



# DIGITAL ACCESS TO SCHOLARSHIP AT HARVARD

## Socially Conscious Decision Making

The Harvard community has made this article openly available.  
[Please share](#) how this access benefits you. Your story matters.

<b>Citation</b>	Glass, Alyssa and Barbara J. Grosz. 2003. Socially conscious decision making. <i>Autonomous Agents and Multi-Agent Systems</i> 6(3): 317-339.
<b>Published Version</b>	<a href="https://doi.org/10.1023/A:1022987709366">doi:10.1023/A:1022987709366</a>
<b>Accessed</b>	September 23, 2017 10:46:28 AM EDT
<b>Citable Link</b>	<a href="http://nrs.harvard.edu/urn-3:HUL.InstRepos:2579646">http://nrs.harvard.edu/urn-3:HUL.InstRepos:2579646</a>
<b>Terms of Use</b>	This article was downloaded from Harvard University's DASH repository, and is made available under the terms and conditions applicable to Other Posted Material, as set forth at <a href="http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#LAA">http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#LAA</a>

*(Article begins on next page)*

# Socially Conscious Decision-Making

Alyssa Glass  
Xerox Palo Alto Research Center  
Palo Alto, CA 94304 USA  
aglass@parc.xerox.com

Barbara Grosz  
Division of Engineering and Applied Sciences  
Harvard University  
Cambridge, MA 02138 USA  
grosz@eecs.harvard.edu

## ABSTRACT

For individually motivated agents to work collaboratively to satisfy shared goals, they must be able to make decisions about actions and intentions in the context of commitments to group activities. This paper examines the role of social consciousness in the process of reconciliation of intentions to do group-related actions with other, conflicting intentions. We define a measure of social consciousness; describe its incorporation into the SPIRE experimental system, a simulation environment that allows the process of intention reconciliation in team contexts to be simulated and studied; and present results of several experiments that investigate the interaction in decision-making of measures of group and individual good. In particular, we investigate the effect of varying levels of social consciousness on the utility of the group and the individuals it comprises. A key finding is that an intermediate level of social consciousness yields better results than an extreme commitment. We suggest preliminary principles for designers of collaborative agents based on the results.

## 1. INTRODUCTION

In many situations, agents must interact with other systems and with people to accomplish tasks, and for many applications, it is necessary to form teams that comprise people and computer agents to work collaboratively to satisfy a shared goal. As rational agents, individual team members must be able to make individually rational decisions about their commitments and plans [13]. They must also, however, be responsible to the teams in which they participate. Agent designers need to be able to construct computer agents that include a sense of group commitment in their reasoning about actions and plans.

Various possibilities arise for externally influencing agents to act in the group's interest, including the imposition of sanctions on agents that default. In this paper, we address

a different possible influence on decision-making, one that is internally motivated. We examine the design of agents that incorporate a notion of *social consciousness* that influences the ways they measure benefit and maximize individual outcome. In informal terms, these socially conscious agents may make decisions that are locally, individually suboptimal, because doing so engenders a society that is globally better off. Castelfranchi has argued for such social consciousness in agents [2]. As we discuss after presenting the model of social consciousness, this notion differs from the notions of benevolence in prior multi-agent work. We aim to specify ways agent designers can create optimally socially conscious, individually rational agents.

We first describe SPIRE (SharedPlans Intention-Reconciliation Experiments), a simulation system that enables investigation of the effectiveness of different decision-making strategies under various environmental conditions [19]. We then provide a method for assigning a measure to social consciousness in collaborative agents and define a model in which non-monetary social factors play a role in how agents make decisions about their intentions. We show how this measure of social consideration can be made compatible with monetary considerations to create agents whose level of social commitment can be altered and studied using SPIRE.

The paper presents the results of several experiments in which the level of social consciousness is varied and the effect on different measures of group income is determined. The experiments also consider environments with different densities of tasks to be accomplished by the group. Contrary to expectations, we find a maximal, intermediate level of social consciousness. We provide suggestions for creating agents that maximize group income while remaining individually rational.

## 2. INTENTION RECONCILIATION

### 2.1 The Problem

The experiments we describe in this paper extend previous work [6, 7, 14, 18, 21] by addressing the need for collaborative agents to manage plans and intentions in multi-agent contexts, reasoning jointly about commitments to individual plans and commitments to group activities. Our investigation focuses on the problem of intention reconciliation which arises because rational agents cannot adopt conflicting intentions [1, 6]. If an agent has adopted an intention to do

some action  $\beta$  and is presented the opportunity to do another action  $\gamma$  that would in some way preclude its being able to do  $\beta$ , then the agent must decide between doing  $\beta$  and doing  $\gamma$ : it must *reconcile* intentions.

