



Automatic analysis of trust over the course of a human-robot interaction using multimodal features and recurrent neural architectures

Karola Pitsch

Alexandre Pauchet

Catherine Pelachaud

Brian Ravenet

Reviewer

Reviewer

Examiner

Examiner

Chloé Clavel

Giovanna Varni

Nicolas Rollet

Thesis supervisor

Thesis co-supervisor

Thesis co-supervisor

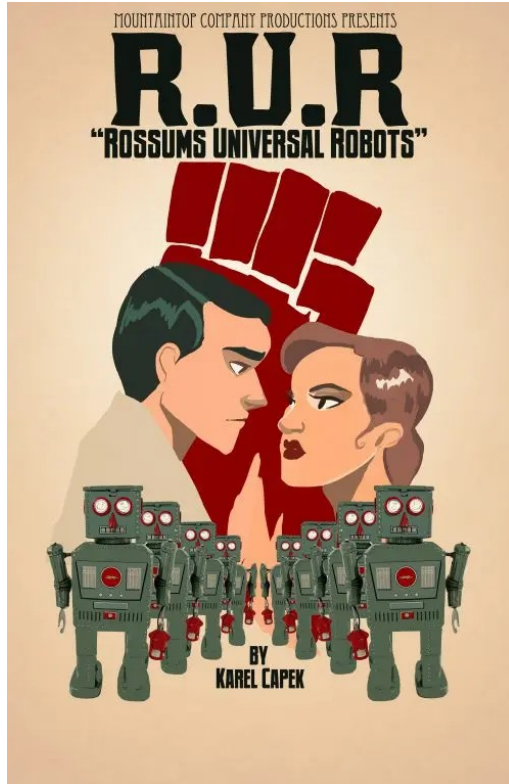
Presentation Plan

- Part I : Introduction
- Part II : State of the Art
- Part III : A New Framework
- Part IV : Computational Models of Trust
- Part V : Conclusion and Perspectives

Part I

Introduction

I - Introduction

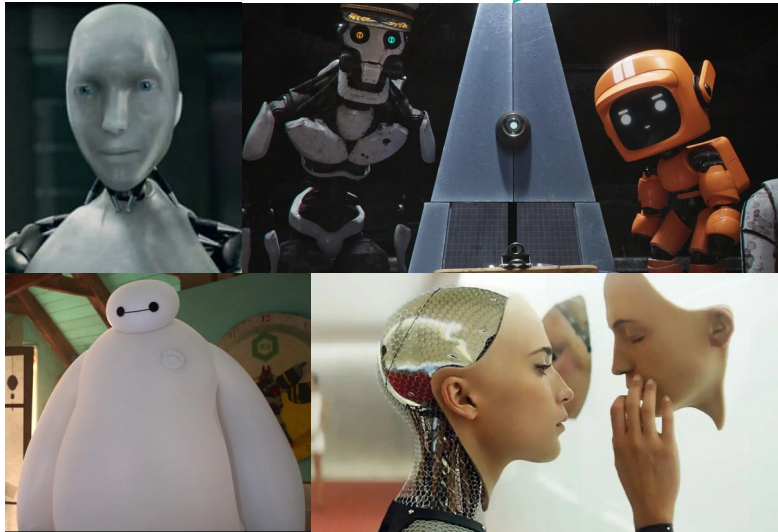


- The term **robot** was first coined in a play by Karel Capek
It comes from the slavic word **robota**
- It was used to designate artifical beings that could be mistaken with humans
- In the play, robots happily work with humans, but eventually revolt and cause the extinction of the human race



I - Introduction

Media depiction



Discrepancies

- Autonomy
- Technical skills
- Social skills
- Emotionnal expressivity



User acceptance ?
Trust ?

Reality



I - Introduction

What is Trust ?

Rousseau's definition :

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

→ Too little trust (**undertrust**) → disuse

→ Too much trust (**overtrust**) → potentially dangerous use

I - Introduction

Trust has an impact on robot's acceptance and task performance



Trust should be monitored during the interaction



How to measure and model trust ?



I - Introduction

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



- Interpersonal Trust Scale (Rotter et al. 1967)
- Negative Attitude towards the Robot Scale (Syrdal et al. 2009)
- Trust Perception Scale - HRI (Schaefer 2016)
- Godspeed questionnaire (Bartneck et al. 2009)

Trust can also be measured through proxy measures (e.g . distance)

I - Introduction

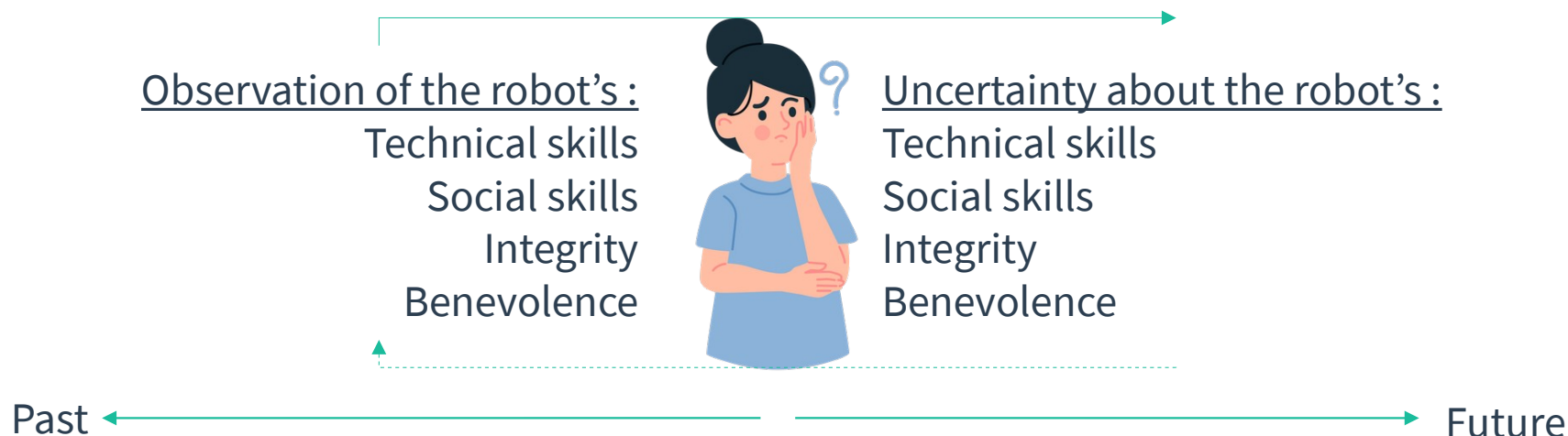
- RQ 1: Which theoretical framework is applicable to perform a multimodal analysis of trust regularly throughout the interaction ?
- RQ 2: Do homogeneous segments of trust arise within the interaction based on observable behavioral cues ?
- RQ 3: How can we discriminate trusting segments from mistrusting ones with tangible behavioral cues ?

