, editor, *Proceedings of Theoretical Aspects of Reasoning about Knowledge*, pages 121–135. Morgan Kaufmann, San Mateo, 1990.

[23] David Poole. Decision-theoretic defaults. In *Proceedings of Canadian Society for Computational Studies of Intelligence Conference*, pages 190–197, Vancouver, 1992.

[24] Austin Tate, James Hendler, and Mark Drummond. A review of AI planning techniques. In J. Allen, J. Hendler, and A. Tate, editors, *Readings in Planning*, pages 26–49. Morgan-Kaufmann, San Mateo, 1990.

[25] Bas C. van Fraassen. The logic of conditional obligation. *Journal of Philosophical Logic*, 1:417–438, 1972.

[26] Georg Henrik von Wright. Deontic logic. *Mind*, 60:1–15, 1951.

[27] Georg Henrik von Wright. A new system of deontic logic. In Risto Hilpinen, editor, *Deontic Logic: Introductory and Systematic Readings*, pages 105–120. D.Reidel, Dordecht, 1964. 1981.

[28] Michael P. Wellman and Jon Doyle. Preferential semantics for goals. In *Proc. of AAAI-91*, pages 698–703, Anaheim, 1991.

# Acknowledgements