

## Dynamic mathematical models of theory of mind for socially assistive robots

Patricio, Maria L.M.; Jamshidnejad, Anahita

**DOI**

[10.1109/ACCESS.2023.3316603](https://doi.org/10.1109/ACCESS.2023.3316603)

**Publication date**

2023

**Document Version**

Final published version

**Published in**

IEEE Access

**Citation (APA)**

Patricio, M. L. M., & Jamshidnejad, A. (2023). Dynamic mathematical models of theory of mind for socially assistive robots. *IEEE Access*, 11, 103956-103975. <https://doi.org/10.1109/ACCESS.2023.3316603>

**Important note**

To cite this publication, please use the final published version (if applicable).  
Please check the document version above.

**Copyright**

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

**Takedown policy**

Please contact us and provide details if you believe this document breaches copyrights.  
We will remove access to the work immediately and investigate your claim.

Received 14 August 2023, accepted 8 September 2023, date of publication 18 September 2023, date of current version 27 September 2023.

Digital Object Identifier 10.1109/ACCESS.2023.3316603

## RESEARCH ARTICLE

# Dynamic Mathematical Models of Theory of Mind for Socially Assistive Robots

MARIA L. M. PATRÍCIO<sup>ID</sup> AND ANAHITA JAMSHIDNEJAD<sup>ID</sup>

Department of Control and Operations, Delft University of Technology, 2629 Delft, The Netherlands

Corresponding authors: Maria L. M. Patrício (M.L.MoraoPatricio@tudelft.nl) and Anahita Jamshidnejad (a.jamshidnejad@tudelft.nl)

This work was supported in part by the Technische Universiteit (TU) Delft AI Labs and Talent Program.

This work involved human subjects or animals in its research. Approval of all ethical and experimental procedures and protocols was granted by the Human Research Ethics Committee of TU Delft under Approval No. 2398.

**ABSTRACT** Interactive machines should establish and maintain meaningful social interactions with humans. Thus, they need to understand and predict the mental states and actions of humans. Based on Theory of Mind (ToM), in order to understand and interact with each other, humans develop cognitive models of one another. Our main goal is to provide a mathematical framework based on ToM to improve the understanding of interactive machines regarding the perception, cognition, and decision-making of humans. Most state-of-the-art models of behavioral theories based on machine learning are focused on input-output black-box representations. Thus, they lack transparency and generalizability, and exhaustive training procedures are needed to personalize them for various humans. Moreover, these models lack dynamics, i.e., they do not mathematically describe the evolution of the mental states and actions of humans in time. Following a systems-and-control-theoretic point-of-view, we represent for the first time the perception, cognition, and decision-making of humans via a dynamic, mathematical framework by introducing a novel formalization and an extension to Fuzzy Cognitive Maps (FCMs). The resulting models are given in a general state-space representation, which can be used by interactive machines within known model-based state estimation and control methods. In a case study, the resulting models were identified and validated for 21 participants, in scenarios where predicting the intentions and behavior of the participants required understanding the dynamics of their mental procedures. The results of these experiments show that our model is capable of incorporating the dynamics to estimate the intentions and predict the behavior of the participants, with an accuracy of, respectively, 81.55% and 66.06%. Moreover, we compared our model with a state-of-the-art formalization of human cognition, which was made dynamic using our introduced FCM framework. Our model, which in addition to the elements of the state-of-the-art model included emotions, personality traits, and biases (thus providing a more transparent insight about the mental procedures of the participants) showed 6.25% and 2.45% more accuracy in, respectively, estimating the intentions and predicting the behavior of the participants.

**INDEX TERMS** Dynamic mathematical models, long-term interactions of rational agents, state-space models of cognition, theory of mind.

## I. INTRODUCTION

Socially assistive robots (SARs) assist humans mentally and physically by socially interacting with them. SARs have

The associate editor coordinating the review of this manuscript and approving it for publication was Gustavo Olague<sup>ID</sup>.

proven successful in boosting the outcomes of therapeutic and educational assistance for humans [1], [2], [3], [4], [5], [6]. Most applications of SARs in literature involve short-term interactions (e.g., a few interactive sessions lasting no more than a month) with humans [7], [8] that exclude in-depth analysis of the perception, cognition, and decision-making

of the individuals [8]. Since intelligent machines, including SARs, are becoming more prominent, it is essential for them to maintain long-term, meaningful and engaging interactions with humans [8], whereas state-of-the-art SARs face serious challenges regarding this [3], [4], [5], [7]. Particularly, rudimentary social skills displayed by SARs in interaction with humans negatively impact the human engagement, and thus the effectiveness of SARs [4], [5], [7], [9]. Although non-verbal cues, e.g., joint attention (i.e., drawing the attention of others to an object or person by looking or pointing at it) [4], eye contact [5], and facial expressions [3], have been implemented for SARs, these robots often fail to recognize the best time to display these cues in social interactions, due to a lack of (deep) understanding of the human cognition [7], [10]. Moreover, personalizing the behavior of SARs to every human is crucial to maintain meaningful interactions [2], [3], [4], [10]. However, without a white-box approach (e.g., transparent mathematical models), personalization of SARs has been task-specific [2], [6], case-based, and ungeneralizable [8]. In fact, personalization for SARs has been limited to learning interactive behaviors that improve the performance of humans in the given tasks [2], [3], [4], [10]. Such task-specific, black-box methods used to control SARs result in interactions that are perceived as less natural and engaging for humans, especially in the long term.

The key to successful social interactions by humans, based on ToM [11], is to build cognitive models of each other and to interact based on these models. In order to interact as humanly as possible, SARs should exhibit similar understanding of rational agents [7], [8], [9], [10], [12]. Thus, we focus on developing dynamic mathematical models for perception, cognition, and decision-making of humans that can be used by intelligent machines in human-machine interactions to estimate the mental states of humans and to predict their behavior. In order to provide generalizable, white-box models with proper mathematical representation and axiomatization that can effectively be used in model-based state estimation and control methods by SARs, we provide the following main contributions:

- We follow systems-and-control-theoretic methods in order to model the perception, cognition, and decision-making of humans via dynamic equation-based models for the first time. We identify and represent the corresponding (controllable and uncontrollable) input, output, auxiliary, and state variables, as well as the identification parameters. To provide accuracy and transparency, we introduce additional state variables (e.g., emotions, intentions) and auxiliary variables (e.g., bias, perceived knowledge) into the perception, cognition, and decision-making models based on the existing behavioral theories from cognitive science literature.
- We propose a novel framework based on an extended version of FCMs. This framework transforms static representations of ToM into dynamic models with a general state-space representation. A main advantage of

representing perception, cognition, and decision-making of humans via equation-based state-space models is that they can directly be embedded within established model-based state estimation and control methods to steer intelligent machines.

- The resulting models are trained and personalized for 21 human participants, and are validated via experiments that simulate scenarios of emergency evacuation in 2-dimensional environments. The participants complete each scenario by controlling a virtual agent in order to select its trajectory and sequence of tasks. The developed models are used to estimate the intentions and actions of each participant. These estimated intentions and actions are then compared to the intentions and actions selected by the participants to assess the accuracy of the models.

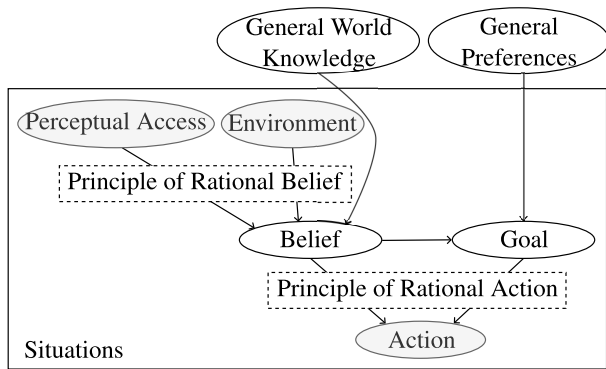
Including the auxiliary variable ‘biases’ and the identifiable parameters, and defining new concepts of *rational intention selection* and *rational action selection*, which are agent-specific, instead of using a universal principle of rational action for all agents, provide our models with the unique potential to be personalized per human.

The rest of the paper is structured as it follows. Section II gives an overview of previous related work and identifies the open challenges. Section III presents the proposed methodologies, including the main motivations and mathematical approaches. In Section IV, the proposed models are implemented and assessed via real-life experiments with 21 volunteer human participants. Finally, Section V concludes the paper and presents topics for future research.

## II. RELATED WORK AND OPEN CHALLENGES

Understanding the (individual) cognitive procedures of humans in their social interactions is important for both cognitive psychology and cognitive computing fields (see [4], [5], [13], [14], [15], [16]).

ToM was first proposed in the neuroscience domain [11], and since then has been used to develop computational frameworks for human cognition [17], [18], [19], [20], [21], [22]. ToM relies on the principle of rationality [23], [24], which implies that humans take actions that, according to their beliefs and goals, maximize their desired outcomes, and/or minimize their losses. This principle is the cornerstone of most ToM-based computational frameworks. While ToM involves three main mental procedures of humans, i.e., perception, cognition, and decision-making, Belief-Desire-Intention (BDI) [25] is a framework that particularly models the decision-making [26], [27], [28]. Based on BDI, the actions of rational agents are a consequence of their beliefs, desires, and intentions (without imposing specific conditions, e.g., rationality) [25]. ToM and BDI can complement one another [27], [29]: In fact, the concept of intentions introduced in BDI can be included in a dynamic, mathematical framework of ToM, in order to account for the interlinks within a sequence of actions of the agent, for achieving a particular intention [26].



**FIGURE 1.** Network representation of human cognition from [30], explaining the actions taken by humans based on the principles of rational action and rational belief.

Baker [30] has proposed a Bayesian ToM model based on the network representation in Figure 1. Using Partially Observable Markov Decision Processes (POMDPs) [31] and Bayesian inference [32], this model makes forward and inverse inferences about, respectively, the actions and the beliefs and goals of rational agents that follow the principle of rationality. A simplified version of the model, excluding general world knowledge and general preferences, was implemented in [33], where virtual rational agents were moving in two-dimensional spaces. The predictions that were made by the model and by human participants about the beliefs and goals of the virtual rational agent were compared, where the results showed close similarities, and acknowledged that goals and beliefs need to be inferred simultaneously. The model and experimental scenarios in [33] did not account for the differences of humans and represented them via idealized rational agents that do not deal with biases, which is not true in real-life scenarios [34], [35]. Additionally, the stochastic nature of the model formulation makes its use limited for model-based predictive methods, since decision trees, which can exponentially grow with the size of the prediction horizon, need to be generated, and using these trees in optimization-based decision-making frameworks is not suited for real-time implementations. Finally, similarly to most computational frameworks of ToM, the model in [33] does not include inferences about the emotions of humans. In fact, based on a recently published review paper [36], the number of literature that model the understanding of emotional states from the perspective of an observer is very limited.

To summarize, the majority of the research on computational frameworks for ToM are focused on representing human behavior via POMDPs [17], [19], [21] or approximating it with Neural Networks (NNs) [18], [19], [20], [21]. The research aiming at using computational models of ToM to perform inverse inferences is even more scarce [37], [38]. Furthermore, most state-of-the-art models are static representations that do not include dynamics [20], [21], [22], i.e., do not mathematically describe the evolution of the mental states of humans in time. The articles that do present

a dynamic model are scarce [19], [39], and none of them proposes a general framework that covers the main aspects of ToM for a general context. For example, [39] dynamically models the emotions but does not focus on any other mental states, or on perception or decision-making, and [19] presents a dynamic model that is tailored to one specific context. To the best of our knowledge, there has been no follow-up research on potential extensions for the model given in [30] in order to address its challenges and to exploit its potential for representing a ToM framework. The core focus of this paper is on how to develop an equation-based model of and mathematically represent the mentioned theories (ToM and BDI), which have before been represented as verbal models, agent-based models, or computer simulations [40]. Thus, the model shown in Figure 1 is our main inspiration for proposing the first dynamic, mathematical framework for ToM. In order to complete this discussion, next we briefly refer to two other behavioral theories that have been proposed to understand human's behavior. Although we do not use these theories in this paper, since the approaches that are proposed, especially the dynamic framework, are not tailored to a specific behavioral theory, they may be adopted for the following theories as well.

The Theory of Planned Behavior (TPB) describes the behavior of rational agents as a consequence of their intentions and perceived behavioral control [41], i.e., the perception of the agent regarding the feasibility of the intended action [42]. The intention is affected by the perceived behavioral control, the attitude (the value of the action according to the agent), and the subjective norm (the social value of the action). Since this theory is mainly used in medical field applications, we have not focused on it in this paper.

Common Model of Mind (CMM) [43] is a framework that has been proposed aiming at integrating the biological, cognitive, and motor procedures for modeling the behavior of rational agents. Due to its novelty when compared with other frameworks that describe human behavior, implementations of CMM are scarce [44], and do not provide significant contributions to the cognitive domain.

### III. METHODOLOGIES

Next, we explain our methodologies for developing a dynamic mathematical framework for ToM, based on an extended version of BDI that includes emotions and biases. In the rest of the paper, a rational<sup>1</sup> agent that makes inferences about the mental states and actions of another rational agent is called an *observer agent*. The other agent is called an *observed agent*. First, we briefly motivate our research.

#### A. MAIN MOTIVATIONS

Mathematical models are used in systems and control theory in order to represent the governing dynamics of systems as

<sup>1</sup>The approaches that are proposed are not limited to the principle of rationality. In fact, by proposing a generalized, white-box model of ToM, it will be possible to model different policies, objectives, and rationality levels.