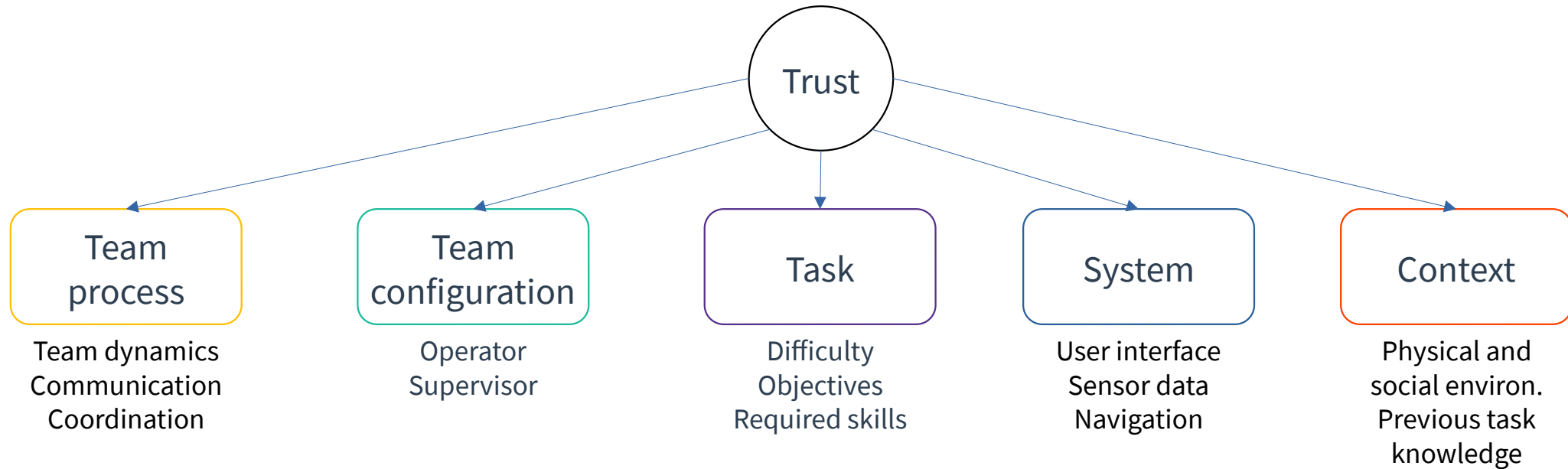
Part II

State of the art

- Bickmore et al. (2011)

« a process of uncertainty reduction, the ultimate goal of which is to reinforce assumptions about a partner's dependability with actual evidence from the partner's behavior »





Trust is mostly measured through questionnaires :

- ~~Interpersonal Trust Scale (ITS)~~ Used for human-human interaction
- Godspeed questionnaire Not a direct trust measure but all items have an impact on trust
- Negative Attitude towards the Robot Scale (NARS)
- Trust Perception Scale – HRI (TPS)

Questionnaires were used to determine trust antecedents and correlates

Trust Perception Scale - HRI

11-Likert scale

« What % of the time will the robot ... »



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	Perform exactly as instructed	Follow directions

* reduced to its 14 items

Each item is a **projection** of one of the robot's capabilities **given past interactions**

Pros

- Direct measure of the participant's perception of the robot
- Comprehensive list of items

Cons

- Only participants themselves can fill the questionnaires
- Subjective interpretation of items
- Time-consuming
- Requires interrupting the interaction
- Discrepancy between what participants think and how they behave

Very few multimodal computational of trust in HRI

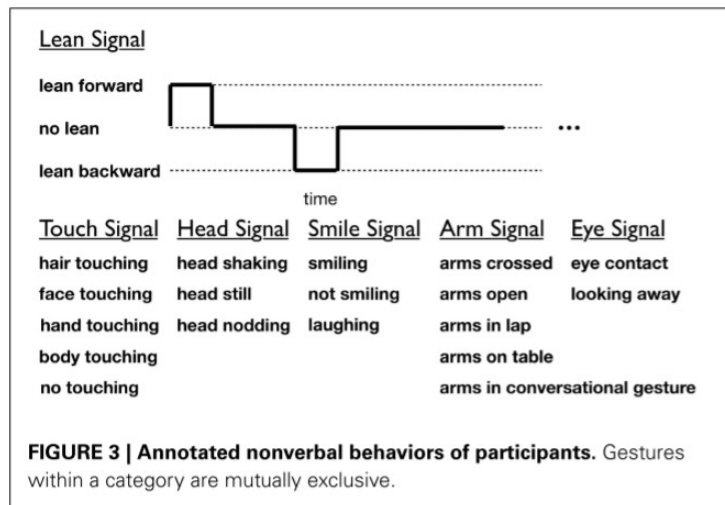
- Previous studies mostly focused how the robot's behavior impacts participants' trust
- Behavioral studies focused on the robot's behaviors and not the participants
- Continuous measures of trust are scarce

II – State of the art

Automatic trust analysis methods



Interaction for 5 minutes then exchange coins.
Can keep or give a certain amount. The other receives double of what is given.

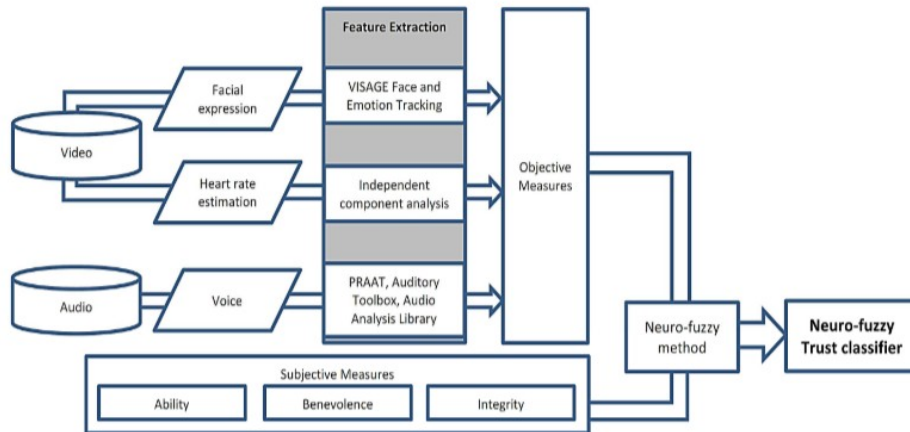
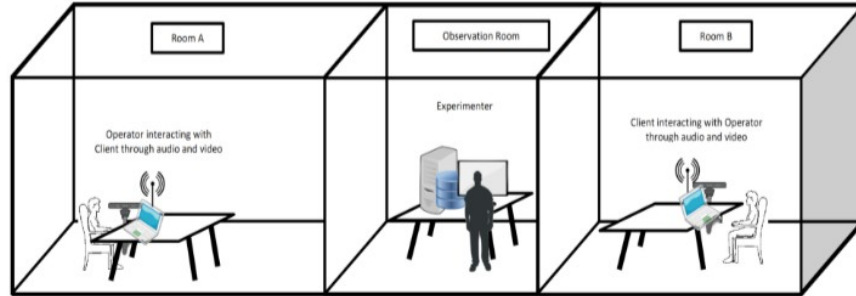


