

# Comparing a mentalist and an interactionist approach for trust analysis in Human-Robot Interaction

Marc Hulcelle - Télécom Paris

Giovanna Varni - University of Trento

Nicolas Rollet - Télécom Paris

Chloé Clavel - Télécom Paris

2  
main  
perspectives

- Trust was mostly studied in Human-Robot Interaction (HRI) from a Psychological *mentalist* perspective
  - Mostly focused on robot-related factors that impact user's trust towards the robot
- Some adopt an *interactionist* perspective as described by Interactionist Sociology theories
  - Entirely focused on the user's display of trust

## Research Question

Can we identify criteria to differentiate both approaches based on their theoretical framework and trust assessment tools' methodologies ?

Trust is defined as a user's mental-state. Definition from Rousseau [1] :

“**psychological state** comprising the intention to **accept vulnerability** based upon **positive expectations** of the intentions or behavior of another”



Trust is therefore a combination of :

- User's mental projection of the robot's capabilities (Cognitive trust)
- An affective response to these (Affective trust) [2]



Trust is mostly measured through questionnaires filled by *users themselves* at the *beginning* and *end* of an interaction :

- Interpersonal Trust Scale [3]
- Negative Attitude towards the Robot Scale [4]
- Robot Trust Scale [5]
- Godspeed questionnaire [6]



Trust can also be measured through physiological sensors (e.g. EEG)

[3] J.B. Rotter. A new scale for the measurement of interpersonal trust. Journal of personality (1967)

[4] D.S. Syrdal, K. Dautenhahn, K. Koay, and M. Walters. The Negative Attitudes towards Robots Scale and Reactions to Robot Behaviour in a Live Human-Robot Interaction Study. Proceedings of AISB09 (2009)

[5] K.E. Schaefer. Measuring trust in human robot interactions: Development of the “trust perception scale-HRI”. Robust Intelligence and Trust in Autonomous Systems (2016)

[6] 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 (2009)

Trust is a result of the **state of the interaction**, and is oriented towards both the **content** and the **format** of the interaction, defined as [7] :

“form of affiliation and credit characterized by **a set of behaviors** that are intentional or not, expressive or propositional”



Observable at different bases :

- the robot's capacity to maintain a fluid and progressive interaction
- its skill in accomplishing a specific action at a given moment
- its knowledge



We conducted a comparative annotation study of two measurement tools

- Dataset: First two phases of all 10 interactions of the Vernissage dataset [8]
- Tools : Robot Trust Scale (Mentalist) / TURIN (Interactionist)
- Procedure : Annotate 10 seconds long segments TOTAL : 180 segments  
5 experts A, B, C, D, and E annotated using one tool

	Interactions 1→ 5	Interactions 6 → 10
Expert A	RTS / TURIN	RTS / TURIN
Expert B	TURIN	
Expert C		TURIN
Expert D	RTS	
Expert E		RTS



We used the Robot Trust Scale reduced to its 14 items.

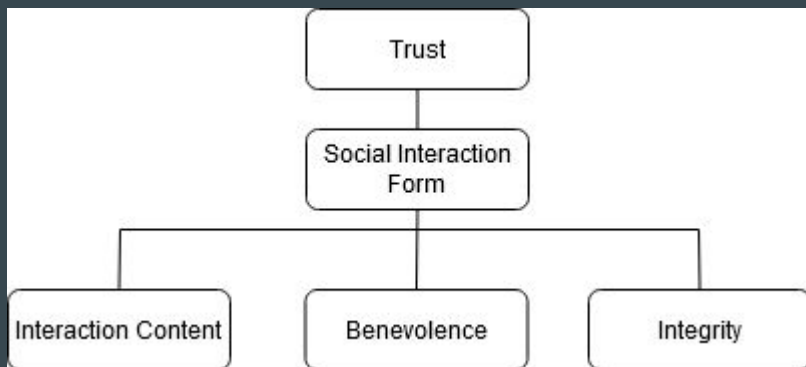
“What % of the time will the robot...” (11 Likert scale)

Have errors	Provide appropriate information	Be unresponsive
Malfunction	Communicate with people	Provide feedback
Function successfully	Act consistently	Be reliable
Be predictable	Be dependable	Meet the needs of the mission

