

# PsychSim: Modeling Theory of Mind with Decision-Theoretic Agents

David V. Pynadath and Stacy C. Marsella

Information Sciences Institute, University of Southern California  
4676 Admiralty Way, Marina del Rey, CA 90292 USA  
{marsella,pynadath}@isi.edu

## Abstract

Agent-based modeling of human social behavior is an increasingly important research area. A key factor in human social interaction is our beliefs about others, a theory of mind. Whether we believe a message depends not only on its content but also on our model of the communicator. How we act depends not only on the immediate effect but also on how we believe others will react. In this paper, we discuss PsychSim, an implemented multiagent-based simulation tool for modeling interactions and influence. While typical approaches to such modeling have used first-order logic, PsychSim agents have their own decision-theoretic model of the world, including beliefs about its environment and recursive models of other agents. Using these quantitative models of uncertainty and preferences, we have translated existing psychological theories into a decision-theoretic semantics that allow the agents to reason about degrees of believability in a novel way. We discuss PsychSim's underlying architecture and describe its application to a school violence scenario for illustration.

## 1 Introduction

People interact within a rich social framework. To better understand people's social interactions, researchers have increasingly relied on computational models [Liebrand *et al.*, 1998; Prietula *et al.*, 1998]. Models of social interaction have also been used to create social training environments where the learner explores high-stress social interactions in the safety of a virtual world [Marsella *et al.*, 2000]. A key factor in human social interaction is our beliefs about others, a *theory of mind* [Whiten, 1991]. Our decisions to act are influenced by how we believe others will react. Whether we believe a message depends not only on its content but also on our model of the communicator. Giving its importance in human social interaction, modeling theory of mind can play a key role in enriching social simulations.

Typical approaches to modeling theory of mind in a computational framework have relied on first-order logic to represent beliefs and goals. However, such representations are often insensitive to the distinctions among conflicting goals

that people must balance in a social interaction. For example, psychological research has identified a range of goals that motivate classroom bullies (e.g., peer approval, sadism, tangible rewards). All bullies share the same goals, but it is the relative priorities that they place on them that leads to the variations in their behavior. Resolving the ambiguity among equally possible, but unequally plausible or preferred, options requires a quantitative model of uncertainty and preference. Unfortunately, more quantitative frameworks, like decision theory and game theory, face their own difficulties in modeling human psychology. Game theoretic frameworks typically rely on concepts of equilibria that people rarely achieve in an unstructured social setting like a classroom. Decision theoretic frameworks typically rely on assumptions of rationality that people constantly violate.

We have developed a social simulation tool, *PsychSim* [Marsella *et al.*, 2004], that operationalizes existing psychological theories as boundedly rational computations to generate more plausibly human behavior. PsychSim allows a user to quickly construct a social scenario where a diverse set of entities, groups or individuals, interact and communicate. Each entity has its own preferences, relationships (e.g., friendship, hostility, authority) with other entities, private beliefs, and mental models about other entities. The simulation tool generates the behavior for these entities and provides explanations of the result in terms of each entity's preferences and beliefs. The richness of the entity models allows one to explore the potential consequences of minor variations on the scenario. A user can play different roles by specifying actions or messages for any entity to perform.

A central aspect of the PsychSim design is that agents have fully specified decision-theoretic models of others. Such quantitative recursive models give PsychSim a powerful mechanism to model a range of factors in a principled way. For instance, we exploit this recursive modeling to allow agents to form complex attributions about others, enrich the messages between agents to include the beliefs and preferences of other agents, model the impact such recursive models have on an agent's own behavior, model the influence observations of another's behavior have on the agent's model of that other, and enrich the explanations provided to the user. The decision-theoretic models in particular give our agents the ability to judge degree of credibility of messages in a subjective fashion that can consider the range of influences that

sway such judgments in humans.

The rest of this paper describes PsychSim’s underlying architecture in more detail, using a school bully scenario for illustration. The agents represent different people and groups in the school setting. The user can analyze the simulated behavior of the students to explore the causes and cures for school violence. One agent represents a bully, and another represents the student who is the target of the bully’s violence. A third agent represents the group of onlookers, who encourage the bully’s exploits by, for example, laughing at the victim as he is beaten up. A final agent represents the class’s teacher trying to maintain control of the classroom, for example by doling out punishment in response to the violence.