Hidden
Markov
Models

- low-trust {  ,  ,  ,  }
- high-trust {  ,  ,  ,  }

II – State of the art

Automatic trust analysis methods



Significant contributors for trust estimation

Facial Expression	Sad, Anger, Fear
Heart Rate	RRV-Min, RRV-Mean, LF/HF ratio, RRI-Range, RRI-Mean, RRV-Median, RRVDistr-Median, RRI-Median, RRVDistr-Range, RRVDistr-Max
Voice	MeanVoicePower, VoicePitch-Range, Formants-Range, StdSpectralCentroid, VoicePitch-Min, VoiceRatio, VoiceFreq-Peak, Formants-Std, VoicePitch-Std, VoicePitch-Max, Formants-Min, Formants-Max, StdByMeanZeroCrossing, TeagerWave-Mean, StdEnergyEntropy, TeagerWaveFreq-Mean, VoicePitch-Mean
Misc. Information	Ethnicity, Attribute

Takebacks

- Trust should be analyzed in a multimodal fashion
- Some specific behaviors are indicative of a certain type of trust
- Facial expressions and vocal descriptors play an important role in trust prediction

Part III

A New Framework

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 :

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



Interactional Trust relates to the ecological validity of the robot as an autonomous agent during a social interaction

It is thus observable at different bases :

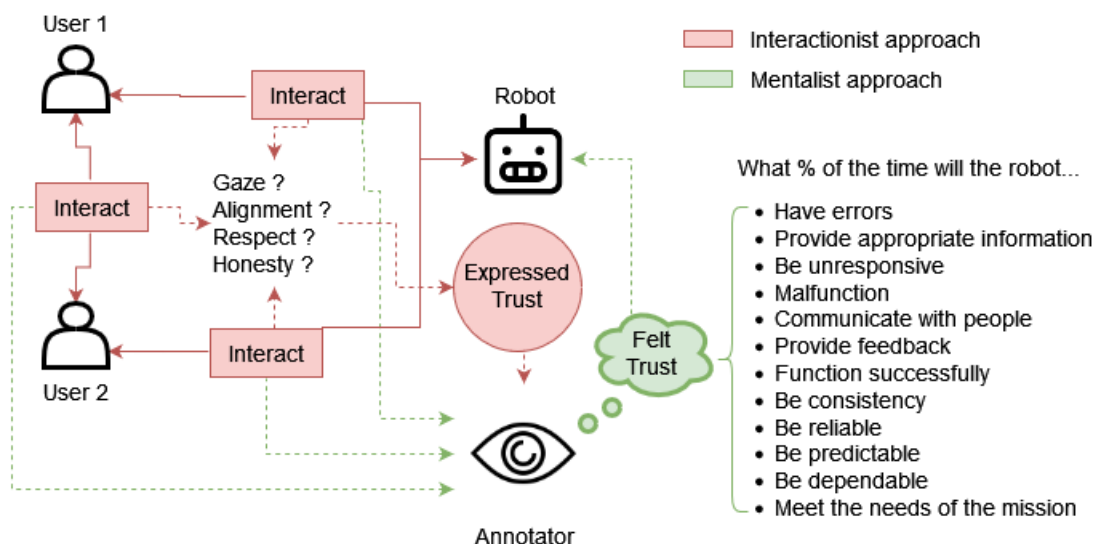
- the robot's capacity to maintain a fluid and progressive interaction
- its skill in accomplishing a specific action at a given moment
- its knowledge



✗ No pre-observation assumption that should be validated with post-analysis statistical tests.

➡ The observer determines the relevance of each participant's action in relation to trust, and to the interaction's history and context.

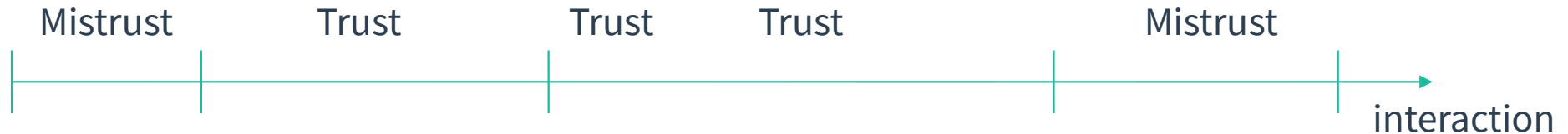
Difference between trust expressed through behaviors and trust felt by participants



Coding scheme for Trust in hUman Robot INteraction (TURIN)

Segmentation situated in the interaction dynamics :

- Start at the single behavioral act
- Acts referring to changes in behavior should be assigned to a trust category
- Consecutive acts of the same category are aggregated



Coding scheme for Trust in hUman Robot Interaction (TURIN)

Trust

Trusting

- Display of naturalness, fluidity of the interaction
- Vulnerability acceptance
- Friendliness display
- Acknowledging the partner's competence

Mistrusting

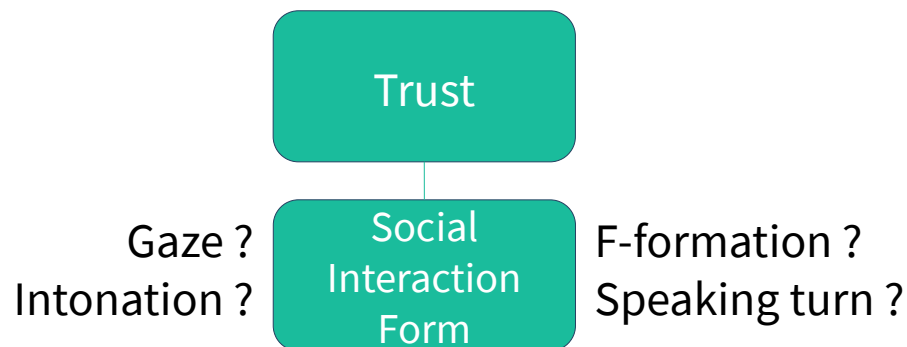
- Uneasiness
- Doubt / confusion
- Aggressiveness
- Unwillingness to cooperate

