# Modeling appraisal in theory of mind reasoning

**Mei Si · Stacy C. Marsella · David V. Pynadath**

**Abstract**    Cognitive appraisal theories, which link human emotional experience to their interpretations of events happening in the environment, are leading approaches to model emotions. Cognitive appraisal theories have often been used both for simulating "real emotions" in virtual characters and for predicting the human user's emotional experience to facilitate human–computer interaction. In this work, we investigate the computational modeling of appraisal in a multi-agent decision-theoretic framework using Partially Observable Markov Decision Process-based (POMDP) agents. Domain-independent approaches are developed for five key appraisal dimensions (motivational relevance, motivation congruence, accountability, control and novelty). We also discuss how the modeling of theory of mind (recursive beliefs about self and others) is realized in the agents and is critical for simulating social emotions. Our model of appraisal is applied to three different scenarios to illustrate its usages. This work not only provides a solution for computationally modeling emotion in POMDP-based agents, but also illustrates the tight relationship between emotion and cognition—the appraisal dimensions are derived from the processes and information required for the agent's decision-making and belief maintenance processes, which suggests a uniform cognitive structure for emotion and cognition.

**Keywords**   Emotion · Decision-making · Appraisal · Multi-agent system

## 1 Introduction

Researchers have increasingly argued that the modeling of human emotion should play an important role in a wide range of intelligent systems and specifically in agent-based systems.

M. Si (✉) · S. C. Marsella · D. V. Pynadath
Institute for Creative Technologies, University of Southern California, Los Angeles, CA, USA
e-mail: meisi@ict.usc.edu

S. C. Marsella
e-mail: marsella@ict.usc.edu

D. V. Pynadath
e-mail: pynadath@ict.usc.edu

Springer

For example, much as emotion plays a role in human–human interaction, the modeling of the user's emotional state has been proposed as a way to facilitate human–computer interactions [13,19]. Similarly, research in embodied conversational agents, virtual characters that can engage in spoken dialog with users, argues that emotion is a key aspect of building virtual characters. To simulate realistic interactions between virtual agents and humans, the virtual characters need to be able to emotionally react to events happening in the environment as well as form expectations about the human user's emotional responses [1,7,15,33].

Computational models of emotion used in agents have often been based on appraisal theory [1,4,5,17,18,21,34], a leading psychological theory of emotion. Appraisal theory argues that a person's subjective assessment of their relationship to the environment, the person-environment relation, determines the person's emotional responses [6,10,11,18,22,24,31, 32]. This assessment occurs along several dimensions, called appraisal variables or checks, such as motivational congruence, accountability, novelty and control. Emotion is decided by the combination of results from these checks. For example, an event that leads to a bad outcome for a person (motivationally incongruent) and is caused by others (accountability) is likely to elicit anger response; but if the event is caused by the person himself/herself, that person is more likely to feel guilt or regret [22]. A real life example of this would be that an employee feels angry if their supervisor has evaluated them unfairly; however, if the employee instead believes that he/she receives a negative evaluation because of his/her own fault, he/she is more likely to feel regret.

The work we report here investigates how to computationally model appraisal within a decision-theoretic framework. Our model is built within the Thespian framework [26–29] for authoring and simulating computer-aided interactive narratives. Computer-aided interactive narratives allow the user to actively participate in the development of a story. The user can play a role in the story and interact with virtual characters realized by software agents. This model of appraisal enables Thespian agents to express emotions as well as to anticipate other's emotions.

We approached the task of incorporating appraisal into the existing Thespian multi-agent framework as a form of thought experiment. We wanted to assess to what extent the processes and representations necessary for modeling appraisal were already incorporated into Thespian's belief revision and decision-making processes. Could we leverage the existing processes and representations to model appraisal? The motivations for this thought experiment were two-fold. We sought to demonstrate how appraisal is in some ways a blueprint, or requirements specification, for intelligent social agents by showing that an existing social agent framework that had not been designed with emotion or appraisal in mind had in fact appraisal-like processes. In addition, we sought a design that was elegant, that reused architectural features to realize new capabilities such as emotion. An alternative approach for creating embodied conversational agents and virtual agents is through integrating modules for emotion, decision-making, dialogue, etc. This can lead to sophisticated but complex architectures [33]. The work here can thus be viewed as part of an alternative minimalist agenda for agent design.

Various computational models for appraisal have been proposed (see Sect. 2.2 for a review.) Key questions in designing computational models of appraisal include how the person-environment relation is represented and how the appraisal processes operate over that representation. Often agent-based models of emotion leverage the agent's decision-making representations to model the person–environment relation [1,7]. For example, EMA [7,15] defines appraisal processes as operations over a uniform plan-based representation, termed a causal interpretation, of the agent's goals and how events impact those goals. Cognitive

processes maintain the causal interpretation and appraisal processes leverage this uniform representation to generate appraisal.

This work is in the spirit of, and closely related to, work on the EMA model of emotion. This work seeks to go further by detailing how the cognitive processes themselves need to realize appraisal as part of decision-making and belief update. Whereas EMA exploits a uniform representation for appraisal and cognition, we seek to identify overlaps not only in representational requirements but also seek to exploit more extensively the overlap in the processes underlying cognition so that appraisal becomes an integral part of the cognitive processes that a social agent must perform to maintain its beliefs about others and to inform its decision-making in a multi-agent social context.

A key distinction between this work and other computational models including EMA is Thespian's modeling of theory of mind, which is a key factor in human social interaction [36], and the role theory of mind plays in decision-making and belief revision. Agents in Thespian possess beliefs about other agents that constitute a fully specified, quantitative model of the other agents' beliefs, policies and goals. In other words, the agents have a theory of mind capability with which they can simulate others. Thespian's representation of agents' subjective beliefs about each other enables the model to better reason about social emotions, in effect agents can reason about other agent's cognitive and emotional processes both from the other agent's and its own perspectives. For example, if an agent's actions hurt another agent's utility, Thespian agent has the capacity to "feel" regret about the situation and at the same time anticipate that the other agent will "feel" angry.

