



HAL
open science

Towards Robots able to Measure in Real-time the Quality of Interaction in HRI Contexts

Amandine Mayima, Aurélie Clodic, Rachid Alami

► **To cite this version:**

Amandine Mayima, Aurélie Clodic, Rachid Alami. Towards Robots able to Measure in Real-time the Quality of Interaction in HRI Contexts. *International Journal of Social Robotics*, 2022, 14, pp.713-731. 10.1007/s12369-021-00814-5 . hal-03326954v2

HAL Id: hal-03326954

<https://laas.hal.science/hal-03326954v2>

Submitted on 23 Jun 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Towards robots able to measure in real-time the Quality of Interaction in HRI contexts

Amandine Mayima* · Aurélie Clodic · Rachid Alami

Abstract When humans interact with each other, collaborating on a shared activity or chatting, they are able to tell whether their interaction is going well or not and if they observe that its quality is deteriorating, they can adapt their behavior or invite their partner to act in order to improve it. A robot endowed with the ability to evaluate the quality of its interaction with its human partners, will have the opportunity to perform better since it will be better informed for its decision making processes. We propose metrics to be integrated in a cognitive and collaborative robot in order to measure in real-time the quality of an interaction (QoI). This permanent evaluation process has been implemented and tested within the high-level controller of an entertainment robot. A first demonstration shows the ability of the scheme to compute QoI for a direction-giving task and exhibit significant differences between its performance in interaction with a fully compliant human, a human confused by the course of action and a non-cooperative one. This paper is an extension and further refinement of work originally reported in [30].

Keywords Quality of Interaction · Human-Robot Interaction · Reasoning in robotic systems · Evaluation Methods · Collaborative tasks

*Corresponding author

Amandine Mayima
LAAS-CNRS, Université de Toulouse, CNRS, Toulouse,
France E-mail: amandine.mayima@laas.fr

Aurélie Clodic · Rachid Alami
LAAS-CNRS, Université de Toulouse, CNRS, Toulouse,
France
Artificial and Natural Intelligence Toulouse Institute (ANITI)
E-mail: surname.name@laas.fr

1 Introduction

Robots dedicated to Human-Robot interactions are not just machines receiving commands and executing them. They should be decisional agents with high-level goals, taking decisions (potentially taking into account social norms), and acting and reacting to not only their actions but those of other agents. Cognitive and interactive robots are becoming more and more capable thanks to the use of human-aware models and algorithms [22, 45], with roboticists endowing them with the ability to execute their share of the work while adapting to contingencies, particularly those caused by human's behaviours and decisions [17, 2, 27]. The decision-making process is based on a range of knowledge about the environment, the interaction, the context... Nevertheless, curiously and interestingly, very little has been done to allow the robot, while performing its collaborative or assistive activity, to permanently evaluate if things are going well or not, as humans do. We name this ability “the measure of the Quality of Interaction from the robot point of view”. We believe that enriching the robot knowledge with a good estimation about how the interaction is going, could enhance its decision-making process and thus, its social behaviour.

For example, if the robot detects that the QoI starts to drop, it can take a decision based on this information and act to try to improve the interaction quality (e.g. it can choose to change some modalities such as the language in which it communicates with the human, the volume of its speakers, or the parameters of its planners). On the contrary, when the QoI is high, the robot can decide to just continue the interaction as planned. Then, endowed with a QoI Evaluator, a robot becomes more adaptive and performs better. Also, a very poor performance all along a task could allow the robot to assess that the human is not really engaged in the in-

teraction, or even is trying to play the robot. In such situation, the robot might perhaps better disengage. Finally, from a methodological point of view, a robot deployed in the wild able to assess interactions, has an asset compared to others as it could reduce the investment in material and human resources to perform user studies. And, a developer might use the logs to improve their design.

In this paper, we only focus on the Quality of Interaction evaluation process and not on how to use its result for decision making. Therefore, we present in the sequel the methods and tools we developed, allowing the robot to evaluate in real-time the quality of the human-robot collaborative activity it is involved in. It is based on a set of metrics we have defined, focused on two concepts: the measure of human engagement and the measure of the effectiveness of collaborative tasks performance. However, this is by no means exhaustive, and other metrics and parameters could (and should) be added later. Our work can be seen as a toolbox among which it is possible to pick the desired metrics according to tasks or contexts. We propose a way to aggregate these metrics, producing the QoI. The evaluation of the QoI is performed at three different levels of abstraction: the interaction session level, the task level and the action level. In further work, this ability could provide additional information to the robot and open the possibility for reconsidering its behaviour in case it estimates that the quality of the interaction is degrading (e.g. changing its plan or the way it is achieving it, informing the human or requesting a change in their behaviour, or even deciding to disengage).

The paper is organized as follows. In the next section, we briefly discuss related work and the main challenges. In section 3 we present the representation of human-robot collaborative activity which we use and its hierarchical decomposition. In sections 4 and 5, we introduce our concept and proposed set of metrics to evaluate the Quality of Interaction. Finally, in section 6, we describe a first implementation and then conclude and discuss future developments.

2 Related work

Inspired from the evaluation methods used in Human-Computer Interactions and User Experience fields, the field of Human-Robot Interaction (HRI) has elaborated its own methods to evaluate robotic systems when they interact with humans. There are various ways to evaluate a human-robot interaction from the human perspective. Bethel *et al.* [7] divided them into five categories: (1) self-assessments, (2) interviews, (3) behavioral measures, (4) psychophysiology measures, and (5) task per-

formance metrics. They reviewed metrics used for each of the categories. They can be grouped into two types: (1) and (2) are subjective metrics and, (3), (4) and (5) are objective ones. Since our aim is to have a robot able to evaluate interactions by itself, human subjective metrics are not usable. Then we focused on the study of existing objective metrics meant to measure how the interaction goes. Steinfeld *et al.* [42] proposed a set of metrics to be used in a wide range of tasks whose goal is to assess the system performance by measuring the task effectiveness (i.e., how well the task is completed) and the task efficiency (i.e., the time required to complete a task). Their work is very thorough and inspiring but does not target the evaluation of the quality of an on-going interaction. Hoffman [16] defined a type of quality of interaction, the *fluency*, pointing out that the notion is not well defined and somewhat vague but can still be assessed and recognized when compared to non-fluent scenario. To measure it, they propose a list of objective metrics, only based on duration measures, designed to be quite general: robot idle time, human idle time, concurrent activity (i.e., active time of both the robot and the human), functional delay (i.e., time difference between the end of one agent’s task and the beginning of the other agent’s task). It is an interesting way to measure the fluency and thus the quality of the human-robot interaction but it only applies to shared workspace tasks and is dedicated to an offline evaluation.

Systems targeting real-time measurements during human-robot interactions, with the purpose to “close the loop” and use the information for decision-making, have been developed. Tanevska *et al.* [44] proposed a framework allowing the robot to perceive with face detection and evaluate in real-time the affective state (i.e. anger, happiness, sadness, surprise, etc) and the engagement state (i.e. whether the person is interested or bored in the interaction) of the people it is interacting with. However, the human affective state measure might not be enough to assess an interaction or a task as an affective state is actually a facial expression which can be misinterpreted (e.g. a smile can be a sign of happiness or embarrassment) and which might be not visible when one of the agent perform an action and looks somewhere else. Moreover, as the notion of engagement is very task specific, it needs further exploration. Real-time engagement measurement has also been investigated by Anzalone *et al.* [1] using metrics such as gaze, head pose, body pose and response times. Their work is interesting and could be an element among others to assess the interaction quality but, it is dedicated to face-to-face interactions.

Cameras are not the only sensor used to assess interactions on-the-fly, some use human physiological responses such as skin conductance and temperature, heart or brain signals. Itoh *et al.* [19], Bekele *et al.* [4] or Kulic *et al.* [24] use them to detect human affective states such as anxiety or liking in real-time. However, physiological measures often imply a lot of sensors which can be invasive for the human. And, as explained by Kulic and Croft [23], physiological signals may be difficult to interpret and there is a large variability in physiological response from person to person. Thus, it can be difficult for a controller to determine which emotional state the subject is in, or whether the response was caused by an action of the system, or by an external stimulus. Moreover, we claim that the human affective state only is not enough to assess the quality of an interaction, a human could be satisfied with an interaction or a task result even though they were stressed during it.

Finally, Bensch *et al.* [6] proposed a formal approach to compute interaction quality in real-time. Their work focused on how to combine metrics together which is in the same line as ours. However, they do not provide implementation examples, remaining at an abstract level.

In summary, while a substantial number of studies have been devoted to the evaluation of collaborative interactions for analysis purposes once the interaction is over, there is a lack of methods allowing the robot to evaluate in real-time the quality of the interaction based on multiple metrics and not only anxiety or engagement. We claim that such an ability is very important and should strongly influence the situation assessment as well as the decisional abilities of interactive and collaborative robots.

3 Representation of a H-R collaborative activity

It is possible to describe and decompose a Human-Robot collaborative activity in various ways. For all the following definitions, we place ourselves in the context of one-to-one human-robot interactions, however we believe that the scheme can be extended to multi-human multi-robot contexts. We draw our inspiration from the literature of sociology and robotics to define a model of interaction with three layered levels: interaction session, tasks and actions; as illustrated in Fig. 1. We chose to represent collaborative tasks and their decomposition using the Hierarchical Task Network (HTN) [13] representation which is often used in cognitive robotics [18, 25] and because it allows to deal with goal-based and situation-based activities at different levels of hierarchy such as task, subtasks and actions and consequently to consider different level of

granularity. In the example of a task with an overall bad QoI, it would be interesting to know that in fact it is only a particular action or subtask ruining it. Indeed, the other parts of the task can be ok, or on the opposite, a particular subtask or action can have performed very well among the others. We need and use this granularity also on three levels defined (interaction session, tasks and actions) to finely evaluate the Quality of Interaction, as a task can be of poor quality but the session is globally going well.