We will use an example from one of our application domains, Systems Administration [19, 8] to illustrate the problem of intention reconciliation in the context of group activities and to motivate our experiments. In this domain, teams of people and computer systems work together to maintain a cluster of workstations. Their overall group activity (which we will refer to as  $\alpha$ ) comprises many subtasks ( $\beta_i$ ) including such actions as upgrading hardware, restoring files from backups, checking system security, and maintaining printers. The need for intention reconciliation would arise, for instance, if an agent that intends to perform an operating system upgrade as part of its commitment to  $\alpha$ , is subsequently offered an opportunity to attend a lecture by a Nobel Prize winner at the same time. The agent must decide whether to remain committed to the group or to renege on its original intention in favor of attending the lecture.

Because the agent is a member of a team, its income, and therefore its utility, depends not just on the completion of its own subtasks but also on the completion of the subtasks of other team members. Similarly, the utility of the other team members depends on the actions of this agent. To reconcile intentions in a group context, an agent must have a method of reasoning about future utility and the influence of social consciousness as well as about current utility. The agent must consider how other group members will view its failure to honor its commitment. It must reason about future utility and consider the costs it may incur as a result of the group’s reaction to its defaulting on a group-related task. For instance, if the agent reneges on its task assignment, it may receive less valuable task assignments in the future, decreasing its overall utility. The focus of this paper is on another component of the agent’s calculations: in addition to monetary utility from income earned, the agent may receive utility simply from being a “good guy” and honoring its commitment to the group. Depending on the level of social consciousness that the agent has, it may assign different weights to this “good guy” factor in its utility calculation.

## 2.2 The SPIRE Framework

The SPIRE system [19] enables manipulation of various agent properties and environmental conditions and the examination of the effect of different decision-making strategies under those conditions. In SPIRE, a team of agents ( $G_1, \dots, G_n$ ) work together on group activities, called *GroupTasks*, each of which consists of doing a set of tasks (task instances). Each task instance is of one of the types  $\beta_1, \dots, \beta_k$  and occurs at one of the times  $T_1, \dots, T_m$ . For example, a GroupTask in the Systems Administration domain might consist of a week’s work (with the times  $T_i$  being the hours of the work week) doing various tasks  $\beta_i$ . Some task types may have only one instance in the week (e.g., printer maintenance); others may have multiple instances (e.g., running backups). Agents receive income for the tasks they do, and this income can be used in determining an agent’s current and future expected utility.

A SPIRE simulation consists of a sequence of GroupTasks. To simplify the simulations and the analysis, the same GroupTask is done repeatedly by the same group, although the individual tasks within the GroupTask will not necessarily be done by the same agent each time. SPIRE considers a given GroupTask to consist of a set of tasks with time constraints on the tasks and capability requirements for agents doing the tasks. To simplify the description that follows, we assume that a GroupTask maps to a weekly task schedule. A simulation then consists of activity over a sequence of weeks.

A *weekly task schedule* (WTS) is a set of pairs  $\langle task_i, time_i \rangle$  where  $task_i$  is to be done at  $time_i$ , and a *weekly task schedule assignment* (WTSA) is a set of triples  $\langle task_i, time_i, agent_i \rangle$  where  $task_i$  is to be done at  $time_i$  by  $agent_i$ . Each agent has a set of task capabilities and a set of available times that constrain the assignment of tasks in the WTS to produce a WTSA. An agent can only be assigned tasks for which it has the needed capabilities and the time availability.

To model the need to reconcile intentions, a sequence of *outside offers* is generated. These offers correspond to actions that an agent might choose to do apart from the GroupTask. Each outside offer  $\gamma$  conflicts with some task  $\beta$  in the WTSA; to accept an outside offer, an agent must default on one of its assigned tasks. The central question we investigate is the ways in which different levels of social consciousness and thus intention reconciliation strategies influence the rates at which agents default and their individual and collective incomes, given a particular group time horizon and configuration of environmental factors.

Each week, agents, chosen randomly, are offered the opportunity to do some  $\gamma$  that conflicts with a task  $\beta$  in the context of doing  $\alpha$  that it has been assigned in the WTSA. The income value of  $\gamma$  is also chosen randomly from a distribution with approximately the same shape as the distribution of task values in the WTS, allowing the distribution itself to be varied without altering the comparative values of group tasks and outside offers. To provide an incentive to default, the distribution of outside offers is shifted so that it has a mean value that exceeds the mean value of the WTS tasks. If the agent chooses the new opportunity, it defaults on the task  $\beta$  with which  $\gamma$  conflicts. If there is an agent that is available and capable of doing  $\beta$ , the task is given to that agent; otherwise,  $\beta$  cannot be completed by the group, and therefore goes undone.

The group as a whole incurs a cost whenever an agent defaults, and this cost is divided equally among the group’s members. The cost of a particular default depends on its impact on the group. At a minimum, it equals a baseline value that represents the cost of finding a replacement agent. If no replacement is available and so the task will not be done, the cost is increased by an amount proportional to the value of the task.