## 2 The Agent Models

We embed PsychSim’s agents within a decision-theoretic framework for quantitative modeling of multiple agents. Each agent maintains its independent beliefs about the world, has its own goals and it owns policies for achieving those goals. The PsychSim framework is an extension to the Com-MTDP model [Pynadath and Tambe, 2002] of agent teamwork. To extend the Com-MTDP framework to social scenarios (where the agents are pursuing their own goals, rather than those of a team), we designed novel agent models for handling belief update and policy application, as described in Section 2.3.

### 2.1 Model of the World

Each agent model starts with a representation of its current state and the Markovian process by which that state evolves over time in response to the actions performed.

#### State

Each agent model includes several features representing its “true” state. This state consists of objective facts about the world, some of which may be hidden from the agent itself. For our example bully domain, we included such state features as `power(agent)`, to represent the strength of an agent. `trust(truster, trustee)` represents the degree of trust that the agent `truster` has in another agent `trustee`’s messages. `support(supporter, supportee)` is the strength of support that an agent `supporter` has for another agent `supportee`. We represent the state as a vector,  $\vec{s}^t$ , where each component corresponds to one of these state features and has a value in the range  $[-1, 1]$ .

#### Actions

Agents have a set of actions that they can choose to change the world. An action consists of an action type (e.g., `punish`), an agent performing the action (i.e., the actor), and possibly another agent who is the object of the action. For example, the action `laugh(onlooker, victim)` represents the laughter of the onlooker directed at the victim.

#### World Dynamics

The state of the world changes in response to the actions performed by the agents. We model these dynamics using a transition probability function,  $T(\vec{s}, \vec{a}, \vec{s}')$ , to capture the possibly uncertain effects of these actions on the subsequent state:

$$\Pr(\vec{s}^{t+1} = \vec{s}' \mid \vec{s}^t = \vec{s}, \vec{a}^t = \vec{a}) = T(\vec{s}, \vec{a}, \vec{s}') \quad (1)$$

For example, the bully’s attack on the victim impacts the power of the bully, the power of the victim, etc. The distribution over the bully’s and victim’s changes in power is a function of the relative powers of the two—e.g., the larger the power gap that the bully enjoys over the victim, the more likely the victim is to suffer a big loss in power.

### 2.2 Preferences

PsychSim’s decision-theoretic framework represents an agent’s incentives for behavior as a reward function that maps the state of the world into a real-valued evaluation of benefit for the agent. We separate components of this reward function into two types of subgoals. A goal of **Minimize/maximize** `feature(agent)` corresponds to a negative/positive reward proportional to the value of the given state feature. For example, an agent can have the goal of maximizing its own power. A goal of **Minimize/maximize** `action(actor, object)` corresponds to a negative/positive reward proportional to the number of matching actions performed. For example, the teacher may have the goal of minimizing the number of times any student teases any other.

We can represent the overall preferences of an agent, as well as the relative priority among them, as a vector of weights,  $\vec{g}$ , so that the product,  $\vec{g} \cdot \vec{s}^t$ , quantifies the degree of satisfaction that the agent receives from the world, as represented by the state vector,  $\vec{s}^t$ . For example, in the school violence simulation, the bully’s reward function consists of goals of maximizing `power(bully)`, minimizing `power(victim)`, and maximizing `laugh(onlookers, victim)`. By modifying the weights on the different goals, we can alter the motivation of the agent and, thus, its behavior in the simulation.

### 2.3 Beliefs about Others

As described by Sections 2.1 and 2.2, the overall decision problem facing a single agent maps easily into a partially observable Markov decision problem (POMDP) [Smallwood and Sondik, 1973]. Software agents can solve such a decision problem using existing algorithms to form their beliefs and then determine the action that maximizes their reward given those beliefs. However, we do not expect people to conform to such optimality in their behavior. Thus, we have taken the POMDP algorithms as our starting point and modified them in a psychologically motivated manner to capture more human-like behavior. This “bounded rationality” better captures the reasoning of people in the real-world, as well as providing the additional benefit of avoiding the computational complexity incurred by an assumption of perfect rationality.