In the work reported here, we focus on five appraisal variables: motivational relevance, motivational congruence, accountability, control and novelty. We demonstrate the application of our model in three different scenarios, including the Little Red Riding Hood fairy tale, small talk between two persons and a firing-squad scenario as described in [14]. The Little Red Riding Hood story will be used as an example to motivate the discussion throughout this paper.

## 2 Related work

In this section we briefly review prevailing cognitive appraisal theories, which provide theoretical background for this work. We also discuss existing computational models of appraisal in comparison to our new model.

### 2.1 Cognitive appraisal theories

Roseman and Smith [23] in their review of cognitive appraisal theories roughly divided recent theories into two categories. One is the "structural models", which concentrate on the content being evaluated—the appraisal dimensions. The other is the "process models", which try to explain the processes that evaluate the content.

Theories in the "structural models" category have significant overlaps on key appraisal dimensions, such as motivational relevance and congruence, causal attribution and coping potential, but also have differences on which dimensions are included and how the dimensions are defined [23]. For example, Roseman [22] proposed five appraisal dimensions for emotion: positive/negative, appetitive/aversive, caused by circumstances/others/self, certain/uncertain, deserved/underserved. Smith and Ellsworth [31] proposed ten dimensions, including pleasantness, certainty, responsibility, control, subjective importance, etc. Similarly, the OCC

model [18], which predicts people's valenced reaction to events based on their goals, evaluates the event's relevance, desirability and causal attribution.

Leventhal and Scherer's [12,24] model can be viewed as a "process model". They view emotions as the outcome of going through a fixed sequence of Stimulus Evaluation Checks (SECs). These checks largely overlap with the appraisal dimensions included in other theories and are grouped into four appraisal objectives. The first one is relevance detection which checks for novelty and goal relevance. The second once is implementation assessment and contains checks such as causal attribution check and goal conduciveness check. The third objective is coping potential determination which evaluates the person's control and power over the situation. Finally, the last objective is normative significance evaluation, which includes both internal standards check and external standards check.

Lazarus et al. [9,10,32] described two types of appraisal, primary appraisal and secondary appraisal. Primary appraisal refers to the significance of the event, which is evaluated by irrelevant encounter, benign-positive encounter and stressful encounter. Each of the counters involves several appraisal dimensions. Secondary appraisal is invoked when the event is appraised as stressful. It evaluates the person's potential for coping. The result of the evaluation will be taken into account by the person's following primary appraisal, and thus form an appraisal-coping-reappraisal loop in people's cognitive/emotion generation process.

The computational model we report here is based on Smith and Lazarus [32]. We demonstrate how the appraisal-coping-reappraisal loop can be flexibly modeled in decision-theoretic goal-based agents. Our model currently includes five appraisal dimensions: motivational relevance, motivational congruence, accountability, control and novelty. We adapted Smith and Lazarus's [32] definitions for modeling motivational relevance, motivational congruence and accountability. Our model of control is roughly equivalent to Smith and Lazarus's [32] definition of problem-focused coping potential, though it is closer to the concept of control in Scherer's [24] theory because it accounts for the overall changeability of the situation and not an individual agent's power to make a change. Finally, novelty is not an appraisal dimension in Smith and Lazarus's [32] theory because they refer the response resulted from a novel stimulus as an affective response rather than an emotional response. The evaluation of novelty is useful for driving virtual characters' non-verbal behaviors and therefore is included in our model. We used Leventhal and Scherer's [12,24] definition of predictability-based novelty to inform our computational model.

## 2.2 Computational models of appraisal

Cognitive appraisal theories have had an increasing impact on the design of virtual agents. Various computational models have been proposed. In FLAME, El Nasr et al. [4] use domain-independent fuzzy logic rules to simulate appraisal. In WILL [17], concern and relevance are evaluated as the discrepancies between the agent's desired state and the current state. Cathexis [34] uses a threshold model to simulate basic variables, which are called "sensors", related to emotion. The OCC model of appraisal [18] has inspired many computational systems. Elliott's [5] Affective Reasoner uses a set of domain-specific rules to appraise events based on the OCC model. Both EM [21] and ParleE [2] deployed the OCC model of emotion over plan-based agents. FearNot! [1] also applied the OCC model for emotion. In FearNot! there are two types of appraisal processes. The reactive appraisal processes directly link the perceptive input and the correspondent memories to generate emotion. The deliberative appraisal processes appraise the environmental perceptional inputs in the light of the agent's current plans, intentions and goals.

Our approach to modeling emotion is inspired by the EMA work [7,15], which follows Smith and Lazarus [32]. In EMA, the cognitive processes for constructing the person-environment relation representation is treated as distinct from appraisal, and appraisal is reduced to simple and fast pattern matching over the representation. Similarly, in Thespian we treat appraisal as leveraging the representations generated by the agent's decision-making and belief revision. We illustrate how key appraisal variables can be straightforwardly extracted from these representations. The difference between this model and EMA is that we seek to go further by arguing that appraisal is an integral part of cognition.

## 3 Example domains

We will demonstrate the application of this domain-independent model of appraisal in three different scenarios, including a simple conversation between two persons, a firing-squad scenario as modeled in [14], and a fairy tale, "the Little Red Riding Hood". The last scenario will be described here as it will be used as an example to motivate the discussion throughout this paper. The details of the other two scenarios are given in Sect. 6.

The story "the Little Red Riding Hood" contains four main characters, Little Red Riding Hood, Granny, the hunter and the wolf. The story starts as Little Red Riding Hood (Red) and the wolf meet each other on the outskirt of a wood while Red is on her way to Granny's house. The wolf has a mind to eat Red, but it dare not because there are some wood-cutters close by. At this point, they can either have a conversation or choose to walk away. The wolf will have a chance to eat Red at other locations where nobody is close by. Moreover, if the wolf heard about Granny from Red, it can even go eat her. Meanwhile, the hunter is searching for the wolf to kill it. Once the wolf is killed, people who got eaten by the wolf can escape.