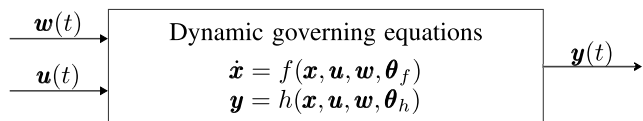


FIGURE 2. Mathematical modeling of dynamical systems.

a set of differential (or difference) equations [45]. Figure 2 shows a system that, in general, receives controlled input vector  $\mathbf{u}$  and uncontrolled input (external disturbances)  $\mathbf{w}$  and generates output vector  $\mathbf{y}$ . [46]. A mathematical dynamic model in continuous time may be represented by the following state space representation<sup>2</sup>

$$\dot{\mathbf{x}}(t) = f(\mathbf{x}(t), \mathbf{u}(t), \mathbf{w}(t), \boldsymbol{\theta}_f) \quad (1)$$

$$\mathbf{y}(t) = h(\mathbf{x}(t), \mathbf{u}(t), \mathbf{w}(t), \boldsymbol{\theta}_h) \quad (2)$$

where  $\mathbf{x}(t)$  is a vector of dimension  $n_x$  that includes all the state variables of the system,  $\mathbf{u}(t)$  is a vector of dimension  $n_u$  that includes all the inputs to the system,  $\mathbf{w}(t)$  is a vector of dimension  $n_w$  that includes all the uncontrollable inputs to the system, and  $\mathbf{y}(t)$  is a vector of dimension  $n_y$  that includes all the outputs of the system. Moreover,  $f(\cdot)$  represents the dynamic model,  $h(\cdot)$  is a function that generates the output. Finally,  $\boldsymbol{\theta}_f$  and  $\boldsymbol{\theta}_h$  are defined as vectors that include all the identification parameters of the corresponding mathematical model. In fact, when an observed system is being analyzed, in addition to determining  $f(\cdot)$  and  $h(\cdot)$ , these parameter vectors should be identified. For instance, for linear time-invariant systems without disturbances, where (1) and (2) are transformed into  $\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t)$  and  $\mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t)$ , with  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$ , and  $\mathbf{D}$  matrices of dimensions, respectively,  $n_x \times n_x$ ,  $n_x \times n_u$ ,  $n_y \times n_x$ , and  $n_y \times n_u$ , [45], [47], the vector  $\boldsymbol{\theta}_f$  includes all elements of  $\mathbf{A}$  and  $\mathbf{B}$ , and similarly the vector  $\boldsymbol{\theta}_h$  includes all elements of  $\mathbf{C}$  and  $\mathbf{D}$ . Assuming that the system is time-invariant,  $\boldsymbol{\theta}_f$  and  $\boldsymbol{\theta}_h$  do not depend on time, and thus the argument time  $t$  is not shown for these two vectors [45]. In most real-life applications, it is common to discretize the model given by (1) and (2) in time, with the corresponding discrete-time model given by

$$\mathbf{x}(k + 1) = f_d(\mathbf{x}(k), \mathbf{u}(k), \mathbf{w}(k), \boldsymbol{\theta}_f) \quad (3)$$

$$\mathbf{y}(k + 1) = h_d(\mathbf{x}(k), \mathbf{u}(k), \mathbf{w}(k), \boldsymbol{\theta}_h) \quad (4)$$

with  $k$  the discrete time step, and  $f_d(\cdot)$  and  $h_d(\cdot)$  the discrete version of  $f(\cdot)$  and  $h(\cdot)$ .

A model given by (1) and (2) (or similarly by (3) and (4) in the discrete time) can be used via model-based state estimation methods in order to provide information about internal mental states of the human (e.g., in a human-machine interaction context) that are not measured directly via the output. Figure 3 illustrates such a model-based state estimator for estimation of the internal mental states of a human (which,

<sup>2</sup>Throughout the paper, for the mathematical notations, italic fonts are used for scalar variables and to distinguish vectors, italic bold fonts are used.

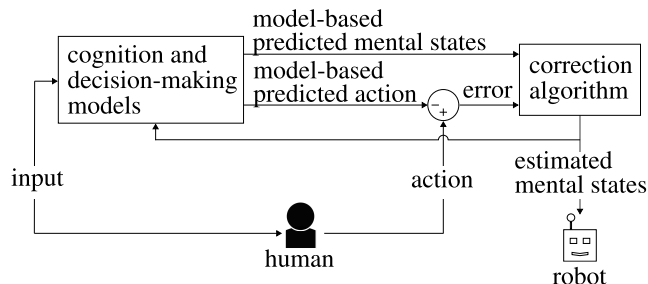
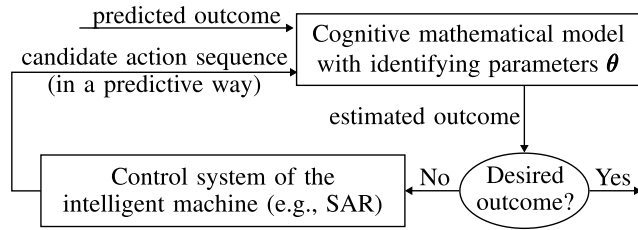


FIGURE 3. Model-based state estimator for making inferences about internal mental states of humans. In this figure, we use the cognitive and decision-making modules of the proposed model to predict internal mental states and action of a human, which correspond to, respectively,  $\mathbf{x}(t)$  and  $\mathbf{y}(t)$  in (1) and (2). A correction algorithm is used to produce the estimated mental states of the human based on the mental states predicted by the model and the error between the predicted output and the real output, i.e., the action of the human. Finally, the estimated mental states of the human determined by this model-based state estimator are fed back to the model (which closes the loop) and are sent to the observer agent (e.g., a social robot).

according to ToM, manifest themselves in the output, i.e., action of the human, and are thus observable according to the systems theory definition [45]). A more advanced version of such a state estimator may also be adaptive in the course of the human-machine interactions. More specifically, re-identification of the parameters of the cognitive and decision-making models (i.e., parameters that are represented via  $\boldsymbol{\theta}_f$  and  $\boldsymbol{\theta}_h$  in (1)-(4) may occur with a frequency lower than the frequency of the operation of the state estimator, i.e., per simulation step. In fact, the adaptive algorithm that performs the re-identification of the parameters may only be activated when necessary (e.g., whenever the estimation error surpasses a certain threshold). Moreover, such models can effectively be used by predictive control methods, e.g., in a loop with an optimizer, to propose a sequence of controlled inputs within a prediction window for the dynamical system that guarantees given requirements on desired performance criteria, and that satisfies imposed and desired constraints.

These two applications, i.e., estimating internal states (see Section III-B1) and making controlled inputs (i.e., decisions) that account for their long-term effect on the performance, are highly relevant for SARs that should understand and interact with humans, and will directly contribute to closing the gaps that have been identified in Section II. Figures 3 and 4 illustrate our main idea for using such mathematical dynamic models in estimating the mental states of humans and in enabling SARs to make future-aware decisions in interactions with humans that guarantee the success of these interactions in long terms.

In order to facilitate the formulation, the proposed ToM network representation has been divided into four submodules (see Figure 5): (1) *cognition module*, including the internal variables that play a role in the dynamic evolution of the fast-dynamics state variables, (2) *perception module*, including the variables and functions that contribute to the perceived knowledge of the rational agent,



**FIGURE 4. Model-based predictive control of intelligent machines that interact with humans.**

(3) *decision-making module*, including rational intention selection and rational action selection, and (4) *world model*, which formulates the influence of the actions of the rational agent on real-life data. The mental states and thus actions and behaviors of an observed agent are influenced by its fixed (i.e., invariant in short terms) characteristics including the general world knowledge, general preferences, and personality traits [30], [48], as well as by its dynamic state variables, including beliefs, goals, and emotions [24], [38]. Moreover, the environmental data of an observed agent is perceived by the agent in a personalized way [34], [35]. Next, we explain how the proposed cognitive models incorporate these characteristics.

**B. NETWORK REPRESENTATION: IDENTIFYING VARIABLES AND INTER-DYNAMICS**

We propose a network representation (see Figure 5) of the perception, cognition, and decision-making procedures of humans that is composed of various elements connected via directed links. These links represent the interdependencies and influences of the elements. The oval-shaped elements are used for variables, and the rectangular elements show mathematical functions. The cognition module in this representation is an extension of the model proposed by Baker [30] (see section II for details), where in addition to beliefs, goals, general world knowledge, and general preferences, we have included emotions and personality traits, which according to the literature (see the discussions in section II as well as [2], [7], [19], [38], [49], [50], [51]) are highly relevant for the cognition and social interactions of humans. Moreover, in order to represent the inter-dynamics of all these variables with more transparency, we have introduced two auxiliary variables, called ‘bias’ and perceived knowledge’.

*Remark 1:* Goals are immediate desires and needs of rational agents, such as finding food or reaching a location. General preferences of rational agents build up in long terms and remain invariant for long, and in order to identify them several interactions with a rational agent are needed. Examples of general preferences include favorite tastes, friends, and hobbies of a rational agent. Beliefs correspond to *temporary* knowledge or interpretations of rational agents from their world, while general world knowledge consists in *persistent* rationally perceived knowledge, which remains unchanged or is rarely updated. For example, a rational agent

believes that a friend who has left an hour ago to fetch a medicine from the drugstore is now in the city center, whereas the exact location of the drugstore is the agent’s general world knowledge.

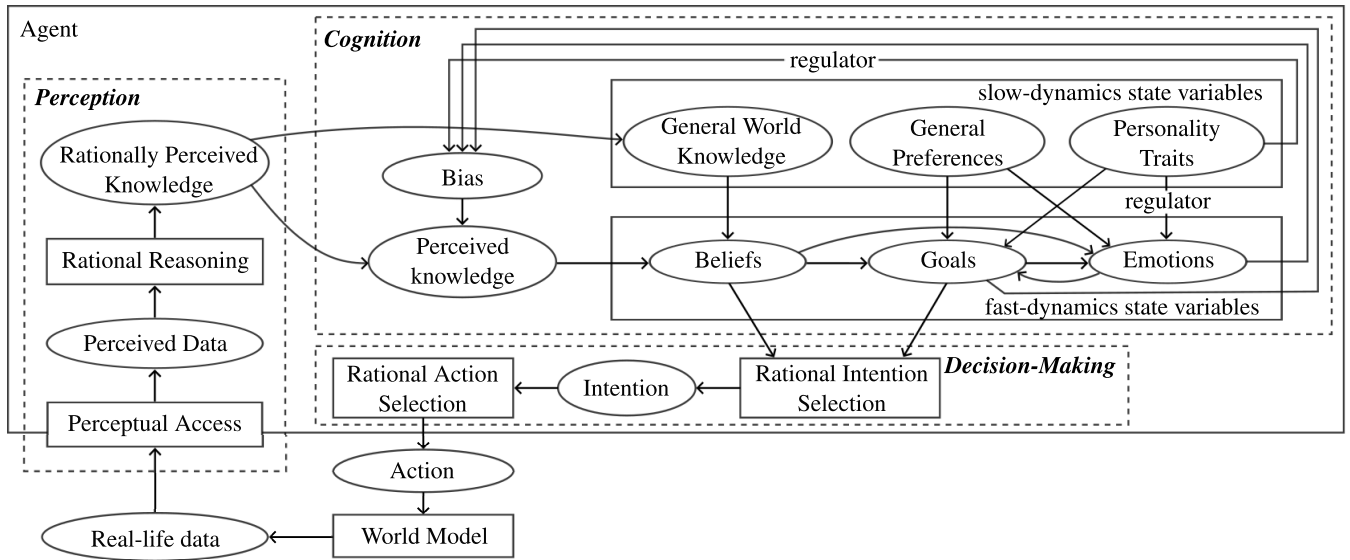
The perception module, which represents the procedure of transitioning real-life data (i.e., the input of the module) into rationally perceived knowledge (i.e., the output of the module), is composed of perceptual access (i.e., a function that models the observation of the agent), perceived data (i.e., the current representation of the real life data according to the observation of the agent), and rational reasoning, which is a function that receives perceived data and generates rationally perceived knowledge. Finally, the decision-making module includes the function “rational intention selection” that outputs the intention(s) of the observed agent, as well as the function “rational action selection”, which generates the action of the agent based on the intention(s). In summary, the main elements represented in Figure 5 correspond to:

- 1) External inputs (e.g., the weather conditions, the actions of a SAR) that may influence the mental states and thus actions and behaviors of observed agents.
- 2) State variables, including the mental states (e.g., beliefs, goals, emotions) and auxiliary variables (e.g., biases) of the observed agent.
- 3) Fixed parameters, including general world knowledge, general preferences, and personality traits of observed agents.
- 4) Dynamic processes, which are functions that receive the fixed parameters, current external inputs, and state variables, and update the next step state variables of the observed agent or predict its behaviors.

Note that when a SAR interacts with the observed agents, the actions of the SAR are the controlled inputs to the model, whereas factors such as the weather conditions are the uncontrollable inputs.

**1) ELEMENTS INTERNAL AND EXTERNAL TO AN OBSERVED AGENT**

While Baker [30] distinguishes the elements in the cognitive network representation of Figure 1 based on whether or not they depend on the current situation, we differentiate the elements of our proposed network representation (Figure 5) considering whether or not they are external to the observed agent. This explicit categorization of internal and external elements provides a model that can be personalized more efficiently than black box approaches. Moreover, since the internal elements of an observed agent are not visible to the observer agent, their inference is in general personalized to an observer agent. Although out of the scope of this paper, our framework is easily expandable for cases where second-order inferences, i.e., inference about the inference of an observer agent about the observed agent [17], are of interest. In addition to fully internal and fully external elements, our representation incorporates partially external elements. In particular, the perceptual access of an observed agent is



**FIGURE 5.** Proposed network representation of human perception, cognition, and decision-making procedures of humans including emotions, personality traits, and biases. Oval-shaped elements show input, output, state, and auxiliary variables, and rectangular elements correspond to processes or functions.

partially external, given that this process is influenced by external inputs, and is partially internal, since it is shaped by internal characteristics of the individual. This dual influence on perceptual access is explained in detail in Section III-B3.

## 2) FAST-DYNAMICS AND SLOW-DYNAMICS STATE VARIABLES

The state variables of the proposed model are distinguished according to their relevance for duration (i.e., short-term or long-term) of interactions between two rational agents and to the frequency of their dynamics. Consequently, two categories of state variables are defined: (i) Fast-dynamics state variables, which may constantly vary (with a timescale in the range of seconds or minutes) as a response to specific situations the observed agent faces. Goals, beliefs, and emotions in Figure 5 are fast-dynamics state variables. (ii) Slow-dynamics state variables, which vary according to large time scales (months or years). General world knowledge, general preferences, and personality traits in Figure 5 are slow-dynamics state variables. Fast-dynamics state variables are more relevant for short-term interactions, while slow-dynamics state variables become more relevant throughout long-term interactions, when the fixed or repetitive patterns of cognitive procedures resulting from these slow-dynamics state variables provide extra information for the observer agent to make more precise estimates and predictions [48].

