

Norms in Game Logics

On the differences between Norms in Game Theory and in
Deontic Logic

Bachelor's thesis
Cognitive Artificial Intelligence
7.5 ECTS

Onno Treep
3245691

Supervisor:
Jan Broersen

February 2012

Contents

1	Introduction	2
2	Introduction in norms	3
2.1	A deontic view on norms	3
2.2	A game theoretic view on norms	4
3	Work of Thomas Ågotnes	5
3.1	Normative system games	5
3.2	Power in normative systems	7
3.3	Examples	7
4	Work of Paolo Turrini	12
4.1	Examples	14
5	Discussion	18
6	Conclusion	20
	References	20

1 Introduction

Norms in daily life are mostly conventions about behavior that tend to make living among other people easier for everyone. They can be either written laws or unwritten promises. An example of the first category would be the norm not to steal, while the norm to shake hands belongs to the second. Most norms are like restrictions to behavior. One has to accept a limitation to its own possibilities, but benefits from the fact that others also have this restriction. A traffic light is a good example: a driver cannot just drive when he wants to, but benefits from the fact that others cannot do that, either.

Norms can be of use in computer science. The first computers that were build, worked entirely separate from other computers. They had their input and output, and no interaction at all with other systems. Nowadays that is completely different: almost all computers are connected in some way. While making a computation, a computer always interacts with other systems, and it often is dependent on these other systems. For computers that need to interact to accomplish their tasks, rules of another kind are needed. Rules only to compute an input are not enough. That is what norms are used for in computer science. They guide the behavior of computers interacting with each other, so that their performance is optimized. It is like optimizing the lives of people who live together in a society.

To define these norms, logicians try to formalize the concept ‘norm’ as we use it. They try to find a suitable definition to use in computer science. As a result of the changing way computers work, there has been a focus shift in agent systems. Rationality of a group has become an important topic apart from single agent rationality.

Deontic logic is the field concerned with prohibitions, obligations and related concepts. Researchers in this area reason about these concepts, and also the formalization of norms is a huge subject. Game theory often plays an important role here. That is not surprising: if the assumption is that every agent is selfish, then accepting a norm only makes sense if it has some good implications for all agents. Most of the time, this also depends on the choices other agents make, so deciding to accept a norm or not can be viewed as a game.

A problem that occurs here is that ideas about norms in game theory seem to differ from those in deontic logic. When a norm in game theory is accepted, it restricts behavior, and violations are not possible. In deontic logic, a norm can be violated. The norm and the system have to be designed in such a way that agents do not benefit from violation.

How are the different views on norms related? Do they contradict each other? Or is it possible to combine them somehow? That will be the subject of this thesis: **how is the game theoretic notion of norms or normative systems related to the deontic notion of norms?**

I will summarize and discuss two approaches on norms in multi agent systems. I’ll start

with two papers by Thomas Ågotnes, in which agents use game theory to optimize their choices. I'll use some examples to illustrate my interpretation of his work. The deontic perspective is based on a paper by Paolo Turrini, who connects coalitional rationality with behavior of individual agents.

Connection to CAI When focussing on Cognitive Artificial Intelligence, this research is relevant in two ways. To understand this, let's first define the two main objectives of CAI: 1) to learn more about human cognition by trying to 'teach' its abilities to computers, and 2) to gain techniques to make smarter machines by studying human cognition. Since computers only compute problems described in logic, for both of these objectives it is needed to formalize concepts about human cognition. Norms, too, can tell something about how human cognition works, so formalizing them is relevant in both objectives of the field of CAI.

2 Introduction in norms

2.1 A deontic view on norms

Most in this section is based on the ideas from [4] and [6]. A distinction can be made between two kinds of norms: written norms (like laws) and unwritten norms (like promises between at least two agents). An intuitive way to describe a norm could be *a description of right and wrong*. Since often it depends on the situation if a certain action is right or wrong, this definition is insufficient, especially in (deontic) logic. The words norm and normative system are mostly used for (sets of) rules that regulate multi-agent behavior. Another distinction to make lies in the origin of the norm. Some norms are rules written by a designer, they come from outside of the system. Other norms are cooperative rules, from which every agent benefits. These norms come from within the system.

When a norm benefits all agents, the situation is simple. An agent, if it has a choice, will comply with the norm. There is of course another category: norms that do not benefit all agents and therefore are not automatically accepted by all agents. An example is found in the prisoner's dilemma. Game theoretic models say agents should not cooperate, while cooperating would be better for the group. So the norm should be to cooperate, even though it is not in the interest of the individuals.

Originally, norms were defined as obligations, based on the distinction between right and wrong. Something was obliged when the opposite was not permitted. Lately, norms and obligations are not defined as the same thing. The two concepts now differ in various ways. Norms are applied to a group of agents over some period of time. They exist independent of an agent. Obligations, on the other hand, are part of the motivational states of single agents. They can influence an agents' choice at a particular moment.

There is another important distinction: the one between ‘ought-to-be’ and ‘ought-to-do’. The first tells what *states* are good to be in, while the second tells what *actions* are good to execute.

While beliefs are descriptive, norms and obligations are prescriptive. Therefore they do not have a truth value, because there is no environment where to test whether a norm or obligation is true or not. So-called norm propositions do have a truth value, however. The logics of these norm propositions differ from the logics of norms. The question whether a logic of norms can exist if norms have no truth value, is called Jorgenson’s dilemma.

2.2 A game theoretic view on norms

For this section, most information comes from Shoham [5], who speaks of social laws. He describes a social law as a restriction on the given strategies of a group of agents. For an agent, a social law presents a tradeoff: it suffers from a loss of freedom, but can benefit from the fact that others lose some freedom as well. A good social law is designed to benefit all agents.