## 4 Thespian

Thespian is a multi-agent framework for authoring and simulating computer-aided interactive narratives. Thespian is built upon PsychSim [16,20], a multi-agent framework for social simulation based on Partially Observable Markov Decision Process (POMDP) [30]. To date, the Thespian framework has been applied to authoring more than thirty interactive narratives in different domains, including both training and entertainment domains.

Thespian's basic architecture uses POMDP-based agents to control each character in the story, with the character's personality and motivations encoded as agent goals. The ability of goal-based agents to decide their actions based on both the environment and their goals makes Thespian agents react to the user and behave with consistent personalities/motivations. In this section we present the basic structure of Thespian agents and their belief revision and decision-making processes.

### 4.1 Thespian agent

Thespian agents are POMDP-based agents built for modeling virtual humans and social groups. Each agent is composed of state, dynamics, goals, beliefs (theory of mind), policy and social relationships.

### 4.1.1 State

State contains information about an agent's current status in the world. An agent's state is defined by a set of state features, such as the name and age of the character, and the relation between that character and other characters (e.g. affinity). Values of state features are represented as real numbers.

### 4.1.2 Dynamics

Dynamics define how actions affect agents' states. For example, we can specify that small talk among a group of agents will increase their affinity with each other by 0.1. The effects of actions can be defined with probabilities. For example, the author may define that when the hunter tries to kill the wolf, the wolf will die only 60% of the time.

### 4.1.3 Goals

We model a character's motivation and personality profile as a set of goals and their relative importance (weight). Goals are expressed as a reward function over the various state features an agent seeks to maximize or minimize. For example, a character can have a goal of maximizing its affinity with another character. The initial value of this state feature can be any value between 0.0 and 1.0; this goal is completely satisfied once the value reaches 1.0. An agent usually has multiple goals with different relative importance (weights). For example, the character may have another goal of knowing another character's name, and this goal may be twice as important to the character as the goal of maximizing affinity. At any moment of time, an agent's utility is simply calculated as $State \times Goals$.

### 4.1.4 Beliefs (theory of mind)

Thespian agents have a "theory of mind". The agent's subjective view of the world includes its beliefs about itself and other agents and *their* subjective views of the world, a form of recursive agent modeling. An agent's subjective view (mental model) of itself or another agent includes every component of that agent, such as state, beliefs, policy, etc.

Each agent has a mental model of self and one or more mental models of other agents. The agent's belief about another agent is a probability distribution over alternative mental models. For example, in the Red Riding Hood story, Red can have two mental models of the wolf—one being that the wolf does not have a goal of eating people and one being otherwise. Initially, Red may believe that there is a 90% chance the first mental model is true and a 10% chance the second mental model is true. This probability distribution will change if Red sees or hears about the wolf eating people.

Within each mental model, an agent's belief about its own or another agent's state is represented as a set of real values with probability distributions. The probability distribution of the possible values of state features indicates the character's beliefs about these values. For example, a character's belief about the amount of money another character has could be {8 with probability of 90%, 0 with probability of 10%}.[1] When reasoning about utilities of actions, the expected value of a state feature is normally used, which is simply calculated as $\sum_{i=0}^{n} value_i \times P(value_i)$. For the simplicity of demonstration, in this paper we only give examples using the expected values.

---

[1] This example only includes one state feature for the simplicity of demonstration. In general, the probability distribution is associated with the values of all state features.
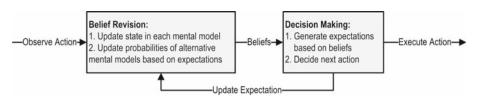
**Fig. 1** Belief revision and decision-making Processes

### 4.1.5 Policy

Policy informs the agent of its best action given the current status of the world—the agent's belief about its own state and other agents' states. By default, all agents use a bounded look-ahead policy to automatically decide their choice of actions during an interaction. The agents project into the future to evaluate the effect of each candidate actions, and choose the one with the highest expected utility (see Sect. 4.2.2 for details).

### 4.1.6 Social relationships

Thespian has a built-in capability of modeling static and dynamic social relationships between agents which in turn can influence the agent's decision-making and belief update. Specifically, Thespian agents maintain a measure of support/affinity for another agent. Support is computed as a running history of their past interactions. An agent increases/decreases its support for another, when the latter selects an action that has a high/low reward, with respect to the preferences of the former.

### 4.2 Belief revision and decision-making processes

Upon observation of an event, each agent updates its beliefs based on the observation and its expectations, and then makes decisions on its next action based on the updated beliefs. The decision-making process also generates new expectations for future events and related self and other agents' states and utilities. These expectations will be used for the following belief update. Figure 1 illustrates this process.

### 4.2.1 Belief revision processes

An agent's beliefs get updated in two ways. One is through dynamics. Upon observation of an event, within each mental model the agent has, the corresponding dynamics are applied and the related state features' values are updated. The other way an agent changes its beliefs is through adjusting the relative probabilities of alternative mental models. Each observation serves as an evidence for the plausibility of alternative mental models, i.e. how consistent the observation is with the predictions from the mental models. Using this information, Bayes' Theorem is applied for updating the probabilities of alternative mental models [8]. During this process, the predictions from alternative mental models are generated by the agent's past decision-making processes. A special case is when the agent only cares about its immediate reward: in this case it performs a zero step lookahead during decision-making and therefore forms no expectations about other's future actions. If the agent needs to adjust the
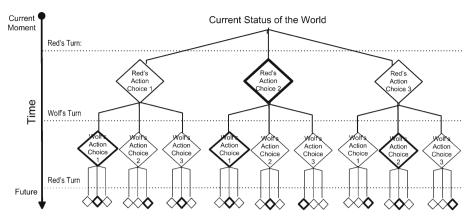
**Fig. 2** Red's lookahead process

probabilities associated with its alternative mental models of others,[2] expectations need to be formed right before the agent's belief revision.