*Remark 2:* Slow-dynamics state variables may influence the evolution of fast-dynamics state variables (see Example 1-Example 4 in Appendix A), while the opposite is not necessarily true (especially in short terms). The main aim of this paper is to formalize and formulate the evolution of fast-dynamics state variables. Modeling the evolution of

slow-dynamics state variables is out of the scope of this paper. Thus, slow-dynamics state variables are mainly considered as fixed parameters.

In the next sections, we explain how the inter-dynamics of the cognition, perception, and decision-making modules of the proposed model, shown in Figure 5, have been developed based on literature and real-life examples given in Appendix A.

## 3) PERCEPTION MODULE

In the model proposed by Baker et al. (see [30] and [33]), the perception procedure that generates the beliefs is represented by a single element called the principle of rational belief (see Figure 1). This simplification was shown to be sufficient to explain the relationship between the environmental inputs and the inferred beliefs in the simple environments and scenarios considered in [33]. In real-life scenarios, however, a more complicated procedure occurs before a belief is developed based on the raw real-life data. More specifically, rational agents may (deliberately or indeliberately) access and perceive only a portion of the real-life data. For instance, in the same environment, different rational agents may notice different types of data [35] (e.g., one may notice a sound, while another agent filters it out). Moreover, rational agents may receive partial data due to external constraints (e.g., missing visual data due to occlusion). Therefore, observed agents may hold false or inaccurate beliefs (c.f. the Sally-Anne experiment [52]), which is essential for ToM-based observer agents to recognize [18], [53], [54]. Furthermore, the same perceived data can be differently interpreted by different rational agents [34], [55].

To address these aspects, in our proposed model, the perception process that transforms real-life data into beliefs

is decomposed into smaller, well-defined sub-processes i.e., *perceptual access* and *rational reasoning* (see Figure 5 and Examples 9 and 10 in Appendix A). Real-life data from the environment is perceived via *perceptual access*, which, as discussed above, is partially personalized and partially depends on the environment. Thus, the rectangular element corresponding to the perceptual access function in Figure 5 is located partially inside the agent box and partially in the environment. The perceived data is then processed via *rational reasoning*, which, as opposed to the universal principle of rational reasoning applied by Baker et al. in [30] and [33], is specific to a rational agent. Accordingly, the rational agent makes a judgment called the *rationally perceived knowledge*, which is then transformed into a belief or a piece of general world knowledge in the cognitive module.

#### 4) COGNITIVE MODULE

The cognitive module receives the rationally perceived knowledge as input. Since, compared to the model of Baker shown in Figure 1, the cognitive module of our model includes emotions and personality traits as well, we first discuss the inter-dynamics of the emotions with the other fast-dynamic state variables. Next, we explain the influence of personality traits on the fast-dynamics state variables.

Emotions of a rational agent are stimulated by its beliefs [56], [57], where goals and general preferences, alongside beliefs, can also stimulate emotions [56]: On the one hand, when a rational agent develops a belief that is in line with the fulfillment of its goal, positive emotions may be stimulated (see Example 6 in Appendix A). On the other hand, when the agent develops a belief that hinders the chances of fulfilling its goal, negative emotions may be stimulated. Similarly, when a general preference is supported by a developed belief, positive emotions can be generated, while beliefs that conflict with the general preferences may result in negative emotions (see Example 7 in Appendix A). We assume that direct influences from general world knowledge on emotions are negligible, since when a piece of general world knowledge is relevant in the current context, it generates a belief (see Figure 5), and this belief can then influence emotions.

Studies show that emotions can affect the goals of rational agents [51], [58], [59], [60], [61]. For instance, gratitude can galvanize rational agents into helping others [60], or anxiety may trigger rational agents to avoid stressful situations [58]. More specifically, emotions may result in the development of a goal that contradicts with general preferences of the rational agent, or in the change of a goal that was previously made by the agent. The influence of emotions on goals is shown in Figure 5 via a directed link (also see Example 8 in Appendix A). While emotions do not directly influence beliefs, they can affect the processes that result in judgements or beliefs of rational agents [55], [58], [59]. Thus, we have introduced the auxiliary variables ‘bias’ and ‘perceived knowledge’ in order

to illustrate this influence. For example, positive emotions may introduce optimistic biases into the process of generation of new beliefs, whereas negative emotions may lead to the formation of overly pessimistic beliefs [60]. Therefore, the influence of emotions on the development of beliefs will be introduced into the proposed cognitive model via bias. Similarly to emotions, goals of rational agents may introduce some bias into their rational reasoning processes [62], where the intensity of the bias may depend on the personality traits of rational agents [55] (see Example 12 in Appendix A). For the purpose of describing the cognition of an observed agent, the bias is introduced into the belief after the rational reasoning process is executed (see Figure 5). This results in the *perceived knowledge* (as opposed to *rationally perceived knowledge*), which generates the biased beliefs.

Regarding the influences of personality traits on fast-dynamics state variables, there is no evidence that beliefs are directly affected by personality traits, whereas goals may directly be affected by personality traits [48] (see Figure 5): For instance, while the goal of an introverted rational agent is to avoid strangers, the goal of an extroverted rational agent is to make new friends. Moreover, while emotions are not generated by personality traits [50], [55], they are regulated by personality traits. In fact, personality traits of rational agents may determine the extent to which certain beliefs affect their emotions [50], [56], [63]. For instance, extroverted individuals are more likely to experience positive and intense emotions when under the same stimuli [55].

*Remark 3:* In summary, beliefs, either alone or in combination with generated goals or general preferences, and boosted or hindered by personality traits, trigger the emotions. For the sake of brevity, we use the following terminology: **emotion trigger 1**, **emotion trigger 2**, and **emotion trigger 3** for, respectively, a combination of beliefs and general preferences, solely beliefs, and a combination of beliefs and goals, triggering emotions.

*Remark 4:* Note that, as it is demonstrated in Figure 5, the rationally perceived knowledge that is the input of the cognitive module can be transformed into a belief, or if this piece of knowledge does not change in the short term, into a general world knowledge ([43], [44]).

#### 5) DECISION-MAKING MODULE

Based on the BDI theory, we introduce the rational intention selection procedure (see the rectangular element in Figure 5), which is an intermediate process in the generation of actions and behaviors of rational agents. This process generates or selects an intended high-level action that maximizes the chance of fulfillment of the goals of the rational agent, given its beliefs. In other words, mathematically, rational intention selection is a function, with beliefs and goals as inputs and intentions as outputs. The generated intention can then be translated into lower-level actions by the rational action selection function. For example, the intention of the agent



may be to drink a glass of water, where one action is to open the cabinet where the glasses are kept.

**C. DYNAMIC MATHEMATICAL REPRESENTATION OF PERCEPTION, COGNITION, AND DECISION-MAKING**

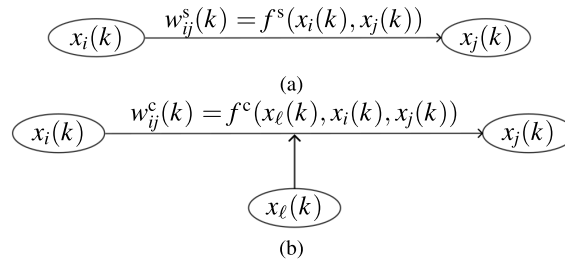
The main aim of this paper is to model the dynamic evolution of the fast-dynamics state variables for the network representation given in Figure 5, and to predict the resulting action. Thus, we need to develop mathematical formulations for the perception, cognition, and decision-making modules given in Figure 5. Note that a mathematical model that represents perception receives real-life data as input and generates rationally perceived knowledge as output. This output is the input to the mathematical model of cognition, which generates two outputs, beliefs and goals. Finally, these two outputs enter the decision-making module as input, and the corresponding model generates action as output.

Previous works [19], [33], [38], [64], [65] mainly use Baye’s theorem to describe human’s cognition. Modeling cognition under this assumption corresponds to representing the cognitive network representation as a Bayesian network [66]. However, this requires knowing beforehand the probability functions of all the variables that are part of the network. Furthermore, the integration of a probabilistic-based formulation of the model with model-based predictive approaches would require the usage of decision-trees, which tend to grow exponentially with the size of the prediction horizon. Therefore, this formulation cannot be integrated with model-based predictive methods in real life applications.

Alternatively, FCMs [67] represent concepts and variables that correspond to complex and/or uncertain systems and their interlinks and interactions. Contrarily to Bayesian networks, FCMs support cyclic connections [68], which is very relevant for modeling cognitive procedures of humans (c.f. Figure 5). Moreover, concepts or variables in an FCM can be represented mathematically by fuzzy variables, which excellently fit the concepts that are involved in human cognitive procedures (e.g., beliefs, goals, emotions, etc.). Therefore, in this paper, based on the idea of FCMs, we propose an extended FCM representation and use it to mathematically formulate the cognitive module of our proposed ToM model.

**1) NOVEL FCM-BASED FRAMEWORK FOR TRANSFORMING STATIC NETWORK REPRESENTATIONS INTO DYNAMIC MATHEMATICAL MODELS**

In an FCM, we denote the vector of all the  $n$  state variables that describe the system at time  $k$  by  $\mathbf{x}(k)$ , and the  $i^{\text{th}}$  state variable by  $x_i(k)$  with  $i = 1\{1, \dots, n\}$ . This vector for the cognitive module of our proposed model includes the fast-dynamics and slow-dynamics state variables, and the auxiliary variables. Mathematically,  $x_i$  (e.g., emotion) can be represented as a fuzzy variable in an FCM, with  $\tilde{X}_{i,j}$  a possible



**FIGURE 6. Simple and complex linkage. (a) Simple Linkage. The horizontal arrow is a simple linkage, shown by  $(i, j)$ . (b) Complex Linkage. The vertical arrow is a side linkage, shown by  $(\ell, i, j)$ .**

realization that is mathematically given by a fuzzy set.<sup>3</sup> Fuzzy variables [69] are used to mathematically represent concepts that are imprecise, vague, and mainly given by human linguistic terms. Fuzzy sets are a generalization of regular (crisp) sets; while an element either belongs to or does not belong to a crisp set (i.e., its membership degree to the crisp set is either 1 or 0), it may partially belong to a fuzzy set, with its membership degree varying within the interval  $[0, 1]$  [70]. In order to represent fuzzy sets and perform mathematical operations on them, their corresponding membership functions are used [70]. The fuzzy sets  $\tilde{X}_{i,j}$  should cover the domain  $X_i$  of  $x_i$ .

The influence of variable  $x_i$  over variable  $x_j$  in an FCM is represented by a directed line called a *linkage* (see Figure 6). In the classic FCM formulation, every linkage between  $x_i$  and  $x_j$  is characterized by a weight  $w_{ij} \in [-1, 1]$  that reflects the level of influence of variable  $x_i$  over variable  $x_j$ . Whenever  $w_{ij}$  is positive (negative), an increase of  $x_i$  implies an increase (a decrease) of  $x_j$ , and the larger the absolute value of  $w_{ij}$ , the larger the influence of  $x_i$  over  $x_j$  (whenever  $w_{ij}$  is null, changes in variable  $x_i$  do not influence variable  $x_j$ ).

In FCMs, the weights  $w_{ij}$  are considered to be constant. However, in order to accurately model most real-world systems with FCMs, variable weights may be required [71], [72]. For instance, in rule-based FCM [71] the values of the weights depend on the realized values of the causing variable  $x_i$ . In our proposed cognitive network representation, due to the mutual influences of the fast-dynamics state variables on one another, the value of weight  $w_{ij}$  may vary in time and, in general, depend on the realized values of the causing variable  $x_i(k)$ , affected variable  $x_j(k)$ , or other intermediate variables  $x_\ell(k)$ . To address this, we consider weights as functions of the causing, affected, or intermediate variables and accordingly define *simple linkages*, *side linkages*, and *complex linkages*, which are explained next.

The linkage  $(i, j)$  that directly connects two variables  $x_i(k)$  and  $x_j(k)$  and is not influenced by any intermediate variables (see Figure 6(a)) is called *simple linkage*. A *side linkage*  $(\ell, i, j)$  (see the vertical arrow in Figure 6(b)) corresponds to the directed influence of an intermediate variable  $x_\ell(k)$  over

<sup>3</sup>In this paper, to distinguish fuzzy sets from crisp sets, we use a tilde symbol.

TABLE 1. Mathematical notations used in (8)-(13).

Mathematical notation	Process or variable
$f_{PA}$	perceptual access
$f_{RR}$	rational reasoning
$d_{RL}$	real life data
$d_p$	perceived data
$\kappa_{RP}$	rationally perceived knowledge
$f_{RIS}$	rational intention selection
$f_{RAS}$	rational action selection
$x_B$	belief
$x_G$	goal
$i$	intention
$a$	action

a linkage  $(i, j)$  that connects variables  $x_i(k)$  and  $x_j(k)$ . The group of a linkage that is influenced by one or several side linkages and all its side linkages is called a *complex linkage* (see Figure 6(b)). The set of all pairs  $(i, j)$  corresponding to simple linkages is given by  $\mathbb{L}$  and the set of all trios<sup>4</sup>  $(\ell, i, j)$  that correspond to complex linkages is given by  $\bar{\mathbb{L}}$ . Note that  $\mathbb{L}$  and  $\bar{\mathbb{L}}$  are complementary sets, and have no common element. The weight  $w_{ij}^s(k)$  of a simple linkage and  $w_{ij}^c(k)$  of a complex linkage for time step  $k$  are computed via, respectively, function  $f^s : \mathbb{X}_i \times \mathbb{X}_j \rightarrow [-1, 1]$  and function  $f^c : \mathbb{X}_i \times \mathbb{X}_j \times \mathbb{X}_\ell \rightarrow [-1, 1]$ . We have

$$w_{ij}^s(k) = f^s(x_i(k), x_j(k)), \quad \forall i, j \text{ for which } (i, j) \in \mathbb{L} \quad (5)$$

$$w_{ij}^c(k) = f^c(x_\ell(k), x_i(k), x_j(k)), \quad \forall i, j \text{ for which } \exists \ell \text{ such that } (\ell, i, j) \in \bar{\mathbb{L}} \quad (6)$$

for all  $k \in \{1, 2, \dots\}$ , and  $x_i(k) \in \mathbb{X}_i$ ,  $x_j(k) \in \mathbb{X}_j$ , and  $x_\ell(k) \in \mathbb{X}_\ell$  are elements of the state vector  $\mathbf{x}(k)$ .