In game theory, each agent has a number of possible strategies. Depending on the strategies selected by each agent, each agent receives a certain payoff. Agents are free to choose their own strategies, based on their guesses about the strategies of other agents. A social law eliminates from a game certain strategies for each of the agents.

When a designer wants to design a system, he cannot always design any system, because often the strategy spaces are given. Still, a designer wants to be able to influence the actions of the agents which is possible using social laws.

Once a social law is imposed, it is not guaranteed that agents will follow it. Three ways to regulate this are the following, using a central (independent) agent to guide the others:

- **Contracts.** A center can propose a contract, which specifies an action for each agent. When all agents accept, the center can fine each agent that did not follow the contract. The downside on this solution is that it requires great effort from the center.
- **Bribes.** The center promises agents a higher payoff in certain outcomes to implement a desired behavior.
- **Mediators.** The center operates as a mediator for the agents. It makes it possible for agents to make their strategy dependent on the strategies of others.

3 Work of Thomas Ågotnes

3.1 Normative system games

The basic principle in [2] is that multi-agent systems are coordinated by social laws, which are constraints on the behavior of agents. By constraining the behavior of the agents in a system, the system will achieve its goal. In this approach preferences of agents are taken into account, an idea that was missing in some older approaches. It is important to take this into account because agents make a choice whether to comply with the normative system or not. Because this means that behavior will now be strategic, a model is needed in which strategy is incorporated.

In the paper, a model is created in which agents have a list of goals. They also have priorities: every goal has some importance relative to the other goals. Based on its goals and priorities, an agent decides to accept or decline a normative system. It is clear that game theory plays a role in formalizing this choice to accept or not.

Kripke structures and CTL A Kripke structure is a directed graph. It has a set of states and a set of transitions between these states. The transitions are directed. In this case transitions are labelled with agents that perform an action and states are labelled with boolean variables that express the properties that are true.

Computation Tree Logic (CTL) is a temporal logic used to express the properties of a Kripke model. A certain computation from some state s is called an *s-path*.

- $K, s \models A \bigcirc \varphi$ means ‘for all s-paths, the second element satisfies φ ’.
- $K, s \models E \bigcirc \varphi$ means ‘for some s-path, the second element satisfies φ ’.
- $K, s \models A(\varphi \mathcal{U} \psi)$ means ‘all elements in all s-paths satisfy φ until some element satisfies ψ ’.
- $K, s \models E(\varphi \mathcal{U} \psi)$ means ‘for some s-path, every element satisfies φ until an element satisfies ψ ’.

Normative systems A normative system tells which transitions in a system are illegal. It is a subset of the set of transitions. When some normative system η is applied to system K , so the illegal transitions are deleted, the resulting system is referred to as $K \dagger \eta$. If C is a subset of the set of agents, $\eta \upharpoonright C$ is the system that only contains the forbidden transitions in η that correspond to actions of agents in C . So in $K \dagger (\eta \upharpoonright C)$ only these transitions are deleted. The same goes for $\eta \downharpoonright C$, but the other way around: it is the same as η but contains only the transitions that do *not* correspond to actions of agents in C .

Goals and utilities The goals of an agent are listed as CTL formulas. The first is the least important one, the last is the most important. If an important goal can be reached, the goals before that one will be ignored. The first formula is always true. A goal in hierarchy γ is true in K if all initial states satisfy this goal. One of the properties of such a hierarchy is monotonicity. If a higher goal is true, then the lower goals are as well. Each agents' hierarchy, together with the Kripke model, forms the multi-agent system $M = \langle K, \gamma_1, \dots, \gamma_n \rangle$.

The utility of a Kripke model for an agent is the index of the highest goal that is always true, which is always at least zero. Because of the definition of utility, it is possible to compare different Kripke models for an agent, but not to compare different agents in a model. Agents mostly do not have the same goals, and the index of a goal does not carry information that can be compared with other agents' goals. $\delta_i(K, K \uparrow \eta)$ (shorter: $\delta_i(K, \eta)$) is the utility of η for agent i . This utility is equal to $u_i(K \uparrow \eta) - u_i(K)$, and is therefore negative if the agent is worse off with the normative system. So if an agent is in a situation where only some universal goal is true, no norm will be a problem for it. The situation cannot be worse than that, so the utility of any normative system for this agent cannot be negative. If for an agent an existential goal is true, no normative system can make its situation better. In that case, the utility of any normative system will be at most zero.

A social system is now defined as $\Sigma = \langle M, \eta \rangle$.

$K_1 \sqsubseteq K_2$ means that K_1 is a subsystem of K_2 . This is the case if $R_1 \subseteq R_2$ and the rest of the properties (S, S^0, A, α, V) are equal.

Normative system games Given an multi-agent system and a normative system, agents get a strategy to accept or defect this normative system. The making of this choice is a process called a normative system game, which takes place before an agent is in the system. A game is defined as follows:

- $G = \langle AG, S_1, \dots, S_n, U_1, \dots, U_n \rangle$
- $AG = \{1, \dots, n\}$ is the set of agents
- S_i is the set strategies for each agent (so the possible actions)
- $U_i : (S_1 \times \dots \times S_n) \rightarrow \mathbb{R}$ is the utility function for an agent

Given an social system $\Sigma = \langle M, \eta \rangle$, every agent chooses a strategy: C (comply) or D (do not comply). For S (a tuple of strategies, one for each agent) and $x \in \{C, D\}$, AG_S^x gives the subset of agents that play x . The utility function then is $U_i(S) = \delta_i(K, \eta \uparrow AG_S^C)$.