### 4.2.2 Decision-making process

In Thespian, all agents use a bounded lookahead policy. When an agent has multiple mental models of others, by default its next action is decided by lookahead reasoning using the most probable mental models, though the expected states/utilities of all alternative mental models are calculated for the purpose of belief revision.

Each agent has a set of candidate actions to choose from when making decisions. When an agent selects its next action, it projects into the future to evaluate the effect of each option on the state and belief of other entities in the story. The agent considers not just the immediate effect, but also the expected responses of other characters and, in turn, the effects of those responses, and its reaction to those responses and so on. The agent evaluates the overall effect with respect to its goals and then chooses the action that has the highest expected value.

Figure 2 lays out the expected states/utilities being calculated when a character performs a one step lookahead with the belief that other characters will also perform a one step look-ahead. The actions in bold square are the actions with the highest expected utilities among all options from the actor's perspective. This lookahead is taking place in the character Red's belief space before she makes a decision. For each of her action options, she anticipates how the action affects each character's state and utility. For example, when Red decides her next action after being stopped by the wolf on her way to Granny's house, the following reasoning happens in her "mind" using her beliefs about the wolf and herself. For each of her action options, e.g. talking to the wolf or walking away, she anticipates how the action directly affects each character's state and utility. Next, Red considers the long term reward/punish-ment. For example, it may be fun to talk to the wolf for a little while (positive immediate reward), but this will delay Granny from getting the cake (long term punishment). To account for long term effects, she needs to predict other agents' responses to her potential actions. For each of her possible action, Red simulates the wolf's lookahead process. Similarly, for each

---

[2] It is often the case that if an agent only performs zero step lookahead for decision-making, it does not care about the probabilities of alternative mental models because the agent does not use the mental models in its decision-making.

of the wolf's possible action choices, Red calculates the immediate expected states/utilities of both the wolf and herself. Next, Red simulates the wolf anticipating her responses. Since the lookahead process only simulates bounded rationality, this recursive reasoning would stop when the maximum number of steps for forward projection is reached. For example, if the number of lookahead steps is set to be one, the wolf will pick the action with highest utility after simulating one step of Red's response rather than several rounds of interaction. Similarly based on the wolf's potential responses in the next step, Red calculates the utilities of her action options—the sum of reward/punishment of the current and all future steps, and chooses the one with the highest utility. Theoretically, each agent can perform lookahead for large enough number of steps until there is no gain for itself and other agents. For performance reasons, we limit the projection to a finite horizon that we determine to be sufficiently realistic without incurring too much computational overhead. For example, three steps of lookahead is simulated for modeling characters in the Little Red Riding Hood story, and for modeling a negotiation scenario, which is not included in this paper, we simulated up to eight steps of lookahead.

## 5 Computational model of appraisal

In this section we illustrate how appraisal dimensions can be derived by leveraging processes involved and information gathered in an agent's belief revision and decision-making processes. We first describe when appraisal happens and where the related information comes from, and then present algorithms for evaluating each appraisal dimension.

5.1 Overview of the model

We model appraisal as a continuous process, that people constantly reevaluate their situations and form a "appraisal-coping-reappraisal" loop, as described in [32].

During decision-making, the lookahead process calculates the agent's belief about what will happen in the future. This information will be kept in the agent's memory as its expectations. The agent will not only keep expectations generated in its last lookahead process, but also those generated in its previous lookahead process, because the evaluations of some appraisal dimensions, e.g. accountability need to trace back more than one step. Note that these expectations not only contain the agent's expected actions of other agents and self in the future, but also the expected states/utilities of every possible action choices of each of the agents, as this information serves as the explanation for why the agent would make the expected choice.

Upon observing a new event—an action performed by an agent or the human user, each agent updates its beliefs and appraises the situation. The calculation of motivational relevance, motivational congruence, novelty and accountability depends only on the agent's beliefs about other agents' and its own utilities in the current step and the previous steps, and therefore can be derived immediately (see Sect. 5.2 for details). Depending on the extent of reasoning the agent performed in the former steps, the agent may or may not have information immediately available regarding its control of the situation. However, when the agent makes its next decision, control will be automatically evaluated and this evaluation will affect the agent's emotion. In fact, at this time the agent may reevaluate along every appraisal dimension as it obtains more updated information about expected states/utilities. In our current model, upon observing an event the agent derives all appraisal dimensions except control, and evaluates control after it makes decision on its next action. In general the appraisal process

could be based on either the expectations formed in previous steps, or the lookahead process being performed at the current step. The agent may also express both emotional responses in sequence.

Thespian agents have mental models of other agents. These mental models enable them to not only have emotional responses to the environment but also form expectations of other agents' emotions. To simulate another agent's appraisal processes, the observing agent's beliefs about the other agent are used for deriving appraisal dimensions. For instance, agent A can use its beliefs about agent B to evaluate the motivational relevance and novelty of an event to agent B, which may be totally different from B's evaluations of those dimensions. If the observing agent has multiple mental models of other agents, currently it uses the mental models with highest probabilities to simulate other agents' appraisals.

## 5.2 Appraisal dimensions

In this section we provide pseudo code for evaluating the five appraisal dimensions (motivational relevance, motivation congruence or incongruence, accountability, control and novelty) using states/utilities calculated during an agent's belief revision and decision-making processes.

### 5.2.1 Motivational relevance & motivational congruence or incongruence

Motivational relevance evaluates the extent to which an encounter touches upon personal goals. Motivational congruence or incongruence measures the extent to which the encounter thwarts or facilitates personal goals [32].