Currently in the SPIRE system, the group’s reaction to an agent defaulting on its commitment is represented by a “social-commitment policy” [19] that constrains assignment

of tasks to agents. Each week, agents are assigned a portion of their tasks based on how responsible they have been thus far in the simulation. Each agent has a rank that reflects the total number of times it has defaulted, with the impact of past weeks' defaults diminishing over time. The higher an agent's relative rank, the more valuable the tasks it receives. Because there is a greater impact on the group when tasks go undone, an agent's rank is reduced by a larger amount if it defaults when no one can replace it.

To assign group tasks to agents, SPIRE makes use of an omniscient scheduler that has total information about all of the agents' ranks and history of defaults.<sup>1</sup> It is important to note, however, that while the central scheduler has complete information about all of the agents, each individual agent does *not* have access to this knowledge. Thus, it must estimate the behavior of other agents based on publicly-available information, as described in the next section.

### 3. DECISION-MAKING IN SPIRE

#### 3.1 Estimating Monetary Utility

In deciding whether to default on a task  $\beta$  and accept an outside offer  $\gamma$ , an agent can weigh the impact of the choice on three factors, the first two of which are essentially monetary: current income ( $CI$ ) and future expected income ( $FEI$ ).  $CI$  only considers the income from the task or outside offer in question, as well as the agent's share of the group cost should it default. For the  $FEI$  calculation, an agent approximates the value of the tasks it will receive the following week, both if it defaults on  $\beta$  and if it does not default. By comparing the value of these two task sets, the agent can approximate the impact that defaulting will have on its income in the following week. The agent then extrapolates beyond the following week to make a more complete estimation, discounting its estimates of subsequent weeks' income by an uncertainty factor  $\delta < 1$ . Mathematical definitions of each of these factors can be found in Sullivan *et al.* [19].

The original  $FEI$  calculation assumed an agent knows the number of weeks the group will continue to work together. In some situations, however, agents will only have indefinite information about how long they will form teams. These situations resemble group activity in an infinite time horizon because, without knowledge of a final week, agents are forced to treat every week as an intermediate week, equally far from the start as the finish [5]. To simulate this infinite horizon, the  $FEI$  calculation can be modified so that it is an infinite sum. If  $F$  is the original estimate of the following week's income, and  $\delta$  is the uncertainty factor,  $FEI$  in an infinite

time horizon is:

$$\begin{aligned} FEI_{infinite}(F) &= \delta F + \delta^2 F + \delta^3 F + \dots \\ &= \left( \frac{\delta}{1 - \delta} \right) F \end{aligned}$$

Once  $CI$  and  $FEI$  have been calculated, they are combined into a total estimated income ( $TEI$ ) to provide a way of measuring total utility from these monetary measures.  $TEI$  in week  $i$  of the simulation in the default and no-default cases is:

$$\begin{aligned} TEI_{def}(\beta, \gamma, i) &= CI_{def}(\beta, \gamma) + FEI(F_{def}, i) \\ TEI_{no-def}(\beta, i) &= CI_{no-def}(\beta) + FEI(F_{no-def}, i) \end{aligned}$$

where  $F_{def}$  and  $F_{no-def}$  are the agent's estimates of its income for the following week if it does and does not default, respectively.

#### 3.2 Utility from Brownie Points

In addition to considering monetary factors, agents may be socially conscious "good guys," that is, willing to sacrifice short-term personal gain for the good of the group. We have developed a model in which monetary factors and social non-monetary utility are weighed during decision-making. In this *brownie point* model, good guy agents who make socially conscious decisions earn brownie points ( $BP$ ) each time they choose not to default. In addition, an agent loses brownie points when it does default. This loss models an agent's disappointment in itself for failing to be the kind of agent it wants on its teams, *i.e.*, failing to live up to its commitment to the group. The number of brownie points an agent gains or loses at each decision point depends on the value of the task that it is considering defaulting on ( $\beta$ ) and the value of the outside offer that it is considering ( $\gamma$ ).

To represent an agent's concern with its historical reputation, its utility from brownie points is also dependent on the total number of brownie points that it has stored up over time. If an agent has remained faithfully committed to the group for a long time, it will not punish itself as much for defaulting, because it has already stored up a lot of brownie points.

We define the  $BP$  value in the default and no-default cases as follows:

$$\begin{aligned} BP_{def}(\beta, \gamma, currentBP) &= currentBP - \frac{value(\beta)^2}{value(\gamma)} \\ BP_{no-def}(\beta, \gamma, currentBP) &= currentBP + \frac{value(\gamma)}{value(\beta)} \end{aligned}$$

where  $currentBP$  is the total number of brownie points that the agent has accumulated thus far in the simulation, based on some initial amount allocated at the beginning of the simulation.

