

Evaluating Empathy in Artificial Agents

Özge Nilay Yalçın

School of Interactive Arts and Technology

Simon Fraser University

Vancouver, BC, CANADA

oyalcin@sfu.ca

Abstract—The novel research area of computational empathy is in its infancy and moving towards developing methods and standards. One major problem is the lack of agreement on the evaluation of empathy in artificial interactive systems. Even though the existence of well-established methods from psychology, psychiatry and neuroscience, the translation between these methods and computational empathy is not straightforward. It requires a collective effort to develop metrics that are more suitable for interactive artificial agents. This paper is aimed as an attempt to initiate the dialogue on this important problem. We examine the evaluation methods for empathy in humans and provide suggestions for the development of better metrics to evaluate empathy in artificial agents. We acknowledge the difficulty of arriving at a single solution in a vast variety of interactive systems and propose a set of systematic approaches that can be used with a variety of applications and systems.

Index Terms—Empathy, Affective Computing, Interactive Agents, Human-Computer Interaction, Evaluation Methods

I. INTRODUCTION

Emerging technologies continue to change the ways in which we interact with computers. Computational systems are evolving from being mere tools to assistants, trainers and companion agents. All of these new roles assigned to these systems highlight the importance of embodying these agents with social and emotional capabilities. The advances in computational interaction techniques allowed for the development of emotionally sensitive, perceptive, socially situated and expressive agents. One of the novel and exciting addition to these behaviors is empathy, as a complex socio-emotional behavior.

Empathy can be defined as the capacity to perceive, understand and respond to others' emotions in a manner that is more suitable to those perceived emotions than one's own [1]. The long history of empathy research with the contribution of many disciplines (e.g. philosophy, psychology, neuroscience, ethology) resulted in a diverse set of definitions and behaviors assigned to empathy (see [2] for the history of the field). Behaviors such as mimicry, affective matching, consolidation and perspective-taking are assigned to empathic ability in humans [3], which are crucial to initiating and maintaining social relationships. Following the centuries-old research that shows the importance of empathy in social interactions, computational empathy emerged as a novel field to equip artificial agents to show empathic behavior during their interactions.

With the recent developments in the capabilities of computational systems, the computational modeling of empathy

has gained increasing attention in the last decade. Research on computational modeling of empathy have shown that empathic capacity in interactive agents lead to more trust [4], [5], increase the length of interaction [6], help coping with stress and frustration [7] and increase engagement [5] (see [8] for a review of the field). These findings suggest that agents with empathy could enhance the social interaction in educational applications, artificial companions, medical assistants and gaming applications. Equipping artificial social agents with empathic capabilities is, therefore, a crucial and yet challenging problem.

Part of this challenge arises from the lack of evaluation methods to measure empathic behavior in artificial agents, which would allow for a systematic assessment of the steps need to be taken to model the components of this complex socio-emotional phenomenon. Although various fields provide well-established evaluation methods to measure empathy in humans, it is not clear how to translate these methods to evaluate computational systems. Empathy research in psychology provides validated methods to measure empathy levels in a person [9]–[11]. These methods often require a first-person report of certain behavioral traits based on subjective questionnaires. This poses a challenge for artificial empathy work in virtual agents, as subjective measurements cannot be used in artificial entities [8]. Moreover, behavioral tests from psychology and neuroscience often rely on physiological signals such as neural activity, heart rate and skin conductance, which cannot be used for machines. Other methods may require the observation of experts in the field, which is hard to automate.

Furthermore, agents that have varying levels of interactive capabilities and application goals can have different effects on the perception of the agent and interaction. The characteristics of the agent (e.g. aesthetics, embodiment) as well as non-functional properties (e.g. fluency, response time) can affect the evaluation of empathy, as much as the empathic functionality. A conversational agent in a text-only environment such as chatbots would require emotion perception and expression using different modalities than an embodied conversational agent. In that sense, every additional capability of the agent would contribute to the evaluation of the interaction and the system. Similarly, the application areas may enforce varying levels of expectations in terms of empathic behavior. A medical assistant may require more sympathy, where a personal trainer would focus on pushing boundaries of the

interaction partner. This diversity highlights the importance of a set of evaluation methods that allow for flexibility, instead of focusing on reaching to a single solution.

This paper aims to provide recommendations on how to systematically evaluate empathy in artificial agents in a variety of contexts and capabilities. We propose to approach this problem by focusing on best practices in both the empathy research in humans and the human-computer interaction (HCI) research. To achieve this, we will examine the methods of evaluation in empathy of humans to draw conclusions on how to adapt the key concepts to computational empathy. We propose to use system-level and feature-level evaluations to systematically list the factors that contribute to the evaluation of empathy. By providing a checklist of these factors, we aim to initiate the discussion towards creating a common ground on the evaluation of empathy in artificial agents.