We dropped “perform exactly as instructed” and “follow directions” items

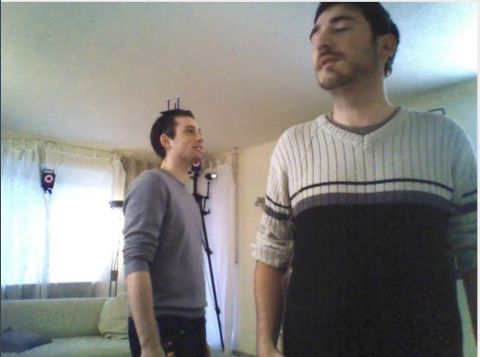

Trust in hUman Robot INteraction (TURIN) coding scheme :

- At most 2 items from “Social Interaction Form”
- At most 2 items from “Interaction Content”
- At most 1 for “Benevolence”
- At most 1 for “Integrity”



Social Interaction Form	Interaction Content
Gaze	Compliance
Facial expression	Cooperation
Nod	Alignment
Gesture	Approval
Phrasing	Out-of-context comment
Intonation	<b>Trusting only</b>
F-formation*	Joke
Speaking turn	<b>Mistrusting only</b>
Repetition	Doubt
Participation status*	
Benevolence	Integrity
Respect	Honesty
Personal info disclosure	Responsibility
Warmth	Promise
	<b>Mistrusting only</b>
	Manipulation

\*Group-specific item

00:00:33.000

Sélection: 00:00:33.000 - 00:00:43.000 10000

☐ Mode de sélection
 ☐ Mode de boucle

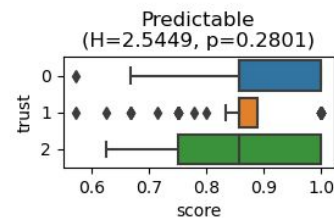
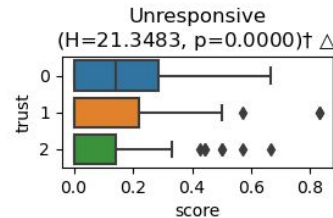
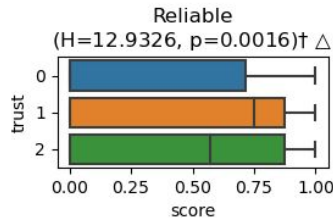
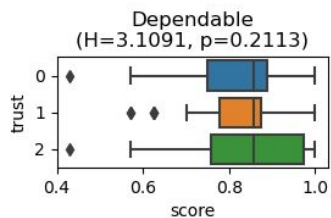
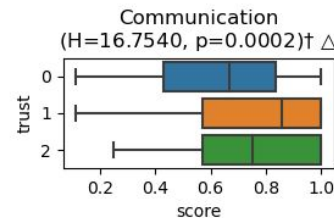
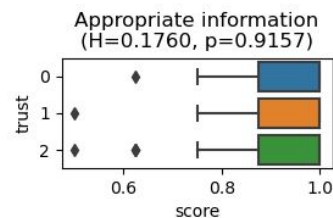
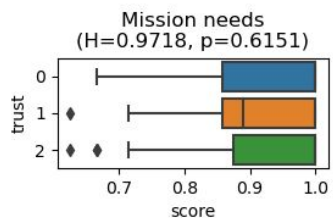
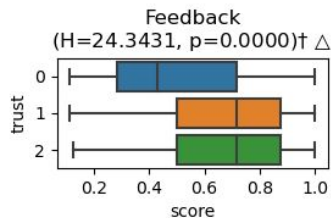
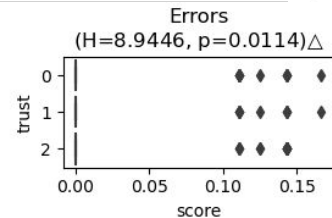
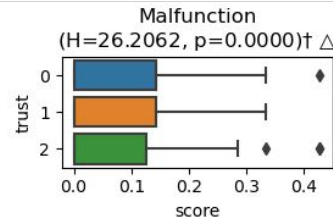
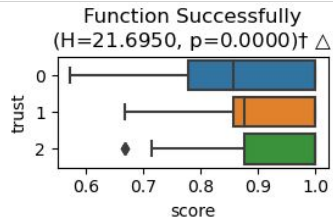
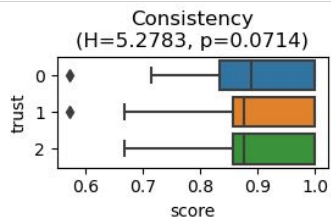
Grille				Texte	Sous-titres	Lexique	Commentaires	Reconnizers	Métadonnées	Contrôles
▼ Trust										
>	Nr	Annotation			Temps de dé...	temps de fin	Durée			
1	Trusting				00:00:33.000	00:00:43.000	00:00:10.000			
2	Neutral				00:00:43.000	00:00:53.000	00:00:10.000			
3	Neutral				00:00:53.000	00:01:03.000	00:00:10.000			
4	Trusting				00:01:03.000	00:01:13.000	00:00:10.000			
5	Trusting				00:01:13.000	00:01:23.000	00:00:10.000			
6	Trusting				00:01:23.000	00:01:33.000	00:00:10.000			
7	Trusting				00:01:33.000	00:01:43.000	00:00:10.000			
8	Trusting				00:01:43.000	00:01:53.000	00:00:10.000			
9	Trusting				00:01:53.000	00:02:03.000	00:00:10.000			
10	Neutral				00:02:03.000	00:02:13.000	00:00:10.000			
11	Trusting				00:02:13.000	00:02:23.000	00:00:10.000			
12	Trusting				00:02:23.000	00:02:33.000	00:00:10.000			
13	Trusting				00:02:33.000	00:02:43.000	00:00:10.000			
14	Neutral				00:02:43.000	00:02:53.000	00:00:10.000			
15	Trusting				00:02:53.000	00:03:03.000	00:00:10.000			
16	Neutral				00:03:03.000	00:03:13.000	00:00:10.000			