A normative system is *individually rational* if each agent would be better off with it. Logically, it is better for each agent if each one of them complies, than if no agent

complies. Individual rationality is necessary, but not sufficient for a norm to succeed (see for example the prisoner’s dilemma).

A normative system is Pareto efficient if there does not exist another system which each agent is better off with.

A social system is a Nash implementation if S_C causes a Nash equilibrium. This means no agent is better off by defecting from the norm, given that other agents will comply.

3.2 Power in normative systems

An issue discussed in [1] is what happens if an agent does not accept a normative system. The goal is to find a way to measure the influence of an agents’ choice to the success of a normative system. The tool used here is the *voting power index*. Power is the ability to let a normative system succeed or not. In the ideal situation this power is equally divided among the agents.

Coalitional games and power A coalition is defined as a subset of the set of agents that accept a normative system. A *swing player* is an agent that can make the difference in letting a normative system succeed or not: if it complies, the system succeeds, otherwise it does not. The *Banzhaf score* of an agent is the number of coalitions of which it is a swing player. The *Banzhaf measure* of an agent is the chance that it is swing player of a random chosen coalition. The *Banzhaf index* of an agent is the proportion of coalitions of which it is a swing player to the number of swings in the game. And last: the Shapley-Shubik index of an agent is the proportion of *permutations* of which it is a swing player to the total number of permutations.

Power in social systems The first thing to do is to show how to associate a coalitional game with a social system. If C , complying with a normative system, will achieve the objective, then it gets the value 1. If not, then it gets the value 0. This way, the relative power (the relative ability to cause a normative system to succeed or fail) of agents in a social system can be measured. The goal could be to evenly divide the power over all of the agents in the system. Coalition monotonicity ensures that if a coalition complies with a certain system, all supersets of that coalition do as well:

$$\forall C : K \uparrow (\eta \uparrow C) \models \varphi \text{ implies } \forall C' \supseteq C : K \uparrow (\eta \uparrow C') \models \varphi$$

3.3 Examples

Example: traffic lights Imagine a crossing of streets. Three cars come from different directions, and the drivers want to cross as soon as possible. There are three time steps at which the drivers can drive (*first, second, third*). If two or more go at the same time,

they crash, which is the worst possible situation for them. The best situation is driving first, while the others wait.

A state in the Kripke structure represents the crossing, with a status for each car: *has crossed* or *has not crossed yet*. The transitions represent the actions of the drivers. These are *crossing* and *waiting*.

When no normative system is applied, every driver will drive first, causing an accident. So each normative system that prevents an accident is individually rational: every agent will be better off with that system.

All of these systems are also Pareto efficient. A system that prevents cars from crashing, does so by telling which driver can go first, which one second, and which one last. There is no system that gives all agents a better position. A driver that has the second or third turn cannot drive earlier, because it will cause an accident. Also the driver that has the first turn obviously cannot go earlier. This last statement is also true for systems in which one driver gets the first turn, but the other two crash. So apart from the systems that prevent an accident, there are more Pareto efficient systems. The set of Pareto efficient systems is a superset of the set of individually rational systems.

Imagine a normative system that forces two or three drivers to drive at the same time. It is better for the drivers to decline this norm and drive at another moment. There is always one moment at which no other driver is driving. So all the systems that cause accidents are not Nash equilibria. All the others are, because no agent can be better off by not accepting the norm, given that others do accept it.

Example: non-shareable resource The goal for both agents is to have the resource as often and as long as possible. The result is that both agents prefer keeping it rather than giving it away. Since transitions in the model correspond to keeping the resource or giving it away, and a normative system contains one or more of these transitions, the norm is preventing one or more of the possible actions. Preventing an agent to give the resource away will make no difference, because he prefers to keep it anyway. So only situations in which one or two agents are forbidden or allowed to keep the resource, are interesting in this example.

None of the possible norms are individually rational. Without a system, the agent that starts with the resource can keep it forever. That is the maximum for this agent, so it will not be better off with any normative system. Therefore no system can be individually rational.

Also none of the systems are Pareto efficient. If keeping the resource is forbidden for only one agent, the other is still able to keep it. The same argument as for individual rationality goes here: one agent get his maximum score, so there is no system with which *both* agents are better off. If the prohibition is for both agents, then it is still impossible to find such a system. One would be better off with a system where keeping the resource is allowed, but the other would not be.

There is no possible Nash implementation in this example. If one agent is prohibited to keep the resource, then it is always better off defecting from the norm. This also is the case if both agents are prohibited to keep it.

Also with three or more agents, this example does not seem to be a good one to explain the notions of individual rationality, Pareto efficiency and Nash implementations. A basic principle in this example is that an agent can only get a better score at the cost of another agents' score. We would have to change the example, in such a way that a better score for both agents can be reached if they cooperate. That is what is missing here.

Example: investment Four agents each have 2 euros. They have the possibility to earn more by investing their money together. Each agent wants to make as much profit as possible. They all have a choice: they can invest 0 or 1 euro. An investment of in total 4 euros is large enough to double the money. The profit will be shared equally, so every agent gets 2 euros. If someone does not invest, there is not enough to double the money, and the agents that gave away their money have 0 euros. So there are three possible outcomes for each agent: 0, 1 or 2 euros.

Let's translate this example to a normative system game like in [2]. The states in the Kripke structure correspond to the different situations in which every agent has invested an amount of money (0 or more). The transitions correspond to the set of investments that are done. As for the goals of the agents: there are certain amounts of money for each agent to end with, these are the goals. Obviously ending with the maximum possible amount of money (in this case 2 euros) is considered the most important goal, ending with 0 euros is the universal goal at the other end of the list. (In this context, *to have x euros* must be read as *to have at least x euros*, so that the condition that higher goals imply lower goals is met.)