These definitions have several important properties. First, if an agent does not default, its  $BP$  value increases as  $value(\beta)$  goes down or  $value(\gamma)$  rises, reflecting an agent's greater pride in turning down highly valued outside offers, and in committing itself to less-valued tasks for the good of the

<sup>1</sup>This central scheduler is used only for convenience. Many domains requiring cooperative agents would most likely not rely on a central scheduler in this way but would instead negotiate each week's schedule based on (possibly incomplete) information about each agent. Since this negotiation is beyond the intended scope of the current SPIRE system, and we wish to study aspects of group-commitment scenarios which come after the initial schedule is made, we simplified this aspect of the problem for these experiments.

group. Second, if an agent defaults on some task  $\beta$  in favor of  $\gamma$ , its  $BP$  value decreases as  $value(\beta)$  goes up or  $value(\gamma)$  goes down. This factor thus reflects an agent’s greater willingness to forgive itself for defaulting in favor of a highly valued outside offer or a less-valued  $\beta$  task, and vice versa. For example, an agent will punish itself more for neglecting to check for security breaks than it would for not putting toner in the printer. Alternatively, an agent will more readily forgive itself for defaulting on its group commitment in order to do a highly-valued activity like attend to a sick family member than it would for defaulting in order to attend a baseball game. In addition,  $BP_{def}(\beta, \gamma, currentBP)$  changes quadratically in  $value(\beta)$ , rather than linearly as in the no-default case. This change guarantees an adequate base deduction in  $BP$  even for small  $value(\beta)$ , since the average  $value(\beta)$  is smaller than the average  $value(\gamma)$ . In addition, a quadratic change in  $value(\beta)$ , as opposed to any type of linear change, provides desirable percent changes in  $BP$  over time, punishing agents for defaulting even at the beginning of a simulation when their initial  $BP$  level is still quite high. Thus, because  $BP$  is not just a constant measure over time but rather is a variable metric based on the situation faced by both the group and the specific agent, it is not just a measure of individual persistence. Furthermore, by taking into account both group benefit and individual gain in the long-term,  $BP$  differs from notions of benevolence that are concerned with helping other agents and increasing their utility [11; 3, inter alia].

### 3.3 Social Consciousness in Decision-Making

For experiments in this paper, an agent’s intention-reconciliation strategy takes into account the three factors just described: current income ( $CI$ ); future expected income ( $FEI$ ); and good guy stature in the community ( $BP$ ).  $CI$  and  $FEI$  have been aggregated to form  $TEI$ . To combine these factors into a single utility function, we use an approach from multi-attribute decision-making [22]. First, since  $TEI$  is in monetary units and  $BP$  is not, both factors are normalized using linear normalization so that they can be combined without unit differences unexpectedly giving more weight to one than the other. In this method of normalization, each factor is divided by the maximum of the two possible values for that factor.

For instance, since  $BP_{no-def}(\beta, \gamma, currentBP) > BP_{def}(\beta, \gamma, currentBP)$  in all cases, normalized  $BP$  in each case is:

$$\begin{aligned} normBP_{no-def}() &= 1 \\ normBP_{def}(\beta, \gamma, currentBP) &= \\ &= \frac{BP_{def}(\beta, \gamma, currentBP)}{BP_{no-def}(\beta, \gamma, currentBP)} \end{aligned}$$

An analogous calculation is used to normalize  $TEI$ , although no similar assumptions can be made concerning whether  $TEI_{def}(\beta, \gamma, i)$  or  $TEI_{no-def}(\beta, i)$  is larger.

Once these factors have been thus normalized, they are weighted relative to each other, allowing varying emphasis to be placed on each factor. The impact of these relative weights is empirically analyzed and discussed in Section 4.

The SPIRE system thus uses the following formulas for the utility an agent receives from defaulting and from not defaulting in week  $i$  of the simulation:

$$\begin{aligned} U_{def}(\beta, \gamma, i, currentBP) &= \\ &= TEIweight \times normTEI_{def}(\beta, \gamma, i) + \\ &= BPweight \times normBP_{def}(\beta, \gamma, currentBP) \\ U_{no-def}(\beta, i) &= \\ &= TEIweight \times normTEI_{no-def}(\beta, i) + BPweight \times 1 \end{aligned}$$

where  $TEIweight$  and  $BPweight$  can be adjusted to create agents with varying levels of social consciousness. Agents default when:

$$U_{def}(\beta, \gamma, i, currentBP) > U_{no-def}(\beta, i).$$

## 4. SOCIALLY CONSCIOUS AGENTS

In earlier work [19] we were able to show that agents default less often and increase their individual and group income when more tasks are assigned based on rank or more weight is given by the agents to future income. We also showed a complex relationship between the number of tasks scheduled concurrently (*task density*) and agent default behavior and income. In this section, we present several experiments that examine various aspects of socially conscious decision-making. Base values for important SPIRE parameters in the majority of experiments are given in Figure 1; any departures from these values are noted in individual experiment descriptions. These experiments all have the simplifying assumption that the agents are homogeneous; we have begun to relax this assumption [20]. To maximize the contrast between socially conscious and socially unconcerned agents, SPIRE parameters for these experiments were set to make a relatively large number of outside offers and to impose relatively large rank deductions and group costs when agents default.