In the following sections, we will explore the evaluation methods on empathy in humans (Section II) and try to apply this know-how into evaluating empathy in artificial agents using system-level and feature-level evaluations (Section III). We will conclude with a discussion of challenges and a call for collective action to work towards new evaluation methods in this new and exciting field.

II. EVALUATION OF EMPATHY IN HUMANS

Empathy research from many disciplines have developed definitions and models for empathy that resulted in a variety of capabilities assigned to empathic behavior [3]. Capabilities such as mimicry, affective matching (emotional contagion), sympathy (empathic concern), altruistic helping (consolidation) and perspective taking are assigned to empathic behavior by scholars [2], [12]. How many of these behaviors should constitute empathic behavior and how they are connected are still highly debated topics in empathy research. Following the differences in the definitions and models of empathy, the evaluation metrics developed to measure empathic behavior may vary dramatically.

Some definitions of empathy separate a number of these capabilities as affective and cognitive empathy [13]. According to this view, affective empathy refers to the relatively automatic emotional responses to other's emotions. Behaviors related to affective empathy can be listed as mimicry, affective matching and empathic concern [14]. On the other hand, cognitive empathy includes behaviors that require the understanding of another's emotional state and behaviors. Behaviors such as consolidation, perspective-taking and altruistic helping are said to originate from the involvement of cognitive mechanisms during the processing of the other's emotional situation [12].

Contrary this dual view of empathy that separates affective and cognitive processes, a unifying view of empathy is gaining attention as an alternative. These recent models and definitions of empathy suggest a more multi-dimensional approach where both affective and cognitive empathy are interconnected with a variety of processes that results in individual differences [1], [15], [16]. One of the most prominent views of empathy

called the Russian Doll Model of Empathy [17] suggests hierarchical levels of affective and cognitive capabilities are connected through evolutionary mechanisms. According to this model, processes such as emotional communication capabilities (recognition and expression), emotion regulation, appraisal processes and theory of mind are considered to be the foundational mechanisms that allow the levels of affective and cognitive empathic behavior [12]. It is also suggested that the individual differences between the empathic responses of people are related to the factors that affect the outcome of these processes [15]. For instance, the recognition of emotions depends on the intensity of the perceived emotion, where the regulation of the emotion would depend on the familiarity between individuals as well as the features of the observer (mood, personality) [18].

The evaluation metrics to measure empathic capacity in humans varies depending on the definition and the capabilities that are assigned to empathic behavior. Evaluation metrics that follow the categorical view of empathy usually assess the affective or the cognitive aspects of empathic behavior as separate constructs. On the other hand, evaluation metrics that follow the multi-dimensional view focuses on defining and evaluating the levels of processes and behaviors that determine the extent of empathic behavior. Although both of these approaches received criticism by some researchers that suggest empathy should only involve the higher level processes [3], it is useful to focus on a broader view of empathy for to arrive at a comprehensive framework to evaluate empathy in artificial agents. Following this notion, we will explore the evaluation methods by focusing their adaptability to computational empathy research.

In this section, we will give an overview of the well-established measurements for empathic behavior, while categorizing them in terms of the focus of empathic behavior and the method of delivery. Evaluation metrics can target a variety of levels of empathic behavior with different levels of granularity and abstraction.

Some of the evaluation methods are designed to measure empathy as a single comprehensive construct and aimed to derive a single value that would indicate the global empathy. Others focus on multiple features or a subset of those features that underlie empathic behavior, such as the affective and cognitive capabilities that are mentioned earlier. Therefore, we will categorize the evaluation metrics according to the level of granularity they are aimed to quantify: global empathy and components of empathy.

Moreover, the method of delivery for these evaluations can be categorized as self-report, observational and physiological approaches. Physiological approaches include measurements of brain activity or autonomic nervous system measures (heart rate, skin conductance, breathing rate). Methodologically, there is no direct way of applying the physiological approaches to virtual agents. Therefore, we will not cover them in our paper (see [19] for a review of these approaches). Self-report measures usually include surveys/questionnaires that rely on the individual's assessment of their behavior. Observational

methods can include the behavioral tests and perceived empathy measures. Behavioral methods rely on performance tests based on experimental stimuli. These methods are often used to assess the components of empathy in humans and aimed to indicate deficiencies. Lastly, perceived empathy metrics are questionnaires that require an observer's assessment of an individual's behaviors. These can include expert observations on the subject's behaviors as well as a second or third person account of a non-expert. In the following sub-sections, we will give detailed examples on the well-established evaluation metrics on each of these methods within global and component-based evaluation metrics.

A. Evaluating Global Empathy

Measures of global empathy are aimed at quantifying a single value that would indicate the strength of empathic capability as a broader concept. Many researchers have attempted to develop self-report measures of empathy based on the definitions and capabilities they assigned to the term. Most of these methods focus on the evaluation of empathy as a whole, while others focus on the specific factors that add up to the global empathic behavior.