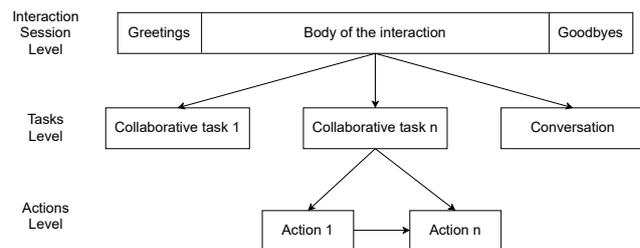


Fig. 1: The hierarchical structure of an interaction session. The highest level is the interaction session. The second level is composed of the tasks. They are included in the body of interaction of the interaction session and, two types of tasks are considered and may overlap, collaborative and conversational tasks. With this representation, a task can be recursively refined as subtasks until reaching the last level, the actions level, which is considered as atomic. Subtasks are not considered as a “real” level of the interaction session, specially to evaluate the QoI, as it may exist or not according to the task.

3.1 Representation of a H-R Interaction Session

We define an **interaction session** as the period during which the robot and a human interact together and are engaged. It is divided in three parts, following the structure proposed by Sidnell and Tanya [35] and the engagement model of Sidner and Lee [40]: the greetings, the body of the interaction and the goodbyes. First, *the greetings* corresponds to the period where an agent starts an interaction by initiating it with another agent. The interaction session lasts as long as the interactants are maintaining the interaction through conversation and collaborative tasks performance which corresponds to the *body of interaction*. Finally it ends when at least one of the interactants is disengaged, either by abruptly ending the interaction or by closing the interaction as described by Schegloff and Sacks [39], it corresponds to “the goodbyes”. For example, for an entertainment

robot in a mall, an *interaction session* starts when a person signals to the robot that they want to be engaged, by greeting it or by approaching it and looking at it. The body of interaction is composed of conversation and eventually direction-giving tasks and, the session lasts until the person says goodbye or leaves. This is the nominal case and, the duty of the robot is to contribute to maintain the session alive until the human decides to close it. However, in some (extreme) cases, the robot might decide to close the interaction by itself.

Social interactions and collaborative tasks involve engagement. There is no unique definition of what it means to be engaged. We chose one that is frequently used and has been proposed by Sidner and Lee [40]: “Engagement is the process by which two (or more) participants establish, maintain and end their perceived connection during interactions they jointly undertake”. The robot must be able to exhibit its engagement and disengagement and also to assess them with respect to its human partner. We defined three states for the body of interaction, corresponding to what is happening during the latter: conversation (i.e. a social chit-chat or a goal negotiation, without any physical action performed except communicative gestures), collaborative task (i.e. both agents executing actions in order to achieve a shared goal) or idle phases (i.e. the agents are not chatting or performing a collaborative task together but remain engaged in the interaction session, it happens in-between active interaction phases). For each of these three states, the way to exhibit the engagement varies (e.g. in a conversation, an agent looking at their partner displays their engagement; during a task, an agent correctly performing their action is a way to demonstrate their engagement). That is why there is a need to define what behavior the robot has to exhibit in each state and what behavior it should expect from the human in each state, as these behaviors are usually very specific (e.g. in a direction-giving task, the robot keeps its head oriented toward its partner’s face to demonstrate its engagement in conversation and idle contexts and when it gives a direction it expects the human to look at the direction it is showing; in a stack task, when the robot gives an instruction it expects the human to take a given cube).

3.2 Collaborative Tasks, Subtasks and Actions

Tasks compose the body of the interaction of an interaction session as shown in Fig. 1. We distinguish conversation (i.e. agents engage in dialogue to exchange ideas, to ask questions, and to resolve differences) from collaborative tasks (i.e. agents work as partners, collaborating to perform tasks and to achieve common goals).

We will not develop more on conversation since it is not the main focus of this paper, assessing the QoI of social dialog being another work.

In collaborative tasks, the robot and the human are committed to achieve a goal together, involving joint actions and shared plans [14]. When a human and a robot perform a task together, as described by Bauer *et al.* [3], we could say that the robot has the intent to help the human, so the human’s intention becomes its own intention. Then, they have the joint intention to reach a common goal and, as shown by Michael and Salice [31], they have a commitment to the joint activity, leading to perform joint actions. Therefore, during its evaluation and decision-making processes, the robot has to take into account that the human and itself should remain engaged all along an interaction session for the tasks to be successful and both have to manage and contribute to maintain expectations about what the other is doing.

The elements composing a *task* are: a goal, a plan and involved agents. A plan is needed to realize a goal. There are many ways to generate a plan. But no matter the way (using a planner to anticipate execution or relying on a reactive planning scheme), a plan is a sequence of **subtasks** which are sequences of actions – *subtasks* are not considered as a “real” level of the interaction session, specially to evaluate the QoI, as it may exist or not according to the task.

Actions are the elementary items of tasks manipulated by the high-level robot supervision controller. They cannot be decomposed further by it (e.g. placement and motion planning are achieved by a lower control system not described here). It is usual to describe an action with its preconditions, its effects and, the agents and entities implied in its execution (e.g. in plans written in PDDL (Planning Domain Definition Language) [12]). We add to this description the notion of expected reactions (which can themselves be actions) from the other agents once the action is executed.

In our model, an agent (human or robot) is a contributor to the task and has a mental state as described by Devin *et al.* [9]. The mental state is a set of facts representing, from the agent point of view, the current world state, the state of the goal and the current task state. Since we are interested here by the robot situation assessment and decisional processes, the mental state of the human is built and managed by the robot as an estimation of the beliefs of the human [33, 15, 43].

4 The Quality of Interaction (QoI)

We believe the real-time assessment of the Quality of Interaction (QoI) with a human partner (i.e. what the

robot “thinks” about how the interaction is going) is a new knowledge that could enhance the robot decision-making process. We define the Quality of Interaction as a measure that indicates how good is the interaction during human-robot collaborative activities. It is computed in real-time based on a set of metrics, at three different levels: the interaction session level, the tasks level and the actions level. The QoI of a given level is computed from selected metrics but also from the QoIs of the level below as shown in Fig. 2.

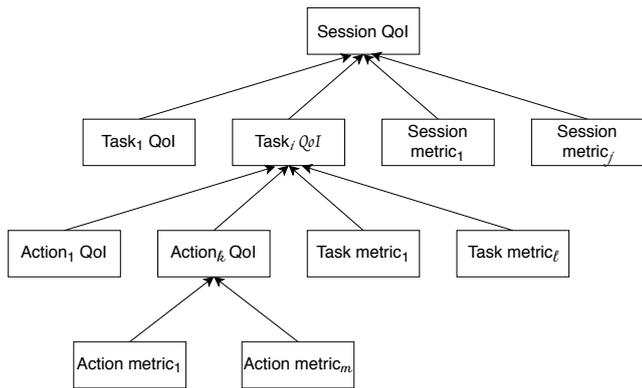


Fig. 2: Representation of the QoI dependencies, with i the number of performed tasks during the interaction session, k the number of performed actions during the task i , j the number of metrics to measure the interaction session QoI, l the number of metrics to measure the task i QoI and m the number of metrics to measure the action k QoI.

The QoI of each level is computed as a score between $[(1)$ for a good quality] and $[(-1)$ for a poor one]. Metrics used to compute the QoI are divided in three categories:

- $Mp \in [0, 1]$ if it can only have a positive effect on the evaluation;
- $Mn \in [-1, 0]$ if a metric can only have a negative effect on the evaluation;
- $M \in [-1, 1]$ if a metric can have a positive or a negative effect.

Defined by the designer according to the needs and context, a metric can belong to one category or another depending on the target application. When needed, metrics values are scaled with the equations presented in Appendix A.

The evaluation of the Quality of Interaction at the level $l \in \{session_f, task_j, action_k\}$ (with f, j and k respectively the identifiers of a given interaction session, task and action), QoI_l , is computed with:

$$QoI_l = \frac{\sum_{i=0}^x W_i * M_i}{\sum_{i=0}^x W_i} + A * \frac{\sum_{i=0}^y Wn_i * Mn_i + \sum_{i=0}^z Wp_i * Mp_i}{\sum_{i=0}^y Wn_i + \sum_{i=0}^z Wp_i} \quad (1)$$

with W_i, Wp_i, Wn_i respectively the corresponding designer-set weights of M_i, Mp_i, Mn_i , A the designer-set weight of the right part of the $+$ sign and x, y, z respectively the number of the metrics M_i, Mp_i, Mn_i .

Equation 1 aggregates the values of the metrics chosen to be indicators of the interaction level quality. As all metrics do not have the same importance in the measure of the QoI, each of them is weighted. Values of these weights are empirically defined. There are two parts in the equation, the left part of the $+$ sign and the right part. The left part of the $+$ sign is a weighted mean of the third category of metrics, the M metrics. The right part is a weighted mean of the metrics seen as bonus (i.e. Mp metrics) or penalty (i.e. Mn metrics). This latter part is weighted with A – whose value is also empirically¹ defined – to be able to adjust its influence on the left part. In such a way, if there are no Mn metrics to compensate for the Mp metrics, it is possible to limit the positive influence of the Mp metrics on the M metrics with A . It is the same if there are no Mp metrics, A can compensate the impact of the Mn metrics on the M metrics. Even though $M, Mp, Mn \in [-1, 1]$, the final result of QoI_l might be less than -1 or greater than 1 because of the addition of the M with the Mn and Mp . If it happens, QoI_l minimal value is set to -1 and its maximal value is set to 1 .