---
**Algorithm 1 Motivational Relevance & Motivation Congruence**
---

\# $preUtility$: utility before the event happens
\# $curUtility$: utility after the event happens

$Motivational\ Relevance$ = abs $\frac{curUtility - preUtility}{preUtility}$

$Motivational\ Congruence$ = $\frac{curUtility - preUtility}{abs(preUtility)}$

---

We model these appraisal dimensions as a product of the agent's utility calculations which are integral to the agent's decision-theoretic reasoning. We use the ratio of the relative utility change and the direction of the utility change to model these two appraisal dimensions. The rationale behind this is that the same amount of utility change will result in different subjective experiences depending on the agent's current utility. For instance, if eating a person increases the wolf's utility by 10, it will be 10 times more relevant and motivationally congruent when the wolf's original utility is 1 (very hungry) than when the wolf's original utility is 10 (less hungry).

Algorithm 1 gives the equations for evaluating motivational relevance and motivational congruence or incongruence. $preUtility$ denotes the agent's expected utility before the other agent takes an action. For agents which perform at least one step of lookahead, this is the expected utility in the future, which is evaluated when the agent made its last decision. For example, in the lookahead reasoning shown in Fig. 2, Red's expected utility before the wolf does any actions is the sum of her utilities over the following sequence of actions: Red's action choice 2 → the wolf's action choice 1 → Red's action choice 2. $curUtility$ denotes

the agent's expected utility after the other agent takes the action. This value is also evaluated when the agent made its last decision. For example, if the wolf does its action choice 1, then *curUtility* is the same as *preUtility* for Red. If the wolf instead chooses action choice 3, then *curUtility* is the sum of Red's utilities over this sequence of actions: Red's action choice 2 → the wolf's action choice 2 → Red's action choice 1. If an agent performs zero step of lookahead, i.e. it only cares about its immediate reward, the value of *curUtility* is not calculated by the agent's previous decision-making process. Rather, it is evaluated when the agent updates its beliefs because the value is associated with the agent's updated state.

The sign of *MotivationalCongruence* indicates whether the event is motivationally congruent or incongruent to the agent. When the value is negative, the event is motivationally incongruent to the extent of *Motivational Relevance*, and otherwise the event is motivationally congruent to the agent.

Thespian agent can have goals regarding other agents' utilities. If taking an action helps a Thespian agent's self-centered goals but hurts a friend that the agent also has a goal to care about, then the agent's happiness (the action's motivational congruence and relevance) is muted accordingly because its overall utility is diminished. For example, Red will feel less satisfied if she eats the cake which she is bringing to Granny than eating other cakes because she also wants Granny to have the cake.

### 5.2.2 Accountability

Accountability characterizes which person deserves credit or blame for a given event [32]. Various theories have been proposed for assigning blame/credit, e.g. [25,35]. The reasoning usually considers factors such as who directly causes the event, does the person foresee the result, does the person intend to do so or is it coerced, etc.
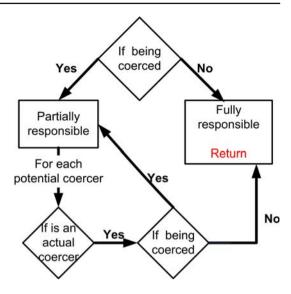
Just as the appraisal of motivational relevance and motivation congruence can be performed as part of the existing Thespian/PsychSim decision-making and belief update processes, we argue here that accountability can be treated as an improvement to Thespian/PsychSim's existing approach to model support/affinity relationships between agents.

In Fig. 3 we use a diagram to illustrate our algorithm for determining accountability. This algorithm first looks at the agent which directly causes the harm/benefit, and judges if the agent is the one who should be fully responsible. The function *If_Coerced()* is called to determine if the agent was coerced to perform the action. If the agent was not coerced, it should be fully responsible and the reasoning stops there. Otherwise, each of the agents that coerced the direct actor will be judged on whether it was coerced by somebody else. If the answer is yes for a coercer, in turn each coercers of that agent will be checked to see if they did the actions voluntarily. The algorithm will trace limited steps back in the history to find out all the responsible agents. While evaluating accountability, we assume that the agent expects others to foresee the effects of their actions. This assumption is correct most of the time because normally a person would expect others to project into the future the same number of steps as what the person will do themselves when making a decision.

Algorithm 2 contains pseudo code for determining if an agent was coerced, and Algorithm 3 finds the coercers if the agent is indeed coerced. We use a qualitative rather than quantitative model to decide coercion. If all action options, other than the action chosen by the agent lead to a drop in its utility (i.e. the agent will be punished if it chooses any other actions), then the agent is coerced by somebody. However, if all of the agent's action options result in utility drops, the agent is regarded as not being coerced. The rationale behind this is that since the agent is going to be punished regardless of what it does, it has the freedom to pick actions which will not hurt the other agent's utility. In this algorithm, *preUtility* is

**Fig. 3** Accountability



---

**Algorithm 2 If_Coerced(*actor*, *pact*)**

# *actor*: the agent being studied
# *pact*: the action performed by *actor*
# *preUtility*: *actor*'s utility before doing *pact*

**for** *action* in *actor.actionOptions*() **do**
  **if** *action* ≠ *pact* **then**
    #if there exists another action which does not hurt *actor*'s own utility
    **if** utility(*action*) ≥ *preUtility* **then**
      Return *F*
**if** utility(*action*) < *preUtility* **then**
  Return *F*
Return *T*

---

defined similarly as in the algorithms for evaluating motivational relevance and motivation congruence. The only difference is that here *preUtility* is the observer (the agent which performs appraisal)'s beliefs about the *actor*'s expected utility. Similarly, utility(*action*) denotes the observer's belief about the *actor*'s utility of alternative option.

To decide who coerced an agent, we treat each agent that acted between the coerced agent's current and last actions as a potential coercer. For each potential coercer, if the coerced agent would not have been coerced in case the potential coercer had made a different choice, then the potential coercer is judged as actually being a coercer. This process is illustrated in Algorithm 3.