The social system is clearly a Nash implementation. To get the maximum amount, all agents need to invest their money. So it is good to invest the money if everyone does, but if someone does not, it is better to keep what you have. The norm to invest the money is also Pareto efficient: all agents get the maximum amount of money. Furthermore, it is individually rational, because every agent is better off with the norm.

$4 \rightarrow 8$	0	1
ind. rat.	no	yes
Pareto	no	yes
Nash	yes	yes

Table 1: Situation in which every agent has 1 euro.

Now, to make the example a bit more interesting, each agent gets 2 euros. They now have three options: give 0, 1 or 2 euros to make an investment together. If the total investment is 4 euros (or more), 8 euros are to be shared. But if 8 euros are invested,

16 come back. If the total amount is 1, 2, or 3, but also if it is 5, 6, or 7, some money is lost.

Investing 1 euro still results in a Nash equilibrium. If all agents give 1, then it is better to also give 1. The same goes for investing 2 euros: if all other agents do that, it is best to follow. Also, both of these norms are individually rational. If everyone invests the same amount of money, then everyone is better off. Then Pareto efficiency: the norm to invest 2 euros results in more money for each agent than the norm to invest 1. So the first is Pareto efficient: each agent gets the highest possible amount. The second is not, because there is a norm with which every agent is better off.

$4 \rightarrow 8, 8 \rightarrow 16$	0	1	2
ind. rat.	no	yes	yes
Pareto	no	no	yes
Nash	yes	yes	yes

Table 2: Situation in which every agent has 2 euros.

Let's change another thing in this example: to get 16 euros, a total investment of 6 is needed instead of 8. The maximum amount an agent can get now changes: an agent will end up with 6 if it keeps his money, while all the others give 2. All investing 1 euro at the same time still results in a Nash equilibrium, all investing 2 does not. It has become profitable to not invest anything in that situation. The norm to invest 2 euros does remain Pareto efficient and individually rational. Each agent is better off with the norm, and no norm gives a better outcome for all agents.

$4 \rightarrow 8, 6 \rightarrow 16$	0	1	2
ind. rat.	no	yes	yes
Pareto	no	no	yes
Nash	yes	yes	no

Table 3: Situation in which investing 6 euros results in 16 euros.

Prisoner's dilemma The previous example is pretty much like a large prisoner's dilemma. It could be interesting to describe the prisoner's dilemma itself in terms of normative systems, and to analyze it like the other examples.

	Comply	Do not comply
Comply	$\langle 4, 4 \rangle$	$\langle 0, 5 \rangle$
Do not comply	$\langle 5, 0 \rangle$	$\langle 1, 1 \rangle$

Table 4: Prisoner's dilemma.

There are two possible norms: to comply and to defect. The first is Pareto efficient and individually rational. It is better to have this norm than no norm at all, and there is no

norm with which both agents are better off. It is not a Nash equilibrium, though. Since it is the default choice for both agents, the second norm is not individually rational nor Pareto efficient, but it is a Nash equilibrium.

The following is always true in games with agents that have equal choices: if some choice is the default choice for each agent, then the norm of choosing that is always a Nash equilibrium and never individually rational. This is because without this ‘norm’, every agent would choose the same. Because the outcome in accepting and defecting is equal, an agent never benefits from accepting or defecting. For a norm to be individually rational or a Nash equilibrium, agents have to benefit making a certain choice, which is impossible in this situation.

Example: voting game (Shoham) There are four parties in a parliament, A, B, C, and D. They have 45, 25, 15 and 15 representatives, respectively. There will be a vote on whether to spend a certain amount of money or not, and how much of that money each party gets to spend. If a coalition of at least 51 representatives votes for the same distribution, that is what will happen. If there is no such majority, no party gets to spend anything. Let’s not worry about the actual distribution of the money, and let’s assume that:

- every representative of the same party votes the same;
- every party in the winning coalition gets a share in proportion to the number of representatives in the coalition;
- because of the assumption above, every coalition wants to be in a winning coalition as small as possible.

In that case there are fifteen possible coalitions. One-party coalitions are never a majority. Two-party coalitions are a majority only if A is one of the two parties in it. Furthermore, every three-party coalition and of course the four-party coalition are majorities. The possible coalitions can function as norms: the representatives (or the parties) in the coalition are agents that all have to vote for the coalition, like for a norm. Let’s see what the notions of individual rationality, Pareto efficiency and Nash implementation can tell us here.

The coalitions that are minorities are neither individually rational, Pareto efficient or a Nash equilibrium. A party wins nothing by voting for such a coalition.

All the winning coalitions are individually rational. The parties in it win by voting for these coalitions, so they are better off with them than without them.

Pareto efficiency is not useful in this example. If you are only taking into account the parties in a certain coalition, then there is never a better coalition for all of the parties in the first one. Such a coalition has to be one that contains at least all the parties of the first coalition. Therefore it is at least as big as the first one, and cannot be better for each party.

About Nash: in the winning two-party coalitions $\{A,B\}$, $\{A,C\}$, and $\{A,D\}$ there is no Nash equilibrium. In the first case, A is better off in a coalition with C or D and therefore should not comply with this coalition. In all three cases, the smallest party is better off in $\{B,C,D\}$, so even if A complies, there is a better situation. Also in the three-party coalitions with A in it, $\{A,B,C\}$, $\{A,B,D\}$, and $\{A,C,D\}$, there are no Nash equilibria. A is better off in a two-party coalition, and the rest will prefer $\{B,C,D\}$. In this last coalition is a Nash equilibrium, since it is the best choice for all three of the parties in it. They all should comply with the insurance the rest does. The four-party coalition is of course the worst winning coalition for all of them, so they are declining this one.