Neutral

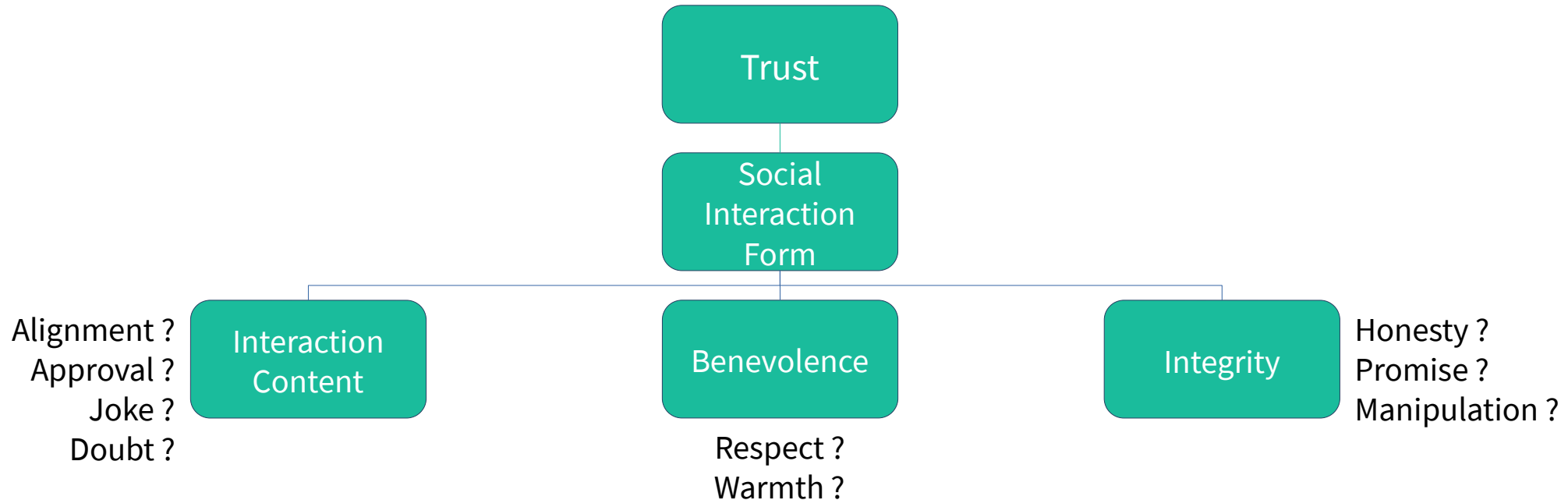
Behaviors that are inconclusive



Coding scheme for Trust in hUman Robot INteraction (TURIN)



Coding scheme for Trust in hUman Robot INteraction (TURIN)



Vernissage corpus

- 10 interactions
 - Quick introduction
 - Vernissage
 - Self-presentation
 - Art quizz
- Recordings
 - Video from 3 different angles
 - Audio for each participant and robot
 - Head pose and rotation



Validation of the corpus :

- Two expert annotators
- Annotated the 1st minute of 3 interactions
- ELAN software



Trusting
example



Mistrusting
example

Peter Wittenburg et al. "ELAN: A professional framework for multimodality research". In: 5th international conference on language resources and evaluation (LREC 2006). 2006, pp. 1556–1559.

Segment category	IRA (κ)	Mean duration (s)	Std (s)
Mistrusting	0.79	4.6	2.2
Trusting	0.64	2.1	1.5
Neutral	0.45	4.7	4.6

- Easier to recognize errors and disfluencies in the interaction
- Neutral segments are considerably longer → more disagreement within
- Segments are short
 - Focus on behavioral changes for the unitizing
 - Beginning of interaction : trust needs to be calibrated

III – A New Framework

RQ 1: Which theoretical framework is applicable to perform a multimodal analysis of trust regularly throughout the interaction ?

- Interactionist Sociology theories through their inductive methods
- We conceived a coding scheme TURIN that can unveil trust dynamics through the observation of participants' behaviors

III – A New Framework

RQ 2: Do homogeneous segments of trust arise within the interaction based on observable behavioral cues ?

We proposed a segmentation method based on the observation of behavioral cues that indicate trust and highlights homogeneous segments of trust.

Part IV

Computational Models of Trust

IV – Computational Models of Trust

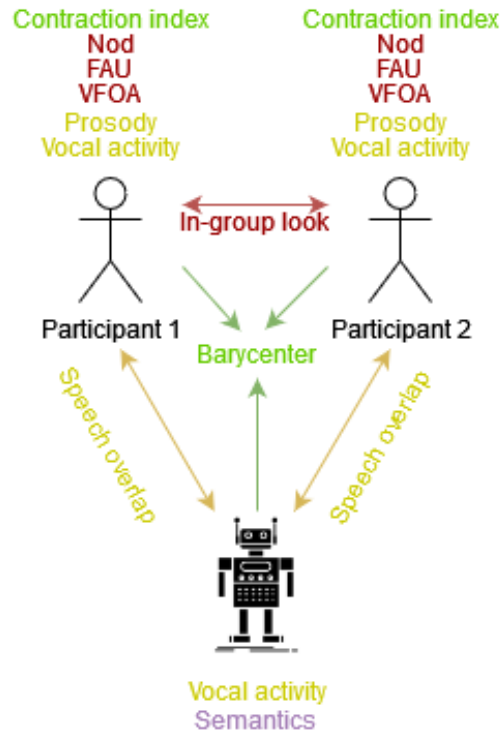
Feature Design and Extraction

Features extracted from 4 different modalities :

Body, Face, Voice, Semantics

- Aggregate within a segment using the mean and standard deviation for continuous values
- Vector length of 222 :
68 for each user, 79 for the robot, 3 for the dyads, and 4 for the triad

 Automatic extraction  Manual extraction



Hulcelle, M., Varni, G., Rollet, N., Clavel, C. (2023). “Computational Multimodal Models of Users’ Interactional Trust in Multiparty Human-Robot Interaction”. In: Rousseau, JJ., Kapralos, B. (eds) Pattern Recognition, Computer Vision, and Image Processing. ICPR 2022 International Workshops and Challenges. ICPR 2022. Lecture Notes in Computer Science, vol 13643. Springer, Cham.

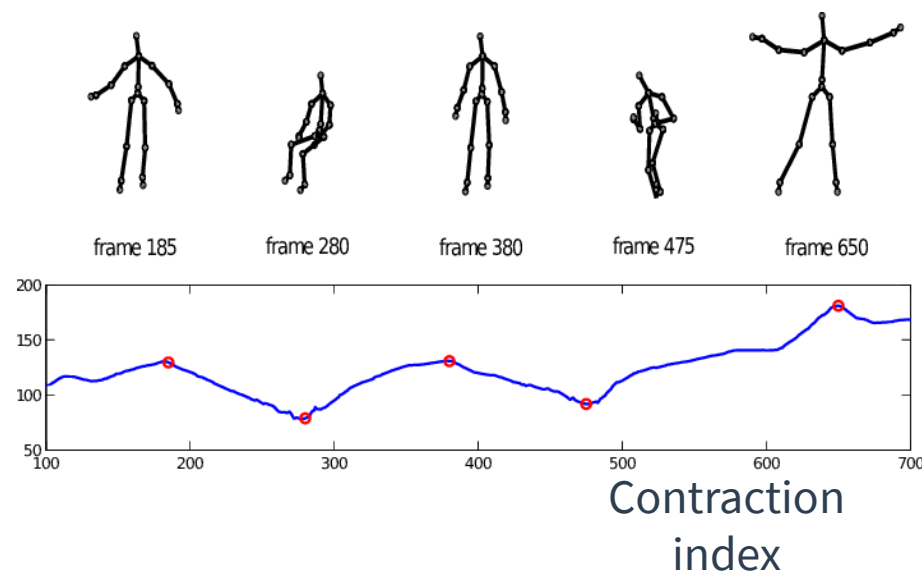
Body modality



Contraction index: ratio between the area of a body's silhouette and its bounding box



Barycenter of the group and kept the 2D point projected on the floor plane



IV – Computational Models of Trust

Feature Design and Extraction

Face Modality

Nod : time % within segment

Facial Action Units : activation value

Visual Focus Of Attention : binary indicator

left painting, central painting, right painting, Nao, other human, other, unclear