5 A set of metrics

In this section, we present a few measures to assess the QoI of an interaction session in Sect. 5.1. Then, we present metrics for the different levels based on engagement in Sect. 5.2 and effectiveness estimations during human-robot joint activities in Sect. 5.3. For example, if the human is engaged and if tasks are performed effectively, the QoI will tend to be high and *vice versa*. Both concepts are difficult to measure, so we do not exactly measure them but we compute their trends from the set of metrics presented in this section. This set is not exhaustive and will be extended in future work but it gave promising results as we show with our implementation in Sect. 6. All metrics are meant to be used

¹ Values are empirically defined given intuition regarding the importance of a given metrics for a given task and a set of testing experiments

for online evaluations of interactions. They are summarized in Table 1.

5.1 Measures to assess the QoI at the interaction session level

According to the context, the duration of an interaction session can be an indicator of the human engagement. Indeed, a human leaving only a few seconds after the beginning of the interaction is probably less engaged than a human staying with the robot several minutes. Also depending on the context, the number of executed tasks is a measure which can be considered as interesting information with respect to the engagement of the human, as well as the ratio of successful tasks. The more the human executes successful tasks with the robot, the higher the session QoI might be. Finally, it can be valuable to take into account how the session has been terminated in the evaluation of the quality of an interaction session. For instance, the fact that the human leaves abruptly in the middle of a task, during an idle time or a conversation without saying goodbye, or only at an appropriate time saying farewell to the robot is significant in terms of social interaction quality.

5.2 Metrics related to human engagement

Michael *et al.* [32] stated that commitments² facilitates “the planning and coordination of joint actions involving multiple agents. Moreover, commitment also facilitates cooperation by making individuals willing to contribute to joint actions to which they would not be willing to contribute if they, and others, were not committed to doing so”. As it is an important element of the joint action, we want to provide the robot with a way to estimate the engagement of its partner during an interaction.

Metrics allowing to state if an agent is engaged or not in an interaction are often specific to the type of interaction. For example, Fan *et al.* [10] implemented their measure of the human engagement as a kind of hysteresis: when the human gaze is on the robot, they are considered as engaged and when the human gaze is somewhere else during more than 3 consecutive seconds, they are considered as not-engaged.

In the same vein, we think that the measure of the engagement for a collaborative activity can be divided

in 2 types of metrics, summed up in Table 1: the Human contribution to the goal and the Fulfilling robot expectations about social interaction.

We define in this section examples of metrics of each types which can be used to estimate the level of engagement of the human partner.

5.2.1 Human contribution to the goal

A good and very promising indicator could be the ability from the robot to evaluate how well the human actions help to the goal progression. We call this indicator *Human contribution to the goal*. To the best of our knowledge, there is no general method to estimate it.

As a first version of the *Human contribution to the goal*, we chose to measure it through the number of times the robot has to repeat an instruction or a question before the human performs correctly, when it expects the human to answer or to perform the action. As, if it needs to repeat, it means that the human is not correctly contributing to the goal, intentionally or not, as they are not performing their part of the HR action as they should. The more the robot needs to repeat because of the human’s bad performance, the less they are contributing to the goal, the more the action QoI should decrease.

5.2.2 Fulfilling robot expectations about social interaction

During a social interaction, agents are expected to behave in a certain way and so the robot has expectations about the human. Then, the robot can monitor the human behavior to check if they are acting as they are expected to. For example, most of the time, when the robot speaks to the human, it will expect them to look at it and so it can monitor if it is the case or not as implemented by Fan *et al.* [10]. Quite similarly, Lemaignan *et al.* [26] developed a way to measure if the human is *with* the robot during their interaction, based on attention assessment, by computing if the human is looking at the desired attentional target or not. This latter metric will be integrated to our framework in future work.

As the works of Lemaignan *et al.* and Fan *et al.*, we estimate the *Fulfilling robot expectations about social interaction* with the human head orientation, in the context of our implementation described in Sect. 6. We compute an attention ratio i.e., the time during which the human is attentive to the robot (i.e. staying close enough and looking at it) when it speaks compared to the total time of the speech:

$$Ar = \frac{duration_{isAttentiveTo(robot)=true}}{duration_{robot_speaks}} \quad (2)$$

² In the robotic domain, it is the word “engagement” and not “commitment” which is often used, unlike in the psychological and philosophical fields.

	Metric names	Measures	Illustration	Session	Task	Action	
Effectiveness	Progress towards goal	Distance-to-Goal	Geometric distance			x	x
		Time-to-Goal	Time			x	x
		Steps-to-Goal	Number of executed actions/subtasks			x	
	Deviation from standard duration	Time			x	x	
Engagement	Fulfilling robot expectations about social interaction	e.g. attention ratio, with-ness,...		x	x	x	
	Human contribution to the goal	e.g. number of repeated instructions, number of successful human actions,...			x	x	

Table 1: The set of metrics presented in Section 5.

5.3 Metrics related to effectiveness

One can elaborate metrics to measure how well a task or an action is achieved. As discussed by Olsen and Goodrich [34], there are a variety of metrics such as time-based metrics which reward the speed of performance or the response times; error metrics which are based on counting retrials, failures, or mistakes; coverage metrics which measure to what extent a goal is achieved, as well as other possible metrics. We use some of them such as counting retrials, however these metrics alone were not enough for our example task as we are in a HRI context.

One can measure for different kinds of tasks, the ratio of successful³ executions to the total number of executions (e.g. $R = \frac{Succ}{Exec}$) or the deviation from the initial plan (distance, cost, trajectory, etc).

We define four metrics, summed up in Table 1, allowing to measure the current task and action effectiveness. Three of them are means to measure how the progress towards the goal of a task or an action varies. Indeed, they are good indicators for the interaction quality as, when executing a task or an action, if the agents are not getting closer from the goal or even diverged from it, it means that something goes wrong. There are three different metrics because the one to use depends on the type of task or action. The fourth

³ Obviously, the success is context and task dependent and should be defined according to the needs

metric allows to compare the current execution duration to the standard execution duration of the task or action, based on durations measured during previous executions.

5.3.1 Metrics to assess the progress towards the goal

We defined three different metrics to assess the progress towards the goal. The first one allows to assess the progress towards the goal of geometric-based actions. The second estimates the progress by using the remaining time to reach the goal. Finally, the last one measures the number of remaining steps (actions or subtasks) before achieving the goal of a task.

Distance-to-Goal When an agent is performing a geometric-based action such as a movement, observing if the agent is getting closer to the target position over time provides a useful information about how well the action is going. Therefore, we introduce the *Distance-to-Goal* ΔDtG metric:

$$\begin{cases} \Delta DtG(t=0) = 0 \\ \Delta DtG(t) = \max(0, \Delta DtG(t-1) - 1) \\ \quad \text{if } path_length(t) < path_length(t-1) \\ \Delta DtG(t) = \Delta DtG(t-1) + 1, \text{ otherwise.} \end{cases} \quad (3)$$

with $path_length(t)$ the length of the path leading the goal at time t (e.g. which can be given by a reactive motion planner [20]). The metric lower bound is 0. If

at time t the agent is closer to its final position than at $t-1$, i.e. progressing towards their goal, the metric is set to decrease or to remain equal to 0. Now, if the agent has not moved or is even further, the metric increases. The closer the metric value is to 0, the better it is, as it means the distance to the goal has decreased over time. We chose to not directly compute the difference between $path_length(t)$ and $path_length(t-1)$ as the results would be very different whether it is an action implying a long path or a short path.

Time-to-Goal This measure is intended to estimate the progress of a given task or action towards its goal based on the estimation of the remaining time to reach it. It compares the current estimated time to goal with the initial estimated time to goal taking into account the current task duration. As so, it is possible to measure the variation compared to the initial plan. We define the *Time-to-Goal* ΔTtG as:

$$\Delta TtG(t) = \max(0, e(t) + TtG(t) - TtG(T_0)) \quad (4)$$

with $e(t) = t - T_0$ the task execution duration (time elapsed since the beginning of the task), $TtG(t)$ the current time to the goal, and $TtG(T_0)$ the initial planned time to goal. In our work, $TtG(t)$ and $TtG(T_0)$ are provided by a reactive motion planner [20] because we used the metric for navigation but it could be provided by other kind of planners.

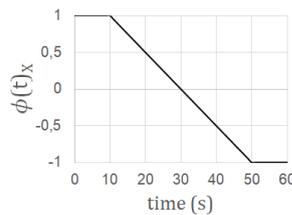
Steps-to-Goal One way to estimate the remaining distance to the goal for a task is to count the number of remaining subtasks or actions (depending on the relevant scale) to perform. In addition, one can add a factor which estimates the weight (or effort needed) of each action or subtask. These weights can be determined by the designer, provided by the planner, etc. Then, the *Steps-to-Goal* \mathcal{D} of a task can be computed as time t :

$$\mathcal{D}(t) = \frac{\sum_{i=1}^c \mathcal{W}_i}{\sum_{i=1}^n \mathcal{W}_i} \quad (5)$$

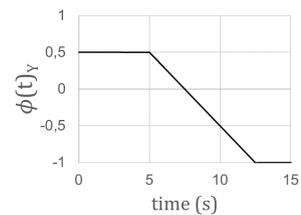
with \mathcal{W}_i the weight of a subtask/action i , c the number of completed subtasks/actions and n the total number of planned subtasks/actions.