One of the earliest self-measure of empathy is Hogan's Empathy Scale (ES) [20] that is mostly used to assess cognitive empathy with 64 true-false statements taken from the standard psychological scales. This questionnaire was intended to examine the relation of empathy with moral and socially appropriate behavior. It was criticized by later works that it is better suited for the evaluation of social skills in a broader sense rather than a specification of empathy [9], [10]. Moreover, the low scores of the validity and reliability of the scale resulted in a continuous decrease in the use of this scale as a valid measure of empathy [21]. However, this attempt encouraged researchers to investigate further and examine the development of a more suitable evaluation of empathy.

A frequently used example is Davis's Interpersonal Reactivity Index (IRI) [9], is a 28-item scale for multi-dimensional measurement of empathy with four sub-scales: perspective-taking, empathic-concern, fantasy and personal distress. However, there have been some discussions around the appropriateness of this scale to measure empathy. Firstly, it was argued that the questionnaire may capture behaviors broader than empathy [10], [22], such as imagination (e.g. item 1 "I daydream and fantasize, with some regularity, about things that might happen to me") and emotional control (e.g. item 10 "I sometimes feel helpless when I am in the middle of a very emotional situation"). An adaptation of IRI to exclude the "fantasy" subscale was later adopted as Feeling and Thinking Scale [23]. It was also suggested that the "personal distress" sub-scale mostly measures anxiety towards distressing situations in general and does not relate to the core functions of empathy. Moreover, some researchers suggested the further refinement of this scale due to the correlation between these sub-scales [22].

The Empathy Quotient (EQ) [10] is one of the most accepted self-report scales that is validated by numerous

studies [11]. Authors define empathy as "the drive to identify another person's emotions and thoughts and to respond to these with appropriate emotion" (p.361). This test is aimed as a clinical screening tool for adults with Autism Spectrum Disorders. In contrast with other self-report questionnaires of empathy, authors did not differentiate between affective and cognitive empathy. They aim to capture empathy in a broader sense where both levels have very interrelated capacities. This questionnaire includes 60 items with 40 empathy-related and 20 filler questions answered with a 4-point Likert scale that scores agreement with the statements. Example questions from the EQ are "I am good at predicting how someone will feel" and "Seeing people cry doesn't really upset me." The questionnaire scores are shown to correlate with autism and gender differences [10].

A recent attempt to further examine and combine these self-report measures uses factor-analysis to reach to a brief and reliable measurement of empathy is called The Toronto Empathy Questionnaire [22] (TEQ). Authors gathered a total of 142 items from several self-report empathy questionnaires such as IRI, ES, BEES, QMEE, AQ as well as empathy questionnaires for specific populations such as Jefferson Scale of Physician Empathy [24], Nursing Empathy Scale [25] and Japanese Adolescent Empathy Scale [26]. Authors used these items to refine a final set of 16 items that are found to be most correlated with Empathy scores compared to other questionnaires. Responses are made with 5-point Likert scale items that show agreeableness of the statements. This questionnaire is a shorter alternative to the EQ with high internal consistency, validity and reliability scores.

A similar approach is taken in the Questionnaire of Cognitive and Affective Empathy (QCAE) [27], which is derived from EQ, ES, IRI and the Impulsiveness-Venturesomeness-Empathy Inventory [28]. Authors finalized a 31-item questionnaire that measures cognitive and affective empathy, as the name suggests. The QCAE consists of five sub-scales, where two sub-scales are related to cognitive empathy (perspective-taking, online simulation), and three of them are related to affective empathy (emotion contagion, proximal responsivity and peripheral responsivity).

Other methods focus on the evaluation of the perception of empathic behavior. These measurements provide a second and third person perspective on an individual's empathy with questionnaires. Jefferson Scale of Physician Empathy is developed to evaluate empathy as a predominantly cognitive attribute [24]. These components are "communication", "understanding" and "cognition", are focused on the cognitive empathy, rather than the affective empathy that was mentioned earlier (see Section II). The scale consists of 20 items with a 7-point Likert type scale on agreement with the statements. Another example for the perceived empathy methods is the Consultation and Relational Empathy (CARE) Questionnaire [29]. CARE was developed to measure "relational empathy" that focuses on the social function of empathy. It consists of 10 statements that start with "How was the doctor at ..." and scored by the patient by using a 5-point likert scale from

“poor” to “excellent”. The items include “Fully understanding your concerns”, “Showing care and compassion” and “Being positive”. Some of these items show significant overlap with self-report questionnaires such as the IRI.

B. Evaluating Components of Empathy

An alternative approach to the evaluation of global empathy is the evaluation of specific components that are required for empathic behavior. These evaluations are usually done by testing behavioral and cognitive abilities to detect the deficits and abnormalities in the behavior. Components such as emotional communication (recognition and expression), emotion regulation, appraisal and perspective taking are usually targeted in these evaluations as critical mechanisms for levels of empathy.