#### Nested Beliefs

The simulation agents have only a *subjective* view of the world, where they form beliefs,  $\vec{b}^t$ , about what they *think* is the state of the world,  $\vec{s}^t$ . Agent *A*’s beliefs about agent *B* have the same structure as the real agent *B*. Thus, our agent belief models follow a recursive structure, similar to previous work on game-theoretic agents [Gmytrasiewicz and Durfee, 1995]. Of course, the nesting of these agent models is potentially unbounded. However, although infinite nesting is required for modeling optimal behavior, people rarely use such

deep models [Taylor *et al.*, 1996]. In our school violence scenario, we found that 2-level nesting was sufficiently rich to generate the desired behavior. Thus, the agents model each other as 1-level agents, who, in turn, model each other as 0-level agents, who do *not* have any beliefs. Thus, there is an inherent loss of precision (but with a gain in computational efficiency) as we move deeper into the belief structure.

For example, an agent’s beliefs may include its subjective view on states of the world: “The bully believes that the teacher is weak”, “The onlookers believe that the teacher supports the victim”, or “The bully believes that he/she is powerful.” These beliefs may also include its subjective view on beliefs of other agents: “The teacher believes that the bully believes the teacher to be weak.” An agent may also have a subjective view of the *preferences* of other agents: “The teacher believes that the bully has a goal to increase his power.” It is important to note that we also separate an agent’s subjective view of itself from the real agent. We can thus represent errors that the agent has in its view of itself (e.g., the bully believes himself to be stronger than he actually is).

Actions affect the beliefs of agents in several ways. For example, the bully’s attack may alter the beliefs that agents have about the state of the world—such as beliefs about the bully’s power. Each agent updates its beliefs according to its subjective beliefs about the world dynamics. It may also alter the beliefs about the bully’s preferences and policy. We discuss the procedure of belief update in Section 2.4.

### Policies of Behavior

Each agent’s policy is a function,  $\pi(\vec{b})$ , that represents the process by which it selects an action or message based on its beliefs. An agent’s policy allows us to model critical psychological distinctions such as reactive vs. deliberative behavior. We model each agent’s real policy as a bounded lookahead procedure that seeks to maximize expected reward simulating the behavior of the other agents and the dynamics of the world in response to the selected action/message. Each agent  $i$  computes a quantitative value,  $V_a(\vec{b}_i^t)$ , of each possible action,  $a$ , given its beliefs,  $\vec{b}_i^t$ .

$$V_a(\vec{b}_i^t) = \vec{g}_i \cdot \vec{b}_i^t + \sum_{\vec{b}^{t+1}} V(\vec{b}^{t+1}) \Pr(\vec{b}^{t+1} | \vec{b}_i^{t+1}, a, \vec{\pi}_{-i}(\vec{b}_i^{t+1})) \quad (2)$$

$$V(\vec{b}^t) = \vec{g}_i \cdot \vec{b}_i^t + \sum_{\vec{b}^{t+1}} V(\vec{b}^{t+1}) \Pr(\vec{b}^{t+1} | \vec{b}_i^{t+1}, \vec{\pi}(\vec{b}_i^{t+1})) \quad (3)$$

Thus, an agent first uses the transition function,  $T$ , to project the immediate effect of the action,  $a$ , and then projects another  $N$  steps into the future, weighing each state against its goals,  $\vec{g}$ . At the first step, agent  $i$  uses its model of the policies of all of the other agents,  $\vec{\pi}_{-i}$ , and, in subsequent steps, it uses its model of the policies of all agents, including itself,  $\vec{\pi}$ . Thus, the agent is seeking to maximize the expected reward of its behavior as in a POMDP. However, PsychSim’s agents are only boundedly rational, given that they are constrained, both by the finite horizon,  $N$ , of their lookahead and the possible error in their belief state,  $\vec{b}$ . By varying  $N$  for different agents, we can model entities who display different degrees of reactive vs. deliberative behavior in their thinking.

### Stereotypical Mental Models

If we applied this full lookahead policy within the nested models of the other agents, the computational complexity of the top-level lookahead would quickly become infeasible as the number of agents grew. To simplify the agents’ reasoning, these mental models are realized as simplified stereotypes of the richer lookahead behavior models of the agents themselves. For our simulation model of a bullying scenario, we have implemented mental models corresponding to *attention-seeking*, *sadistic*, *dominance-seeking*, etc. For example, a model of an attention-seeking bully specifies a high priority on increasing the approval (i.e., support) that the other agents have for it, a dominance-seeking bully specifies a high priority on increasing its power as paramount, and a bully agent specifies a high priority on hurting others.