Then, the dynamic equation for updating variable  $x_j(k)$  that evolves per time step  $k$  within the proposed extended FCM is formulated by

$$x_j(k + 1) = \sum_{(i,j) \in \mathbb{L}} f^s(x_i(k), x_j(k))x_i(k) + \sum_{(i,j,\ell) \in \bar{\mathbb{L}}} f^c(x_\ell(k), x_i(k), x_j(k))x_i(k) + f_j(x_j(k)) \quad (7)$$

where function  $f_j(\cdot)$  determines the influence of variable  $x_j(k)$  on its evolved value  $x_j(k + 1)$ .

Next, we describe the mathematical representations of the perception, cognition, and decision-making modules. Since, for the sake of compactness of the mathematical notations, we use a number of abbreviations in the mathematical notations of the following sections, Table 1 represents the meaning of these mathematical notations.

<sup>4</sup>In our cognitive network representation given in Figure 5, complex linkages include no more than one side linkage. Thus, for the brevity of notations, we talk about pairs and trios only.

## 2) MATHEMATICAL REPRESENTATION OF PERCEPTION

The perception procedure includes two functions (see Figure 5), i.e., perceptual access (shown by  $f_{PA}(\cdot)$ ) and rational reasoning (shown by  $f_{RR}(\cdot)$ ). Perceptual access receives real life data (which we show by  $d_{RL}$ ) and generates perceived data (which we show by  $d_p$ ). Thus, at time step  $k$  we have

$$d_p(k) = f_{PA}(d_{RL}(k)) \quad (8)$$

Similarly, rational reasoning receives perceived data and generates rationally perceived knowledge (shown by  $\kappa_{RP}$ ). Thus, at time step  $k$  we have

$$\kappa_{RP}(k) = f_{RR}(d_p(k)) \quad (9)$$

Hence, from (8) and (9), the mathematical representation for the perception module is given by

$$\kappa_{RP}(k) = f_{RR}(f_{PA}(d_{RL}(k))) \quad (10)$$

## 3) MATHEMATICAL REPRESENTATION OF COGNITION

For the cognitive module of our proposed model, (7) may be used with  $x_j(k) \in \{\text{current beliefs, current goals, current emotions, current biases}\}$ ,  $x_i(k) \in \{\text{current beliefs, current goals, current emotions, current biases, general world knowledge, general preferences, personality traits, current perceived knowledge}\}$ , and  $x_\ell(k) \in \{\text{current goals, general preferences, personality traits}\}$ . In fact, using (7) will provide us with the desired state space representation (3), where the current ‘perceived knowledge’ acts as input  $\mathbf{u}(k)$ , and the weights computed via  $f^s(\cdot)$  and  $f^c(\cdot)$ , as well as the ‘general world knowledge’, the ‘general preferences’, and the ‘personality traits’, act as the parameters in  $\theta_f$ . Note that in this case there are no particular external disturbances  $\mathbf{w}(k)$  and  $f_d(\cdot)$  in (3) has been represented as additive functions in (7).

*Remark 5:* In order to define functions  $f^s(\cdot)$ ,  $f^c(\cdot)$ , and  $f_j(\cdot)$  for the proposed cognitive models, different approaches may be used, such as using explicit mathematical representations or describing these functions via fuzzy inference systems (see, e.g., [71]). When we use explicit expressions for  $f^s(\cdot)$ ,  $f^c(\cdot)$ , and  $f_j(\cdot)$ , in order to use fuzzy variables for the extended FCM, we should in general use Zadeh’s extension principle [69]. In case this is (computationally or analytically) not possible, these functions can be modeled via fuzzy inference systems [71].

## 4) MATHEMATICAL REPRESENTATION OF DECISION-MAKING

The decision-making procedure is composed of two functions, rational intention selection  $f_{RIS}(\cdot)$ , and rational action selection  $f_{RAS}(\cdot)$ . The intention selection receives the belief  $x_B(k + 1)$  and goal  $x_G(k + 1)$  predicted for the next step and generates the intention of the rational agent  $i(k + 1)$  that is predicted also for the next step. Thus, for time step  $k + 1$  we have

$$i(k + 1) = f_{RIS}(x_B(k + 1), x_G(k + 1)) \quad (11)$$

Similarly, the rational action selection receives the intention that is predicted for the next time step, and generates the predicted action for the next time step  $k + 1$ . We have

$$a(k + 1) = f_{RAS}(i(k + 1)) \quad (12)$$

Thus, the decision-making is mathematically given by

$$a(k + 1) = f_{RAS}(f_{RIS}(x_B(k + 1), x_G(k + 1))) \quad (13)$$

where  $a(k + 1)$  is the predicted action of the rational agent for time step  $k + 1$  based on the observations, measurements, and estimations at time step  $k$ . Moreover,  $x_B(k + 1)$  and  $x_G(k + 1)$  are the predicted next belief and next goal.

*Remark 6:* For a predictive decision-making by, e.g., a SAR with a prediction horizon bigger than 1, the influence of the actions on real life data should also be modeled (see Figure 5). This, however, is out of the scope of this paper, since it concerns a modeling of the environment of the rational agent.

#### IV. CASE STUDY

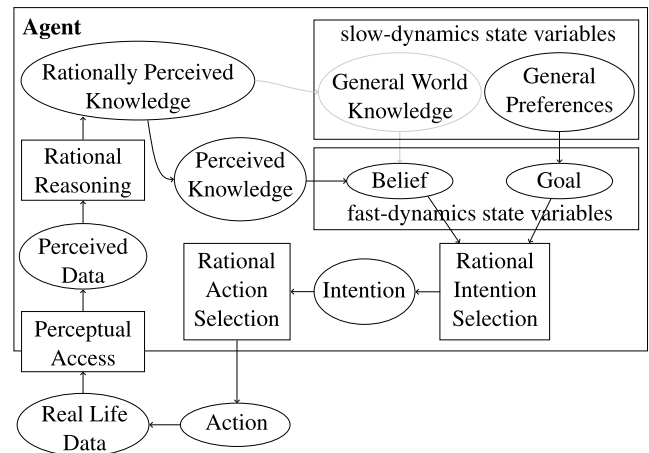
The ToM model proposed in this paper (see Figure 5) was implemented via Python and was used to estimate the expected perception, cognition, and decision-making procedures of human participants in virtual scenarios of emergency evacuation. The experiment were designed to assess the following research questions:

- 1) How well does the proposed model predict the behavior and estimate the mental states of rational agents, particularly humans, in a 2-dimensional environment (similarly to the work in [18] and [33])?
- 2) How much do the cognitive elements (emotions, biases, and personality traits) that were proposed in this paper contribute to improving such predictions?
- 3) How much do the models that were identified for each participant provide a relevant insight about the general preferences of the participants?

In particular, we answer these research questions for our model (which, in the rest of the case study, we call ToM-I) in comparison with a model that did not include the elements that were proposed in this paper (which we call ToM-II and is illustrated in Figure 7). Since the personality trait is an adaptable parameter, ToM-II has one less parameter than ToM-I. To assess the contribution of the emotions and biases independently of the number of adaptable parameters, a third model, ToM-III, that is similar to ToM-I but does not include the personality traits, was also considered.

#### A. EXPERIMENT SETUP

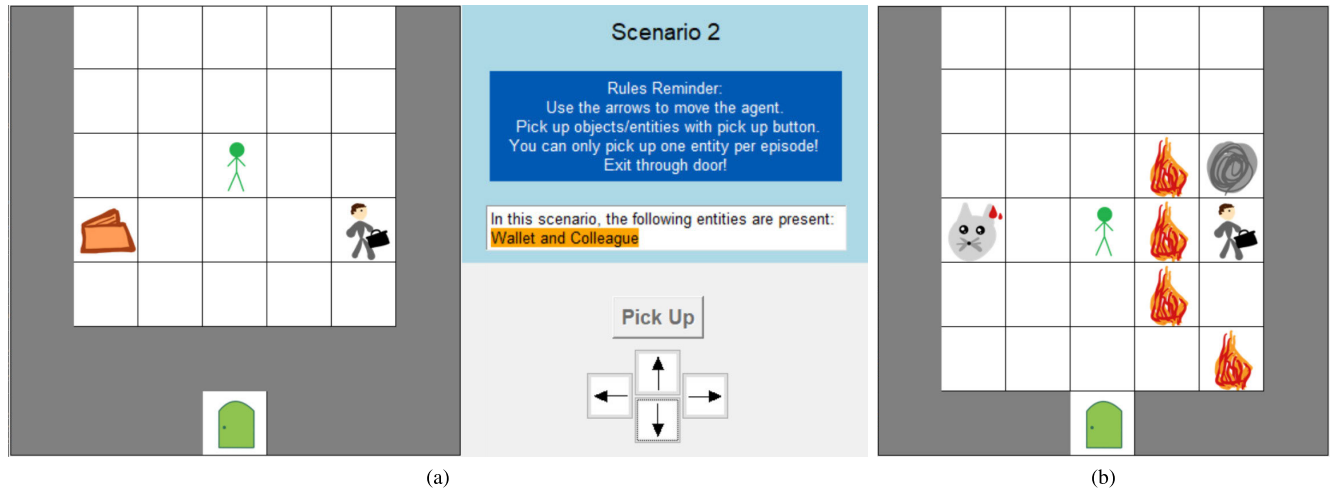
While Baker et al. [33] and Rabinowitz et al. [18] artificially generated the behavior of the observed agents to use as ground truth, we used the observed behavior of human participants in order to generate the ground truth behavior of the observed agents. This is because our main objective for answering research question (1) is to assess how well the proposed model estimates the mental states and behavior



**FIGURE 7.** Baseline model used for comparison in the experiments, ToM-II. This baseline model is similar to the model represented in Figure 5, but it does not include the emotions, biases, and personality traits.

of real humans. Hence, the performance of our ToM model, ToM-I, was assessed based on how accurately it estimated the goals and predicted the behavior of the participants. Moreover, in order to answer research question (2), i.e., to assess the particular contribution of the elements emotions, biases, and personality traits that have been added to our model, ToM-I was compared with the model proposed by Baker [30]. For a fair comparison, the three models were formulated as described in section III-C1.

The experiments were performed via a computer simulation, with 21 human participants. The participants were aged between 14 and 32 years old. Every participant received an executable of a set of simulations that they could run on their own computer. In this set of simulations, participants were able to control the movements, actions, and decisions of a virtual rational agent in 16 simulated scenarios. The decisions taken by each participant were recorded into an Excel file that was sent back to the researchers. The aforementioned scenarios consisted of emergency evacuation scenarios that included a closed space that was populated by living beings (e.g., pets, humans) and objects (e.g., personal belongings). The participants were able to move the rational agent around and decide which entity they wanted to save in each scenario (see 8). To increase the diversity in the general preferences, the scenarios included three types of entities: a dependent being (e.g., a child or a pet), a human colleague, and a personal object (e.g., a wallet or a set of keys). In order to make the scenarios more realistic and to increase the chance that the general preferences had a broad span of values, at the start of the experiment, each participant chose which entity they related to most in each category from a predefined set of options. During the simulation, the participants were able to save an entity by picking it up (see Figure 8(a)) and going to the exit (shown as a green gate in Figure 8(a)). Alternatively, the participants could go to the exit without saving any objects or beings. Once the rational agent reached



**FIGURE 8.** Examples of emergency evacuation scenarios used in our case study. (a) A screenshot of a simulated scenario that a participant can perform: On the left-hand side of Figure 8, the participant sees the simulated emergency evacuation scenario, including various entities (e.g., a colleague and a wallet). On the top right-hand side of Figure 8, the participant sees reminders about the main rules of the simulation. At the bottom right-hand side of Figure 8, the participant can use the buttons to move the rational agent and to pick up objects or beings. (b) An example of a scenario where one of the entities (the colleague) is surrounded by fire and another entity (the pet) is injured.

the green door, the scenario terminated. In some scenarios, the entities were injured or surrounded by fire (see Figure 8(b)). These variations were added to induce emotions of fear in the participants for the entities that were at risk. Note that at every simulation step (i.e., after each action of the rational agent), the participants could only observe the closest positions to the rational agent and the previously visited positions. In order to decrease the exploring movements by the participants, the position of the exit was always visible.

## B. IMPLEMENTATION SETUP

For the perception module, the real-life data contained the inputs from the two-dimensional environment, including the positions of all entities, the rational agent, and the obstacles, as well as the characteristics of the entities (e.g., if they are injured, or if they are a living being or an object). The perceptual access was modeled based on the field of view of the rational agent per simulation step. The field of view of the rational agent was assumed to cover an area of five by five cells. The cells that fell within the current field of view of the rational agent, as well as the cells that were previously in the field of view of the rational agent, were visible to both the participant and the corresponding ToM model (as it is also shown in 8). The perceived data was based on the real-life data corresponding to the visible positions, but did not include any data from the unseen positions. The rationally perceived knowledge was deduced according to logical propositions and based on the perceived data. For example, the location of the fire and the current position of the rational agent were used to deduce a rationally perceived knowledge about whether or not the rational agent was near the fire.

Regarding the cognitive module, since humans express their mental states mainly via quantified expressions

(e.g., relatively happy), the fast-dynamics state variables of the cognitive module (e.g., beliefs, goals, emotions) were represented by fuzzy variables, i.e., verbal terms were associated with the realizations of each variable. These variables in the experiments corresponded, for example, to the belief of the rational agent about an entity being injured, or to the goal of the rational agent to save one entity, or to the emotion of fear. As for the slow-dynamics state variables, we considered the preference towards each of the entities (i.e., the general preferences), and the altruism of the agent (i.e., a personality trait). The full list of all the fuzzy variables, and the corresponding linguistic terms and range values, can be found in [73]. The linkages that connected the elements of the cognitive module were mathematically represented by fuzzy inference systems (see Remark 5 for details).

As for the decision-making module, the intention of the rational agent was computed based on the goal fired with the highest value at each simulation step. The action of the rational agent was computed as the movement of the rational agent along the shortest path (according to the Dijkstra's Algorithm) to the entity that was associated with the current intention of the rational agent, except for when the position of the rational agent and the position of the entity were the same: in that case, the action would be to pick up the object or the being.