---

52 weeks per simulation
12 agents
20 task types (values = 5, 10, 15, ..., 100)
40 time slots per week
10 tasks per time slot = 400 tasks per week
10 tasks per agent per week, assigned based on the agent’s rank; the rest assigned randomly
250-350 outside offers per week
$\delta$ weighting factor for $FEI = 0.8$
$TEIweight = 0.5$
$BPweight = 0.5$

---

**Figure 1: Default SPIRE settings. Departures from these values are noted in experiment descriptions.**

The results presented below are averages of 30 runs that used the same parameter settings but had different, randomly-chosen starting configurations (the values of the tasks in the WTS, and the possible values of the outside offers). In each run, the first ten weeks serve to put the system into a state in which the agents have different ranks; these weeks are not included in the statistics that SPIRE gathers.

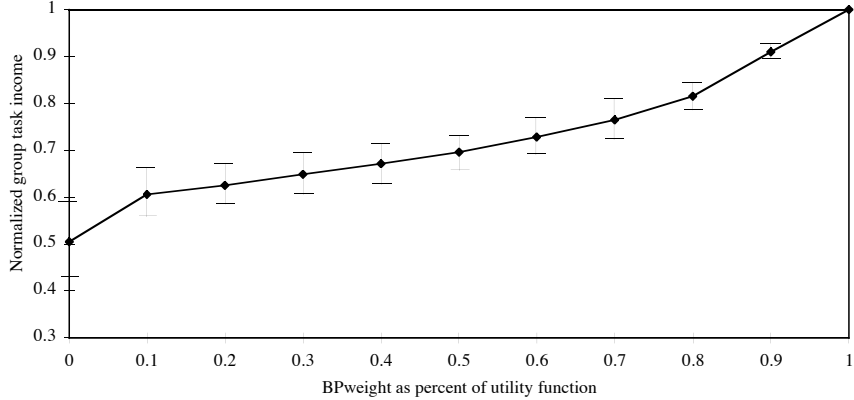


Figure 2: Effect of  $BPweight$  on mean group task income.

#### 4.1 Varying Time Horizons

To provide a basis for further experiments, we investigated the effect of different time horizons on the average number of defaults over the course of a MWS. We compared two sets of runs: one with a known finite horizon of 52 weeks and one with an infinite horizon. The two differed in their  $FEI$  calculations: the finite horizon runs used  $FEI_{finite}$  with  $M = 52$ ; the infinite horizon runs used  $FEI_{infinite}$ .

For the early weeks of an MWS, when there are many weeks left,  $FEI_{finite}$  and  $FEI_{infinite}$  are very close in value. The values diverge as the finite horizon approaches, because  $FEI_{finite}$  becomes much smaller than  $FEI_{infinite}$ . This change reflects the fact that agents have less to lose by defaulting in later weeks, when there is little time left during which they can be punished for their defaulting on their commitment to the group. As a result, in our experiments, during the first several weeks of the MWS the two default rates coincide, but later they diverge.

In particular, for our experiments, the number of defaults in both the finite and infinite cases appears to generally increase until about week 25, when it reaches a plateau. After this initial upward trend, the infinite case appears to reach a steady number of defaults, while the finite case continues with an increasingly upward trend. The major divergence in number of defaults occurs at week 41 (*i.e.*, about 10 weeks from the end of the simulation). By the last week, the number of defaults in the finite case is nearly three times higher than in the infinite horizon.

The plateau phenomenon may be explained in terms of an adjustment of brownie point level. At the beginning of the simulation, agents are initialized with a starting  $BP$  level. Different values for  $BPweight$ , however, dictate a different  $BP$  equilibrium level for the agent. The initial rise in defaults indicates a period of time in which the agents are gradually moving toward this stable  $BP$  level. The plateau occurs when this stable level is reached.

To provide more background data on this plateau effect, we

also investigated the effect of increasing  $BPweight$  in a finite horizon. We found that, for high values of  $BPweight$ , the number of defaults does not “plateau” as it does for lower values, and as was noted above. As  $BPweight$  decreases, this plateau slowly becomes more and more pronounced, starting earlier in the simulation and continuing for more weeks.

This result reflects the relative weight of different factors in the utility function over time and the resulting effect on default equilibrium. Thus, in addition to altering average number of defaults overall, manipulating  $BPweight$  allows us to alter agent behavior over time in a finite horizon environment.

#### 4.2 Optimal Group Commitments

While altering behavior with respect to number of defaults is useful to agent designers in some domains, in many other domains the number of defaults will be of less interest than the actual income earned by the agents. Our next set of experiments examined this income effect in the infinite time horizon case, where  $BPweight$  has a more steady influence. The results in Figure 2 show how average group-only income varies as  $BPweight$  is increased. In this figure, and subsequently in the paper, incomes are normalized with respect to the income that would have been earned if the originally assigned tasks had all been completed. *Group-only income* is the income earned from  $\beta$ -tasks assigned by the group, minus the penalties incurred by the group from defaults. It does not include income earned by agents for any outside offers ( $\gamma$ ) that they complete. As expected, the results show that group-only income strictly increases as  $BPweight$  increases and agents thus default less often.