These simplified mental models also include potentially erroneous beliefs about the policies of other agents. Although the real agents use lookahead exclusively when choosing their own actions (as described in Section 2.3), the agents *believe* that the other agents follow much more reactive policies as part of their mental models of each other. PsychSim models reactive policies as a table of “Condition  $\Rightarrow$  Action” rules. The left-hand side conditions may trigger on an *observation* of some action or a *belief* of some agent (e.g., such as the bully believing himself as powerful). The conditions may also be more complicated combinations of these basic triggers (e.g., a *conjunction* of conditions that matches when each and every individual condition matches).

The use of these more reactive policies in the mental models that agents have of each other achieves two desirable results. First, from a human modeling perspective, the agents perform a shallower reasoning that provides a more accurate model of the real-world entities they represent. Second, from a computational perspective, the direct action rules are cheap to execute, so the agents gain significant efficiency in their reasoning.

## 2.4 Modeling Influence and Belief Change

### Messages

Messages are attempts by one agent to influence the beliefs of another. Messages have four components: source, recipients, subject, and content. For example, the teacher (source) could tell the bully (recipient) that the principal (subject of the message) will punish violence by the bully (content). Messages can refer to beliefs, preferences, policies, or any other aspect of other agents. Thus, a message may make a claim about a state feature of the subject (“the principal is powerful”), the beliefs of the subject (“the principal believes that he is powerful”), the preferences of the subject (“the bully wants to increase his power”), the policy of the subject (“if the bully thinks the victim is weak, he will pick on him”), or the stereotypical model of the subject (“the bully is selfish”).

### Influence Factors

A challenge in creating a social simulation is addressing how groups or individuals influence each other, how they update their beliefs and alter behavior based on any partial observation of, as well as messages from, others. Although many psychological results and theories must inform the modeling

of such influence (e.g., [Cialdini, 2001; Abelson *et al.*, 1968; Petty and Cacioppo, 1986]) they often suffer from two shortcomings from a computational perspective. First, they identify factors that affect influence but do not operationalize those factors. Second, they are rarely comprehensive and do not address the details of how various factors relate to each other or can be composed. To provide a sufficient basis for our computational models, our approach has been to distill key psychological factors and map those factors into our simulation framework. Here, our decision-theoretic models are helpful in quantifying the impact of factors in such a way that they can be composed. Specifically, a survey of the social psychology literature identified the following key factors:

**Consistency:** People expect, prefer, and are driven to maintain consistency, and avoid cognitive dissonance, between beliefs and behaviors.

**Self-interest:** The inferences we draw are biased by self-interest (e.g., motivated inference) and how deeply we analyze information in general is biased by self-interest.

**Speaker’s Self-interest:** If the sender of a message benefits greatly if the recipient believes it, there is often a tendency to be more critical and for influence to fail.

**Trust, Likability, Affinity:** The relation to the source of the message, whether we trust, like or have some group affinity for him, all impact whether we are influenced by the message.

### Computational Model of Influence

To model such factors in the simulation, one could specify them exogenously and make them explicit, user-specified factors for a message. This tactic is often employed in social simulations where massive numbers of simpler, often identical, agents are used to explore emergent social properties. However, providing each agent with quantitative models of itself and, more importantly, of other agents gives us a powerful mechanism to model this range of factors in a principled way. We model these factors by a few simple mechanisms in the simulation: *consistency*, *self-interest*, and *bias*. We can render each as a quantitative function of beliefs that allows an agent to compare alternate candidate belief states (e.g., an agent’s original  $\vec{b}$  vs. the  $\vec{b}'$  implied by a message).

**Consistency** is an evaluation of the degree to which a potential belief agreed with prior observations. In effect, the agent asks itself, “If this belief holds, would it better explain the past better than my current beliefs?”. We use a Bayesian definition of consistency based on the relative likelihood of past observations given the two candidate sets of beliefs (e.g., my current beliefs with and without believing the message). An agent assesses the quality of the competing explanations by a re-simulation of the past history. In other words, it starts at time 0 with the two worlds implied by the two candidate sets of beliefs, projects each world forward up to the current point of time, and computes the probability of the observation it received. In particular, the consistency of a sequence of observations,  $\omega^0, \omega^1, \dots$ , with a given belief state,  $\vec{b}$ , cor-

responds to:

$$\begin{aligned} & \text{consistency}(\vec{b}^t, [\omega^0, \omega^1, \dots, \omega^{t-1}]) \\ &= \Pr \left( [\omega^0, \omega^1, \dots, \omega^{t-1}] \mid \vec{b}^t \right) \\ &\propto \sum_{\tau=0}^{t-1} \sum_{a \in A} V_a(\vec{b}^\tau) \Pr(\omega^\tau \mid a, \vec{b}^\tau) \end{aligned} \quad (4)$$

The value function,  $V$ , computed is with respect to the agent performing the action at time  $\tau$ . We are summing the value of the observed action to the acting agent, given the set of beliefs under consideration. The higher the value, the more likely that agent is to have chosen the observed action, and, thus, the higher the degree of consistency.