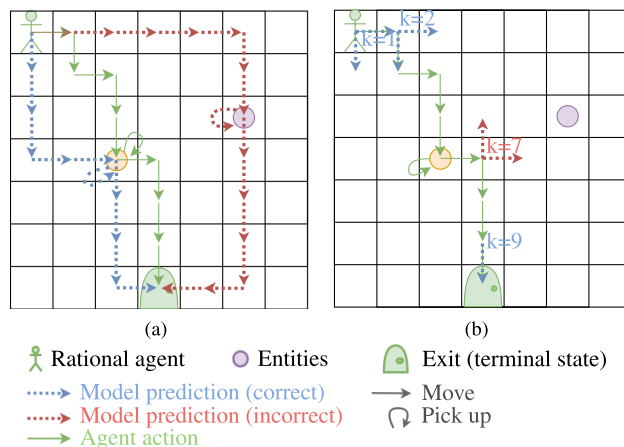
## C. TRAINING PHASE FOR THE MODELS

A proportion of the data that was collected from each participant was used to personalize the model to that participant. In fact, the data that was collected in 75% of the scenarios was used to train the ToM model, and the remaining was used to test the accuracy of the model in predicting the behavior of the agent in unseen scenarios.

The personalization consisted in identifying the values of the slow-dynamic state variables (i.e., the general preferences towards each of the entities for ToM-I, ToM-II, and ToM-II, as well as the personality trait altruism for ToM-I) for each participant. The slow-dynamic state variables are also called the model parameters due to their very low update frequency. These parameters  $\theta$  were optimized using a grid search algorithm with a step of 0.1. The number of scenarios per participant was only 16, where this number was selected to provide sufficient data for training the models, while preventing the participants from getting distracted or tired, resulting in inaccurate answers. Hence, in order to increase the reliability of the results, the combination of scenarios that were used for training and for testing were shuffled in order to create 4 different combinations. Running one iteration of the model using an Intel CPU (i7-1185G7 3.0GHz, single thread implementation) took 0.20554 s on average. Furthermore, around 8 hours per participant were required to train the ToM-I model on a super computer with 48 cores (2x Intel XEON E5-6248R 24C 3.0GHz).

The action of the participants per simulation step included moving the rational agent to one of its (allowed) neighboring cells or picking up/rescuing an object/being. This action depended on two decision levels: (1) Whether or not the rational agent decided to rescue a particular being or pick up a particular object (2) Which path the rational agent followed to approach the corresponding being/object or the exit. Therefore, for identifying the models per participant, for the training sets, two independent optimizations were performed in order to train and assess the models for two types of predictions: the prediction about the entity that the rational agent would save in each scenario (i.e., predicting the overall intention of the rational agent per scenario), as well as the prediction of the action that would be taken by the rational agent per simulation step. In order to train and assess the models regarding the first prediction, per simulation step the models were run, and the predicted actions of the rational agent were executed consecutively, until the exit was reached. The saved entity, as predicted by each model, was compared with the entity that was really saved by the participant who steered the rational agent in that scenario (see Figure 9(a)). In order to train and assess the models regarding the second prediction, the actions that were chosen by the participant in each simulation step were compared with the actions that were predicted via each model for that simulation step (see Figure 9(b)).

Equation (14a) shows the cost function used to train and assess the models with respect to the prediction of the overall intention for participant  $p$  throughout a set  $\mathcal{S}$  of training or test scenarios, where  $i_{p,s}$  is the intention of participant  $p$  for scenario  $s$  that is predicted via the model, and  $i_{p,s}^*$  is the intention of participant  $p$  in scenario  $s$ . Both  $i_{p,s}$  and  $i_{p,s}^*$  take one out of the four following realizations: save one of the three entities (i.e., the dependent being, the colleague, or the personal belonging), or directly move to the exit. The cost function used to train and assess the prediction of the



**FIGURE 9.** Graphic representation of the two types of predictions made by the models. (a) Predictions by the models regarding the overall intention of the rational agent, where these predictions are considered correct if the model predicted a trajectory in which the agent chose the same intentions as the participants (i.e., either to save a specific entity or to go directly to the exit). The dashed arrows in blue represent an accurate prediction of the overall intention (i.e., the rational agent rescued the orange entity, which was also rescued by the participant in this scenario, and then went to the exit), while the dashed red arrows represent a wrong prediction about the intention of the participant via the models (i.e., the rational agent rescued the purple entity, which was different from the orange entity that was in reality rescued by the participant in this scenario, and then went to the exit). (b) Predictions by the models regarding the action(s) of the rational agent for different simulation steps. In the first two simulation steps (specified by  $k = 1$  and  $k = 2$ ), correct predictions are shown by the dashed blue arrows. Since the model estimates that the rational agent wants to save the orange entity, any of the two predicted actions shown in the figure for the first simulation step (i.e., moving downwards or moving to the right) are equally likely to be chosen by the agent. An incorrect prediction is shown via dashed red arrows for the simulation step  $k = 7$ , where the model estimates the optimal actions to be moving towards the purple entity, while the participant wants to go to the exit.

step-wise actions of the rational agent for participant  $p$  is given by (14b), where  $c_{k,p,s}(\theta)$  is 0 if the action taken by the participant  $p$  in scenario  $s$  at simulation step  $k$  is one of the predicted actions by the model for the same simulation step, and is 1 otherwise. Moreover,  $n_{steps,p,s}$  is the total number of the simulation steps taken by participant  $p$  in scenario  $s$ . For  $p = 1, \dots, n_p$  (with the number of participants  $n_p = 21$  in our case study), we have

$$J_{int,p}(\theta) = \sum_{s \in \mathcal{S}} J_{int,p,s} \quad \text{with } J_{int,p,s}(\theta) = \begin{cases} 0, & \text{if } i_{p,s} = i_{p,s}^* \\ 1, & \text{otherwise} \end{cases} \quad (14a)$$

$$J_{act,p}(\theta) = \sum_{s \in \mathcal{S}} J_{act,p,s} \quad \text{with } J_{act,p,s}(\theta) = \frac{1}{n_{steps,p,s}} \sum_{k=1}^{n_{steps,p,s}} c_{k,p,s}(\theta) \quad (14b)$$

Note that for both optimization procedures,  $\theta$  is the optimization variable, which includes all the identification parameters of the corresponding model.

**TABLE 2.** Average minimum cost achieved by each model when predicting the intentions of each participant. The smallest this minimum cost value, the more satisfactory the performance of the corresponding model.

Model	Training phase			Test phase		
	ToM-I	ToM-II	ToM-III	ToM-I	ToM-II	ToM-III
Combination 1	7.14	15.48	10.32	20.24	34.52	22.62
Combination 2	6.75	13.49	10.71	16.67	19.05	17.86
Combination 3	5.56	15.08	7.94	21.43	25.00	23.81
Combination 4	8.33	18.25	11.11	15.48	20.24	17.86
Average	6.94	15.57	10.02	18.45	24.70	20.54

**TABLE 3.** Average minimum cost achieved by each model when predicting the actions of each participant.

Model	Training Phase			Test Phase		
	ToM-I	ToM-II	ToM-III	ToM-I	ToM-II	ToM-III
Combination 1	31.81	34.97	33.44	34.56	36.52	35.13
Combination 2	32.90	34.70	34.34	32.15	37.26	32.95
Combination 3	31.84	35.38	33.32	33.90	36.16	34.50
Combination 4	32.06	35.31	33.37	35.14	35.61	35.68
Average	32.15	35.09	33.62	33.94	36.39	34.57

**D. RESULTS AND DISCUSSIONS**

In order to assess the performance of our model and the influence of our contributions, the average optimal costs achieved by each one of the three models for the test phase were compared. The optimal costs  $J_{int,p}^*$  and  $J_{act,p}^*$  that are achieved by each model for participant  $p$  are obtained by evaluations equations (14a) and (14b) using the parameters  $\theta^*$  that have been obtained by minimizing these cost functions in the training phase. The average optimal costs  $\bar{J}_{int}^*$  and  $\bar{J}_{act}^*$  for each trained model are an average of the aforementioned costs  $J_{int,p}^*$  and  $J_{act,p}^*$  across all the  $n_p$  participants and are divided by the number of the scenarios, as it is given by (15) for the average intention cost  $\bar{J}_{int}^*$ . Given that the costs computed via (14) are normalized, the costs computed via (15) are also normalized, and we report them in percentage.

$$\bar{J}_{int}^* = \frac{1}{n_s} \cdot \frac{1}{n_p} \sum_{p=1}^{n_p} J_{int,p}^*(\theta^*) \tag{15}$$

The average optimal cost  $\bar{J}_{act}^*$  for the action prediction is calculated similarly to  $\bar{J}_{int}^*$ . The average optimal costs obtained both in training and in test conditions when the models estimated the intention of the participants and predicted their actions are shown, respectively, in Tables 2 and 3 for the four combinations of training and test scenarios.

Regarding the estimations of the intentions of the participants, we first focus on the comparison between ToM-I and ToM-II, and we subsequently focus on comparing ToM-I with ToM-III. ToM-I obtained the best performance, i.e., the lowest estimation error compared to ToM-II in both the training and the test phases for all the four combinations of training and test scenarios (compare the average costs given in the last row of Table 2 for TOM-I (see column 2 and column 5) with those for TOM-II). On average, ToM-I achieved an error of 6.94% on the training scenarios and of 18.45% on the test scenarios, representing an improvement of

8.63% in training conditions and of 6.25% in test conditions when compared to the performance of ToM-II. Although ToM-I had one more parameter to be optimized compared to ToM-II, it was possible to establish that the improvement of the performance of ToM-I was not uniquely prompted by the extra parameter by comparing the performances of ToM-II and ToM-III: In fact, ToM-III outperformed ToM-II for all combinations of training and test sets in both the training and the test phase, where on average, ToM-III had an error of 5.55% less in estimating the intentions of the participants than ToM-II in the training phase and of 4.16% less in the test phase. Finally, compared to ToM-III, ToM-I showed a better accuracy in estimating the intentions of the participants both in the training and in the test phases, which confirms that including and personalizing the parameter that corresponds to the personality trait of altruism for each participant improves the performance of the model.

Regarding the predictions of the actions that each rational agent took per simulation step, ToM-I made more accurate predictions than ToM-II for all combinations of the training and test sets in both training and test phases. On average, compared to ToM-II, ToM-I was 2.94% more accurate in the training phase and 2.45% more accurate in the test phase. In order to show that the improved performance of ToM-I with respect to ToM-II was not only due to the extra parameter that was identified in ToM-I, but also that the additional cognitive elements that have been added to ToM-I have influenced its performance, the performances of ToM-II and ToM-III were also compared. ToM-III was more accurate than ToM-II in the training phase for all the combinations of training and test scenarios. ToM-III also outperformed ToM-II for three out of the four combinations of training and test scenarios in the test phase. On average, ToM-III was 1.47% more accurate than ToM-II in the training phase and 1.82% more accurate in the test phase. Nonetheless, including and identifying the altruism parameter allows ToM-I to improve

its prediction capabilities, which is shown by the fact that the performance of ToM-I was better than the performance of ToM-III in all cases.

Furthermore, the parameter vector  $\theta$  that was identified per participant has the advantage of providing transparency and an opportunity for more detailed analysis, due to assigning real (rather than abstract) meaning to the elements of this vector, compared to the existing black-box approaches, e.g., machine learning methods. More specifically, for example, the three parameters that are used in all the three models ToM-I, ToM-II, and ToM-III, correspond to the preferences of each participant towards each one of the three entities in the scenarios. Thus, in order to assess how much the identified parameters provided a correct insight about the preferences of the participants, at the end of each experiment, the participants were asked to evaluate in a scale from 0 to 10 how much they cared for each one of the entities. The self-assessed preferences  $\theta_p^s$  of each participant  $p$  were compared to the parameter vector  $\theta_p^*$  that was identified by each model for that participant.<sup>5</sup> However, the scale of the parameters in the parameter vector  $\theta_p^*$  can be different from the scale used by the participants to self-assess their preferences, since the former depends on the fuzzy systems that were defined to represent the cognitive module. This occurs because the choice of the membership functions and the ranges of the variables, as well as the definition of the rules, was done manually and intuitively by the researchers as to represent the relationships between the variables in Fig. 5 in a way that is general (i.e., not personalized) to all rational agents. These design choices can introduce a shift (i.e., an offset) in the range of the parameters, where this offset was naturally corrected during the identification procedure. For example, the set of parameters identified for a participant  $p$  (e.g.,  $\theta_p = [0.9, 0.2, 0.5]$ ) when using a certain fuzzy system to represent the cognitive module, could be slightly different (e.g.,  $\theta_p = [0.8, 0.1, 0.4]$ ) if the fuzzy system had been projected following different design choices (e.g., using different shapes of the membership functions, or a different number of linguistic terms per variable). Nevertheless, when we compare the values of one parameter amongst participants (or of the parameters identified for one participant  $p$ ), it is still possible to obtain relevant insights about the preferences of the participants. For this reason, there may be an offset between the subjective scales used by the participants and the scale of the parameters, which should be corrected. Similarly, the ranges of ToM-I and ToM-II can have different offsets, since the fuzzy rules that are defined for these models are different. Furthermore, the range of the parameters that corresponded to the preferences of the participants for the living entities (i.e., the dependent being and the colleague) was different from the range of the parameter that corresponded to the preference of the

<sup>5</sup>For this purpose, the answers that were given by the participants were divided by a factor of 10, as to be in the same scale as the parameters identified by the models.

participants for the personal belonging. This is because the fuzzy rules that defined the influence of the parameters corresponding to the preference for the living beings over the fast-dynamics state variables were different from the fuzzy rules that described the influence of the parameter corresponding to the preference for the personal belonging over these state variables (the entire rule base of these fuzzy systems has been made available via [73]). For example, the variable “fear for an entity” (emotion) is only considered when the entity is a living being. Thus, there are fuzzy rules that describe the influence of the parameters corresponding to general preferences towards the living beings over emotions, while there are no fuzzy rules describing the influence of the general preference towards the objects over emotions. In practice, for each model, we computed two offsets: the offset of the parameters  $\theta_{j=1}$  and  $\theta_{j=2}$  (for the preferences towards the living beings), and the offset of the parameter  $\theta_{j=3}$  (for the preference towards the objects). The offset of each group of parameters was computed as the average error between the identified parameter  $\theta_{j,p}^*$ , and the self-assessed preference  $\theta_{j,p}^s$  across all the participants and all the parameters of the group. The computation of the offsets is shown in (16).