*5.2.3 Control*

The appraisal of control evaluates the extent to which an event or its outcome can be influenced or controlled by people [24]. It captures not only the individual's own ability to control the situation but also the potential for seeking instrumental social support from other people. Different from the evaluations of motivational relevance, motivational congruence and accountability in which the most probable mental models of other agents are

---

**Algorithm 3 Is_Coercer_For(*agent*, *actor*, *agent_pact*, *actor_pact*)**

---

```
# check if agent coerced actor
# agent_pact: the action performed by agent
# actor_pact: the action performed by actor

for action in agent.actionOptions() do
  if action ≠ agent_pact then
    Simulate action agent_pact
    if If_Coerced(actor,actor_pact)== F then
      Return T
Return F
```

---

used for reasoning, here we factor in the probabilities of the mental models because the degree of control is affected by the estimation of how likely certain events will happen in the future.

---

**Algorithm 4 Control(*preUtility*)**

---

```
# preUtility: utility before the event happens

control ← 0
for m1 in mental_models_about_agent1 do
  for m2 in mental_models_about_agent2 do
    for m3 in mental_models_about_self do
      #project limited steps into the future using this set of mental models
      lookahead(m1,m2,m3)
      #curUtility: utility after the lookahead process
      if curUtility ≥ preUtility then
        control ← control + p(m1) * p(m2) * p(m3)
Return control
```

---

Algorithm 4 gives the pseudo code for evaluating control. This algorithm first simulates future steps of the interaction using each possible combination of mental models of self and others, and checks whether the utility drop will be recovered. The algorithm then considers the probabilities of the mental models to be correct, and therefore the event, if being predicted, will actually happen in the future. For example, assume Granny has two mental models of the wolf. In the first mental model, the wolf will always die after being shot by the hunter. In the second mental model, the wolf will never die even after being shot. Granny believes that there is a 60% possibility that the first mental model is true. Next assume Granny has two mental models regarding the hunter. One mental model indicates that the hunter is close by and this mental model has a 50% chance to be true. The other mental model indicates that the hunter is far away. After Granny is eaten by the wolf, the only event that can help her is that the wolf is killed by the hunter. Therefore, she would evaluate her control as: $60\% \times 50\% = 30\%$.

Algorithm 4 contains pseudo code for the three-agent interaction case. It is straightforward to configure the algorithm to be applied when more or less agents are in the interaction. In this algorithm, *preUtility* is defined the same way as in the algorithms for evaluating motivational relevance and motivation congruence. *curUtility* denotes the agent's utility associated with its state after the lookahead projection.

*5.2.4 Novelty*

In this work, we adapt Leventhal and Scherer's definition of "novelty at the conceptual level"—novelty describes whether the event is expected from the agent's past beliefs[3] [12,24].

In our model, novelty appraisal is treated as a byproduct of an agent's belief maintenance. Specifically, in a multi-agent context the novelty of an agent's behavior is viewed as the opposite of the agent's motivational consistency, i.e. the more consistent the event is with the agent's motivations, the less novel. Of course, this evaluation is performed from the observing agent's perspective and using the observing agent's beliefs, and there can be discrepancies between what the observing agent feels and what the agent who did the action feels. Computationally, we define novelty as $1 - consistency$, where $consistency$ is calculated using one of the methods proposed by Ito et al. [8] for deciding motivational consistencies of actions for POMDP-based agents.

$$consistency(a_j) = \frac{e^{rank(a_j)}}{\sum_j e^{rank(a_j)}} \tag{1}$$

*Consistency* is calculated based on the most probable mental model of the actor. The algorithm first ranks the utilities of the actor's alternative actions in reversed order ($rank\left(a_j\right)$). The higher an action's utility ranks compared to other alternatives, the lower consistency it has with the observing agent's expectation about the actor, and hence the higher novelty if the action happens. For example, if from Red's perspective the wolf did an action which has the second highest utility among the wolf's five alternative actions, the amount of novelty Red will feel if seeing that action is calculated as $1 - \frac{e^3}{\sum_{j=0-4} e^j} = 0.37$.

# 6 Sample results

All the previous examples of our new appraisal model are derived from a Thespian implementation of the Little Red Riding Hood story. In this section we provide two additional scenarios to illustrate the usage of our computational model of appraisal in modeling social interactions. In particular, in Scenario 1 we demonstrate the tight relationship between emotion and cognitive decision-making by showing how appraisal is affected by the depth of reasoning in decision-making. In Scenario 2 we provide a complex situation for accountability reasoning and show that the result of our model is consistent with another validated computational model of social attribution.

6.1 Scenario 1: small talk

To reveal the tight relationship between cognitive processes and emotion in our model, we implemented an abstract domain of two persons (A and B) taking turns talking to each other. Both of them have these goals: to be talkative and to obey social norms. In fact, just the norm following behavior itself is an incentive to them—they will be rewarded whenever they do an action that is consistent with social norms. Table 1 contains the two persons' appraisals

---

[3] Leventhal and Scherer have also defined novelty at sensory-motor level and schematic level. We did not model them because they are mainly related to people's low level perceptual processes rather than cognitive processes.

**Table 1** Small talk between two persons

| Step | Action | Perspective | Lookahead steps | Motivational relevance |
|---|---|---|---|---|
| 1 | A greets B | B | 1 | 0 |
|   |   | B | 2 | 100 |
| 2 | B greets A | A | 1 | 0 |
|   |   | A | 2 | 0.99 |
| 3 | A asks B a question | B | 1 | 0 |
|   |   | B | 2 | 0.99 |
| 4 | B answers the question | A | 1 | 0 |
|   |   | A | 2 | 0.49 |

of motivational relevance regarding each other's actions. We did not include results of other appraisal dimensions as they are less interesting in this scenario.