**Self-interest** is similar to consistency, in that the agent compares two sets of beliefs, one which accepts the message and one which rejects it. However, while consistency evaluates the past, we compute self-interest by evaluating the future using Equation 3. An agent can perform an analogous computation using its beliefs about the sender’s preferences to compute the sender’s self-interest in sending the message.

**Bias** factors represent subjective views of the message sender that influence the receiver’s acceptance/rejection of the message. We treat support (or affinity) and trust as such a bias on message acceptance. Agents compute their support and trust levels as a running history of their past interactions. In particular, one agent increases (decreases) its trust in another, when the second sends a message that the first decides to accept (reject). Similarly, an agent increases (decreases) its support for another, when the second selects an action that has a high (low) reward, with respect to the preferences of the first. In other words, if an agent selects an action  $a$ , then the other agents modify their support level for that agent by a value proportional to  $\vec{g} \cdot \vec{b}$ , where  $\vec{g}$  corresponds to the goals and  $\vec{b}$  the new beliefs of the agent modifying its support.

Upon receiving any information (whether message or observation), an agent must consider all of these various factors in deciding whether to accept it and how to alter its beliefs (including its mental models of the other agents). For a message, the agent determines acceptance using a weighted sum of the five components: consistency, self-interest, speaker self-interest, trust and support. Whenever an agent observes an action by another, it checks whether the observation is consistent with its current beliefs (including mental models). If so, no belief change is necessary. If not, the agent evaluates alternate mental models as possible new beliefs to adopt in light of this inconsistent behavior. Agents evaluate these possible belief changes using the same weighted sum as for messages.

Each agent’s decision-making procedure is sensitive to these changes that its actions may trigger in the beliefs of others. Each agent accounts for the others’ belief update when doing its lookahead, as Equations 2 and 3 project the future beliefs of the other agents in response to an agent’s selected action. Similar to work by [de Rosi *et al.*, 2003] this mechanism provides PsychSim agents with a potential incentive to deceive, if doing so leads the other agents to perform actions that lead to a better state for the deceiving agent.

We see the computation of these factors as a toolkit for the user to explore the system’s behavior under existing theories, which we can encode in PsychSim. For example, the elaboration likelihood model (ELM) [Petty and Cacioppo, 1986] argues that the way messages are processed differs according to the relevance of the message to the receiver. High relevance or importance would lead to a deeper assessment of the message, which is consistent with the self-interest calculations our model performs. PsychSim’s linear combination of factors is roughly in keeping with ELM because self-interest values of high magnitude would tend to dominate.

### 3 PsychSim in Operation

The research literature on childhood aggression provides interesting insight into the role that theory of mind plays in human behavior. Investigations of bullying and victimization [Schwartz, 2000] have identified four types of children; we focus here on *nonvictimized aggressors*, those who display proactive aggression due to positive outcome expectancies for aggression. Children develop expectations on the likely outcomes of aggression based on past experiences (e.g., did past acts of aggression lead to rewards or punishment). This section describes the results of our exploration of the space of different nonvictimized aggressors and the effectiveness of possible intervention strategies in dealing with them.

#### 3.1 Scenario Setup

The user sets up a simulation in PsychSim by selecting generic agent models that will play the roles of the various groups or individuals to be simulated and specializing those models as needed. In our bullying scenario, we constructed generic bully models that compute outcome expectancies as the expected value of actions ( $V_a$  from Equation 2). Thus, when considering possible aggression, the agents consider the immediate effect of an act of violence, as well as the possible consequences, including the change in the beliefs of the other agents. In our example scenario, a bully has three subgoals that provide incentives to perform an act of aggression: (1) to change the power dynamic in the class by making himself stronger, (2) to change the power dynamic by weakening his victim, and (3) to earn the approval of his peers (as demonstrated by their response of laughter at the victim). Our bully agent models the first incentive as a goal of maximizing `power(bully)` and the second as minimizing `power(victim)`, both coupled with a belief that an act of aggression will increase the former and decrease the latter. The third incentive seeks to maximize the `laugh` actions directed at the victim, so it must consider the actions that the other agents may take in response.