In-group look : time % within segment

Upper Face Action Units					
AU 1	AU 2	AU 4	AU 5	AU 6	AU 7
					
Inner Brow Raiser	Outer Brow Raiser	Brow Lowerer	Upper Lid Raiser	Cheek Raiser	Lid Tightener
*AU 41	*AU 42	*AU 43	AU 44	AU 45	AU 46
					
Lid Droop	Slit	Eyes Closed	Squint	Blink	Wink
Lower Face Action Units					
AU 9	AU 10	AU 11	AU 12	AU 13	AU 14
					
Nose Wrinkler	Upper Lip Raiser	Nasolabial Deepener	Lip Corner Puller	Cheek Puffer	Dimpler
AU 15	AU 16	AU 17	AU 18	AU 20	AU 22
					
Lip Corner Depressor	Lower Lip Depressor	Chin Raiser	Lip Puckerer	Lip Stretcher	Lip Funneler
AU 23	AU 24	*AU 25	*AU 26	*AU 27	AU 28
					
Lip Tightener	Lip Pressor	Lips Part	Jaw Drop	Mouth Stretch	Lip Suck

Voice modality



Vocal activity: binary indicator
speech, silence, laughter



Prosody:
F0, F0', loudness, jitter, shimmer, spectral flux, first 4 MFCCs and their derivative



Speech overlap: time % within segment

Semantics

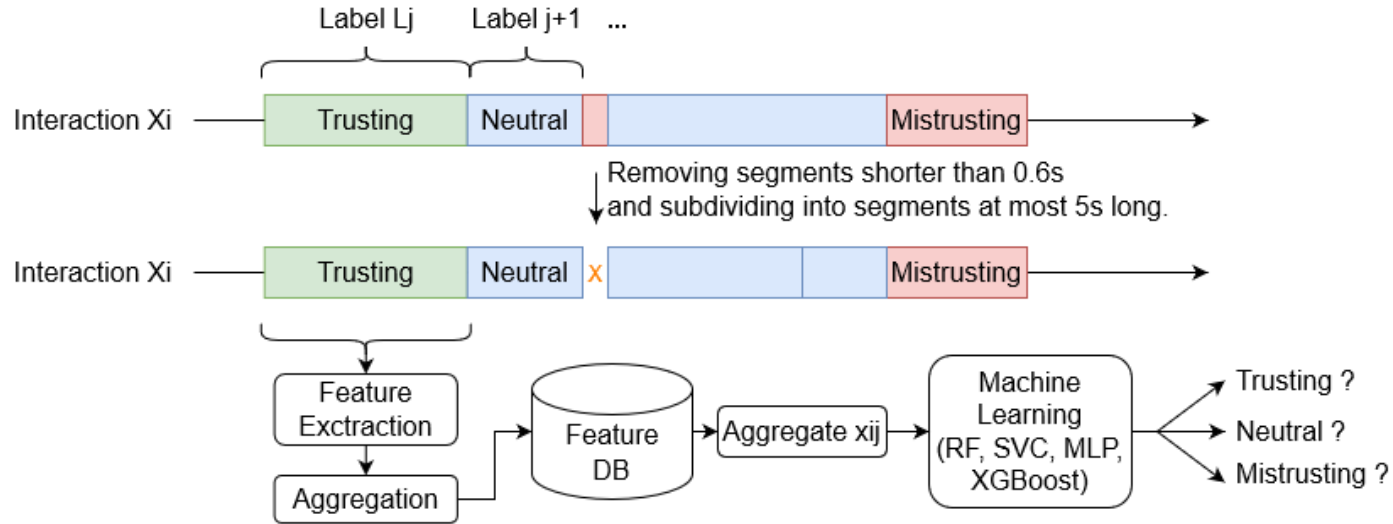


Extracted with a TinyBERT : vector size 312

- PCA to reduce to vector size of 50 (99 % explainable variance)
- Propagate semantics in silent segments from previous one

IV – Computational Models of Trust

ML Models

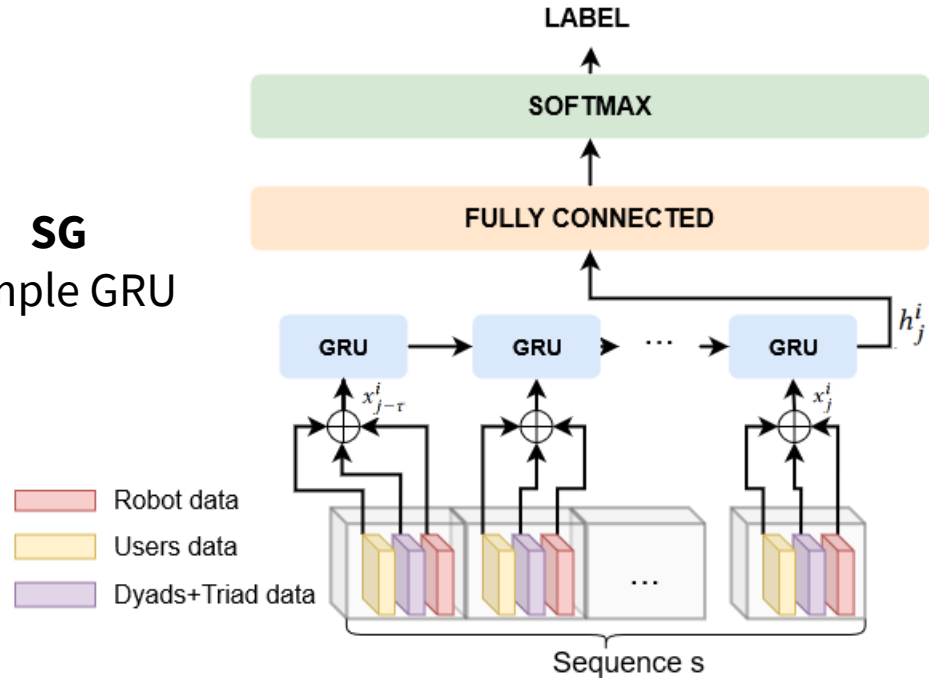


- 4 models : Ridge Classifier, Random Forest (RF), SVM-Classifer, Multi-Layer Perceptron (MLP)
- Two classification tasks
 - One-Vs-Rest
 - Multi-class
- Two fusion mechanisms :
 - Early-fusion
 - Late-fusion

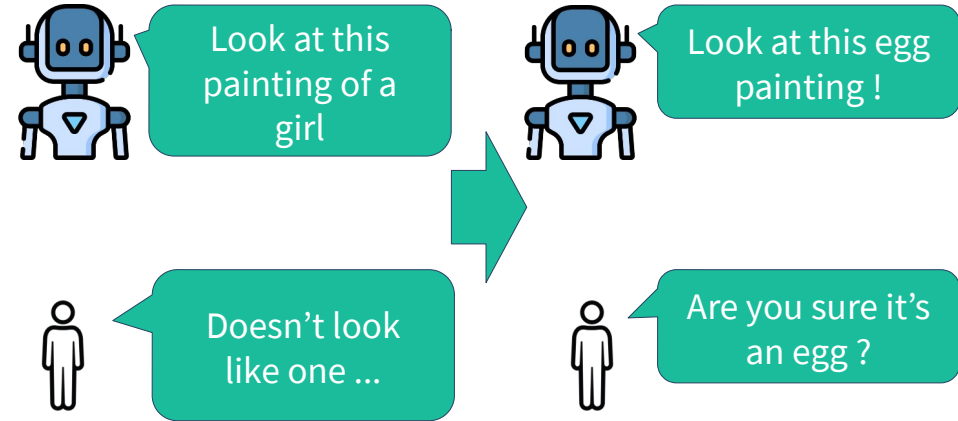
IV – Computational Models of Trust

ML Models

SG Simple GRU



A1 [RNN]: participants' actions are relevant within the sequence of previous behaviors of all participants, and produced in response to another's speaking turn