$$\delta_{1,2} = \frac{1}{2} \cdot \frac{1}{n_p} \sum_{j=1}^2 \sum_{p=1}^{n_p} (\theta_{j,p}^* - \theta_{j,p}^s) \quad (16a)$$

$$\delta_3 = \frac{1}{n_p} \sum_{p=1}^{n_p} (\theta_{3,p}^* - \theta_{3,p}^s) \quad (16b)$$

Afterwards, each of these parameters can be corrected per participant using the offsets, where the corrected parameter is represented by  $\hat{\theta}_{p,j}$  for  $j = 1, 2, 3$ , as it is shown in (17).

$$\hat{\theta}_{j,p} = \theta_{j,p}^* - \delta_{1,2}, \quad j = 1, 2 \quad (17a)$$

$$\hat{\theta}_{3,p} = \theta_{3,p}^* - \delta_3 \quad (17b)$$

The computed offsets for ToM-I and for ToM-II in the contexts of intention estimation and action prediction can be found in Table 4 and Table 5, respectively. Note that the optimal parameters of ToM-I and ToM-III are the same, and consequently the optimal parameters of ToM-III are not represented. Finally, to assess the degree to which the identified parameters corresponded to the self-assessed preferences of each participant  $p$ , an average of the absolute errors between the corrected parameters  $\hat{\theta}_{j,p}$  identified by the model and the self-assessed preferences  $\theta_{j,p}^s$  across all the participants and across the three parameters was computed as it is shown in (18). The errors obtained by ToM-I and ToM-II are also given in Table 4 and Table 5.

$$e_\theta = \frac{1}{3} \cdot \frac{1}{n_p} \sum_{j=1}^3 \sum_{p=1}^{n_p} |\hat{\theta}_{j,p} - \theta_{j,p}^s| \quad (18)$$

The average errors  $e_\theta$  in the parameters that were identified by ToM-I and by ToM-II when estimating the intentions of the participants were 18.4% and 23.4%, respectively. When

**TABLE 4. Parameter offsets and errors between parameters and self-assessed preferences achieved by ToM-I and ToM-II when estimating the intentions of the participants.**

Model	ToM-I			ToM-II		
	$\delta_{\theta_{1,2}}$	$\delta_{\theta_3}$	$e_{\theta}(\%)$	$\delta_{\theta_{1,2}}$	$\delta_{\theta_3}$	$e_{\theta}(\%)$
Combination 1	-0.569	-0.281	17.8	-0.313	-0.270	26.2
Combination 2	-0.462	-0.276	21.3	-0.233	-0.285	17.8
Combination 3	-0.576	-0.300	17.3	-0.370	-0.265	25.2
Combination 4	-0.602	-0.305	17.0	-0.335	-0.350	24.3
Average	-	-	18.4	-	-	23.4

**TABLE 5. Parameter offsets and errors between parameters and self-assessed preferences achieved by ToM-I and ToM-II when predicting the actions of the participants.**

Model	ToM-I			ToM-II		
	$\delta_{\theta_{1,2}}$	$\delta_{\theta_3}$	$e_{\theta}(\%)$	$\delta_{\theta_{1,2}}$	$\delta_{\theta_3}$	$e_{\theta}(\%)$
Combination 1	-0.602	-0.271	24.1	-0.350	-0.315	26.0
Combination 2	-0.669	-0.310	20.3	-0.328	-0.330	24.1
Combination 3	-0.631	-0.319	20.0	-0.443	-0.300	21.9
Combination 4	-0.657	-0.333	21.0	-0.303	-0.365	24.8
Average	-	-	21.4	-	-	24.2

predicting the actions of the participants, the average errors  $e_{\theta}$  were 21.4% for ToM-I and 24.2% for ToM-II. In both cases, the parameters identified by ToM-I are closer to the self-identified preferences of the participants, compared to the parameters that were identified by ToM-II. It is important to mention that an error of 0 would be impossible to reach in practice, since the self-assessment of the participants might not correspond to the ground truth, as they can be biased in the cognitive procedure of assessing their own preferences. Moreover, the interpretation of the scale is subjective and can diverge per participant. Nonetheless, the achieved results show a strong relationship between the corrected identified parameters and the opinions of the participants about their own preferences. We can then conclude that these parameters give a relevant practical insight about the preferences of the participants.

All in all, by assessing the performance of the three ToM-based models when predicting the actions and estimating the intentions of the participants, the experiments acknowledged the benefits of including the emotions, biases, and slow-dynamics state-variables such as personality traits and general preferences in the ToM-based models. They also showed that the model can make dynamic predictions and that the identified parameters provide a realistic insight into the meaning and contribution of each parameter, contrarily to the parameters that are often identified in black-box approaches.

## V. CONCLUSION AND TOPICS FOR FUTURE RESEARCH

The main aim of this paper is to provide interactive machines, e.g., socially assistive robots, with the capability of exhibiting theory of mind in their interactions with humans. Theory of mind states that in their interactions, humans create mental models of each other's cognitive procedures, in order to estimate the mental states of one another. Therefore, we proposed a novel formalization for the perception, cognition, and decision-making procedures of humans using

a network representation. We transformed this network representation into a dynamic formulation by introducing an extended version of FCMs, and generated a state-space representation for the mental procedures of humans, in order to model the theory of mind for interactive machines.

Compared to the state-of-the-art representations of human cognition, which mainly include beliefs, goals, and intentions of humans, we also included emotions, personality traits, and biases in our ToM-based models, which have been shown to be important elements in cognition and decision-making procedures of humans (see, e.g., [2], [7], [19], [38], [49], [50]).

The resulting ToM models were identified and validated in experiments that included computer-based simulations with 21 human participants. We performed several analyses about the performance and accuracy of the resulting models, where the results showed both the success of the model in predicting the intentions and actions of the participants, which are not possible unless the model is dynamic, as well as the accuracy of the model compared to a model that does not include emotions, personality traits, and biases. In the future, a more extensive validation with a larger number of participants and longer-term interaction per participant can be carried out for various contexts of human-machine interaction, so that the model can be validated and used for different real-life applications. Furthermore, the linkages of our mathematical representation can be represented by parameterized polynomial functions rather than by fuzzy inference systems (FIS) in case the training needs to be accomplished faster. In fact, this change of the mathematical representation leads to an improvement of more than 400 times in the computation time. Moreover, for online applications, the identification of the parameters can be done using a gradient descent algorithm instead of a grid search, which can significantly decrease the computational time.

In the future, the proposed cognitive models will be used for interactive robots for two purposes: (1) to estimate the internal mental states of humans, using model-based state estimation methods, and (2) to predict the evolution of the mental states and actions of humans, and let the robot behave accordingly, in order to improve the quality of the human-robot interactions.

## APPENDIX A EXAMPLES THAT ILLUSTRATE THE IMPORTANCE OF THE ELEMENTS OF THE PERCEPTION, COGNITION, AND DECISION-MAKING MODEL

This appendix represents a number of examples according to real-life scenarios, where these examples motivate the elements and/or linkages that are introduced in the proposed cognitive model.

*Role of Slow-Dynamics State Variables:* The following four examples show the influence of the slow-dynamics state variables on the dynamic evolution of fast-dynamics state variables. In these examples, the observer agent (described by first-person pronouns) makes inferences about the mental



state of the observed agent (referred to by a specific name in the examples) based on their observed actions (*inverse inference*).

*Example 1:* Ana and I are both in the library at 5:30 PM. Ana picks up her wallet and walks towards the door (*action of the observed agent noticed by the observer agent*). The coffee house nearby has a late opening hour until 6:00 PM, but I do not know whether or not Ana knows about this (no access to the general world knowledge of the observed agent). I guess Ana knows about the opening hour of the coffee house and believes that it is still open (*guessing the general world knowledge and inferring about the belief of the observed agent*). I infer that Ana is going to buy a cup of coffee (*goal of the observed agent inferred by the observer agent*).

In this example, since the observer agent does not have access to the general world knowledge of the observed agent, the inference involves an intermediate procedure, i.e., inference about the belief of the observed agent based on a guess rather than facts.

Suppose that according to Ana's general world knowledge the opening hour of the nearby coffee house is until 5:00 PM. Thus Ana does not go at 5:30 PM to the coffee house, which she supposes to be closed. In case I had access to Ana's general world knowledge, then I would not conclude that Ana is going to grab a coffee. Therefore, my inference about Ana's goal was more reliable.

*Example 2:* Consider the scenario of Example 1, but this time I have heard from Ana before that she knows that the coffee house nearby is open until 6:00 PM (*access to the general world knowledge of the observed agent*). This time I infer with more certainty that Ana is going to buy a cup of coffee.

Although this example shows that when the general world knowledge of the observed agent is known by the observer agent, an inverse inference about the goals of the observed agent are less prone to uncertainties, the next example shows that knowing the general world knowledge of the observed agent alone may not suffice to make precise inverse inferences.

*Example 3:* Suppose that in Example 2, in addition to being aware of Ana's general world knowledge, I know that she needs to drink coffee when studying late (*access to the general preference of the observed agent*). Then I infer that Ana is going to buy coffee with a much higher certainty than in Example 2. On the contrary, if I know that Ana never drinks coffee in the afternoon (*general preference of the observed agent*) I will not infer that her beliefs and goals are related to grabbing a coffee.

This example shows that the access of an observer agent to the general world knowledge and general preferences of an observed agent significantly improves the reliability and level of certainty of the inferred fast-dynamics state variables. Moreover, having access to only one of these slow-dynamics state variables may still result in inaccurate or erroneous inferences. In particular, having access to the general preferences of the observed agent in addition to the

general world knowledge in Example 3 supports the certainty of the inferences or prevents the observer agent from making erroneous inferences.

*Example 4:* Now consider Example 3, where I am aware that Ana never drinks coffee in the afternoon. While the combination of Ana's general world knowledge and general preference prevents me from making a wrong inference, they do not provide me with a chance either to make an inference about what Ana's beliefs and goals at the moment are. Now suppose that I know Ana for long enough to be aware that she is an introvert (i.e., she needs some personal time after long interactions) with higher levels of neuroticism<sup>6</sup> (i.e., she often feels worried). These personality traits of Ana together with her actions help me to infer that Ana believes she needs some personal time (thus she leaves me to be by herself). I also infer that Ana believes that, if she leaves her wallet unattended, someone may steal it (thus takes her wallet with her).

This example shows the importance of incorporating the specific personality traits of rational agents in cognitive models for achieving more precision and reliability with the resulting inferences.

*Elements that Influence Emotions:* The following three examples illustrate the role of various state variables in triggering the emotions of an observed agent (referred to by a specific name).

*Example 5:* While walking on the street, Elisa's wallet falls out of her purse (*real-life data*). Later on in a shop, Elisa reaches for her wallet and realizes that it is not in her purse (*perceptual access*). She reasons that she has lost the wallet (*rationally perceived knowledge*). She then supposes that she has lost her wallet (*inference of a belief based on the rationally perceived knowledge*). This belief makes her anxious (*stimulation of emotions*).

In the given example, before Elisa notices that her wallet is missing (i.e., *without perceptual access*) and before she reasons that she has lost the wallet (i.e., *without rational reasoning*), she was not anxious (*no stimulation of emotions*). In a different situation, for the same perceptual access that causes the same perceived data, i.e., a missing wallet, Elisa may reason and believe that she has left her wallet on the dining table at home (*different rational reasoning and hence different rationally perceived knowledge*). Therefore, Elisa will not be anxious (*no stimulation of emotions*). In summary, independent of what the real-life data is (e.g., the wallet has fallen on the street or is at home) the emotions of a rational agent may be moderated by the perceptual access of the agent to that data and by the reasoning the agent applies to the perceived data. In other words, the emotions of a rational agent depend on its beliefs rather than on real-life data directly.

*Example 6:* Frank is exploring a new city for the first time and wants to buy an ice cream (*goal*). While walking, he notices a few people across the street who are eating

<sup>6</sup>We use the terms introversion and neuroticism according to the categorization introduced by the Big Five Personality Traits [74].

ice cream (*perceived data*). Correspondingly, he reasons and believes that there should be an ice cream shop close by (*rationally perceived knowledge transformed into a belief*), which makes him feel satisfied (*stimulated emotions*).

This example shows a case where a belief by itself does not stimulate emotions, but the belief together with a goal does. In other words, if Frank did not want to eat an ice cream, the belief that an ice cream shop is nearby would not influence his emotional status.

*Example 7:* Grace is afraid of dogs (*general preference*). While walking in a park, she notices the footprints of a dog (*perceptual access*) and correspondingly reasons and believes that there should be a dog nearby (*rationally perceived knowledge transformed into a belief*). This belief makes her anxious (*stimulated emotion*).

This example illustrates how general preferences alongside beliefs may directly stimulate emotions in rational agents. Note that in this case, if Grace was not afraid of dogs, the belief that a dog is nearby by itself would not stimulate particular emotions in her.

*Elements That are Influenced by Emotions:* The following example illustrates how emotions of rational agents may influence their goals.

*Example 8:* Hailey has planned to go to a party tonight (*original goal*). In the afternoon, Hailey receives some bad news that make her deeply sad (*emotion*). As a consequence, she decides not to go to the party anymore (*change in the goal due to the emotions*).

This example shows that emotions may change the goals of rational agents. Although Hailey's general preference may be to participate in such parties, and she may in general be an extrovert with high levels of conscientiousness (i.e., self-discipline and tendency to follow her schedules), due to her sadness she may make a goal (i.e., skipping the party) that contradicts her initial goal that corresponded to her general preferences, personality traits, etc.

*Importance of Personalizing the Perceptual Access and the Rational Reasoning Processes:* The first example below illustrates the importance of personalizing the perception procedure of various rational agents. In the given example, the tourist is the observed agent and the tour leader and the tourist's close friend are the observer agents. The second example below demonstrates the importance of decomposing the process that yields the rationally perceived knowledge into the sub-processes that are explained in subsection III-B3.