5.3.2 Deviation from standard duration

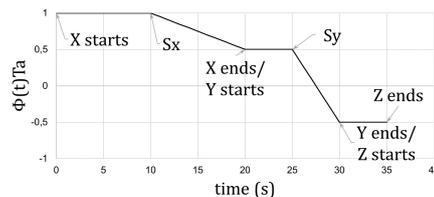
We introduce here a metric to measure the deviation from standard execution duration, the *Deviation from standard duration* ϕ for subtasks/actions and the *Deviation from standard duration* Φ for a whole task. This measure is intended to represent the degradation of the



(a) Plot of $\phi(t)_X$ of the sub-task X lasting 60 seconds, with $SD_X = 10sec$, $V_X = 0.5$ and $\alpha = 1$



(b) Plot of $\phi(t)_Y$ of the sub-task Y lasting 15 seconds, with $SD_Y = 5sec$, $V_Y = 1$ and $\alpha = 0.5$



(c) Plot of $\Phi(t)_{Ta}$ for a task composed of a sequence of three subtasks X, Y, Z : the duration of X exceeded $SD_X = 10s$ and reached $20s$, the duration of Y exceeded $SD_Y = 5s$ and reached $10s$, finally the duration of Z was less than $SD_Z = 10s$

Fig. 3: Examples of plots of the ϕ and Φ functions

quality of execution of a HR task when its duration exceeds a certain time.

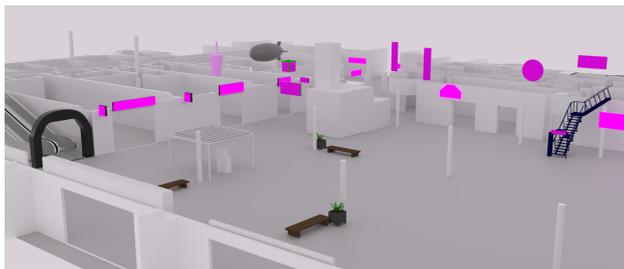
To each subtask/action a_i , we associate two attributes whose values are defined by the designer: a soft deadline SD_i and a decreasing quality speed V_i . If, at time t , the execution duration $e(t) = t - T_0$ of a subtask or action a_i which has started at T_0 exceeds SD_i , the quality will decrease over time at speed V_i :

$$\phi(t)_i = \max \left(V_i * \frac{-\max(e(t) - SD_i, 0)}{SD_i} + \alpha, -1 \right) \quad (6)$$

where α is the value initial value and the upper bound (as at $t = 0$, $\max(e(t) - SD_i, 0) = 0$) of ϕ_i , when the subtask/action a_i starts.

Then, we define a metric Φ for a task. It is an aggregation of the ϕ_i computed for each performed subtask/action a_i of the task. At any moment, Φ can be seen as a memory of the previous steps, so the initial value α of a_i is equal to the final value of ϕ_{i-1} of the previous subtask/action a_{i-1} , $\alpha = \phi(T_{final})_{i-1}$.

We can notice that it is not possible for this metric to increase over time since it memorizes the values of the previous actions. However, the total computed QoI can get higher thanks to the other metrics. Moreover, ϕ can be used independently of Φ . In such a case, the initial of value α of ϕ can be set to 1.



(a) 3D model of the Ideapark mall in Finland. The pink elements are storefronts.



(b) 3D model of the copy of the Ideapark mall in our lab. The red elements are the signs of fake storefronts.

Fig. 4: We have built a mockup of the Finnish mall environment in our lab in order to be able to test and debug the direction-giving task in our lab. This environment comprises a two-level area with corridors, “shops”, passages, stairs, open central space and consequently allowed us to run realistic guiding scenarios and finalize the QoI Evaluator proof-of-concept.

Three examples are given in Fig. 3. Fig. 3a and 3b represent $\phi(t)_X$ and $\phi(t)_Y$ for two independent subtasks X and Y . Fig. 3c is a plot of $\Phi(t)_{Ta}$ for the task Ta composed of the subtasks X, Y, Z with $SD_X = 10s$, $V_X = 0.5$, $SD_Y = 5s$, $V_Y = 1$, $SD_Z = 10s$ and $V_Z = 1$.

6 Implementation of the QoI Evaluator for a direction-giving task

As a proof-of-concept, we have implemented the Quality of Interaction Evaluator as part of a fully integrated robotic system [11] developed in the context of the MuMMER European project⁴. This project led to the deployment in a mall of an autonomous robot based on a Pepper platform. There, the robot reacted by starting an interaction session with any person willing to interact with it. The person had the possibility to have a conversation on several topics with the robot and/or to ask how to reach a shop or location in the mall or where to buy a given item. This latter possibility is one of the robot core tasks. It consists in giving guidance to the customers to reach locations in the mall, by pointing at places and explaining the route to the desired location. In order to guide as best as possible, the robot was allowed to move in a limited area to place itself and to invite the customer to move in order to reach a configuration where the landmarks it has chosen to indicate are visible to the human and to itself.

The direction-giving task is run by a robotic architecture presented in Fig. 5 which has been inspired from [28]. The architecture by itself is not a contribution to this paper but it is necessary to briefly describe it to understand how the direction-giving task is achieved. Thus, we give hereafter a rapid overview

of the role of the different components and their articulation within the architecture. Two of them were not developed in our lab and are described in [11]: the *Dialog* and the *Visual Perception*, identifying and tracking humans the robot is interacting with. Based on the data provided continuously by the *Visual Perception* and on a pre-defined 3D model of the mall (Fig. 4a), a *Situation Assessment* module, based on Underworlds [36], computes continuously symbolic relations between agents, and between agents and objects in the environment such as $[isSpeakingTo(X, Y)]$ when X speaks and looks at Y or $[isLookingAt(X, Y)]$. The *Route Handler* [38], based on a semantic spatial representation, handles the search for the best route to get to the destination – based on criteria such as accessibility or ease of explanation – and the verbalization of the chosen route. The *Human-Aware Navigation* of the robot is implemented using a reactive navigation planner inspired from [21]. This algorithm is able to plan and continuously adapt robot motion close to humans while respecting social constraints [41]. It is the role of the *Shared Visual Perspective Planner* [46] to try to find a position where the human will have to be in order to see an element of the environment such as a passage, a staircase or a store. It computes a new position for the robot as well, to form a triangle whose vertices are the planned robot position, the planned human position and the targeted landmark. Finally, the *Supervisor* handles the direction-giving task execution through reactive plans, coordinating all the components described above. Throughout the task, it supervises the execution, adapts the robot’s responses to human actions and to contingencies and computes the Quality of Interaction thanks to a process which we coined the Quality of Interaction Evaluator.

⁴ <http://mummer-project.eu/>

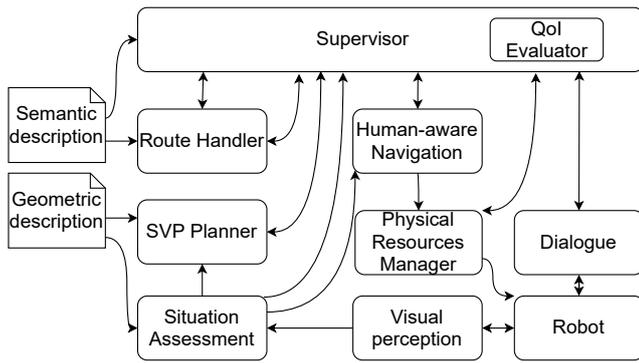


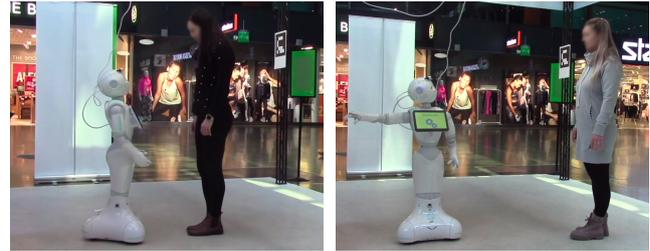
Fig. 5: The general architecture of the system running the direction-giving task.

For clarity purposes, we summarize here below the chronicle of the various development and deployment steps that have been followed to run the direction-giving task and to refine the QoI Evaluator:

1. March 2018: beginning of the design and implementation of the direction-giving task
2. September 2018: First tests of the task on the field, i.e., in a Finnish mall
3. June 2019 and September 2019: New tests of the direction-giving task on the field
4. From September to December 2019 (project formal end): The robot autonomously ran three days a week in the mall (with only remote monitoring of the robot performance by our team for debugging and tuning)
 - (a) November 2019: Integration in the *Supervisor* of a preliminary version of Quality of Interaction Evaluator implementing the model described in [29] \Rightarrow version 1 of the QoI Evaluator
 - (b) From November 2019 to December 2019: Around 350 direction-giving tasks were performed with usual mall customers. Bug corrections and tuning of the direction-giving task. This allowed us to improve the QoI Evaluator thanks to: (i) data collection of task failures and standard durations of the subtasks executions (ii) lessons drawn about metric definitions and choices. \Rightarrow version 2 of the QoI Evaluator
5. March 2020: Refinement of the QoI Evaluator, i.e., improvement of the metric functions and tuning of their parameters. In the lab, with the same direction-giving task than the one used in the mall, comparison of the QoI computed by the robot when it is dealing with an “ideal” human, a “confused” human and a “non-compliant” human. \Rightarrow version 3 of the QoI Evaluator



(a) A customer listening to Pepper after re-positioning (b) A customer listening and Pepper pointing to a corridor



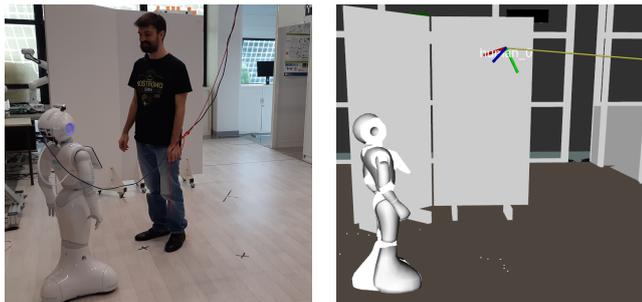
(c) A customer answering to Pepper (d) A customer listening and Pepper pointing to a shop

Fig. 6: MuMMER robot engaged in direction-giving tasks. Around 350 trials with customers in the mall allowed us to gather empirical data to select the metrics and tune the measuring functions parameters.