Charles Goodwin. "Conversational organization". In: Interaction between speakers and hearers (1981)

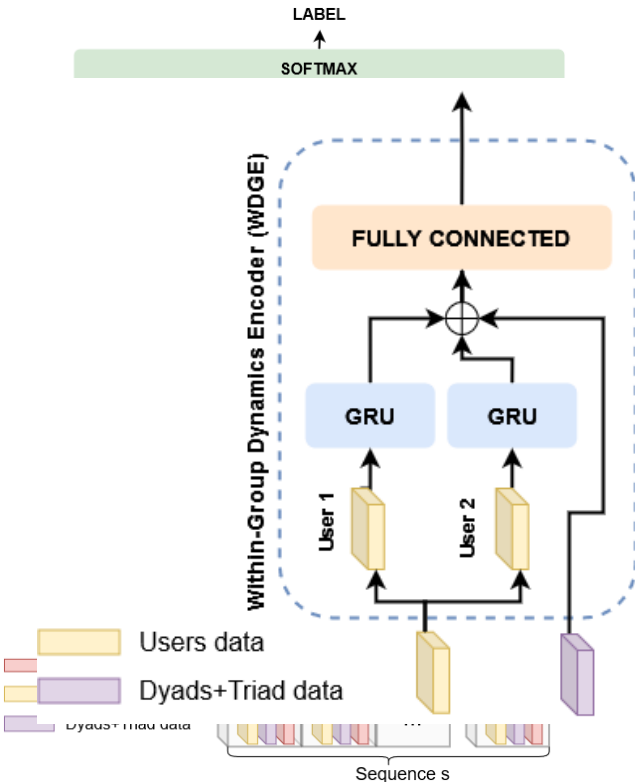
Erving Goffman. Forms of talk. University of Pennsylvania Press, 1981

Charles Goodwin et al. "Restarts, pauses, and the achievement of a state of mutual gaze at turn beginning". In: Sociological inquiry 50.3-4 (1980), pp. 272-302.

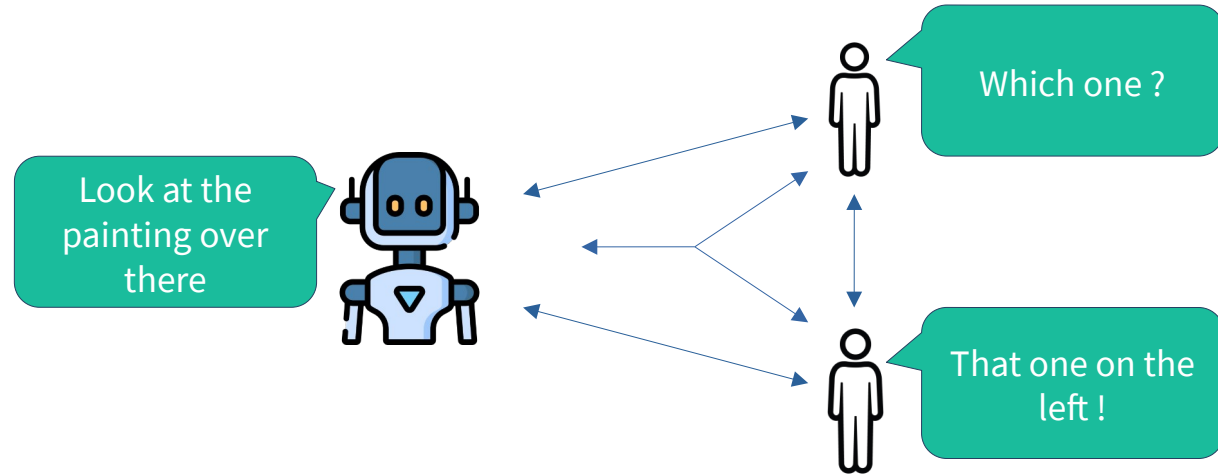
IV – Computational Models of Trust

ML Models

WGDE - SG



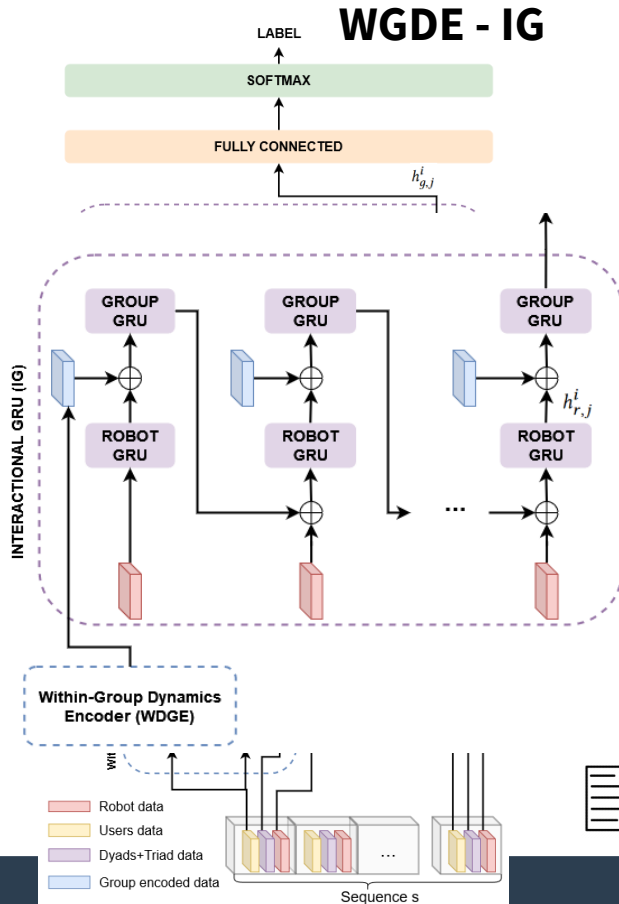
A2 [\[WGDE\]](#): participants can either be speakers - addressing the whole group or a part of it - or be listeners - by actively or passively being engaged. It is necessary to analyze the interaction between all participants to fully understand the group dynamics