For example, a bully motivated by the approval of his classmates would use his mental model of them to predict whether they would laugh along with him. We implemented two possible mental models of the bully’s classmates: *encouraging*, where the students will laugh at the victim, and *scared*, where the students will laugh only if the teacher did not punish them for laughing last time. Similarly, the bully would use his mental model of the teacher to predict whether he will be punished or not. We provide the bully with three possible mental

models of the teacher: *normal*, where the teacher will punish the bully in response to an act of violence; *severe*, where the teacher will more harshly punish the bully than in the *normal* model; and *weak*, where the teacher never punishes the bully.

The relative priorities of these subgoals within the bully’s overall reward function provide a large space of possible behavior. When creating a model of a specific bully, PsychSim uses a fitting algorithm to automatically determine the appropriate weights for these goals to match observed behavior. For example, if the user wants the bully to initially attack a victim and the teacher to threaten the bully with punishment, then the user specifies those behaviors and the model parameters are fitted accordingly [Pynadath and Marsella, 2004]. This degree of automation significantly simplifies simulation setup. In this experiment, we selected three specific bully models from the overall space: (1) *dominance-seeking*, (2) *sadistic*, and (3) *attention-seeking*, each corresponding to a goal weighting that favors the corresponding subgoal.

#### 3.2 Experimental Results

PsychSim allows one to explore multiple tactics for dealing with a social issue and see the potential consequences. Here, we examine a decision point for the teacher after the bully has attacked the victim, followed by laughter by the rest of the class. At this point, the teacher can punish the bully, punish the whole class (including the victim), or do nothing. We explore the impact of different types of proactive aggression by varying the type of the bully, the teacher’s decision to punish the bully, the whole class, or no one, and the mental models that the bully has of the other students and the teacher.

A successful outcome is when the bully does not choose to act out violently toward the victim the next time around. By examining the outcomes under these combinations, we can see the effects of intervention over the space of possible classroom settings. Table 1 shows all of the outcomes, where we use the “\*” wildcard symbol to collapse rows where the outcome was the same. Similarly, a row with “¬severe” in the **Teacher** row spans the cases where the bully’s mental model of the teacher is either *normal* or *weak*.

We first see that the PsychSim bully models meets our intuitive expectations. For example, we see from Table 1 that if the bully thinks that the teacher is too weak to ever punish, then no immediate action by the teacher will change the bully from picking on the victim. Thus, it is critical for the teacher to avoid behavior that leads the bully to form such mental models. Similarly, if the bully is of the *attention-seeking* variety, then punishment directed at solely himself will not dissuade him, as he will still expect to gain peer approval. In such cases, the teacher is better off punishing the whole class.

We can see more interesting cases as we delve deeper. For example, if we look at the case of a *sadistic* bully when the teacher punishes the whole class, we see that bully can be dissuaded only if he thinks that the other students will *approve* of his act of violence. This outcome may seem counter-intuitive at first, but the *sadistic* bully is primarily concerned with causing suffering for the victim, and thus does not mind being punished if the victim is punished as well. However, if the bully thinks that the rest of the class is *encouraging*, then the teacher’s punishment of the whole class costs him peer

Bully Type	Punish Whom?	Model of Students	Model of Teacher	Success?	
Sadistic	bully	*	$\neg$ severe	N	
			severe	Y	
	class	scared	*	N	
			encouraging	$\neg$ severe	N
	no one	*	$\neg$ severe	N	
			severe	Y	
Attention Seeking	bully	*	*	N	
			class	scared	weak
	normal	Y			
	severe	Y			
	no one	*	encouraging	*	N
			*	*	N
Dominance Seeking	*	*	weak	N	
			normal	Y	
			severe	Y	

Table 1: Outcomes of intervention strategies

approval. On the other hand, if the bully thinks that the rest of the class is already *scared*, so that they will not approve of his violence, then he has no peer approval to lose.