Analyzing the notions in [1], we can say the following: A is a swing player for six of the seven winning coalition that contain A. This means its Banzhaf score is 6 and its Banzhaf measure is $\frac{6}{15}$. It has a Banzhaf index of $\frac{6}{12}$. The Banzhaf score of the rest of the parties is 2, so they all have a Banzhaf measure of $\frac{2}{15}$ and a Banzhaf index of $\frac{2}{12}$.

4 Work of Paolo Turrini

Introduction The aim of [3] is to formally capture the notion of coalitional rationality. To achieve this, it is needed to combine two views on the representation of group decisions: one is to represent the preferences of the whole group and the other is to represent the preferences of the individuals. Therefore policies are studied where the desirable properties to be achieved by a coalition reflect the preferences of some superset of that coalition. In this way the different views can be accounted for as particular cases. In multi-agent systems there are situations in which individual preferences are not compatible and agents' capabilities affect the realization of other agents' preferences. In this paper, enactment of norms as aimed at the regulation of such interactions is studied.

The goal is to isolate the notions of *betterness* (to compare players' possibilities), *choice restriction* and *interest*.

Preliminaries A *dynamic effectivity function* $E : W \rightarrow (2^{Agt} \rightarrow 2^{2^W})$ gives a choice set for a certain coalition in a certain state. If a set X is member of $E(w)(C)$, the coalition is able to force that the next state after w will be some member of X .

The following properties are considered:

- **Regularity:** if a coalition is able to force the outcome of an interaction to belong to a particular set, then no possible combinations of moves by the other agents can prevent this to happen.

- **Outcome monotonicity:** if a coalition is able to force the outcome of an interaction to belong to a particular set, then it is also able to force the outcome to belong to all of its supersets.
- **Inability of the empty coalition:** the empty coalition cannot bring about non-trivial consequences.

To reason about effectivity functions, *coalition logic* is used. The language of this logic is $\phi ::= p \mid \neg\phi \mid \phi \wedge \psi \mid [C]\phi$. A *coalition model* is a triple (W, E, V) . W is a set of states, E is a effectivity function and V is a valuation function that associates to every state a set of atomic propositions. In a model, in a certain state, $[C]\phi$ is true iff $\phi^M \in E(w)(C)$ (where $\phi^M = \{w \in W \mid M, w \models \phi\}$).

The relation \succeq_i shows a preference ordering. In $v \succeq_i w$, state v is as least as nice to be in as w , from the perspective of agent i . The relations \preceq_i , \succ_i , and \prec_i work in a similar way. A coalition model extended by a preference relation, (W, E, \preceq_i, V) , is called a *cooperative game model*. $M, w \models \diamond_i^{\preceq} \phi$ now means that for some w' with $w \preceq w'$, $M, w' \models \phi$. $L^{[C], \preceq}$ is the language of coalition logic extended with the \diamond_i^{\preceq} modality.

Regulating strategic decisions Effectivity functions do not say anything about what coalitions would force if they had a choice, and preference relations are defined on outcomes and not on sets. So there needs to be a way to lift preferences over outcomes to effectivity functions and a way to take choices of others into account.

$E(w)(C) \cap Y$ gives the choices of C in w , intersected with some choice Y of the opponent. The possible outcomes are restricted with Y . With this notion of choice restriction, it is possible to define a notion of undomination. An undominated choice remains Pareto optimal for all possible answers of the opponents. A possible choice X is undominated iff $(X \cap Y)$ is Pareto optimal in the Y -choice restriction of $E(w)(C)$ (for all $Y \in E(w)(\overline{C})$).

A deontic logic for strategic interactions The ingredients for a deontic logic for this purpose are a coalition logic and some way to express coalitional rationality. The operator for rationality:

$$M, w \models [rational_C]\phi \text{ iff } \phi^M \triangleright_{C,w}$$

With this, it is possible to say which coalitional choices are rational.

Also needed is a way to express what actions should be performed by coalitions, and the link with rationality. The following operators are used:

- $F(C : \phi)$
Coalition C is forbidden to choose ϕ . It makes no difference whether ϕ is available as a choice.

- $P(C : \phi)$
Coalition C is permitted to choose ϕ . This means it is not forbidden to choose ϕ .
- $O(C : \phi)$
Coalition C is obliged to choose ϕ . It is forbidden to choose $\neg\phi$.

Stating that a choice is rational for a coalition itself, is the same as stating that this choice is permitted in the interest of the coalition.

The following statements are true in the described system:

- If something is permitted for a coalition, it is not obliged to do the opposite.
- For a choice to be permitted for a smaller coalition, it needs to be rational for the grand coalition.
- If ϕ is permitted or ψ is permitted, then $\phi \vee \psi$ is permitted.
- If a possible choice of C is obligated, then it is rational for the grand coalition, and therefore also obligated for other smaller coalitions.
- A rational choice of a coalition is forbidden if it is in conflict with a rational choice of the grand coalition.
- If ϕ is forbidden and ψ is forbidden, then $\phi \wedge \psi$ is forbidden.

The following statements are not true in the system:

- If ϕ is obliged and ψ is obliged, then $\phi \wedge \psi$ is obliged.
- If $\phi \vee \psi$ is permitted, then ϕ is permitted or ψ is permitted.
- If ϕ is obliged, then $\neg\phi$ is not obliged.
- If a choice is rational for a coalition, it is not always rational for the grand coalition (and the other way around).
- If ϕ is obliged, then ϕ is permitted.
- If ϕ is obliged for C , then C can choose ϕ .