The “reading the mind in the eyes” test [30] was one of the first examples of these behavioral tests. This test was aimed to be used as a screening test for adults or children with Aspergers Syndrome, who are considered to have a deficit in empathy. The revised version of this test consists of 36 photographs of the eye-region of the face that shows different emotional expressions [31]. The participants are presented with these photographs and are asked the most appropriate word to describe “what the person in the photograph is thinking or feeling” (p.241) among four words that are presented. The target terms include words that show mental states such as “thoughtful”, “interested” or “fantasizing”, as well as words that relate to the emotional state such as “upset”, “nervous” or “hostile”. The test results indicate the ability to perceive social and emotional cues where a lower score is associated with a broader set of phenomena than just measuring empathy.

Similarly “reading the mind in the voice” test [32] and “reading the mind in films” test [33] are aimed to measure the ability to detect socio-emotional cues in voice and movie stimuli respectively. The “voice” test uses segments of dialogue taken from dramatic performances, where the “films” test uses audio-visual recordings from movies that shows complex situations. These behavioral tests target the Theory of Mind (ToM), that is the ability to attribute mental states (beliefs, desires, intentions and emotions) to others that are distinct to one’s own. ToM being a crucial part of cognitive empathy, these simple tests allow to spot deficiencies while targeting necessary perceptual abilities.

Understanding appraisals [34], [35] and intuitive physics [36] can also be used to determine the capacity to understand cause and effect relationships, which is based on the higher level cognitive mechanisms. The Picture-Stories task [34] consists of a series of pictures that show the cause and effect relationship in social situations when appropriately sequenced. Similarly, the Social Stories Questionnaire (SSQ) [35] consists of 10 short stories that may involve situations where one character could upset the other character in the story. Participants are asked to whether a selected utterance from the story contains an upsetting utterance and whether the behavior of one character could have upset the other character.

The number of correct answers defines the SSQ score in this test.

III. EVALUATION OF EMPATHY IN INTERACTIVE AGENTS

Being a novel field, empathy studies in artificial intelligence (AI) has no strong standardization and validated methods to measure empathy in artificial agents. In the previous section, we laid out some of the most accepted evaluation methods to evaluate empathy in humans. Although these methods are well-established and agreed upon in the academic community, applying them in the context of artificial agents is not straightforward. Most of these tools rely on self-measurement which cannot be applied to computational systems or the assessment of an expert that is difficult to automatize. Moreover, the differences in the capabilities of the agents and the application context restrict the usage of general behavioral measurements. These issues made it challenging to use this know-how to the evaluation of empathy in artificial interactive agents.

Empathy measurements in psychology literature include the evaluation of specific cognitive and behavioral capabilities as well as an overall evaluation of empathy. Specific features include evaluations of emotion recognition [31], perspective taking [15] and empathic concern [9]. Understanding the user’s emotion depends on the correct recognition of the facial expressions and the performance of the emotion classifier. The perception of empathic behavior depends on the successful expression of the intended empathic emotion. Overall evaluation of empathy should take these feature’s performance along with the system-level evaluation of empathy.

Similarly, the performance of computational systems highly depends on the performance and accuracy of the individual components as well as the integration of these components at the system-level. Due to the complexity and multi-component nature of interactive agents, scholars suggested [37], [38] to provide feature-level and system-level evaluations separately. System-level evaluations focus on the behavior of the agent as a whole, where feature-level evaluations are aimed to isolate individual components of the system separately.

Following this notion, we propose to combine the best practices in the HCI research with the traditional methods of evaluating empathy. In the following sub-sections, we will focus on how to evaluate empathy using system-level and feature-level evaluation methods. We will systematically list the factors that contribute to the evaluation of empathy in artificial agents to initiate the discussion towards creating a common ground.

A. System-Level Evaluation

System-level evaluation is focused on the measurement of the behavior of the system in a broader sense. Similar to the self-report and perceived empathy evaluations that are aimed at capturing the global empathic behavior, system-level evaluations in artificial agents focus on the overall perception of empathy of the agent. In these type of evaluations, the participants interact with the complete system according to the interaction context, and a set of subjective and objective

evaluations are used to compare the different versions of the system or with human behavior.

Previous studies in artificial empathy often focus on the second person or third person perception of empathy of the systems by using empathy-related terminology such as “feeling with” [39], “feeling with” [40], “emotion matching” [41], “compassionate” [42] or “caring” [4]. However, these terms only focus on one specific aspect of empathic behavior or related constructs. An interesting approach was used to train and evaluate the CARE framework [41] by comparing the behavior of the agent with human behavior in a goal-directed environment. This approach can be automated but requires additional data-collection and evaluation steps of human behavior in a similar context to allow for a direct comparison.

As computational empathy research gaining more attention, researchers are beginning to raise awareness on the importance of using more suitable metrics. Recently scholars [8] suggested using a variation of the IRI questionnaire [9] by adapting the first-person evaluation to a perceived-empathy survey. This idea was applied as a part of the EMOTE project [?] authors assessed the perceived empathy of a social robot using the IRI questionnaire. Similarly, Toronto empathy questionnaire [22] was used as a perceived empathy metric by converting the self-report questionnaire into a second or third person evaluation [43]. These evaluation methods can be used to evaluate the system by the interaction partner using a questionnaire. Moreover, the evaluation can be done by a third-person after watching the live or recorded interaction between the system and a participant. Although these methods provide an evaluation that is aligned with the related research on empathy, they were not validated.