*Example 9:* Suppose that a tourist tells her tour leader that she has already been to the historical city center of the city they are visiting. The tour leader may assume that the tourist has a perfect knowledge of the real-life data, including the location of the church, the old building of the City Hall, and all the souvenir shops (*general world knowledge of the observed agent according to the observer agent*), while in her previous visit, the tourist has overlooked the old building of the City Hall. Then the general world knowledge considered by the tour leader for the tourist is inaccurate. Now suppose

that the tourist tells her close friend (*an observer agent who is aware of the personalized perception of the observed agent*) that she has once been to the historical city center. The friend assumes - knowing the personalized perception procedure of the tourist - that she might have overlooked the old building of the City Hall.

*Example 10:* Brian, Charlie, and Diana are inside a shopping mall. Although they cannot see the outside, before they entered the shopping mall, it was sunny. Someone soaked in water enters the shopping mall. Brian does not notice this person (*no updated perceptual access*) and thus, keeps the *belief* that outside is sunny. Charlie and Diana notice this person (*updated perceptual access*). Charlie reasons and accordingly *believes* that it must be raining now, while Diana reasons that this person has fallen into a ditch (*different [personalized] rational reasoning*) and keeps the *belief* that it is still sunny outside.

Note that in the above example, if an observer agent infers about the beliefs of Brian, Charlie, and Diana (all as observed agents), if their personalized perceptual access and rational reasoning processes are excluded, the observer agent may infer that all three observed agents hold the same belief that, e.g., it is now raining outside.

*Elements That Bias the Generation of Beliefs:* The two following examples show how the generation of beliefs for rational agents may be biased by their emotions and goals, respectively. Igor and Jane (in the first example), and Kevin and his father (in the second example) are the observed agents.

*Example 11:* Igor and Jane are having a walk together. While walking, they both see a dog (*real-life data*, which after perceptual access and rational reasoning, for both observed agents results in the belief that there is a dog nearby). Since Igor is afraid of dogs (*general preference*), he feels afraid (*emotion resulting from emotion trigger 1*, see Remark 3), and starts to believe that the dog might harm him (*belief biased by emotion*). Jane, however, does not feel any fear and thus believes there is no threat from the dog.

In the above example, although both observed agents initially had the same belief (i.e., there is a dog nearby) one of them develops a biased new belief about the dog because of his triggered emotions.

*Example 12:* Kevin is a football fan (*general preference*). The team he supports is currently in the second place in the championship. When all evidences are studied by an objective analyst, they conclude that - although not impossible yet - the chances that Kevin's favorite team wins the championship are very small (*unbiased belief*). Since Kevin wants his team to win (Kevin's *desire* or *goal*), he believes that his team will win (*belief biased by a goal*). Now suppose that Kevin's father has the same general preference and goal as Kevin, while he has a much lower level of conscientiousness (a personality trait that implies Kevin's father is generally less stubborn and more flexible regarding various situations). Consequently, Kevin's father develops a much less strong belief than his son about their team winning the championship.

The above example illustrates that goals can bias the beliefs of rational agents, whereas personality traits - while not generating a bias by themselves - may regulate (boost or hinder) this influence.

*Inverse Inference of Emotions from Actions:* The following two examples illustrate the process of inverse inference of emotions based on the observed actions of rational agents, which is based on updating the belief and the goal of the observed agent (due to the influence of emotions), respectively. The observer agent is specified by the first-person pronouns, and the other person in the examples is the observed agent.

*Example 13:* I see Anabel and smile at her. Anabel and I are friends, so I expect Anabel to smile back (*expected action*). Anabel, however, looks unfriendly instead and turns her back at me (*observed action*). Thus, I conclude that Anabel may feel negative emotions about me at the moment (*inferred emotion*).

In the above example, the observer agent initially estimates the belief of the observed agent to be “we are good to each other”. The observed action of the observed agent, however, implies that this belief is wrong. Therefore, the observer agent updates the belief of the observed agent to “Anabel is not good to me”, and deduces that she holds negative emotions. Now, in order for the observer agent to infer precisely about the emotion of the observed agent (e.g., whether Anabel is angry or sad), the observer agent should be aware of the belief(s) and goal(s) of the observed agent in their previous interaction(s).

Suppose that this morning Anabel showed me a picture of a dress. Her goal was to wear it for her sister’s wedding (Anabel’s goal in the previous interaction). I said the dress may not look nice on her. Anabel believed that I was being mean to her (Anabel’s belief in the previous interaction). This made her feel **angry** at me, and thus in our next interaction she updated her belief from “we are good to each other” to “we are not friends anymore”.

*Example 14:* Last week, I had a good discussion with Lewis about his projects. While walking on the campus, I come across Lewis. I know he has a vacancy for his new project. I tell to Lewis that I am seeking a new project to join. I suppose that Lewis has developed the goal of establishing a collaboration with me after our previous talk, so I expect him to invite me to join his new project (*expected action*). Instead, Lewis wishes me good luck and leaves (*observed action*). Clearly my inference about Lewis’ goal was wrong. In this example, the observer agent expects a particular action from the observed agent based on inferring the observed agent’s belief to be “the observer agent has very high qualifications”. Although this belief is correctly deduced, the goal is wrongly inferred, because the observer agent has not considered the beliefs and goals of the observer agent during their last interaction, i.e., their talk (Lewis believed that he should always be better than his employees and while talking his goal was to show he knows the best), and thus the emotions that were triggered (Lewis feels threatened by the

qualifications of the observer agent and thus develops the goal to avoid collaborating with him).

## REFERENCES

- [1] D. Feil-Seifer and M. J. Mataric, “Defining socially assistive robotics,” in *Proc. 9th Int. Conf. Rehabil. Robot. (ICORR)*, Chicago, IL, USA, Jul. 2005, pp. 465–468, doi: [10.1109/ICORR.2005.1501143](https://doi.org/10.1109/ICORR.2005.1501143).
- [2] A. Tapus and M. J. Mataric, “Socially assistive robots: The link between personality, empathy, physiological signals, and task performance,” in *Proc. AAAI Spring Symp.*, Stanford, CA, USA, Jan. 2008, pp. 133–140.
- [3] C. Clabaugh, K. Mahajan, S. Jain, R. Pakkar, D. Becerra, Z. Shi, E. Deng, R. Lee, G. Ragusa, and M. Mataric, “Long-term personalization of an in-home socially assistive robot for children with autism spectrum disorders,” *Frontiers Robot. AI*, vol. 6, no. 110, pp. 1–18, Nov. 2019, doi: [10.3389/frobt.2019.00110](https://doi.org/10.3389/frobt.2019.00110).
- [4] B. Scassellati, L. Boccanfuso, C.-M. Huang, M. Mademtzi, M. Qin, N. Salomons, P. Ventola, and F. Shic, “Improving social skills in children with ASD using a long-term, in-home social robot,” *Sci. Robot.*, vol. 3, no. 21, Aug. 2018, Art. no. eaat7544, doi: [10.1126/scirobotics.aat7544](https://doi.org/10.1126/scirobotics.aat7544).
- [5] C. D. Kidd and C. Breazeal, “Robots at home: Understanding long-term human-robot interaction,” in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, Nice, France, Sep. 2008, pp. 3230–3235, doi: [10.1109/IROS.2008.4651113](https://doi.org/10.1109/IROS.2008.4651113).
- [6] E. Martinez-Martin and M. Cazorla, “A socially assistive robot for elderly exercise promotion,” *IEEE Access*, vol. 7, pp. 75515–75529, 2019, doi: [10.1109/ACCESS.2019.2921257](https://doi.org/10.1109/ACCESS.2019.2921257).
- [7] I. Leite, C. Martinho, and A. Paiva, “Social robots for long-term interaction: A survey,” *Int. J. Social Robot.*, vol. 5, no. 2, pp. 291–308, Apr. 2013, doi: [10.1007/s12369-013-0178-y](https://doi.org/10.1007/s12369-013-0178-y).
- [8] M. J. Mataric and B. Scassellati, “Socially assistive robotics,” in *Springer Handbook of Robotics*, 2nd ed. Berlin, Germany: Springer, 2016, pp. 1973–1993.
- [9] M. K. Al-Qaderi and A. B. Rad, “A brain-inspired multi-modal perceptual system for social robots: An experimental realization,” *IEEE Access*, vol. 6, pp. 35402–35424, 2018, doi: [10.1109/ACCESS.2018.2851841](https://doi.org/10.1109/ACCESS.2018.2851841).
- [10] A. Tapus, M. J. Mataric, and B. Scassellati, “Socially assistive robotics [grand challenges of robotics],” *IEEE Robot. Automat. Mag.*, vol. 14, no. 1, pp. 35–42, Mar. 2007, doi: [10.1109/MRA.2007.339605](https://doi.org/10.1109/MRA.2007.339605).
- [11] D. Premack and G. Woodruff, “Does the chimpanzee have a theory of mind?” *Behav. Brain Sci.*, vol. 1, no. 4, pp. 515–526, Dec. 1978, doi: [10.1017/S0140525X00076512](https://doi.org/10.1017/S0140525X00076512).
- [12] B. Scassellati, “Theory of mind for a humanoid robot,” *Auto. Robots*, vol. 12, no. 1, pp. 13–24, Jan. 2002, doi: [10.1023/A:1013298507114](https://doi.org/10.1023/A:1013298507114).
- [13] O. Guest and A. E. Martin, “How computational modeling can force theory building in psychological science,” *Perspect. Psychol. Sci.*, vol. 16, no. 4, pp. 789–802, Jan. 2021, doi: [10.1177/1745691620970585](https://doi.org/10.1177/1745691620970585).
- [14] S. Farrell and S. Lewandowsky, *Computational Modeling of Cognition and Behavior*. Cambridge, U.K.: Cambridge Univ. Press, 2018.
- [15] D. Mareschal and M. S. C. Thomas, “Computational modeling in developmental psychology,” *IEEE Trans. Evol. Comput.*, vol. 11, no. 2, pp. 137–150, Apr. 2007, doi: [10.1109/TEVC.2006.890232](https://doi.org/10.1109/TEVC.2006.890232).
- [16] S. Farrell and S. Lewandowsky, *Computational Modeling in Cognition*. New York, NY, USA: SAGE Publications, 2011.
- [17] C. Baker, N. D. Goodman, and J. B. Tenenbaum, “Theory-based social goal inference,” in *Proc. Annu. Meeting Cogn. Sci. Soc.*, Washington, DC, USA, Jan. 2008, pp. 1447–1452.
- [18] N. Rabinowitz, F. Perbet, F. Song, C. Zhang, S. M. A. Eslami, and M. Botvinick, “Machine theory of mind,” in *Proc. 35th Int. Conf. Mach. Learn.*, Stockholm, Sweden, Jul. 2018, pp. 4218–4227.
- [19] J. J. Lee, F. Sha, and C. Breazeal, “A Bayesian theory of mind approach to nonverbal communication,” in *Proc. 14th ACM/IEEE Int. Conf. Human-Robot Interact. (HRI)*, Mar. 2019, pp. 487–496, doi: [10.1109/HRI.2019.8673023](https://doi.org/10.1109/HRI.2019.8673023).
- [20] Y. Zeng, Y. Zhao, T. Zhang, D. Zhao, F. Zhao, and E. Lu, “A brain-inspired model of theory of mind,” *Frontiers Neurobotics*, vol. 14, pp. 1–17, Aug. 2020, doi: [10.3389/fnbot.2020.00060](https://doi.org/10.3389/fnbot.2020.00060).
- [21] L. Dissing and T. Bolander, “Implementing theory of mind on a robot using dynamic epistemic logic,” in *Proc. 29th Int. Joint Conf. Artif. Intell.*, Yokohama, Japan, Jul. 2020, pp. 1615–1621, doi: [10.24963/ijcai.2020/224](https://doi.org/10.24963/ijcai.2020/224).