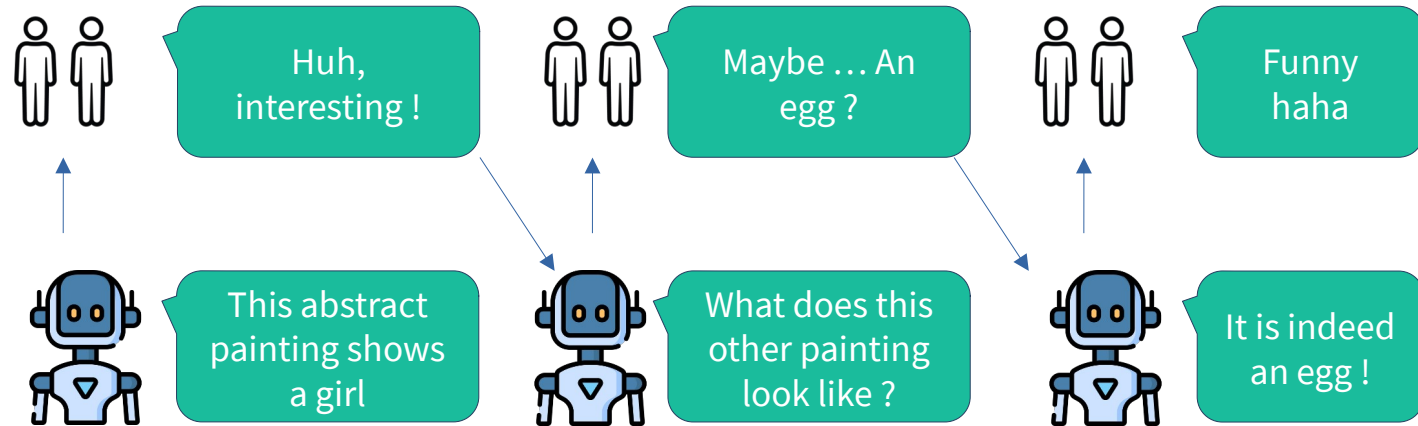
IV – Computational Models of Trust

ML Models

WGDE - IG



A3 [IG]: participants continuously exchange social signals shaping the interactional context which other participants use to build their answer, and hence renewing the context at each speaking turn



M. Hulcelle, L. Hemamou, G. Varni, N. Rollet and C. Clavel, "Leveraging Interactional Sociology for Trust Analysis in Multiparty Human-Robot Interaction," International Conference on Human-Agent Interaction (HAI '23), 2023, Gothenburg, Sweden

Training without interaction history

- Leave-Three-Groups-Out cross-validation
- Augmentation with SMOTE to obtain a balanced training set
- ROC for the OVR task, F1 score for the Multiclass
- All modalities except semantics

IV – Computational Models of Trust

ML Models

One-Vs-Rest

Early-fusion	RF	MLP	SVM-C
Trusting-vs-rest	0.72 ± 0.04	0.70 ± 0.04	0.74 ± 0.04
Neutral-vs-rest	0.77 ± 0.04	0.74 ± 0.04	0.75 ± 0.04
Mistrusting-vs-rest	0.59 ± 0.06	0.54 ± 0.07	0.58 ± 0.06
Late-fusion			
Trusting-vs-rest	0.67 ± 0.04	0.60 ± 0.04	0.66 ± 0.04
Neutral-vs-rest	0.74 ± 0.04	0.65 ± 0.04	0.70 ± 0.03
Mistrusting-vs-rest	0.54 ± 0.08	0.48 ± 0.08	0.49 ± 0.10

Multiclass

Fusion	Rand.	Majority	RF	MLP	SVM-C
<i>Early</i>	0.38 \pm 0.03	0.52 \pm 0.05	0.66 \pm 0.04	0.65 \pm 0.04	0.60 \pm 0.04
		Body	Face	Voice	
	Trusting-vs-rest	0.54 \pm 0.05	0.62 \pm 0.04	0.65 \pm 0.05	
	Neutral-vs-rest	0.57 \pm 0.03	0.65 \pm 0.03	0.73 \pm 0.05	
	Mistrusting-vs-rest	0.46 \pm 0.08	0.51 \pm 0.06	0.60 \pm 0.07	

- RF performs better
- Early-fusion mechanism works best
- Voice modality is more important

- **Training with history**
- Leave-One-Group-Out cross-validation
- Augmentation 4x the dataset with Gaussian noise $\sigma = 2 \cdot 10^{-3}$
- Weighted random sampling during training with weights equal to inverse class proportions
- Regularization term with $\lambda = 1 \cdot 10^{-2}$
- All modalities

F1 score for all sequential models

	τ	1	2	3	4	5	6	7	8
Simple GRU	SG	0.733 $\pm.120$	0.739 $\pm.119$	0.733 $\pm.118$	0.733 $\pm.127$	0.735 $\pm.123$	0.735 $\pm.125$	0.734 $\pm.125$	0.735 $\pm.124$
Simple GRU (only human data)	SG (no robot)	0.621 $\pm.058$	0.613 $\pm.081$	0.621 $\pm.080$	0.604 $\pm.095$	0.597 $\pm.084$	0.598 $\pm.092$	0.603 $\pm.085$	0.605 $\pm.087$
Within-Group Dynamics Encoder + Simple GRU	WGDE-SG	0.726 $\pm.116$	0.732 $\pm.119$	0.723 $\pm.144$	0.730 $\pm.141$	0.724 $\pm.138$	0.725 $\pm.146$	0.723 $\pm.146$	0.731 $\pm.148$
Interactional GRU	IG [†]	0.730 $\pm.113$	0.717 $\pm.105$	0.695 $\pm.120$	0.698 $\pm.163$	0.694 $\pm.182$	0.710 $\pm.145$	0.689 $\pm.175$	0.694 $\pm.188$
Within-Group Dynamics Encoder + Interactional GRU	WGDE-IG	0.730 $\pm.102$	0.730 $\pm.098$	0.715 $\pm.143$	0.736 $\pm.124$	0.735 $\pm.135$	0.745 $\pm.110$	0.730 $\pm.146$	0.714 $\pm.137$

- Full model performs best
- No optimal history size (maybe due to the lack of data)

- **SHAP value analysis without history**
- Trusting:
Nod more often, closer to the robot, tighten their lid more, more variations of contraction index, lower F0
- Mistrusting:
Lower their brows more, more changes in VFOA, further away from the robot, talk more, higher speech overlap

- **Most common errors with history**
- Segments annotated « Alignment »
- Trusting:
Gaze, Facial expression, F formation
- Mistrusting:
Gaze, Facial expression, Intonation

Part V

Conclusion and Perspectives

V – Conclusion and Perspectives

RQ 1: Which theoretical framework is applicable to perform a multimodal analysis of trust regularly throughout the interaction ?

- Interactionist Sociology theories through their inductive methods
- We conceived a coding scheme TURIN
 - It can unveil trust dynamics through the observation of users' behaviors
 - Can be used as ground truth for computational models

V – Conclusion and Perspectives

RQ 2: Do homogeneous segments of trust arise within the interaction based on observable behavioral cues ?

We proposed a segmentation method based on the observation of behavioral cues that indicate trust and highlights homogeneous segments of trust.

V – Conclusion and Perspectives

RQ 3: How can we discriminate trusting segments from mistrusting ones with tangible behavioral cues ?

- A set of multimodal features. Some features carry more weight than others, and that some were more specific to a certain trust category.
- Early-fusion mechanism leads to quite optimistic performance in binary classification with traditional ML techniques.
- We designed a neuronal architecture to model interactional dynamics within the user-group, as well as between the robot and users, which led to increased performance.