person1.wav

00:00:10.000 00:00:20.000 00:00:30.000 00:00:40.000 00:00:50.000 00:01:00.000 00:01:10.000 00:01:20.000 00:01:30.000 00:01:40.000 00:01:50.000 00:02:00.000

	00:00:10.000	00:00:20.000	00:00:30.000	00:00:40.000	00:00:50.000	00:01:00.000	00:01:10.000	00:01:20.000	00:01:30.000	00:01:40.000	00:01:50.000	00:02:00.000
Trust [10]	Trusting	Neutral	Neutral	Trusting	Trusting	Trusting	Trusting	Trusting	Trusting	Trusting	Trusting	Neutral
Interaction_Form_0 [10]	Intonation	Intonation	Intonation	Facial expression	Gaze	Gaze	Gaze	Gaze	Facial expression	Participation status	Gaze	
Interaction_Form_1 [10]	Gaze		Facial expression	Gesture			Nod	Nod	Nod			
Interaction_Content_0 [10]	Approval	Doubt	Approval	Compliance	Alignment	Alignment	Alignment	Alignment	Approval			
Interaction_Content_1 [10]	Compliance			Approval			Compliance	Approval				
Benevolence [10]	Respect	Disclose personal	Disclose personal	Warmth			Respect	Warmth				
Integrity [10]												

- We compute the global trust score for each segment by averaging all items' score
- For each annotator, we rescale the global trust score through a min-max scaling
- We aggregate segments according to their assigned TURIN label
- We compare the RTS global score distribution between different TURIN labels



We identified 4 criteria that differentiates both approaches to analyze trust:

- Time-framing
- Orientation
- Generalization
- Scalability

	Time-framing			Orientation		Generalization		Scalability	
	BU	ST	EI	Data-driven	Theoretical-framework-driven	Specific	Generic	Individual	Group
<b>Mentalist</b>			X		X		X	X	X
<b>Interactionist</b>	X	X		X		X		X	X

**Table 1: Summary of the comparison of the mentalist and interactionist approach based on 4 criteria.**

**BU: Behavioral Unity. ST: Speaking Turns. EI: Entire interaction.**

QUESTIONS ?