More specifically, this implementation of the Quality of Interaction Evaluator measured the interaction quality at the direction-giving task level and at the elementary actions level, omitting the interaction session level as this latter was not our focus in the MuMMER project. The QoI Evaluator was integrated into the *Supervisor* which is programmed using Jason [8], an agent-oriented framework. The supervision system handles the HR collaborative task execution through Jason reactive plans. The QoI Evaluator is implemented into a Jason function (the reasoning cycle) which is invoked periodically. After multiple testings, we reached the conclusion that it was pertinent, at least in the context of the direction-giving task, to have the Evaluator computing the QoI every second for both levels. Therefore, every second, the system computes the value of each metric and then outputs a value for QoI_{task} and QoI_{action} .

As mentioned in the step 4b of the chronicle, the robot interacted in the wild with dozens of usual customers (Fig. 6), executing around 350 direction-giving tasks. This allowed us to improve the performance of the direction-giving task, to gather standard durations of the subtasks executions and to draw lessons about metric definitions and choices (e.g. we realized it was not relevant to measure the human visual attention to-

wards the robot when it was giving the route explanation as humans look around at this moment). Unfortunately, the practical conditions of the project deployments did not offer us the possibility to evaluate the QoI Evaluator based on a study in the mall with real customers. So, we demonstrated – after improvements of the metrics equations such as the Distance-to-Goal one, and manual tuning of their parameters based on the experience in the mall – our finalized concept through tests in our lab (step 5). This is shown in Sect. 6.3 where we present and discuss, a comparison of the QoI computed by the robot when it is dealing with an “ideal” human, a “confused” human and a “non-compliant” human during a direction-giving task, performed in the lab. Before that, we present in Sect. 6.1 and Sect. 6.2 how the QoI is evaluated at both task and action levels for the direction-giving task.



(a) Initial positions of the human and the robot. The human asked the robot direction to reach a place behind him.



(b) The robot and the human are in their final positions as planned by the robot. The blue spheres are the computed position for the human by the robot. The robot is pointing at the place direction.

Fig. 7: Initial and final positions of a direction-giving task in the lab context. On the left are pictures and on the right screenshots of Rviz.

6.1 QoI Evaluation at the task level

The direction-giving task is triggered when a human asks for a location. In order to guide, the robot will choose and then point and give an explanation to reach the desired location. Specifically, it first computes two positions: one for itself and one for the human. These planned positions, once reached, will allow the human to properly see what the robot will be pointing [46]. Then, the robot navigates to its planned position and waits for the human to move in front. Since the human may not reach exactly the position expected by the planner, the robot checks again the visibility by the human after she/he has moved of the landmark it wants to point. In case the visibility is too low, the robot verbally asks them to adjust their position (i.e., “come closer”, “move back”). Fig. 7 illustrates the initial and final positions of both agents, in the lab context. Finally, it provides verbal explanation about how to reach their destination, in accordance with protocols identified thanks to a human-human guidance study conducted in the very same environment [5].

This task can be represented as a succession of sub-tasks, as shown in Fig. 8. This figure also exhibits the incremental refinement of the task into a sequence of HR interactive actions which are described in Sect. 6.2.

In the context of the direction-giving, we have selected two metrics to evaluate the QoI at the task level: a metric defined in the Sect. 5, the *Deviation from standard duration* and, the aggregation over time of the actions QoIs. Following the process of Fig. 2, we measure the QoI of the $\text{Task}_i = \text{direction-giving_task}$, based on the QoI of all task actions and $\text{Task metric}_1 = \text{Deviation from standard duration}$.

The *Deviation from standard duration* is used to measure the QoI at the task level as the task is a sequence of subtasks. Indeed, if the subtask lasts longer than expected, the QoI should decrease. Then, as needed for the metric computation we have determined the values of the soft deadlines SD_i for each subtask $a_i, i \in [0, 4]$, using the empirical data we gathered as explained in the introduction of the Sect. 6. Specifically, we have computed the average time execution of each subtask, after removing the cases for which the execution of the subtask was annotated as not smooth. These soft deadlines are presented in table 2. Finally, we chose $V_i = 0.5$ for all the subtasks.

The task QoI is also dependent on the actions QoI values (their computation is described in Sect. 6.2). Indeed, the actions QoIs should be reflected on the task QoI as, if a majority of the actions have a low QoI, the task QoI cannot remain high. That is why, besides the *Deviation from standard duration*, we take into account

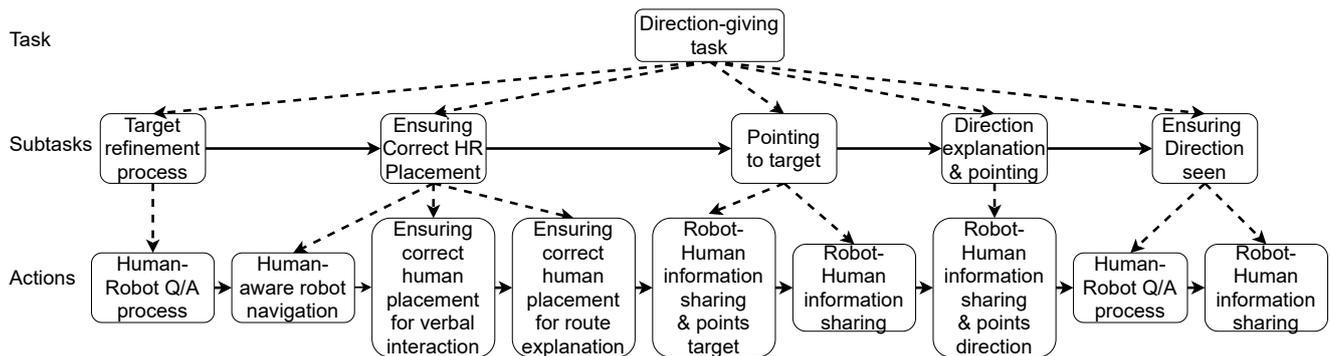


Fig. 8: The representation of the direction-giving task as a hierarchical task network with task, subtasks and actions levels. All the horizontal arrows are sequential links and the rest are decomposition ones.

Subtasks	soft deadline (s)
Target refinement process	30
Ensuring Correct HR Placement	30
Ensuring target seen	20
Direction explanation and pointing	30
Ensuring Direction Seen	20

Table 2: Soft deadlines SD_i for each subtask of the direction-giving task

the average of the action QoIs of the actions already executed or still running.

Then, the task QoI is computed using Equation (1) presented in Sect. 4. After various trials we have empirically chosen the weights W_i for each metric M_i , $i \in [0, 1]$. The final equation to compute the task QoI is:

$$QoI_{dir-giv.task}(t) = \frac{\Phi_{dir-giv.task}(t) + 3 * \overline{QoI}_{actions}}{4}$$

6.2 QoI Evaluation at the action level

As mentioned earlier, each subtask of the direction-giving task can be decomposed into actions. These actions involve several turn-taking steps, the robot asking complementary information, informing the human or expecting an action or reaction from them. We need to measure the QoI during the execution of each action. To do so, we have chosen one or more metrics for each action.

For each action of the following list, we explain which metrics M of Table 3 we have used and scaling functions of Appendix A and then, how we compute the action QoI.

- (a) *Robot-Human information sharing*: The robot speaks to the human, shares information such as the route direction and announces the next steps of the plan. The robot expects that they are paying attention

to it. Therefore, we use the *Fulfilling robot expectations about social interaction* M_{Exp-SI} based on the attention ratio. Two parameters need to be defined for the scaling function, the bounds b_1 and b_2 . As the minimum value for the metric, a ratio, is 0 and the maximum value is 1, then $b_1 = 0$ and $b_2 = 1$. The QoI of the action is computed with this only metric.

- (b) *Human-Robot Q/A process*: The robot asks a question to the human. As for the previous action, the robot expects the human to pay attention to it so we compute the QoI with M_{Exp-SI} . It also expects the human to give an appropriate answer. If it does not happen, it will ask the human to repeat, specifying that the answer has not been understood. We have limited the possible number of attempts to 3. After 3 attempts, the robot ends the task, as it cannot carry on with the task without an answer. So, we use *Human contribution to the goal* $M_{H-contrib}$, the number of times the robot repeats. Because the maximal number of repetitions is 3, we set for the scaling function $b_1 = 3$ and $b_2 = 0$.

The QoI is computed with the two metrics: *Fulfilling robot expectations about social interaction* and *Human contribution to the goal*. The trials showed that the action QoI results were satisfying with the weights $W_i = 1, i \in [0, 1]$ as applying the Equation (1).