In Thespian, we explicitly model the depth of reasoning in agents as the number of steps they project into the future. In this example we provide a comparison of appraisal results when the person's previous reasoning process takes a different number of steps. It can be observed in Table 1 that different depths of reasoning lead to different appraisals. A person appraises another person's initiatives as irrelevant when performing shallow reasoning (lookahead steps = 1). In this case, even though the person has predicted the other person's action, because the action does not bring him/her immediate reward, the person can not see the relevance of the action. Once the person reasons one step further, he/she finds out that by opening up a topic the other person actually provides him/her a chance to engage in further conversation and perform a norm following action, the person will then appraise the other person's action as relevant.

### 6.2 Scenario 2: firing-squad

We implemented the Firing-squad scenario from [14] to illustrate accountability reasoning in which agents are coerced and have only partial responsibility. The scenario goes like this:

*In a firing-squad, the commander orders the marksmen to shoot a prisoner. The marksmen refuse the order. The commander insists that the marksmen shoot. They shoot the prisoner and he dies.*

We modeled the commander as an agent with an explicit goal of killing the prisoner, and the marksmen as having no goals related to the prisoner but will be punished if they do not follow the commander's order. Using our appraisal model, from the prisoner's perspective, the marksmen hold responsibility for his/her death because they are the persons who directly perform the action. Further, the prisoner simulates the decision-making process of the marksmen and finds out that the marksmen are coerced because their utilities will be hurt if they perform any action other than shooting. The commander acts right before the marksmen in the scenario and therefore is identified as a potential coercer for the marksmen. Using Algorithm 3, the prisoner can see that if the commander chose a different action, the marksmen are not coerced to shoot. Assuming the prisoner does not find a coercer for the commander, he/she will now believe that the commander holds full responsibility for his/her death. This prediction is consistent with the prediction from Mao's model of social attribution and the data collected from human subjects to validate that model [14].

## 7 Discussion

In the Thespian system, comparison among expected utilities plays the central role in decision-making and mental model update. Comparison of expected utilities also plays a central role for deriving appraisal dimensions. Our algorithms for deriving appraisal dimensions demonstrate that no additional calculation of utilities or states other than what has already been performed in the Thespian agent's existing decision-making and belief revision processes is required for appraisal.

Compared to other computational models of appraisal, the main advantage of this model is that the agents explicitly model other agents' goals, states and beliefs (theory of mind). Modeling theory of mind makes this model particularly suitable for simulating emotions in social interactions in two ways. First, appraisals are strongly embedded in the social context. For example, novelty is not simply evaluated as whether the physical stimulus is unexpected, but whether the other agents behave as expected. Second, appraisals that are explicitly relevant to social interaction and derivation of social emotion, such as accountability, have to leverage theory of mind.

Further, the fact that Thespian agents have a theory of mind capability enables them to simulate others' emotions. This ability allows us to simulate an agent's potential mis-expectations about other agents' emotional states. For example, if Granny believes that the hunter can always kill the wolf successfully and the hunter believes that he can only kill the wolf successfully 60% of the time, Granny's control when being eaten by the wolf will be evaluated differently from Granny's and the hunter's perspectives.

The appraisal model can not only simulate ego-centric emotions, but also can simulate emotions that take social relationship into account. Thespian agent can have goals regarding other agents' utilities and emotions (emotion can be modeled as a feature of an agent's state). Therefore, an agent's emotion can be related to other agents' utility changes and emotions. For example, we can simulate an agent having goals of facilitating another agent's goals, or even more specifically having goals of making the other agent feel happy. This agent will act deliberately to help the other agent, and "feel" bad if it hurts the other agent's utility or emotional state.

Finally, we explicitly model the depth of reasoning in agents as the number of steps they project into the future. As shown in Scenario 1, different depths of reasoning lead to different appraisals. Though we have only demonstrated this effect using one appraisal dimension—motivational relevance, this effect is general. Different steps of projection lead to different predictions of future events; and a character's prediction about the future affects the character's reasoning about whether an event is novel, whether the effect of the event is changeable and who caused the event.

## 8 Future work

Our future work involves improving our model of appraisal and performing further validation/examination of our hypothesis on the relationship between appraisal and cognition.

In particular, the future work is planned in three directions. First, we want to factor in the probabilities associated with alternative mental models when evaluating utilities. Currently, utilities are calculated based on the agent's most probable mental models about others. For the stories we have built so far, most of the times the less probable mental models only have very small chance to be true, and therefore using only the most probable mental models for evaluating utility is sufficient. However, if the probabilities of alternative mental models are

similar to each other, the character will not be confident on using a single mental model to predict other characters' behaviors, and therefore multiple mental models need to be considered when evaluating utilities. Further, we want to study how changing the probabilities of alternative mental models affect an agent's utility and therefore its appraisal. For example, is there a cost for updating the probabilities of alternative mental models? Is there a delay before the new mental model affects the person's emotions? Secondly, we want to add additional emotion-cognition interaction to Thespian agents by modeling how emotion affects the agent's decision-making and belief update processes. It has been argued that affect plays an important role in how people process information and make decisions [3]. For example, positive affect often leads to more heuristic processing and negative affect signals a problematic situation and therefore leads people to perform more systematic processing. Finally, we want to enrich the current model with a more complete set of appraisal dimensions. In particular, we are interested in modeling emotion-focused coping potential, which is a useful appraisal dimension for social agents. This extension will not only make our model of appraisal more complete, but also provide further validation/examination of whether appraisal dimensions can always be derived from information generated in the agent's existing cognitive processes.

## 9 Conclusion

In this paper, we provide a computational model of appraisal for POMDP-based agents, implemented in the Thespian framework for interactive narratives. The focus is on five key appraisal dimensions for virtual agents: motivational relevance, motivational congruence, accountability, control and novelty. The approach demonstrates that appraisal is an integral part of an agent's cognitive processes.

We demonstrate the model on three different scenarios. Compared to other computational models of appraisal, this model is particularly suitable for simulating emotions in social interactions because it models theory of mind.