In contrast to group-only income, *total income* is an agent’s income from group-assigned tasks and outside offers accepted. When we examine total income, we find rather surprising results (Figure 3). Instead of strictly increasing with increased  $BPweight$ , as with group-only income, total income reaches a maximum at  $BPweight = 0.3$  (although differences between this  $BPweight$  value and a few surrounding values

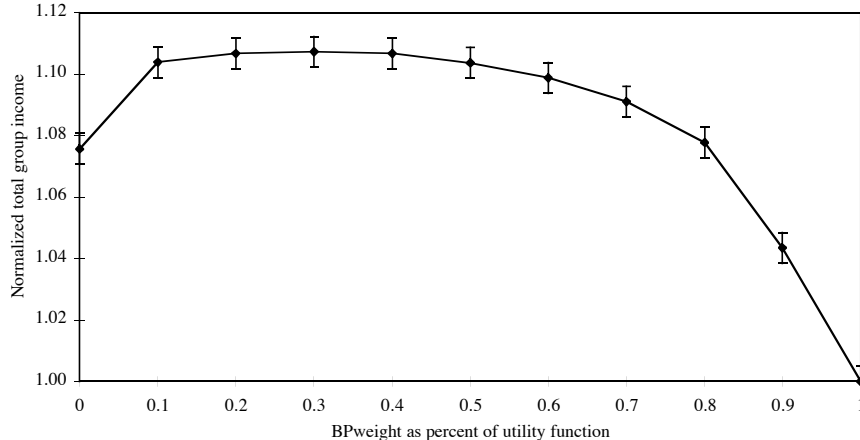


Figure 3: Effect of  $BPweight$  on mean total group income.

are within error) and then begins to drop. This result contradicts our original hypothesis, which was that the optimal number of defaults for maximizing total income would correspond with a high  $BPweight$ , near 1. Instead, this result shows that the optimal number of defaults actually occurs when  $BPweight$  is less than 1, which corresponds to approximately 61 defaults each week.

On further analysis, this result can be explained based on task density and the range of values of outside offers available in a given WTS. As explained in Section 2.2, the mean value of the outside offers exceeds the mean value of the group-assigned tasks. Thus, when deciding between  $\beta$  and  $\gamma$ , it is possible that an agent may be deciding between jobs that vary widely in value. When  $BPweight$  is low, an agent’s decisions are driven mainly by  $TEI$ . For comparatively high outside offers, then, the current income ( $CI$ ) gained by defaulting will be much bigger than the potential loss in future income ( $FEI$ ). Thus, the agent will default.

When  $BPweight$  is too high and social factors are being weighted very heavily in the agent’s utility function, the brownie point factor in the agent’s utility function may overcome this potential gain in income. As a result, the agent may give up some very lucrative outside offers in order to stay committed to the group. In these cases,  $value(\gamma)$  is actually high enough to offset any group penalties suffered from the default, but the agent is, in a way, blinded to this fact by its unusually strong conscience.

In addition, with a very high density of group tasks assigned each week, the effect on future income of defaulting is decreased [19]. In these situations, when defaulting has a very small effect, brownie point considerations may tend to over-emphasize this effect. When combined with a high  $BPweight$ , then, the agent may not default even when it is truly beneficial to do so. Thus, it seems that “good guys” do not finish first. Rather, communities made up of *less* socially conscious, more balanced agents actually do better.

### 4.3 Optimality Across Different Environments

Our final set of experiments extends this finding by varying an environmental factor and observing the effect on the income-maximizing number of defaults. In these experiments, the number of tasks scheduled in each time slot (task density) was varied. As was shown in previous work [19] varying task density has a complex relationship with the number of defaults. To overcome the problems discussed in that paper, as we varied the task density in these experiments, we also varied the number of tasks assigned based on rank so that the *percent* of rank-based tasks was constant across different task densities. In these experiments, the number of rank-based tasks was kept to approximately 30 percent of the total group tasks, as was also the case in the other experiments. The results are shown in Figure 4.

Given the results from the previous section, these results are unexpected. While at the highest task density studied in this paper<sup>2</sup> (10) there is a local maximum for total group income for a  $BPweight$  value greater than 0, this result does not hold for lower task densities. Instead, a pattern emerges in which a large range of values for  $BPweight$  greater than 0 provides a total group income which is approximately equal for a given task density. For example, for a task density of 2, the total group incomes associated with the range of  $BPweight$  values from 0.1 to 0.7 are all roughly the same, within error ranges.

This result points to an interesting conclusion. While we have shown that  $BPweight$  has a large and consistent effect on the average number of defaults and on group task income, Figure 4 shows that it does not have this same effect on total income. Thus, when creating cooperative agents, designers do not have to be as concerned about total group income as they are about these other factors. Instead, within a fairly large range of  $BPweight$  values, agents can be de-

<sup>2</sup>In previous experiments [19] the highest possible task density, in which all agents are busy all of the time, proved to be a special case, and thus was not experimented with here.