Such exploration can offer the user an understanding of the potential pitfalls in implementing an intervention strategy. Rather than providing a simple prediction of whether a strategy will succeed or not, PsychSim maps out the key conditions, in terms of the bully’s preferences and beliefs, on which a strategy’s success depends. PsychSim provides a rich space of possible models that we can systematically explore to understand the social behavior that arises out of different configurations of student psychologies. We are continuing to investigate more class configurations and the effects of possible interventions as we expand our models to cover all of the factors in school aggression identified in the literature.

## 4 Conclusion

We have discussed PsychSim, an environment for multi-agent simulation of human social interaction that employs a formal decision-theoretic approach using recursive models. This approach allows us to model phenomena rarely if at all addressed in simulated worlds. Within PsychSim, we have developed a range of technology to simplify the task of setting up the models, exploring the simulation, and analyzing results. This includes new algorithms for fitting multi-agent simulations. There is also an ontology for modeling communications about theory of mind. Finally, there is an analysis/perturbation capability that supports the iterative refinement of the simulation by reporting sensitivities in the results, as well as potentially interesting perturbations to the scenario. We have exploited the recursive models to provide a psychologically motivated computational model of how agents influence each other’s beliefs. We believe PsychSim has a range of innovative applications, including computational social science and the model of social training environments. Our current goals are to expand the exploration already begun in the school violence scenario and begin evaluating the application

of PsychSim there and in these other areas.

## References

- [Abelson *et al.*, 1968] R.P. Abelson, E. Aronson, W.J. McGuire, T.M. Newcomb, M.J. Rosenberg, and P.H. Tannenbaum, eds. *Theories of Cognitive Consistency: A Sourcebook*. Rand McNally, Chicago, IL, 1968.
- [Cialdini, 2001] R. Cialdini. *Influence: Science and Practice*. Allyn and Bacon, Boston, MA, 2001.
- [de Rosiis *et al.*, 2003] F. de Rosiis, C. Castelfranchi, V. Carofiglio, and G. Grassano. Can computers deliberately deceive? A simulation tool and its application to Turing’s imitation game. *Computational Intelligence*, 19(3):253–263, 2003.
- [Gmytrasiewicz and Durfee, 1995] P.J. Gmytrasiewicz and E.H. Durfee. A rigorous, operational formalization of recursive modeling. In *ICMAS*, pp. 125–132, 1995.
- [Liebrand *et al.*, 1998] W. Liebrand, A. Nowak, and R. Hegselmann, editors. *Computer Modeling of Social Processes*. Sage, London, UK, 1998.
- [Marsella *et al.*, 2000] S.C. Marsella, W.L. Johnson, and C. LaBore. Interactive pedagogical drama. In *Agents*, pp. 301–308, 2000.
- [Marsella *et al.*, 2004] S.C. Marsella, D.V. Pynadath, and S.J. Read. PsychSim: Agent-based modeling of social interactions and influence. In *ICCM*, pp. 243–248, 2004.
- [Petty and Cacioppo, 1986] R. Petty and J. Cacioppo. *Communication and persuasion: Central and peripheral routes to attitude change*. Springer, New York, NY, 1986.
- [Prietula *et al.*, 1998] M. Prietula, K. Carley, and L. Gasser, eds. *Simulating Organizations: Computational Models of Institutions and Groups*. AAAI Press, Menlo Park, CA, 1998.
- [Pynadath and Marsella, 2004] D.V. Pynadath and S.C. Marsella. Fitting and compilation of multiagent models through piecewise linear functions. In *AAMAS*, pp. 1197–1204, 2004.
- [Pynadath and Tambe, 2002] D.V. Pynadath and M. Tambe. Multiagent teamwork: Analyzing the optimality and complexity of key theories and models. In *AAMAS*, pp. 873–880, 2002.
- [Schwartz, 2000] D. Schwartz. Subtypes of victims and aggressors in children’s peer groups. *Journal of Abnormal Child Psychology*, 28:181–192, 2000.
- [Smallwood and Sondik, 1973] R.D. Smallwood and E.J. Sondik. The optimal control of partially observable Markov processes over a finite horizon. *Operations Research*, 21:1071–1088, 1973.
- [Taylor *et al.*, 1996] J. Taylor, J. Carletta, and C. Mellish. Requirements for belief models in cooperative dialogue. *User Modeling and User-Adapted Interaction*, 6:23–68, 1996.
- [Whiten, 1991] A. Whiten, editor. *Natural Theories of Mind*. Basil Blackwell, Oxford, UK, 1991.