- (c) *Ensuring that Human moves aside*: This action is used if, for pointing, the robot decides to place itself in a position which is very close to where the human is currently standing. In this case, the robot asks the human to step aside to the right or left, depending on the human's future position. Then, we want to measure the progress of the human going further from the planned robot position. In order to do this, we use the *Distance-to-Goal* M_{DtG} but with the condition of the ΔDtG equation adapted,

Metric id	Metric name	Metric equation – with Equations of Section 5	Scaled metric – with functions of Appendix A
$M_{H_contrib}$	Human contribution to the goal	nb_R_repet	$n_1(nb_R_repet) = 2 * \frac{nb_R_repet - 3}{-3} - 1$
M_{Exp_SI}	Fulfilling robot expectations about social interaction	$Ar = \frac{duration_{isAttentiveTo(robot)=true}}{duration_{robot_speaks}}$	$n_1(Ar) = 2 * Ar - 1$
M_{DtG}	Distance-to-Goal	$\begin{cases} \Delta DtG(t=0) = 0 \\ \Delta DtG(t) = \max(0, \Delta DtG(t-1) - 1) \\ \quad \text{if } path_length(t) < path_length(t-1) \\ \Delta DtG(t) = \Delta DtG(t-1) + 1, \text{ otherwise.} \end{cases}$	$-s_1(DtG(t)) = -1 + 2 \exp\left(-\ln(2) \left(\frac{DtG(t)}{5}\right)^{1.5}\right)$
M_{TtG}	Time-To-Goal	$\Delta TtG(t) = \max(0, e(t) + TtG(t) - TtG(T_0))$	$-s_1(TtG(t)) = -1 + 2 \exp\left(-\ln(2) \left(\frac{TtG(t)}{5}\right)^{1.5}\right)$

Table 3: Metrics used in the implementation presented in Section 6.

being if $path_length(t) > path_length(t-1)$ instead of if $path_length(t) < path_length(t-1)$. We scale the metric with $-s_1$, the additive inverse of the scaling function and not directly s_1 as the closer to 0 ΔDtG is, the better it is in terms of goal completion. From trials, we set $-s_1$ parameters values with $th = 5$ and $k = 1.5$.

If the human does not move or does not go far enough from the robot position, the robot will ask again with a limit of 3 trials (if the robot cannot move, it will carry on the task from their current positions). So, we use $M_{H_contrib}$ as for the previous action.

- (d) *Human-aware robot navigation*: The robot has to move from its initial position to its computed one. It navigates while respecting social constraints and its path may change as it adapts according to what the human is doing. At execution time, to measure the robot progress towards its goal, we use the *Time-to-goal* M_{TtG} , with the same scaling function than M_{DtG} . The QoI of the action is computed with this only metric.
- (e) *Ensuring correct human placement for verbal interaction*: After it has moved, the robot asks the human to come in front of it. If the human is not perceived after a few seconds, the robot will ask again and so on in a maximum of 3 trials. If after these 3 times the human is still not perceived, the robot ends the task.

The QoI of this action is computed with $M_{H_contrib}$ – we do not use M_{Exp_SI} as the human is not in the field of view when the robot is calling them.

- (f) *Ensuring correct human placement for route explanation*: Once the human is in the robot field of

view after the HR motion, they may not be at the right place to properly see what the robot has to point at. In this case, the robot will ask the human to move forward or backward according to what it has computed about the human perspective (e.g. this is to avoid that an object occludes the view for the human). Then, we want to measure the human progress towards the position the robot has computed for them. In order to do this, we use the *Distance-to-Goal* M_{DtG} .

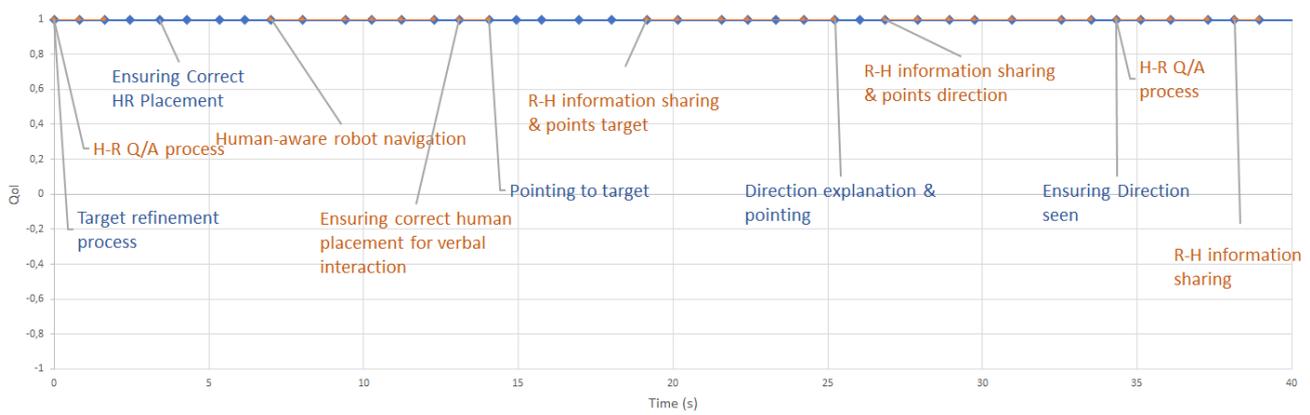
The robot stops giving instructions if it computes that the position of the human allows them to see the target, or after 3 trials, so we use $M_{H_contrib}$. After 3 trials, if the human cannot see the target, still, the robot will carry on the task taking this into account.

Mall elements	Mockup mall	Real mall
Shops	19	140
Doors, stairs, elevators	10	50
Corridors	11	41
Levels	2	2

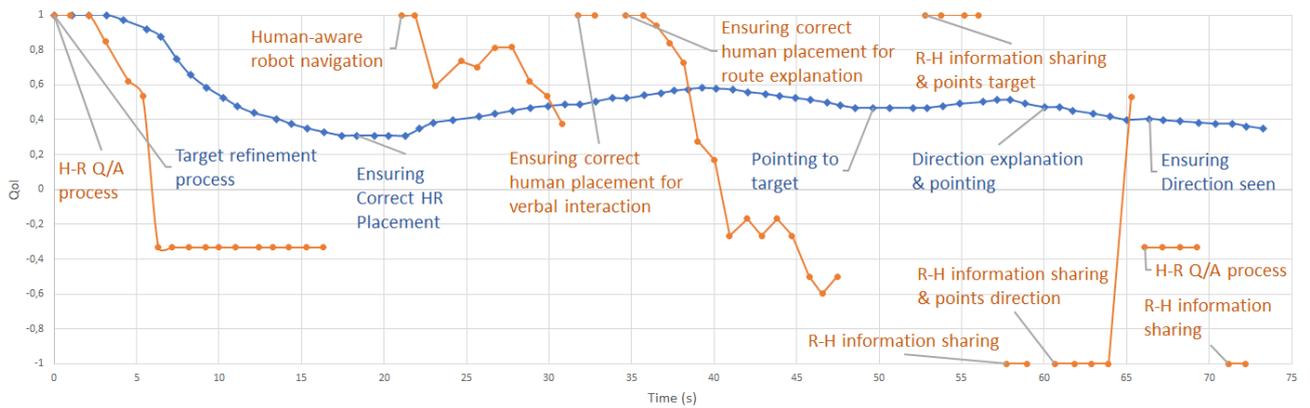
Table 4: Number of elements described in the mockup and real malls (geometric, topologic and semantic models in Fig. 4).

6.3 A first proof of concept

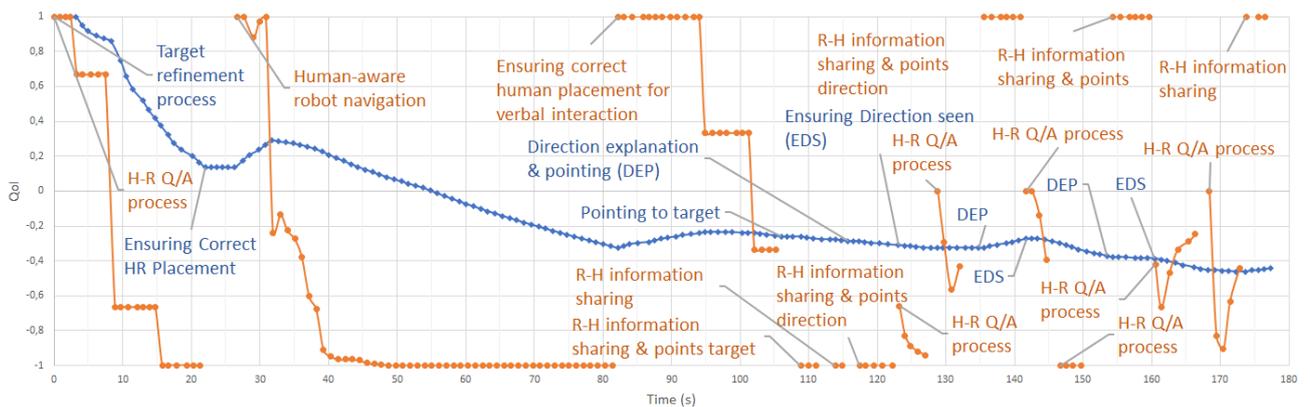
This section reports on an effective implementation of the approach as an illustrative proof of concept. We



(a) Evolution over time of the measured QoI for the 'ideal' human. Both action and task QoIs remain at 1 as the task is proceeding smoothly.



(b) Evolution over time of the measured QoI for the "confused" human. They took time to answer the first robot question and to move forward but the task QoI does not drop too much because the robot was able to give the route explanation without any issue even though the human was not very attentive.



(c) Evolution over time of the measured QoI for the non-compliant human. Several times the human did not give the expected answer to the robot during the target refinement process. Then, they blocked the robot path. After that, the robot had to ask twice the human to come in front of it. Finally, the robot repeated the route direction three times but still the human kept saying that they did not understand. Therefore, the task QoI decreases all along the task.