However, these perceptual evaluations of empathy can be affected by several factors that should be taken into consideration while applying these system-level evaluations. These factors can be categorized as user-related factors, context-related factors and system-related factors.

1) *User-related Factors*: Research on empathy shown that humans empathize with each other on different levels depending on factors such as their gender, mood, personality, similarity and social capabilities [15], [16], [18]. These findings highlight the importance of controlling for these factors in a comparative evaluation of the agent behavior. Moreover, individual traits such as culture, socio-economic background and computer experience might affect the evaluation of the system as an interactive tool [44].

2) *Context-related Factors*: Relationship and context related factors would impact the strength and expression of empathic behavior. The context, the appraisal of the situation or the social role of the empathizer are suggested to influence the regulation of emotions [13]. Systems that act as companions as opposed to trainers are expected to be more friendly. This user expectancy based on the role of the agent and the context can effect the perception of empathy, where people tend to be more empathic towards in-group members such as friends and family members [16], [18]. Moreover, goal-directed factors

that show the quality of experience such as effectiveness, efficiency, user-satisfaction, utility and acceptability can influence the overall perception of the system [38].

3) *System-related Factors*: Factors related to the system behavior that are not directly linked to its empathic capacity can also impact the evaluation of empathy. Studies have shown that aesthetic characteristics of the interaction partner have a dramatic influence on the perception of empathy in humans [45]. These aesthetic considerations might translate into the factors related to the looks, human-likeness, fluency of movement and believability of the agents [46], [47]. HCI research has developed evaluation metrics to control the effect of these factors such as anthropomorphism, animacy, likability, perceived intelligence and perceived safety [48].

Computational empathy research has already been measuring some of these factors as control variables as well as additional metrics for the overall success of their system [5], [39], [42]. Although the effects of the factors related to empathic behavior are examined in detail in empathy research, the relationship between these factors and the perception of empathy is yet to be examined.

B. Feature-Level Evaluation

In addition to the system-level evaluation, the evaluation of individual aspects of the system is necessary to assess the empathic capabilities of an interactive system. Feature-level evaluations can provide an incremental assessment of each component and capability of the agent. This allows for capturing the propagation of errors in empathic behavior, similar to the behavioral evaluations in empathy research that focuses on capturing deficits in empathic capacity.

In complex interactive systems, the evaluation methodologies usually include the metrics from various sub-fields, such as speech recognition, emotion recognition and speech synthesis [38]. The performance of the implementation is depended on the success of the separate features of the system, as each component affects the evaluation of other components. Therefore, the deficits in one capability might drastically influence the other. For example, the appraisal mechanism could be effected by simply a poorly performing emotion recognition component. Similarly, according to the empathic capabilities implemented to an agent, the features of every capability should be evaluated systematically at every stage of development.

For the evaluation of empathy as a broader concept, we will use the categorization of empathy features based on the evolutionary approaches [12], [17] as we discussed in Section II. According to these approaches, the empathic capacity can be categorized into three hierarchical mechanisms: emotional communication, emotion regulation and cognitive processes. Similar components have been proposed by other researchers in empathy [8], [49] and emotional intelligence research [50]. However, it should be noted that different types of definitions, models of empathy, as well as the capabilities and goals of interactive agents would require the evaluation of different subsets of these capabilities.

1) *Emotional Communication*: Emotional communication capacity forms the foundation of affective behaviors including empathy [12], [50]. This capacity can be further categorized as emotion recognition and emotion expression components. The successful detection and recognition of the input emotions would directly impact the empathic behavior of every level, hence the perception of empathy of the agent. Similarly, as the empathic behavior is essentially an emotional response to the stimuli, the emotional expression ability of the agent would directly influence the empathic behavior and the evaluation of the behavior. Therefore, it is crucial to include the individual evaluations of emotional communication capacity to assess the empathic capabilities of an interactive agent.

The evaluation of the emotion recognition ability may include a variety of well-established tests depending on the input modalities of the agent. For example, a text-based conversational agent's emotional communication capability can only be tested via the text-based linguistic emotional recognition and expression, where an embodied conversational agent should also be evaluated according to its speech, body gestures and facial expressions. Similarly, the success of the emotion expression behavior should be evaluated depending on the output modalities of the agent that are going to be used for expressing the empathic emotions. Following the behavioral metrics for empathy that are designed to evaluate the emotion recognition from pictures [31], voice [32] and complex emotions from movies [33], the metrics for the recognition of agents should include the evaluation of each modality. Affective computing research provides well-established evaluation metrics for emotion recognition in computational systems [51].