4.1 Examples

Voting game The voting game example used before might be a good example here, as well. With four parties, there are fifteen possible coalitions (the empty one not counted). All of the parties have the choice to vote for a certain coalition. If the coalition is a majority, then it gets money to spend. The set of agents is the set of parties, renamed to P, Q, R and S. They consist of 45, 25, 15 and 15 members respectively.

$$A = \{P, Q, R, S\}$$

The set of worlds consists of one world in which nobody has money (n), and for each coalition a possible world in which that coalition gets the money (let's name these worlds after the coalitions). So the number of possible worlds is equal to the number of possible coalitions plus one, which is sixteen.

$$W = \{n, p, q, r, s, pq, pr, ps, qr, qs, rs, pqr, pqs, prs, qrs, pqrs\}$$

Choices: if a set X is part of choice set $E(w)(C)$, then C is able to force that the world after w is an element of X . This is where the difference between winning coalitions and the rest becomes clear. The winning coalition $\{Q, R, S\}$ for example is able to force that the money goes to them: $\{qrs\} \in E(n)(\{Q, R, S\})$.¹ But $\{Q, R\}$, which is not a winning coalition, is unable to do so. It is unable to force any world after n , so $E(n)(\{Q, R\}) = \emptyset$. To take a look at choice sets $E(w)(C)$ in which $w \neq n$ is not very interesting, so that will not be discussed.

Investment Imagine the same situation as before: four players all have 1 euro. They have a choice to invest, and if they all do, everyone gets twice as much back. Because agents with the same amount of money have the same possible actions, they are considered equal in this example. Therefore, there are four different coalitions: a coalition with one agent (C_1), with two agents (C_2), with three (C_3), or with four (C_4). The possible worlds differ in the amount of money the agents have. Let's call world w_{2222} the world in which all agents have 2 euros, w_{1100} the world in which two agents have 1 and two agents have 0, etc. This notation only gives the number of agents that have a certain amount of money; it doesn't specify which agent has which amount. Since agents with the same amount of money are considered equal, worlds with an equal number of agents with the same amount of money are considered the same. So, for example, w_{1100} , w_{1001} and w_{0101} all have two agents with 1 euro and two agents with 0 euros, and are therefore considered as the same state. The possible worlds then are:

$$W = \{w_{2222}, w_{1111}, w_{1110}, w_{1100}, w_{1000}\}$$

The choices we want to discuss are all made in w_{1111} , since agents always have 1 euro to start with. So w_{1111} is the starting state. As said before: only if all four of the agents invest, then they all get 2 euros. So only a coalition of four agents is able to force w_{2222} as the next world. Formally: $\{w_{2222}\} \in E(w_{1111})(C)$ if and only if $C = C_4$. But also w_{1111} is a possibility for this coalition, in case the agents choose not to invest. Therefore $\{w_{1111}\} \in E(w_{1111})(C_4)$. You could say C_4 is also able to force the other coalitions, but then agents need to make decisions that are not equal, which would be a strange thing to do. I will not discuss that option here and assume that agents in a coalition take the same action.

$$E(w_{1111})(C_4) = \{\{w_{2222}\}, \{w_{1111}\}\}$$

¹In the original example, a winning coalition got to decide which party could spend which amount of money. In that case, all worlds could be forced by the coalition, so $\{p\} \in E(n)(\{Q, R, S\})$, $\{q\} \in E(n)(\{Q, R, S\})$, etc.

A three agent coalition is unable to force a set of one world. When three agents choose to invest, there are two possible outcomes: the last agent also invests, and everyone gets 2 euros, or it does not, in which case it is the only one that keeps 1. So: $\{w_{2222}, w_{0001}\} \in E(w_{1111})(C_3)$. When the coalition chooses not to invest, there are also two possibilities: when the agent invests, it loses its money, and when it does not invest, it keeps 1. So, in this case: $\{w_{1110}, w_{1111}\} \in E(w_{1111})(C_3)$. Together this results in:

$$E(w_{1111})(C_3) = \{\{w_{2222}, w_{0001}\}, \{w_{1110}, w_{1111}\}\}$$

In a similar way, we can describe the effectivity functions of coalitions C_2 and C_1 :

$$E(w_{1111})(C_2) = \{\{w_{2222}, w_{0010}, w_{0011}\}, \{w_{1100}, w_{1110}, w_{1111}\}\}$$

$$E(w_{1111})(C_1) = \{\{w_{2222}, w_{0001}, w_{0011}, w_{0111}\}, \{w_{1000}, w_{1100}, w_{1110}, w_{1111}\}\}$$

The following is an example of a valuation function, which is also needed for a complete coalition model. In this example, the propositions that are assigned to the states contain information about the amount of money each agent has (for example, $p =$ ‘agent x has 1 euro’). A modality $[C]\phi$ in this context could mean that C is able to force that a certain agent has 1 euro and another has 0.

The coalition model is further extended with a preference relation, in which different worlds are compared by the propositions that are assigned to them. With this, it is possible to express which worlds individual agents prefer. In this example, agents probably prefer worlds in which they have more money. w_{2222} is preferred by all of the agents, because they have the maximum amount of money. But only the individual outcome counts: a hypothetical w_{3000} would be preferred by the agent that has 3 euros in that world.

Now the connection is made between the abilities of the coalitions and the preferences of the agents. Suppose that for an agent w_{2222} is at least as good as w_{1110} , which in turn is at least as good as w_{1000} :

$$w_{2222} \succeq_i w_{1110} \succeq_i w_{1000}$$

Then by lifting these preferences to sets, it is possible to say what choice sets are preferable over another. For example:

$$\{w_{2222}\} \succeq_i \{w_{1110}\} \succeq_i \{w_{1000}\}$$

$$\{w_{2222}, w_{1110}\} \succeq_i \{w_{1110}, w_{1000}\}$$

A choice X for a coalition C is Pareto optimal if there is no choice Y (in the same state) that is strictly better than X for all agents in the coalition. So in our example, for C_4 the Pareto optimal choice is $\{w_{2222}\}$. For the others, no particular set is strictly better than another, so all choices are Pareto optimal.