- [22] M. Patacchiola and A. Cangelosi, "A developmental cognitive architecture for trust and theory of mind in humanoid robots," *IEEE Trans. Cybern.*, vol. 52, no. 3, pp. 1947–1959, Mar. 2022, doi: [10.1109/TCYB.2020.3002892](https://doi.org/10.1109/TCYB.2020.3002892).
- [23] K. Popper, *The Myth of the Framework*. Thames, U.K.: Routledge, 1994.
- [24] D. C. Dennett, *The Intentional Stance*, 6th ed. Cambridge, MA, USA: MIT Press, 1987.
- [25] M. Bratman, *Intention, Plans, and Practical Reason*, 1st ed. Cambridge, MA, USA: Harvard Univ. Press, 1987.
- [26] A. S. Rao, "Decision procedures for propositional linear-time belief-desire-intention logics," in *Proc. Int. Joint Conf. Artif. Intell., Agent Theories, Archit., Lang. Workshop*, Montreal, QC, Canada, Aug. 1995, pp. 33–48.
- [27] T. Bosse, Z. A. Memon, and J. Treur, "A recursive BDI agent model for theory of mind and its applications," *Appl. Artif. Intell.*, vol. 25, no. 1, pp. 1–44, Jan. 2011, doi: [10.1080/08839514.2010.529259](https://doi.org/10.1080/08839514.2010.529259).
- [28] J. Broersen, M. Dastani, J. Hulstijn, Z. Huang, and L. van der Torre, "The BOLD architecture: Conflicts between beliefs, obligations, intentions and desires," in *Proc. 5th Int. Conf. Auto. agents*, New York, NY, USA, May 2001, pp. 9–16, doi: [10.1145/375735.375766](https://doi.org/10.1145/375735.375766).
- [29] P. Vossen, S. Baez, L. Bajčetić, and B. Kraaijeveld, "Leolani: A reference machine with a theory of mind for social communication," in *Proc. Int. Conf. Text, Speech, Dialogue*, Brno, Czech Republic, Sep. 2018, pp. 15–25, doi: [10.1007/978-3-030-00794-2\\_2](https://doi.org/10.1007/978-3-030-00794-2_2).
- [30] C. L. Baker, "Bayesian theory of mind : Modeling human reasoning about beliefs, desires, goals, and social relations," Ph.D. dissertation, Massachusetts Inst. Technol., Cambridge, U.K., 2012.
- [31] K. J. Åström, "Optimal control of Markov processes with incomplete state information," *J. Math. Anal. Appl.*, vol. 10, no. 1, pp. 174–205, Feb. 1965, doi: [10.1016/0022-247X\(65\)90154-X](https://doi.org/10.1016/0022-247X(65)90154-X).
- [32] S. M. Stigler, "Thomas Bayes's Bayesian inference," *J. Roy. Stat. Society. Ser. A, Gen.*, vol. 145, no. 2, p. 250, 1982, doi: [10.2307/2981538](https://doi.org/10.2307/2981538).
- [33] C. Baker, R. Saxe, and J. Tenenbaum, "Bayesian theory of mind: Modeling joint belief-desire attribution," in *Proc. Annu. Meeting Cogn. Sci. Soc.*, Boston, MA, USA, Jul. 2011, pp. 2469–2474.
- [34] J. Block and D. C. Funder, "Social roles and social perception: Individual differences in attribution and error," *J. Personality Social Psychol.*, vol. 51, no. 6, pp. 1200–1207, 1986, doi: [10.1037/0022-3514.51.6.1200](https://doi.org/10.1037/0022-3514.51.6.1200).
- [35] H. A. Witkin, "The nature and importance of individual differences in perception," *J. Personality*, vol. 18, no. 2, pp. 145–170, Dec. 1949, doi: [10.1111/J.1467-6494.1949.TB01237.x](https://doi.org/10.1111/J.1467-6494.1949.TB01237.x).
- [36] D. C. Ong, J. Zaki, and N. D. Goodman, "Computational models of emotion inference in theory of mind: A review and roadmap," *Topics Cognit. Sci.*, vol. 11, no. 2, pp. 338–357, Apr. 2019, doi: [10.1111/tops.12371](https://doi.org/10.1111/tops.12371).
- [37] J. Jara-Ettinger, "Theory of mind as inverse reinforcement learning," *Current Opinion Behav. Sci.*, vol. 29, pp. 105–110, Oct. 2019, doi: [10.1016/J.COBEHA.2019.04.010](https://doi.org/10.1016/J.COBEHA.2019.04.010).
- [38] R. Saxe and S. D. Houlihan, "Formalizing emotion concepts within a Bayesian model of theory of mind," *Current Opinion Psychol.*, vol. 17, pp. 15–21, Oct. 2017, doi: [10.1016/j.copsyc.2017.04.019](https://doi.org/10.1016/j.copsyc.2017.04.019).
- [39] E. B. Sönmez, H. Han, O. Karadeniz, T. Dalyan, and B. Sarioglu, "EMRES: A new EMotional RESpondent robot," *IEEE Trans. Cognit. Develop. Syst.*, vol. 14, no. 2, pp. 772–780, Jun. 2022, doi: [10.1109/TCDS.2021.3120562](https://doi.org/10.1109/TCDS.2021.3120562).
- [40] P. Smaldino, *Modeling Social Behavior: Mathematical and Agent-Based Models of Social Dynamics and Cultural Evolution*, 1st ed. Princeton, NJ, USA: Princeton Univ. Press, 2023.
- [41] I. Ajzen, "The theory of planned behavior," *Organizational Behav. Hum. Decis. Processes*, vol. 50, no. 2, pp. 179–211, Dec. 1991, doi: [10.1016/0749-5978\(91\)90020-T](https://doi.org/10.1016/0749-5978(91)90020-T).
- [42] J.-E. Navarro-Barrientos, D. E. Rivera, and L. M. Collins, "A dynamical model for describing behavioural interventions for weight loss and body composition change," *Math. Comput. Model. Dyn. Syst.*, vol. 17, no. 2, pp. 183–203, Mar. 2011, doi: [10.1080/13873954.2010.520409](https://doi.org/10.1080/13873954.2010.520409).
- [43] J. E. Laird, C. Lebiere, and P. S. Rosenbloom, "A standard model of the mind: Toward a common computational framework across artificial intelligence, cognitive science, neuroscience, and robotics," *AI Mag.*, vol. 38, no. 4, pp. 13–26, Dec. 2017, doi: [10.1609/AIMAG.V38I4.2744](https://doi.org/10.1609/AIMAG.V38I4.2744).
- [44] A. Chella and A. Pipitone, "A cognitive architecture for inner speech," *Cognit. Syst. Res.*, vol. 59, pp. 287–292, Jan. 2020, doi: [10.1016/j.cogsys.2019.09.010](https://doi.org/10.1016/j.cogsys.2019.09.010).
- [45] K. Ogata, *Modern Control Engineering*, 4th ed. New Jersey, NJ, USA: Prentice-Hall, 2002.
- [46] J. Matz, A. Birouche, B. Mourllion, F. Bouziani, and M. Basset, "Parameter identification for nonlinear models from a state-space approach," *IFAC-PapersOnLine*, vol. 53, no. 2, pp. 13910–13915, 2020, doi: [10.1016/j.ifacol.2020.12.905](https://doi.org/10.1016/j.ifacol.2020.12.905).
- [47] S. H. Zak, *Systems and Control*, vol. 198. New York, NY, USA: Oxford Univ. Press, 2003.
- [48] M. Snyder and W. Ickes, "Personality and social behavior," in *Handbook of Social Psychology*, vol. 2, 3rd ed. New York, NY, USA: Random House, 1985, pp. 883–947.
- [49] D.-S. Kwon, M. J. Chung, J. C. Park, C. D. Yoo, E.-S. Jee, K.-S. Park, Y.-M. Kim, H.-R. Kim, J.-C. Park, H.-J. Min, J. W. Park, S. Yun, and K.-W. Lee, "Emotional exchange of a socially interactive robot," in *Proc. 17th IFAC World Congr.*, Seoul, South Korea, Jul. 2008, pp. 4330–4335, doi: [10.3182/20080706-5-kr-1001.00729](https://doi.org/10.3182/20080706-5-kr-1001.00729).
- [50] J. E. Bono and M. A. Vey, "Personality and emotional performance: Extraversion, neuroticism, and self-monitoring," *J. Occupational Health Psychol.*, vol. 12, no. 2, pp. 177–192, 2007, doi: [10.1037/1076-8998.12.2.177](https://doi.org/10.1037/1076-8998.12.2.177).
- [51] J. Heredia, E. Lopes-Silva, Y. Cardinale, J. Diaz-Amado, I. Dongo, W. Graterol, and A. Aguilera, "Adaptive multimodal emotion detection architecture for social robots," *IEEE Access*, vol. 10, pp. 20727–20744, 2022, doi: [10.1109/ACCESS.2022.3149214](https://doi.org/10.1109/ACCESS.2022.3149214).
- [52] S. Baron-Cohen, A. M. Leslie, and U. Frith, "Does the autistic child have a 'theory of mind?'" *Cognition*, vol. 21, no. 1, pp. 37–46, Oct. 1985, doi: [10.1016/0010-0277\(85\)90022-8](https://doi.org/10.1016/0010-0277(85)90022-8).
- [53] H. M. Wellman, D. Cross, and J. Watson, "Meta-analysis of theory-of-mind development: The truth about false belief," *Child Develop.*, vol. 72, no. 3, pp. 655–684, May 2001, doi: [10.1111/1467-8624.00304](https://doi.org/10.1111/1467-8624.00304).
- [54] A. Fumero, C. Santamaría, and P. Johnson-Laird, "The effect of personality on reasoning," *Nature Precedings*, vol. 3, pp. 1–9, Jul. 2008, doi: [10.1038/NPRE.2008.2099.1](https://doi.org/10.1038/NPRE.2008.2099.1).
- [55] J. Zelenski, "The role of personality in emotion, judgment, and decision making," in *Do Emotions Help or Hurt Decisionmaking? A Hedgefoxian Perspective*. New York, NY, USA: Russell Sage Foundation, 2007, pp. 117–132.
- [56] K. R. Scherer, "The dynamic architecture of emotion: Evidence for the component process model," *Cognition Emotion*, vol. 23, no. 7, pp. 1307–1351, Nov. 2009, doi: [10.1080/02699930902928969](https://doi.org/10.1080/02699930902928969).
- [57] P. Ekman and D. Cordaro, "What is meant by calling emotions basic," *Emotion Rev.*, vol. 3, no. 4, pp. 364–370, Sep. 2011, doi: [10.1177/1754073911410740](https://doi.org/10.1177/1754073911410740).
- [58] R. Raghunathan and M. T. Pham, "All negative moods are not equal: Motivational influences of anxiety and sadness on decision making," *Organizational Behav. Human Decis. Processes*, vol. 79, no. 1, pp. 56–77, Jul. 1999, doi: [10.1006/obhd.1999.2838](https://doi.org/10.1006/obhd.1999.2838).
- [59] E. B. Andrade and D. Ariely, "The enduring impact of transient emotions on decision making," *Organizational Behav. Human Decis. Processes*, vol. 109, no. 1, pp. 1–8, May 2009, doi: [10.1016/j.obhdp.2009.02.003](https://doi.org/10.1016/j.obhdp.2009.02.003).
- [60] J. S. Lerner, Y. Li, P. Valdesolo, and K. S. Kassam, "Emotion and decision making," *Annu. Rev. Psychol.*, vol. 66, no. 1, pp. 799–823, Jan. 2015, doi: [10.1146/annurev-psych-010213-115043](https://doi.org/10.1146/annurev-psych-010213-115043).
- [61] J. M. George and E. Dane, "Affect, emotion, and decision making," *Organizational Behav. Human Decis. Processes*, vol. 136, pp. 47–55, Sep. 2016, doi: [10.1016/j.obhdp.2016.06.004](https://doi.org/10.1016/j.obhdp.2016.06.004).
- [62] A. R. Mele, "Real self-deception," *Behav. Brain Sci.*, vol. 20, no. 1, pp. 91–136, Mar. 1997, doi: [10.1017/S0140525X97420034](https://doi.org/10.1017/S0140525X97420034).
- [63] T. Canli, Z. Zhao, J. E. Desmond, E. Kang, J. Gross, and J. D. E. Gabrieli, "An fMRI study of personality influences on brain reactivity to emotional stimuli," *Behav. Neurosci.*, vol. 115, no. 1, pp. 33–42, 2001, doi: [10.1037/0735-7044.115.1.33](https://doi.org/10.1037/0735-7044.115.1.33).
- [64] C. L. Baker, J. Jara-Ettinger, R. Saxe, and J. B. Tenenbaum, "Rational quantitative attribution of beliefs, desires and percepts in human mentalizing," *Nature Hum. Behaviour*, vol. 1, no. 4, pp. 1–10, Mar. 2017, doi: [10.1038/s41562-017-0064](https://doi.org/10.1038/s41562-017-0064).
- [65] J. Jara-Ettinger, H. Gweon, L. E. Schulz, and J. B. Tenenbaum, "The Naïve utility calculus: Computational principles underlying commonsense psychology," *Trends Cognit. Sci.*, vol. 20, no. 8, pp. 589–604, Aug. 2016, doi: [10.1016/j.tics.2016.05.011](https://doi.org/10.1016/j.tics.2016.05.011).

- [66] D. Heckerman, "A tutorial on learning with Bayesian networks," in *Innovations in Bayesian Networks: Theory and Applications*, vol. 156, 1st ed. Berlin, Germany: Springer, 2008, pp. 33–82.
- [67] B. Kosko, "Fuzzy cognitive maps," *Int. J. Man-Mach. Stud.*, vol. 24, no. 1, pp. 65–75, Jan. 1986, doi: [10.1016/S0020-7373\(86\)80040-2](https://doi.org/10.1016/S0020-7373(86)80040-2).
- [68] C. D. Stylios and P. P. Groumos, "Modeling complex systems using fuzzy cognitive maps," *IEEE Trans. Syst., Man, Cybern. A, Syst. Humans*, vol. 34, no. 1, pp. 155–162, Jan. 2004, doi: [10.1109/TSMCA.2003.818878](https://doi.org/10.1109/TSMCA.2003.818878).
- [69] L. A. Zadeh, "The concept of a linguistic variable and its application to approximate reasoning-I," *Inf. Sci.*, vol. 8, no. 3, pp. 199–249, 1975, doi: [10.1016/0020-0255\(75\)90036-5](https://doi.org/10.1016/0020-0255(75)90036-5).
- [70] L. A. Zadeh, "Fuzzy logic," *Computer*, vol. 21, no. 4, pp. 83–93, Apr. 1988, doi: [10.1109/2.53](https://doi.org/10.1109/2.53).
- [71] J. P. Carvalho and J. A. B. Tome, "Rule based fuzzy cognitive maps—Expressing time in qualitative system dynamics," in *Proc. 10th IEEE Int. Conf. Fuzzy Syst.*, Melbourne, VIC, Australia, Dec. 2001, pp. 280–283, doi: [10.1109/FUZZ.2001.1007303](https://doi.org/10.1109/FUZZ.2001.1007303).
- [72] A. Mourhir, T. Rachidi, E. I. Papageorgiou, M. Karim, and F. S. Alaoui, "A cognitive map framework to support integrated environmental assessment," *Environ. Model. Softw.*, vol. 77, pp. 81–94, Mar. 2016, doi: [10.1016/J.ENVSOF.2015.11.018](https://doi.org/10.1016/J.ENVSOF.2015.11.018).
- [73] A. Jamshidnejad and M. M. Patrício. (2023). *Supplementary Data to the Paper: Dynamic Mathematical Models of Theory of Mind for Socially Assistive Robots*. [Online]. Available: <https://doi.org/10.4121/0cc4d150-1c9d-4b54-9e75-fbe7dd3efacf.v1>
- [74] O. John, L. P. Naumann, and C. J. Soto, "Paradigm shift to the integrative big five trait taxonomy: History, measurement, and conceptual issues," in *Handbook of Personality: Theory and Research*. New York, NY, USA: Guilford Press, 2008, pp. 114–158.



**MARIA L. M. PATRÍCIO** received the B.Sc. degree in aerospace engineering from Instituto Superior Técnico, Portugal, in 2019, and the M.Sc. degree (cum laude) in control and operations from the Faculty of Aerospace Engineering, Delft University of Technology, The Netherlands, in 2021, where she is currently pursuing the Ph.D. degree. Her main research interests include AI-based control, autonomous control, and cognitive robotics.



**ANAHITA JAMSHIDNEJAD** received the Ph.D. degree (cum laude) from the Delft University of Technology, The Netherlands, in 2017. She is currently an Assistant Professor of mathematical decision-making with the Department of Control and Operations, Delft University of Technology. She is also the Founder and the Co-Director of the AI\*MAN Laboratory, Delft University of Technology. Her main research interests include optimization theory in engineering problems, AI-based control, integrated control methods, and real-time model predictive control, with applications in autonomous and cognitive socially assistive robots.

...