2) *Emotion Regulation*: A variety of models on empathy and emotional intelligence assign central importance in the ability to regulate emotions based on a variety of dynamics [12], [50]. Emotion regulation capacity can be based on personality and mood of the individual that allows for automatic regulation [51]. Humans are found to automatically assign attributes such as personality, gender and mood to interactive systems [44]. Personality metrics such as the Big Five are widely used in affective computing research [52]. Subjective evaluation metrics for emotional control have been proposed [53], [54]. However, these approaches have not been used by the empathic computing research and may require adjustments.

3) *Cognitive Processes*: The higher level of empathic capacity is suggested to include the cognitive processes such as appraisal, re-appraisal, self-oriented perspective taking and other-oriented perspective taking behaviors [12]. These cognitive processes would also control the emotion regulation abilities that allow for suppression or enhancement of emotions based on the context [51]. As we covered in Section II, behavioral measures such as understanding appraisals [34], the picture-stories task [34] and the Social Stories Questionnaire (SSQ) [35] are used to assess the deficiencies in the cognitive empathy. However, there are no standardized method to evaluate these capabilities in artificial agents. Moreover, the domain dependence and the problem of scalability for the cognitive capabilities makes it problematic to perform these

tests to interactive agents with various capabilities.

Even though we suggested solutions for the feature-level evaluations to adopt the existing metrics, most of them needs further adjustments and validations to be applied in artificial agents.

IV. CONCLUDING REMARKS

Empathy as a complex socio-emotional phenomena where the variety of definitions and models in the research community makes it problematic to evaluate and compare the implementation of the behavior in interactive artificial agents. This article is aimed to describe in detail the methods have been developed in the empathy research to evaluate empathic behavior, that can be translated into the emerging computational empathy research. We attempted to provide a systematic approach to the evaluation of this complex by suggesting the approach the evaluation on system-level and feature-level. As we acknowledge the difficulties of establishing a common ground in a diverse set application areas and capabilities of agents, we believe the importance of specifying the broader picture in the evaluation of empathy.

Our goal was to provide a guide on how we can evaluate the empathic behavior in artificial agents. We proposed system-level and feature-level evaluations for computational empathy systems to approach the issue systematically. We further provided a list of factors and components that can be used as a road-map to create individual evaluations for empathic systems in various application areas. We propose that the extensive body of work in the evaluation of empathy in humans, and the evaluation methods from affective and social computing can be used for computational empathy research. We hope to initiate the discussion towards creating a common ground to evaluate and compare computational empathy methods with this paper. We believe that a collective effort is required to develop specific measures and evaluation frameworks of empathy for interactive artificial agents.

REFERENCES

- [1] S. D. Preston and F. B. De Waal, "Empathy: Its ultimate and proximate bases," *Behavioral and brain sciences*, vol. 25, no. 1, pp. 1–20, 2002.
- [2] A. Coplan and P. Goldie, *Empathy: Philosophical and psychological perspectives*. Oxford University Press, 2011.
- [3] A. Coplan, "Understanding empathy: Its features and effects." in *Empathy: Philosophical and psychological perspectives*, A. Coplan and P. Goldie, Eds. Oxford University Press, 2011, pp. 3–18.
- [4] S. Brave, C. Nass, and K. Hutchinson, "Computers that care: investigating the effects of orientation of emotion exhibited by an embodied computer agent," *International journal of human-computer studies*, vol. 62, no. 2, pp. 161–178, 2005.
- [5] I. Leite, G. Castellano, A. Pereira, C. Martinho, and A. Paiva, "Empathic robots for long-term interaction," *International Journal of Social Robotics*, vol. 6, no. 3, pp. 329–341, 2014.
- [6] T. W. Bickmore and R. W. Picard, "Establishing and maintaining long-term human-computer relationships," *ACM Transactions on Computer-Human Interaction (TOCHI)*, vol. 12, no. 2, pp. 293–327, 2005.
- [7] H. Prendinger, J. Mori, and M. Ishizuka, "Using human physiology to evaluate subtle expressivity of a virtual quizmaster in a mathematical game," *International journal of human-computer studies*, vol. 62, no. 2, pp. 231–245, 2005.
- [8] A. Paiva, I. Leite, H. Boukricha, and I. Wachsmuth, "Empathy in virtual agents and robots: a survey," *ACM Transactions on Interactive Intelligent Systems (TiIS)*, vol. 7, no. 3, p. 11, 2017.