The next thing to illustrate is the notion of undominated choices (and therefore choice restriction). Suppose a coalition of three C_3 has to decide what to invest. Then its choice set, as said before, contains the following sets:

$$E(w_{1111})(C_3) = \{\{w_{2222}, w_{0001}\}, \{w_{1110}, w_{1111}\}\}$$

Let's say $INV_3 = \{w_{2222}, w_{0001}\}$ and $KEEP_3 = \{w_{1110}, w_{1111}\}$. Now, assuming that the agents in C_3 are taking the same action, the choice set of $E(w_{1111})(\overline{C_3})$ contains the sets $\{w_{2222}, w_{0111}\}$ (which I will refer to as INV_1) and $\{w_{1000}, w_{1111}\}$ (called $KEEP_1$).

For a choice of C_3 to be undominated, it needs to be Pareto optimal, no matter what the choice of $\overline{C_3}$ is. See for example INV_3 . Two things need to be true for it to be undominated:

- $INV_3 \cap INV_1$ needs to be Pareto optimal in $E(w_{1111})(C_3) \sqcap INV_1$
- $INV_3 \cap KEEP_1$ needs to be Pareto optimal in $E(w_{1111})(C_3) \sqcap KEEP_1$

The first condition is true:

$$E(w_{1111})(C_3) \sqcap INV_1 = \{\{w_{2222}\}, \{w_{1110}\}\}$$

So $INV_3 \cap INV_1 = \{w_{2222}\}$ is Pareto optimal for C_3 . The second condition, however, is not:

$$E(w_{1111})(C_3) \sqcap KEEP_1 = \{\{w_{0001}\}, \{w_{1111}\}\}$$

And $INV_3 \cap KEEP_1 = \{w_{0001}\}$ is not Pareto optimal. w_{1111} is better for all agents in C_3 . The conclusion is that INV_3 is not undominated for C_3 in w_{1111} . The same argumentation goes for $KEEP_3$, which is also not undominated.

It is rational for C in w to try to achieve ϕ , if and only if the set of possible worlds that satisfy a formula ϕ is equal to an undominated choice for a coalition C in w . Formally:

$$M, w \models [rational_C]\phi \text{ iff } \phi^M \triangleright_{C,w}$$

For example, if ϕ means that every agent has 2 euros, then $[rational_{C_4}]\phi$ is true, because $\{w_{2222}\}$ is an undominated choice for this coalition. It is not an undominated choice for C_3 , though, since it is not part of the choice set of C_3 at all. So we cannot say $[rational_{C_3}]w_{2222}$. The sets that *are* part of the choice set of C_3 , are not undominated, so there is nothing that is rational to force.

But coalitional rationality is not as much used for expressing what is rational, as for expressing what is *not* rational. For example: if C_3 chooses to keep what it has, it forces that the three agents in the coalition keep 1 euro. Let's call this situation ϕ . Now, ϕ^M contains two worlds: w_{1110} and w_{1111} . This set is of course equal to the earlier mentioned choice $KEEP_3$. From the fact that this choice is not undominated, we can derive that it is not rational for C_3 to force ϕ .

The forbiddance operator $F(C : \phi, \overline{C} : \psi)$ means that C forcing ϕ and C' forcing ψ imply that $\phi \wedge \psi$ is not rational for C and C' together. In the case of C_3 : if C_3 and $\overline{C_3}$ both

keep their money, then the resulting world is w_{1111} . For the two coalitions together, however, this is not a rational choice, because w_{2222} is a better outcome. Therefore, if C_3 and $\overline{C_3}$ are forming a coalition together, the action to not invest should be forbidden for both of them. In this stage, the notion of interest becomes important: the action is forbidden *in the interest of the greater coalition*.

5 Discussion

Now that I have tried to illustrate the work of Thomas Ågotnes ([2],[1]) and Paolo Turrini ([3]) by the same examples, it may be easier to compare the two approaches.

In [2] and [1], notions of norms are simply used to describe possible situations. Let's return to the main example. When speaking of the norm *to invest 1 euro*, the situation in which all agents invest some amount of money is described. We can reason about whether or not this situation is possible to achieve, by checking if it is profitable for each agent. If some agent loses money, then it will not comply with the norm, and the situation we are talking about will never happen. In other words, the state or world in which the agents invest that amount will not be reached. The norms that *are* possible, are the ones that are profitable for all of the agents.

Norms are constraining the behavior of agents. This is, as we have seen, illustrated by the deletion of transitions in the Kripke structures that represent the possible worlds and transitions. But which transitions exactly are deleted? An agent will only accept a norm if it does not delete the transition the agent wanted to use. So a norm that is accepted by all agents is only making actions impossible which the agents would not perform anyway. The profit lies in the fact that agents have knowledge about what others do. The moment a norm is accepted, agents have knowledge about what others will not do. Given this information, their choice of action might change. So what this approach does is providing a strategy based on game theory to make better choices, but it is still focussed on individual agents. The agents remain selfish, and if cooperation takes place, then it is to make individual outcomes better.

The approach in [3] uses a notion of interest to define coalitional rationality apart from individual rationality. With this, it is possible to tell what an agent should do to optimize the outcome for the coalition it is participating in. This can very well be another choice than the agent would make in its own interest. To compare the two approaches, we could say that in [2], agents act in their own interest, while in [3], a way to act in the interest of a coalition is provided. In this approach, a norm is a guideline, translating interests of higher levels to lower levels.