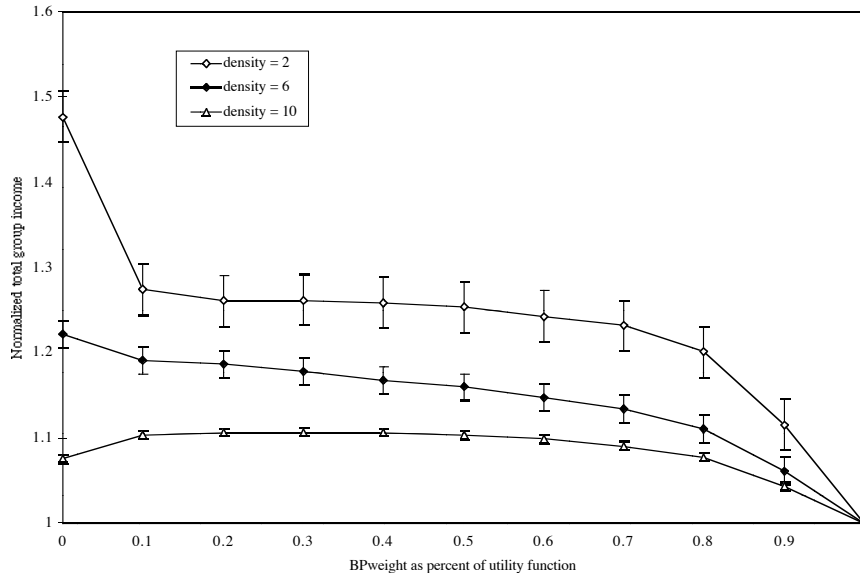


Figure 4: Effect of  $BPweight$  on mean total group income across various task densities.

signed with utility functions that optimize for number of defaults or group task income, with the assurance that total group income will remain near its maximum.

This conclusion does not hold, however, if one considers total income when  $BPweight$  is 0. Figure 4 appears to show that income is maximized at this lowest value for  $BPweight$ . It is unlikely, however, that this is the behavior that would actually be desired in collaborative agents. Agents with low  $BPweight$  default much more frequently than more socially conscious agents. Furthermore, Figure 2 shows that group task income is at its lowest for this value of  $BPweight$ . Since most agent designers would attempt to optimize all of these factors, keeping defaults low and all measures of income high, this extreme behavior would be undesirable despite the apparent payoff as measured by total group income.

The seemingly inconsistent change in total income as  $BPweight$  initially rises from 0 to 0.1 remains an open question. A similar “jump” over this same interval in an otherwise smooth curve can also be observed in Figure 2, although to a lesser extent. While we suspect that this behavior reveals the large effect that adding brownie points at *any* weight has on agent behavior, the exact cause has not been isolated. This result seems to indicate that agents with no social consciousness (*i.e.*, brownie points do not factor at all into their utility functions) have radically different behavior from even mildly socially conscious agents. Analysis into this behavior is an area for future investigation.

## 5. RELATED WORK

Kalenka and Jennings [12] propose several “socially responsible” decision-making principles and examine their effects in the context of a warehouse loading scenario. Our work dif-

fers from theirs in that their policies are domain-dependent and not decision-theoretic and they do not look at conflicting intentions but rather at whether or not agents choose to help each other. Sen [16] considers decision-making strategies that encourage cooperation among self-interested agents, but his work focuses on interactions between pairs of individual agents. Sen’s “philanthropic” agents take the good of the group into account, but do not always necessarily do what is best for the group.

There is a significant body of economics literature on rational choice and intention reconciliation. Iannaccone [10] examines social policies that alter individual utility functions to encourage group commitment. While these policies are similar in spirit to the social-commitment policies that SPIRE incorporates, they are aimed at group formation, not at conflicting intentions. Additionally, that approach is not applicable to agents that face multiple decision points over time. Höllander [9] studies incentives for encouraging group commitment and cooperation under a more limited definition of cooperation, in which an agent is required to incur a personal cost in order to cooperate. His model considers “emotional” cooperation within this limited definition, but assumes a rigid standard shared by all players, a requirement that we relax.

The social-commitment policies in SPIRE also differ from Shoham and Tennenholtz’s [17] social laws. Social laws constrain the ways agents perform actions whereas social-commitment policies constrain decision-making. Social laws are by their nature domain specific; they constrain domain actions. In contrast, social-commitment policies affect decision-making across domains and tasks. The conventions Rosen-schein and Zlotkin [15] present play a role in negotiation sim-



ilar to the role social-commitment policies play in SPIRE.

Cooper *et al.* [4] examine social consciousness in an approach similar to the brownie point model, referred to as the “warm glow” model, in which agents receive a constant amount of added utility when they do the “right thing.” Our approach differs from this warm glow model in that an agent’s utility from brownie points is task-value dependent and depends not just on the number of points gained for a given task but also on the total number of brownie points that it has stored up over time.