- [9] M. H. Davis, "Measuring individual differences in empathy: Evidence for a multidimensional approach." *Journal of personality and social psychology*, vol. 44, no. 1, p. 113, 1983.
- [10] S. Baron-Cohen and S. Wheelwright, "The empathy quotient: an investigation of adults with asperger syndrome or high functioning autism, and normal sex differences," *Journal of autism and developmental disorders*, vol. 34, no. 2, pp. 163–175, 2004.
- [11] E. J. Lawrence, P. Shaw, D. Baker, S. Baron-Cohen, and A. S. David, "Measuring empathy: reliability and validity of the empathy quotient," *Psychological medicine*, vol. 34, no. 5, pp. 911–920, 2004.
- [12] F. B. de Waal and S. D. Preston, "Mammalian empathy: behavioural manifestations and neural basis," *Nature Reviews Neuroscience*, vol. 18, no. 8, p. 498, 2017.
- [13] B. L. Omdahl, *Cognitive appraisal, emotion, and empathy*. Psychology Press, 2014.
- [14] N. Eisenberg and J. Strayer, *Empathy and its development*, ser. Cambridge studies in social and emotional development, 1987.
- [15] M. H. Davis *et al.*, "A multidimensional approach to individual differences in empathy," 1980.
- [16] M. L. Hoffman, *Empathy and moral development: Implications for caring and justice*. Cambridge University Press, 2001.
- [17] F. B. De Waal, "The russian dollmodel of empathy and imitation," *On being moved: From mirror neurons to empathy*, pp. 35–48, 2007.
- [18] F. De Vignemont and T. Singer, "The empathic brain: how, when and why?" *Trends in cognitive sciences*, vol. 10, no. 10, pp. 435–441, 2006.
- [19] D. L. Neumann, R. C. Chan, G. J. Boyle, Y. Wang, and H. Rae Westbury, "Measures of empathy: Self-report, behavioral, and neuroscientific approaches," in *Measures of Personality and Social Psychological Constructs*, 2015, pp. 257–289.
- [20] R. Hogan, "Development of an empathy scale." *Journal of consulting and clinical psychology*, vol. 33, no. 3, p. 307, 1969.
- [21] R. D. Froman and S. M. Peloquin, "Rethinking the use of the hogan empathy scale: A critical psychometric analysis," *The American Journal of Occupational Therapy*, vol. 55, no. 5, pp. 566–572, 2001.
- [22] R. N. Spreng*, M. C. McKinnon*, R. A. Mar, and B. Levine, "The toronto empathy questionnaire: Scale development and initial validation of a factor-analytic solution to multiple empathy measures," *Journal of personality assessment*, vol. 91, no. 1, pp. 62–71, 2009.
- [23] A. F. Garton and E. Gringart, "The development of a scale to measure empathy in 8-and 9-year old children." *Australian Journal of Educational & Developmental Psychology*, vol. 5, pp. 17–25, 2005.
- [24] M. Hojat, S. Mangione, J. S. Gonnella, T. Nasca, J. J. Veloski, and G. Kane, "Empathy in medical education and patient care," *Academic Medicine*, vol. 76, no. 7, p. 669, 2001.
- [25] W. J. Reynolds, *The measurement and development of empathy in nursing*. Routledge, 2017.
- [26] H. Hashimoto and K. Shiomi, "The structure of empathy in japanese adolescents: Construction and examination of an empathy scale," *Social Behavior and Personality: an international journal*, vol. 30, no. 6, pp. 593–601, 2002.
- [27] R. L. Reniers, R. Corcoran, R. Drake, N. M. Shryane, and B. A. Völlm, "The qcae: A questionnaire of cognitive and affective empathy," *Journal of personality assessment*, vol. 93, no. 1, pp. 84–95, 2011.
- [28] S. B. Eysenck and H. J. Eysenck, "Impulsiveness and venturesomeness: Their position in a dimensional system of personality description," *Psychological reports*, vol. 43, no. 3_suppl, pp. 1247–1255, 1978.
- [29] S. W. Mercer, M. Maxwell, D. Heaney, and G. Watt, "The consultation and relational empathy (care) measure: development and preliminary validation and reliability of an empathy-based consultation process measure," *Family practice*, vol. 21, no. 6, pp. 699–705, 2004.
- [30] S. Baron-Cohen, T. Jolliffe, C. Mortimore, and M. Robertson, "Another advanced test of theory of mind: Evidence from very high functioning adults with autism or asperger syndrome," *Journal of Child psychology and Psychiatry*, vol. 38, no. 7, pp. 813–822, 1997.
- [31] S. Baron-Cohen, S. Wheelwright, J. Hill, Y. Raste, and I. Plumb, "The reading the mind in the eyes test revised version: a study with normal adults, and adults with asperger syndrome or high-functioning autism," *The Journal of Child Psychology and Psychiatry and Allied Disciplines*, vol. 42, no. 2, pp. 241–251, 2001.
- [32] O. Golan, S. Baron-Cohen, J. J. Hill, and M. Rutherford, "The reading the mind in the voicetest-revised: a study of complex emotion recognition in adults with and without autism spectrum conditions," *Journal of autism and developmental disorders*, vol. 37, pp. 1096–1106, 2007.