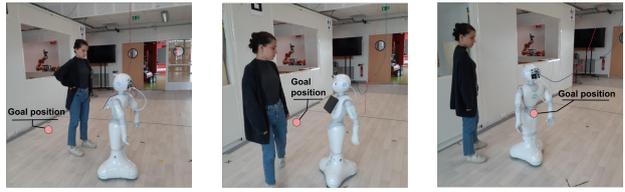
Fig. 9: Evolution over time of the measured QoI for the route guidance task with three different human behaviors. The QoI for the task is drawn in blue, and the QoI for the actions is drawn in orange.

Action	QoI formula (metric aggregation)
Robot-Human information sharing	$M_{Exp_SI}(t)$
Human-Robot Q/A process	$\frac{M_{Exp_SI}(t) + M_{H_contrib}(t)}{2}$
Ensuring that Human moves aside	$\frac{M_{DtG}(t) + M_{H_contrib}(t)}{2}$
Human-aware robot navigation	$M_{TtG}(t)$
Ensuring correct human placement for verbal interaction	$M_{H_contrib}(t)$
Ensuring correct human placement for route explanation	$\frac{M_{DtG}(t) + M_{H_contrib}(t)}{2}$

Table 5: QoI computation for each action as an aggregation of metrics

show the ability of the robot to conduct an interactive task, to assess in real-time the QoI and to track its evolution during three direction-giving task executions where the same human displayed a different way of behaving. In the three cases, the task was conducted until its end, in our lab where we reproduces the mall environment (Fig. 4b, Table 4). The computed QoI for each way is presented in Fig. 9. The three different ways of behaving are described in the following list:

- A human executed perfectly the expected actions and was not disturbing the robot when it navigated (i.e. the “ideal” human from the robot point of view).
- A bit “confused” human tried to contribute to the task success but did not execute everything well. The human was, from time to time, not very attentive, as looking around. Also, they gave an answer to the first question that the robot did not understand, and then they took their time before answering again. Then, they prevented a bit the robot to move as it had planned and once the robot reached its position, they took time to come as close as the robot wanted.
- A human wanted to disturb the robot during the task. They gave three incomprehensible answers to the first question, blocked multiple times the robot in its move, waited for the robot to ask twice to come in front of it and finally asked the robot to point and explain the route three times.



(a) Human who put themselves on the robot path, preventing the robot to navigate towards its goal position
 (b) Human who put themselves on the robot path after it computed a new path to reach its goal position
 (c) Human finally getting outside of the robot path, allowing it to reach its goal position

Fig. 10: A human disturbing the robot during *Human-aware navigation*, preventing it to reach its goal position as planned.

Now, if we take a look at the QoI outputs of Fig. 9, we can see that their three shapes are very different. In Fig. 9a, we can observe that the task and actions QoIs remain with the highest value 1 all along. A graph as this one allows us to infer that everything went very smoothly during this direction-giving task. Then, we can guess that it corresponds to the execution performed with the ‘ideal’ human.

In Fig. 9b, we note that each subtask was executed in respect of the standard duration. If the QoI of *Target refinement process* drops it is because of the action QoI as the QoI of the *H-R Q/A process* drops because the robot had to repeat the question and the human was not looking at it. From 21 seconds to 40 seconds, we can see the task QoI getting higher as the QoIs of *Human-aware robot navigation*, *Ensuring correct human placement for verbal interaction* and the beginning of *Ensuring correct human placement for route explanation* are quite high. Next, seeing the shape of the computed QoI of the action *Ensuring human placement for route explanation*, we can infer that the human was not moving as the robot wanted. Indeed, they took 10 seconds to make one step forward (they had 1 meter to cross). Because of that, the task QoI started to decrease again. In the final part of the task, the human was time to time attentive to the robot answered quickly to the last question, so the task QoI remained rather equal with its final value being 0.34 which is above 0 so meaning a correct interaction.

Finally, we can see in Fig. 9c that the final QoI of the task is -0.44 which allows us to infer that the task was not executed smoothly. And indeed, when we look at the shape of the task QoI, it only went down (or almost) all along the task. It is explained by some subtasks that took more time than they should have and also by some actions QoIs that are very low, especially the one of

Human-aware robot navigation. At the beginning of the robot navigation, the estimated time to goal returned by the planner was 6 seconds but the robot actually took 50 seconds to reach its goal then the action QoI computed with $= M_{T+G}(t)$ was -1 for 40 seconds. And indeed, all along its navigation, the human was blocking the robot until they got tired of this game, as visible on Fig. 10.

In this example, we showed the QoI evaluation process integrated to a complete robotic architecture. The robot was able to assess the QoI in real-time while interacting with a human.

7 Discussion

While a number of evaluation methods has been proposed to evaluate a human-robot interaction from the human perspective and often for analysis after performance, our choice to let the robot evaluate, on its own and in real-time the quality of its interaction with a human is quite new and original. To endow the robot with such an ability, we designed, implemented and tested a number of metrics and a method to aggregate them.

The work of Steinfeld *et al.* [42] was very helpful to design a first set of metrics and as an inspiration about what could be used. From there, we have elaborated and proposed a set of metrics which are meant to estimate of the quality of an ongoing interaction and not once it is over. The work of Hoffman [16] regarding the *fluency* definition and how to measure it was also inspiring. In a way, we extended his work by giving a meaning to the fluency measurement on the robot side, and in real-time – while their work applies to offline evaluation of shared workspace tasks. In Sect. 2, we mentioned systems measuring human affective states in real-time such as the framework developed by Tanevaska *et al* [44]. Although we think such metric could be an interesting additional information to assess if an interaction is going well, we believe that these measurements do not offer an accuracy that would lead to objective measurement of the quality of interaction, thus, we did not introduce them in our set for now. However, this could be done since our framework is designed to be open to new metrics. As for contributions, like the one proposed by Anzalone *et al.* [1], based on metrics such as gaze, head pose, body pose and response times to measure real-time engagement, we took them into account to some extent. However, the measure of the engagement that we propose should be refined depending on the inputs available online to the robot. Moreover, we will investigate how their work could be used in a more general way (e.g. depending on the action that should be done and its context, human head pose and body posture could be

a good indicator of effectiveness and not only engagement).

Our intention, when we developed the idea of the Quality of Interaction Evaluation, was to use such computation to feed the decision-making process of the robot and this is what we intend to do in the future. However, such framework can also be used to compare interactions between different humans and/or robots, eventually as a benchmark similarly to the work of Sanchez-Matilla [37] or as a way for developers to detect repetitive interaction issues with an unsupervised robot in a real-world environment.

As a proof-of-concept, we implemented and deployed a first version of a QoI Evaluator assessing task and actions QoIs. We tested it on an interactive robot dedicated to provide route guidance to customers in a large mall. The approach gave satisfactory results. It showed the potential ability of a robot to detect momentary decreases of the Quality of Interaction and also more serious degradation of it which may need drastic change of behavior for the robot. This is only a first step and it should be validated with a study where we will ask humans to evaluate the quality of their interaction with the robot in a similar manner. The goal will be to analyse and compare this to the evaluation of the interaction quality estimated by our robot and, based on that, investigate potential improvements.

Finally, we do not claim to have a perfect measure of the Quality of Interaction. However, although the concept of Quality of Interaction is quite abstract, Movellan *et al.* showed that when it is measured by human observers, the inter-observer reliability of the concept is quite high. Therefore, we believe we can endow the robot with an effective and pertinent ability aiming at measuring the quality of an interaction. We are aware that the set of metrics we proposed to do so is not exhaustive but the framework is designed to be easily extended with new metrics.

8 Conclusion and future work

We claim that the robot could enhance its decision-making process by estimating if an interaction is going well or not. To endow it with this ability, we have proposed an original framework where the Quality of Interaction is measured from the robot point of view in real-time during its collaborative activities. We proposed in this paper a set of metrics and a method to aggregate them.

The evaluation of the QoI relies on the model of interaction, considered at three levels: the interaction session level, the tasks level and the actions level. In

future work, this granularity will allow the robot to know precisely on what level it needs to act when a low QoI is computed. Indeed, for instance, a task can be of poor quality while the session can be considered as going well.

Therefore, we intend to exploit this QoI evaluation process in order to allow the robot to “close the loop” and smoothly adapt its decisions and execution modalities and also to detect if the human partner is trying to “pull the robot strings”. Then, our next step will be to refine our set of metrics and to expand it. For example, we plan to investigate the possibility to elaborate plan-based algorithms in order to track the evolution of such the *Human contribution to the goal* over time. Finally, we will test and improve the metrics dedicated to the interaction session level.

A Appendix: Scaling functions for the metrics

As the metrics are aggregated to compute the QoI, their values need to be on the same scale. In order to do this, we use scaling functions rescaling metrics into a range of $[-1, 1]$, as the QoI bounds. As all the metrics does not have the same properties, they have to be scaled by using different functions. The two properties to check to choose which function to apply to which metric are the following ones:

- does the metric already have a bounded value ?
- what value of the metric should make the QoI decrease, increase or remain the same ?

Therefore, we designed three functions to be used with metrics having bounded values and three functions for metrics that do not have upper bounds. Then, among these two sets of functions, it is possible to choose the one to use according to the positive, neutral or negative impact a value should have on the QoI.