## 6. CONCLUSIONS

The brownie point model provides a framework in which non-monetary social factors can be effectively considered alongside monetary ones in agent decision-making. By considering not just the default itself but also task value and agent history in the calculation of brownie points, this model realistically emulates expected behavior. Additionally, because this system measures social consciousness over time and allows it to play a variable role in utility functions, the effect of social consciousness on the behavior of the group as a whole can be studied. Current work is extending these findings to heterogeneous societies in which agents have different levels of social consciousness [20].

## 7. ACKNOWLEDGMENTS

This research has been supported by National Science Foundation grants IRI-95-25915, IRI-96-18848, and CDA-94-01024. Sarit Kraus, David Sullivan, and Mike Epstein participated in the development of the SPIRE framework. This paper describes research that formed part of the first author’s undergraduate honors thesis at Harvard University. Luke Hunsberger and David Sullivan provided helpful comments on earlier drafts.

## 8. REFERENCES

- [1] Bratman, M.E., Israel, D.J., and Pollack, M.E. 1988. Plans and resource-bounded practical reasoning. *Computational Intelligence*, 4(4):349-355.
- [2] Castelfranchi, C. 1998. Modelling social action for AI agents. *Artificial Intelligence*, 103:157-182.
- [3] Conte, R. and Castelfranchi, C. 1995. *Cognitive and Social Action*. UCL Press, London, UK.
- [4] Cooper, R., DeJong, D.V., Forsythe, R., Ross, T.W. 1996. Cooperation without reputation: experimental evidence from prisoner’s dilemma games. *Games and Economic Behavior*, 12(1):187-218.
- [5] Gibbons, R. 1992. *Game Theory for Applied Economists*. Princeton University Press.
- [6] Grosz, B.J. and Kraus, S. 1996. Collaborative plans for complex group action. *Artificial Intelligence*, 86(2):269-357.
- [7] Grosz, B.J. and Kraus, S. 1999. The evolution of SharedPlans. Rao, A. and Wooldridge, M., editors, *Foundations of Rational Agency*. Kluwer Academic Press, 227-262.
- [8] Grosz, B.J., Hunsberger, L., and Kraus, S. 1999. Planning and Acting Together. *AI Magazine*, 20(4):23-34.
- [9] Höllander, H. 1990. A social exchange approach to voluntary cooperation. *American Economic Review*, 80(5):1157-1167.
- [10] Iannaccone, L.R. 1992. Sacrifice and stigma: reducing free-riding in cults, communes, and other collectives. *Journal of Political Economy*, 100(2):271-291.
- [11] Jennings, N. and Wooldridge, M. 1995. Agent theories, architecture, and languages: a survey. Jennings, N. and Wooldridge, M., editors, *Intelligent Agent: ECAI-94 Workshop on Agent Theories, Architectures, and Languages*. Springer Verlag.
- [12] Kalenka, S. and Jennings, N.R. 1999. Socially responsible decision making in autonomous agents. Korta, K., *et al.*, editors, *Cognition, Agency and Rationality*. Dordrecht: Kluwer, 79:135-149.
- [13] Keeney, R. and Raiffa, H. 1976. *Decisions with Multiple Objectives: Preferences and Value Tradeoffs*. John Wiley and Sons, New York.
- [14] Levesque, H., Cohen, P., and Nunes, J. 1990. On acting together. *Proceedings of AAAI-90*, 94-99.
- [15] Rosenschein, J.S. and Zlotkin, G. 1994. *Rules of Encounter: Designing Conventions for Automated Negotiation among Computers*. MIT Press.
- [16] Sen, S. 1996. Reciprocity: a foundational principle for promoting cooperative behavior among self-interested agents. *Proceedings of ICMAS-96*, 322-329.
- [17] Shoham, Y. and Tennenholtz, M. 1992. On the synthesis of useful social laws for artificial agent societies. *Proceedings of AAAI-92*, 276-281.
- [18] Sonenberg, E., Tidhar, G., Werner, E., Kinny, D., Ljungberg, M., and Rao, A.S. 1994. Planned team activity. Castelfranchi, C. and Werner, E., editors, *Artificial Social Systems, Lecture Notes in Artificial Intelligence (LNAI-830)*. Springer Verlag, 227-256.
- [19] Sullivan, D.G., Glass, A., Grosz, B.J., and Kraus, S. 1999. Intention reconciliation in the context of teamwork: an initial empirical investigation. Klusch, M., Shehory, O., and Weiss, G., editors, *Cooperative Information Agents III, Lecture Notes on Artificial Intelligence (LNAI-1652)*. Springer Verlag, 149-162.
- [20] Sullivan, D.G., Grosz, B.J., and Kraus, S. 2000. Intention reconciliation by collaborative agents. *Proceedings of ICMAS-2000* (to appear).
- [21] Tambe, M. 1997. Towards flexible teamwork. *Journal of Artificial Intelligence Research*, 7:83-124.
- [22] Yoon, K.P. and Hwang, C.-L. 1995. *Multiple Attribute Decision-Making: An Introduction*. Sage University papers, Series on Quantitative Applications in the Social Science, vol. 104.