- [33] O. Golan, S. Baron-Cohen, J. J. Hill, and Y. Golan, "The reading the mind in films task: complex emotion recognition in adults with and without autism spectrum conditions," *Social Neuroscience*, vol. 1, no. 2, pp. 111–123, 2006.
- [34] S. Baron-Cohen, A. M. Leslie, and U. Frith, "Mechanical, behavioural and intentional understanding of picture stories in autistic children," *British Journal of developmental psychology*, vol. 4, no. 2, pp. 113–125, 1986.
- [35] J. Lawson, S. Baron-Cohen, and S. Wheelwright, "Empathising and systemising in adults with and without asperger syndrome," *Journal of autism and developmental disorders*, vol. 34, no. 3, pp. 301–310, 2004.
- [36] S. Baron-Cohen, S. Wheelwright, A. Spong, V. Scahill, J. Lawson *et al.*, "Are intuitive physics and intuitive psychology independent? a test with children with asperger syndrome," *Journal of Developmental and Learning Disorders*, vol. 5, no. 1, pp. 47–78, 2001.
- [37] L. Dybkjaer, N. O. Bernsen, and W. Minker, "Evaluation and usability of multimodal spoken language dialogue systems," *Speech Communication*, vol. 43, no. 1-2, pp. 33–54, 2004.
- [38] Z. Ruttkay, C. Dormann, and H. Noot, "Evaluating ecas-what, how and why?" in *Dagstuhl Seminar Proceedings*. Schloss Dagstuhl-Leibniz-Zentrum für Informatik, 2006.
- [39] S. H. Rodrigues, S. Mascarenhas, J. Dias, and A. Paiva, "A process model of empathy for virtual agents," *Interacting with Computers*, vol. 27, no. 4, pp. 371–391, 2015.
- [40] H. Boukricha, I. Wachsmuth, M. N. Carminati, and P. Knoeferle, "A computational model of empathy: Empirical evaluation," in *2013 Humaine Association Conference on Affective Computing and Intelligent Interaction*. IEEE, 2013, pp. 1–6.
- [41] S. W. McQuiggan, J. L. Robison, R. Phillips, and J. C. Lester, "Modeling parallel and reactive empathy in virtual agents: An inductive approach," in *Proceedings of the 7th international joint conference on Autonomous agents and multiagent systems-Volume 1*. International Foundation for Autonomous Agents and Multiagent Systems, 2008, pp. 167–174.
- [42] M. Ochs, D. Sadek, and C. Pelachaud, "A formal model of emotions for an empathic rational dialog agent," *Autonomous Agents and Multi-Agent Systems*, vol. 24, no. 3, pp. 410–440, 2012.
- [43] Ö. N. Yalçın and S. DiPaola, "Evaluating levels of emotional contagion with an embodied conversational agent," in *Proceedings of the 41st Annual Conference of the Cognitive Science Society*, 2019.
- [44] B. Reeves and C. I. Nass, *The media equation: How people treat computers, television, and new media like real people and places*. Cambridge university press, 1996.
- [45] B. C. Müller, M. L. Van Leeuwen, R. B. Van Baaren, H. Bekkering, and A. Dijksterhuis, "Empathy is a beautiful thing: Empathy predicts imitation only for attractive others," *Scandinavian journal of psychology*, vol. 54, no. 5, pp. 401–406, 2013.
- [46] C. Misselhorn, "Empathy with inanimate objects and the uncanny valley," *Minds and Machines*, vol. 19, no. 3, p. 345, 2009.
- [47] A. B. Loyall, "Believable agents: Building interactive personalities." Carnegie-Mellon Uni Pittsburg PA, Tech. Rep., 1997.
- [48] C. Bartneck, D. Kulić, E. Croft, and S. Zoghbi, "Measurement instruments for the anthropomorphism, animacy, likeability, perceived intelligence, and perceived safety of robots," *International journal of social robotics*, vol. 1, no. 1, pp. 71–81, 2009.
- [49] Ö. N. Yalçın and S. DiPaola, "A computational model of empathy for interactive agents," *Biologically inspired cognitive architectures*, vol. 26, pp. 20–25, 2018.
- [50] K. R. Scherer, "Componential emotion theory can inform models of emotional competence." in *The Science of emotional intelligence : knowns and unknowns / edited by Gerald Matthews, Moshe Zeidner, and Richard D. Roberts.*, ser. Series in affective science. Oxford University Press, 2007, pp. 101–126.
- [51] K. R. Scherer, T. Bänziger, and E. Roesch, *A Blueprint for Affective Computing: A sourcebook and manual*. Oxford University Press, 2010.
- [52] A. Vinciarelli and G. Mohammadi, "A survey of personality computing," *IEEE Transactions on Affective Computing*, vol. 5, no. 3, pp. 273–291, 2014.
- [53] E. Gullone and J. Taffe, "The emotion regulation questionnaire for children and adolescents: A psychometric evaluation." *Psychological assessment*, vol. 24, no. 2, p. 409, 2012.
- [54] D. A. Preece, R. Becerra, K. Robinson, J. Dandy, and A. Allan, "Measuring emotion regulation ability across negative and positive emotions: The perth emotion regulation competency inventory," *Personality and Individual Differences*, vol. 135, pp. 229–241, 2018.