V – Conclusion and Perspectives

Other contributions

- All annotations collected on the Vernissage dataset made public
- Code of the sequential model soon available on Github

V – Conclusion and Perspectives

Perspectives

- Online trust detection
 - Different features ?
 - Computational cost ?
- Collecting data with a trust-specific scenario

V – Conclusion and Perspectives

Perspectives

- Refining TURIN and links with other social phenomena
 - Alignment ?
 - Engagement ?
- Improving the models of interactional dynamics for trust
 - Different modules ?
 - Hierarchical architecture ?



Thank you for
listening !

QUESTIONS

Engagement

“the process by which two (or more) participants establish, maintain, and end their perceived connection to one another” C. L. Sidner, C. Lee, C. D. Kidd, N. Lesh, and C. Rich, “Explorations in Engagement for Humans and Robots”, *Artificial Intelligence*, pp. 140–164, 2005

MHHRI Dataset to study this link ?

Dyadic interactions VS group interactions ?

Feature category	User's mode	
	Speaker	Listener
Distance (front sonar, face distance, head position, engagement zone)	User	User
Gaze (direction, is looking at robot)	User	User
Head & Face (head angles, 17 face AUs)	User	User
Speech (voicing probability, F0 loudness, log-energy, 12 MFCCs, is robot speaking, speech duration)	User	Robot

Table 1: *Multimodal features used to detect user's SED per category. Columns 2 & 3 show whether the features characterize the user or the robot depending on the user's mode.*

QUESTIONS

Coding schemes we took inspiration from

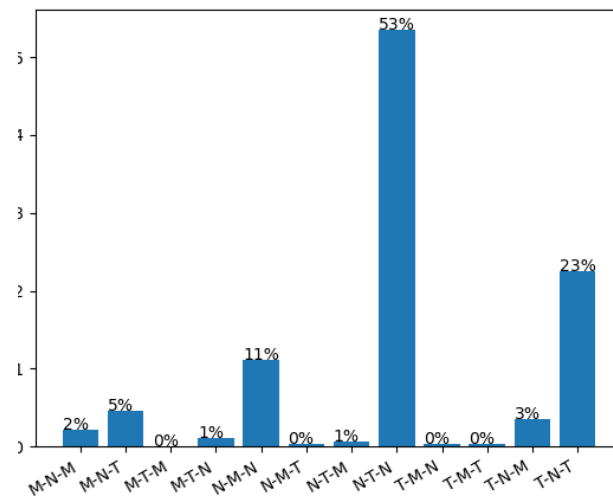
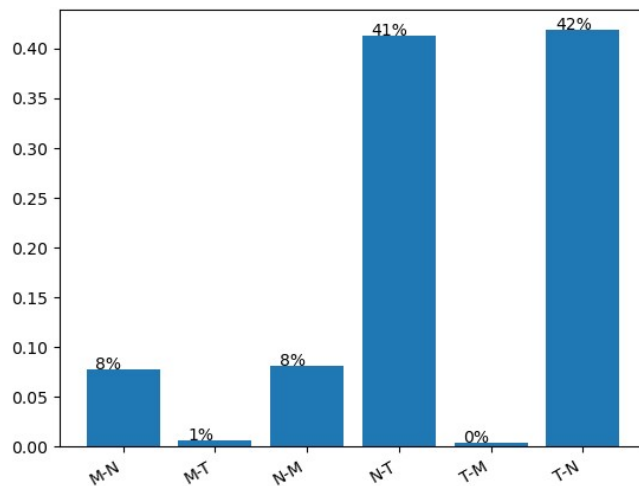
- Act4teams
 - Thought unit
 - Exclusive coding
- DCS
 - Sentence level
 - Who speaks to whom ?
 - Function of the state in the interaction process
 - Affiliation/Dominance valence

Problem-focused statements	Procedural statements	Socio-emotional statements	Action-oriented statements
Problem	Positive:	Positive:	Positive, proactive:
Describing a problem	Goal orientation	Encouraging	Expressing positivity
Connections	Clarifying	participation	Taking responsibility
with problems	Procedural suggestion	Providing support	Action planning
Defining the objective	Procedural question	Active listening	Negative,
Solution	Prioritizing	Reasoned	counterproductive:
Describing a solution	Time management	disagreement	No interest in change
Problem with a	Task distribution	Giving feedback	Complaining
solution	Visualization	Humor	Seeking someone to blame
Arguing for a solution	Summarizing	Separating opinions	Denying responsibility
Organizational	Negative:	from facts	Empty talk
knowledge	Losing the train of	Expressing feelings	Ending the discussion early
Knowing who	thought	Offering praise	
Question	(running off	Negative:	
	topic)	Criticizing/backbiting	
		Interrupting	
		Side conversations	
		Self-promotion	
Additional codes:			
Pause; laughter; incomprehensible; incomplete sentence (cut off by an interruption); other (not fitting any of the above codes)			

QUESTIONS

Class distribution

LABEL	COUNT	AVERAGE LENGTH
Mistrusting	78	2.6s ($\pm 1.6s$)
Neutral	604	4.1s ($\pm 1.4s$)
Trusting	240	3.0s ($\pm 1.5s$)



QUESTIONS

MPR Dataset « A Multimodal Multiparty Human-Robot Dialogue Corpus for Real World Interaction »

- Multiparty Dialogue Management
- Boredom detection
- Surprise detection
- Repair detection

MPR2012	MPR2016
English 1st phase (WOZ) 20-Questions 2nd phase (automatic) Mimic gestures	English and Japanese 1st phase (manual, participants know) 20-Questions 2nd phase (automatic) 20-Questions Random irrelevant and irrational messages from the robot to induce surprised reactions

QUESTIONS

Steven Marsh computational model

Description	Representation	Value Range
Situations	α, β, \dots	
Agents	a, b, c, \dots	
Set of agents	\mathcal{A}	
Societies of agents	$\mathcal{S}_1, \mathcal{S}_2 \dots$ $\mathcal{S}_n \in \mathcal{A}$	
Knowledge (e.g., x knows y)	$K_x(y)$	True/False
Importance (e.g., of α to x)	$I_x(\alpha)$	$[0, +1]$
Utility (e.g., of α to x)	$U_x(\alpha)$	$[-1, +1]$
Basic Trust (e.g., of x)	T_x	$[-1, +1]$
General Trust (e.g., of x in y)	$T_x(y)$	$[-1, +1]$
Situational Trust (e.g., of x in y for α)	$T_x(y, \alpha)$	$[-1, +1]$

$$T_x(y, \alpha) = U_x(\alpha) \times I_x(\alpha) \times \widehat{T_x(y)}$$

$$\widehat{T_x(y)} = \frac{1}{|A|} \sum_{\alpha \in A} T_x(y)$$

$$T_x(y, \alpha) > \text{Cooperation Threshold}_x(\alpha) \\ \Rightarrow \text{Will_Cooperate}(x, y, \alpha)$$

$$\text{Cooperation Threshold}_x(y, \alpha) = \frac{\text{Perceived Risk}_x(\alpha)}{\text{Perceived Competence}_x(y, \alpha) + \widehat{T_x(y)}} \times I_x(\alpha)$$