Having made that distinction, let's see how these approaches can be useful for the purpose they have: improving the interaction and cooperation of computers in a system. Recall the distinction made earlier, between norms from within a group and norms 'from above'. For a system that is guided by the latter, the approach in [3] is very relevant.

A designer of such a system can define goals that it has to reach. These goals can be translated to behavior constraints for subsystems, then for subsystems of these systems and so on. In the end, this will result in constraints for the smallest parts of the system, that still, indirectly, optimize the outcome for the entire system. In a top-down way, each part of the original system gets these limitations that are needed to optimize the outcome on the highest level. In the opposite situation, if these norms from above are no part of the system, then a game theoretic approach like in [2] might be better. The smallest parts of the system will act selfish, and thereby determine the behavior of subsystems on higher levels. The outcome for the entire system depends, indirectly, on the selfish agents on the lowest level. This is a bottom-up way of working. When it comes to applying in computer science, both notions of norms can be relevant and important, but in different situations.

To return to my question which I started with: how are the two notions of norms related? At first, they seem to exclude each other. But by examining the two approaches further, we can see that different assumptions are made about how agents behave and how the accepting of a norm works. An agent is selfish in one approach, but acting in the interest of the group in the second. One approach is focussed on optimizing scores of individual agents, and explains cooperation as letting each other know what actions will not be taken. That is what norms do, in that case: providing information, which is the best way of letting selfish agents cooperate. The other approach is focussed on optimizing scores for the group, while individual agents are not important. Here, norms are the translation of higher level interests to lower level behavior constraints. With these constraints, the score of the group is optimized. So, in both cases, norms help agents to cooperate, but with different goals, and therefore in different ways. If the goal is optimizing the score of the group, it is not clear why agents should cooperate from their own perspective. So apart from the norm, there have to be rules about what happens if an agent does not cooperate. In the game theoretical approach, this is not needed, since the norm does not exclude individually preferred actions. Every agent has to accept a norm before it is applied to their system, so there is no doubt they will follow it. So the deontic norms are explaining what should be achieved with the use of prohibitions, while the game theoretical norms tell us what can be achieved without prohibitions.

Are these both relevant in computer science? Let's recall the two main fields of research in Cognitive Artificial Intelligence I mentioned in the introduction: 1) learning more about human cognition by trying to 'teach' its abilities to computers, and 2) gaining techniques to make smarter machines by studying human cognition. Since a human takes its own decisions and acts selfish, in a way, norms as in game theory can explain a lot about human behavior. The second kind of norms are less relevant in this research. Things like making laws for citizens to rule a country are not part of Cognitive Science.

Focussing on the other field, where the goal is to make smart systems, it is very relevant to have a way to tell what individual agents in a system should do. Individual scores do not matter. So the game theoretical approach seems to answer a question which is

not relevant here. The deontic approach is very relevant, but misses something. In the introduction, I mentioned a few ways in game theory to force agents to perform a certain action, from [5]. Being able to force agents to perform or not perform certain actions is essential, so these elements should have a place in the approach.

6 Conclusion

To discuss the difference between game theory and deontic logic on the subject of norms, I have summarized [2], [1] and [3]. I have used the first two as an example for the game theoretic view on norms, while the last one was an example of the deontic approach. After summarizing the relevant parts of these papers, I have tried to illustrate both approaches with the same examples, to be able to compare them. Comparing the two approaches, I saw that they are making different assumptions on the behavior of agents and the goal to be reached. I concluded that they were solving two different problems, and that the notions of norm are related but not the same. The norms in the game theoretical approach have a more descriptive purpose, while the norms in deontic logic prescribe what agents should do. Therefore the game theoretical approach is not solving the problem it is focussing on. It can be more relevant in the other field of CAI. The deontic approach is a sufficient one if the goal is making intelligent systems, but aspects of forcing agents to behave in a certain way should have a place in it.

References

- [1] T. Ågotnes, W. van der Hoek, M. Tennenholtz, and M. Wooldridge. Power in normative systems. In Decker, Sichman, Sierra, and Castelfranchi, editors, *Proceedings of the Eighth International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2009)*, pages 145–152, Budapest, Hungary, May 2009. IFAMAAS/ACM DL.
- [2] Thomas Ågotnes, Michael Wooldridge, and Wiebe van der Hoek. Normative system games. In M. Huhns and O. Shehory, editors, *Proceedings of the Sixth International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2007)*, pages 876–883. IFAMAAS, May 2007.
- [3] Jan Broersen, Rosja Mastop, John-Jules Ch. Meyer, and Paolo Turrini. A deontic logic for socially optimal norms. In L.W.N. van der Torre and R. van der Meyden, editors, *Proceedings 9th International Workshop on Deontic Logic in Computer Science (DEON'08)*, volume 5076 of *Lecture Notes in Computer Science*, pages 218–232. Springer, 2008.
- [4] J.M. Broersen. Issues in designing logical models for norm change. In George Vouros, Alexander Artikis, Kostas Stathis, and Jeremy Pitt, editors, *Organized Adaption in*

Multi-Agent Systems, First International Workshop, OAMAS 2008, Estoril, Portugal, May 13, 2008. Revised and Invited Papers, volume 5368 of *Lecture Notes in Computer Science*, pages 1–17, 2009.

- [5] Yoav Shoham and Kevin Leyton-Brown. *Multiagent Systems: Algorithmic, Game-Theoretic, and Logical Foundations*. Cambridge University Press, New York, 2009.
- [6] Leendert W. N. van der Torre. Violation games: a new foundation for deontic logic. *Journal of Applied Non-Classical Logics*, 20(4):457–477, 2010.