## References

1. Aylett, R., Dias, J., & Paiva, A. (2006). An affectively-driven planner for synthetic characters. In *ICAPS*.
2. Bui, T. D., Heylen, D., Poel, M., & Nijholt, A. (2002). Parlee: An adaptive plan based event appraisal model of emotions. Parlevink internal report.
3. Clore, G. L., & Storbeck, J. (2006). Affect as information about liking, efficacy, and importance. In J. P. Forgas (Ed.), *Hearts and minds: Affective influences on social cognition and behaviour*. New York: Psychology Press.
4. El Nasr, M. S., Yen, J., & Ioerger, T. (2000). Flame: Fuzzy logic adaptive model of emotions. *Autonomous Agents and Multi-Agent Systems, 3*(3), 219–257.
5. Elliott, C. (1992). *The affective reasoner: A process model of emotions in a multi-agent system.* PhD thesis, Northwestern University Institute for the Learning Sciences.
6. Frijda, N. (1987). *The emotions*. Cambridge University Press.
7. Gratch. J., & Marsella, S. (2004). A domain-independent framework for modeling emotion. *Cognitive Systems Research, 5*(4), 269–306.
8. Ito, J. Y., Pynadath, D. V., & Marsella, S. C. (2007). A decision-theoretic approach to evaluating posterior probabilities of mental models. In *AAAI-07 workshop on plan, activity, and intent recognition*.
9. Lazarus, R. (1984). On the primacy of cognition. *American Psychologist, 39*, 124–129.

10. Lazarus, R. S. (1991). *Emotion & adaptation*. New York: Oxford University Press.
11. Lazzaro, N. (2004). Why we play games: Four keys to more emotion in player experiences. In *Game developers conference*.
12. Leventhal, H., & Scherer, K. R. (1987). The relationship of emotion and cognition: A functional approach to a semantic controversy. *Cognition and Emotion, 1*, 3–28.
13. Lisetti, C., & Nasoz, F. (2005, July). Affective intelligent car interfaces with emotion recognition. In *11th international conference on human computer interaction*, Las Vegas, USA.
14. Mao, W., & Gratch, J. (2005). Social causality and responsibility: Modeling and evaluation. In *IVA*, Kos, Greece.
15. Marsella, S. C. & Gratch, J. (2009). Ema: A model of emotional dynamics. *Cognitive Systems Research, 10*(1), 70–90.
16. Marsella, S. C., Pynadath, D. V., & Read, S. J. (2004). PsychSim: Agent-based modeling of social interactions and influence. In *Proceedings of the international conference on cognitive modeling* (pp. 243–248).
17. Moffat, D., & Frijda, N. (1995). Where there's a will there's an agent. In *ECAI-94 workshop on agent theories, architectures, and languages*, Amsterdam, The Netherlands.
18. Ortony, A., Clore, G. L., & Collins, A. (1998). *The cognitive structure of emotions*. Cambridge: Cambridge University Press.
19. Picard, R. W. (1997). *Affective computing*. Cambridge: MIT Press.
20. Pynadath, D. V., & Marsella, S. C. (2005). Psychsim: Modeling theory of mind with decision-theoretic agents. In *IJCAI* (pp. 1181–1186).
21. Reilly, W. S., & Bates, J. (1992). Building emotional agents. Technical Report CMU-CS-92-143, Carnegie Mellon University.
22. Roseman, I. (1984). Cognitive determinants of emotion: A structural theory. *Review of Personality and Social Psychology, 2*, 11–36.
23. Roseman, I. J., & Smith, C. A. (2001). Appraisal theory: Overview, assumptions, varieties, controversies. In K. Scherer, A. Schorr, & T. Johnstone (Eds.), *Appraisal processes in emotion: Theory, methods*. Oxford: Oxford University Press.
24. Scherer, K. (2001). Appraisal considered as a process of multilevel sequential checking. In K. Scherer, A. Schorr, & T. Johnstone (Eds.), *Appraisal processes in emotion: Theory, methods*. Oxford: Oxford University Press.
25. Shaver, K. G. (1985). *The Attribution theory of blame: Causality, responsibility and blameworthiness*. Springer.
26. Si, M., Marsella, S. C., & Pynadath, D. V. (2005). Thespian: An architecture for interactive pedagogical drama. In *AIED*.
27. Si, M., Marsella, S. C., & Pynadath, D. V. (2005). Thespian: Using multi-agent fitting to craft interactive drama. In *AAMAS* (pp. 21–28).
28. Si, M., Marsella, S. C., & Pynadath, D. V. (2006). Thespian: Modeling socially normative behavior in a decision-theoretic framework. In *IVA*.
29. Si, M., Marsella, S. C., & Pynadath, D. V. (2007). Proactive authoring for interactive drama: An author's assistant. In *IVA*, Paris, France.
30. Smallwood, R. D., & Sondik, E. J. (1973). The optimal control of partially observable Markov processes over a finite horizon. *Operations Research, 21*, 1071–1088.
31. Smith, C. A., & Ellsworth, P. C. (1987). Patterns of appraisal and emotion related to taking an exam. *Personality and Social Psychology, 52*, 475–488.
32. Smith, C. A., & Lazarus, R. S. (1990). Emotion and adaptation. In L. A. Pervin (Ed.), *Handbook of personality: Theory and research*. New York: Guilford.
33. Swartout, W., Hill, R., Gratch, J., Johnson, W., Kyriakakis, C., LaBore, C., Lindheim, R., Marsella, S.C., Miraglia, D., Moore, B., Morie, J., Rickel, J., Thiúbaux, M., Tuch, L., Whitney, R., & Douglas, J. (2001). Toward the holodeck: Integrating graphics, sound, character and story. In *Agents* (pp. 409–416).
34. Velasquez, J. (1997). Modeling emotions and other motivations in synthetic agents. In *AAAI*.
35. Weiner, B. (1995). *The judgment of responsibility: A foundation for a theory of social conduct*. Guilford.
36. Whiten, A. (Ed.). (1991). *Natural theories of mind*. Basil Blackwell, Oxford.