A.1 Scaling of bounded metrics: Min-Max Normalization

We defined three min-max normalization functions, illustrated in Fig. 11. They were designed to be used for metrics whose values belong to a bounded set, i.e., metrics for which the minimum and maximum values are known. The first function is to apply in cases for which a measure approaching the bound value b_1 has a negative impact on the quality evaluation whereas a measure approaching b_2 has a positive one. It allows to scale a measure x between -1 and 1:

$$n_1(x) = 2 * \frac{x - b_1}{b_2 - b_1} - 1 \quad (7)$$

The second function is intended to be applied in cases for which a measure approaching the bound value b_1 has a neutral impact on the quality evaluation whereas a measure approaching b_2 has a positive one. It allows to scale a measure x between 0 and 1:

$$n_2(x) = \frac{x - b_1}{b_2 - b_1} \quad (8)$$

Finally, the last function is to apply in cases for which a measure approaching the bound value b_1 has a negative impact on the quality evaluation whereas a measure approaching b_2 has a neutral one. It allows to scale a measure x between -1 and 0:

$$n_3(x) = \frac{x - b_2}{b_2 - b_1} \quad (9)$$

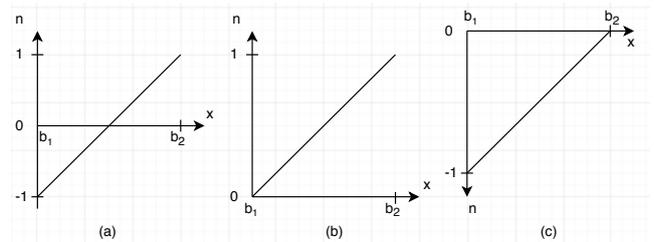


Fig. 11: (a), (b) and (c) respectively represent the min-max normalization functions (7), (8) and (9)

A.2 Scaling of unbounded metrics: Sigmoid Normalization

We defined three sigmoid-like functions to scale and squash values of metrics without an upper bound. As for the min-max normalization, there is one function to scale the metrics values between -1 and 1, another one to scale between 0 and 1 and the last one to scale between -1 and 0.

The first function allows to scale between -1 and 1 the values of a metric, for a metric whose values are between 0 and $+\infty$ (e.g. a duration whose final value is unknown during the execution). The function is defined as:

$$s_1(x) = 1 - 2 \exp\left(-\ln(2) \left(\frac{x}{th}\right)^k\right), x > 0 \quad (10)$$

with $s_1(x) \in [-1, 1]$, th the value of the sigmoid’s midpoint (i.e., $s_1(th) = 0$) and, k setting the shape of the function curve. k and th values are set off-line by the designer and they allow to define the shape of the metric scaling.

The second function is designed for metric which cannot have a negative impact on the QoI as it scales the value between 0 and 1 (and with $x \in [0, +\infty]$ as well):

$$s_2(x) = 1 - \exp\left(-\ln(2) \left(\frac{x}{th}\right)^k\right), x > 0 \quad (11)$$

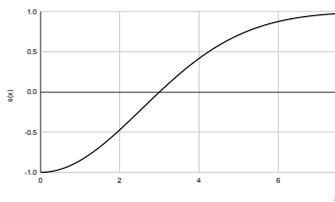
with $s_2(x) \in [0, 1]$, th the value of the sigmoid’s midpoint (i.e., $s_2(th) = 0.5$) and, k setting the shape of the function curve.

The third function is designed for metric which cannot have a positive impact on the QoI as it scales the value between -1 and 0 (and with $x \in [0, +\infty]$ as well):

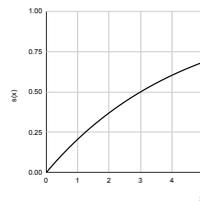
$$s_3(x) = -1 + \exp\left(-\ln(2) \left(\frac{x}{th}\right)^k\right), x > 0 \quad (12)$$

with $s_3(x) \in [-1, 0]$, th the value of the sigmoid’s midpoint (i.e., $s_3(th) = -0.5$) and, k setting the shape of the function curve.

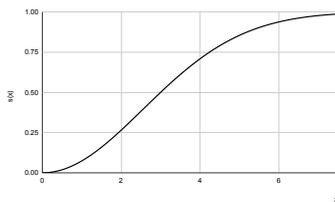
The functions $s_1(x)$ and $s_2(x)$ are illustrated in Fig. 12 with four examples.



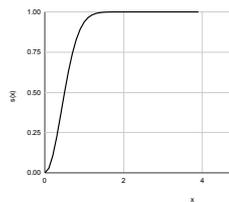
(a) Plot of $s_1(x)$ with $th = 3$ and $k = 2$



(b) Plot of $s_2(x)$ with $th = 3$ and $k = 1$



(c) Plot of $s_2(x)$ with $th = 3$ and $k = 2$



(d) Plot of $s_2(x)$ with $th = 0.5$ and $k = 2$

Fig. 12: Plots of the sigmoid-like functions $s_1(x)$ and $s_2(x)$ with different parameters values

Acknowledgements Many thanks to Michaël Mayer for his technical expertise and his help on the mathematical formalization.

Declarations

This paper or a similar version is not currently under review by a journal or conference. This paper is void of plagiarism or self-plagiarism as defined by the Committee on Publication Ethics and Springer Guidelines.

Funding

This work has been supported by the European Union's Horizon 2020 research and innovation program under grant agreement No. 688147 (MuMMER project), and by the French National Research Agency (ANR) under grant references ANR-16-CE33-0017 (JointAction4HRI project), and ANR-19-PI3A-0004 (ANITI).

Conflicts of interest/Competing interests

The authors have no conflicts of interest to declare that are relevant to the content of this article.

Availability of data and material

Not applicable.

Code availability

The code of the Quality of Interaction Evaluator is available at https://github.com/amdia/guiding_task.

References

1. Anzalone SM, Boucenna S, Ivaldi S, Chetouani M (2015) Evaluating the engagement with social robots. *International Journal of Social Robotics* vol. 7(4):pp. 465–478, DOI 10.1007/s12369-015-0298-7
2. Baraglia J, Cakmak M, Nagai Y, Rao RPN, Asada M (2017) Efficient human-robot collaboration: When should a robot take initiative? *International Journal of Robotics Research* vol. 36(5-7):pp. 563–579
3. Bauer A, Wollherr D, Buss M (2008) Human-robot collaboration: A survey. *International Journal of Humanoid Robotics* vol. 5(01):pp. 47–66, DOI 10.1142/S0219843608001303, <https://doi.org/10.1142/S0219843608001303>
4. Bekele E, Sarkar N (2014) Psychophysiological feedback for adaptive human-robot interaction (hri). In: Fairclough SH, Gilleade K (eds) *Advances in Physiological Computing*, Springer London, pp pp. 141–167, DOI 10.1007/978-1-4471-6392-3_7
5. Belhassein K, Clodic A, Cochet H, Niemelä M, Heikkilä P, Lammi H, Tammela A (2017) Human-Human Guidance Study. Hal-01719730
6. Bensch S, Jevtić A, Hellström T (2017) On interaction quality in human-robot interaction. In: *Proceedings of the 9th International Conference on Agents and Artificial Intelligence (ICAART)*, pp pp. 182–189, DOI 10.5220/0006191601820189
7. Bethel CL, Murphy RR (2010) Review of human studies methods in hri and recommendations. *International Journal of Social Robotics* vol. 2(4):pp. 347–359, DOI 10.1007/s12369-010-0064-9
8. Bordini RH, Hübner JF, Wooldridge M (2007) *Programming Multi-Agent Systems in AgentSpeak Using Jason* (Wiley Series in Agent Technology). John Wiley & Sons, Inc.
9. Devin S, Alami R (2016) An implemented theory of mind to improve human-robot shared plans execution. In: *The Eleventh ACM/IEEE International Conference on Human Robot Interaction (HRI)*, Christchurch, New Zealand, pp pp. 319–326
10. Fan J, Bian D, Zheng Z, Beuscher L, Newhouse PA, Mion LC, Sarkar N (2017) A robotic coach architecture for elder care (rocare) based on multi-user engagement models. *IEEE Transactions on Neural Systems and Rehabilitation Engineering* vol. 25(8):pp. 1153–1163
11. Foster ME, Craenen B, Deshmukh A, Lemon O, Bastianelli E, Dondrup C, et al (2019) Mummer: Socially intelligent human-robot interaction in public spaces. In: *AAAI 2019 Fall Symposium Series*, Arlington, United States, 1909.06749
12. Ghallab M, Knoblock C, Wilkins D, Barrett A, Christianson D, Friedman M, et al (1998) PDDL - The Planning Domain Definition Language
13. Ghallab M, Nau DS, Traverso P (2016) *Automated Planning and Acting*. Cambridge University Press
14. Grosz BJ, Kraus S (1996) Collaborative plans for complex group action. *Artificial Intelligence* vol. 86(2):pp. 269–357, DOI [https://doi.org/10.1016/0004-3702\(95\)00103-4](https://doi.org/10.1016/0004-3702(95)00103-4)

15. Hiatt LM, Narber C, Bekele E, Khemlani SS, Trafton JG (2017) Human modeling for human-robot collaboration. *International Journal of Robotics Research* vol. 36(5-7):pp. 580–596, DOI 10.1177/0278364917690592
16. Hoffman G (2019) Evaluating fluency in human-robot collaboration. *IEEE Transactions on Human-Machine Systems* 49(3):pp. 209–218
17. Hoffman G, Breazeal C (2007) Cost-based anticipatory action selection for human-robot fluency. *IEEE Transactions on Robotics* vol. 23(5):pp. 952–961
18. Ingrand F, Ghallab M (2017) Deliberation for autonomous robots: A survey. *Artificial Intelligence* vol. 247:pp. 10–44, DOI 10.1016/j.artint.2014.11.003
19. Itoh K, Miwa H, Nukariya Y, Zecca M, Takanobu H, Roccella S, et al (2006) Development of a bioinstrumentation system in the interaction between a human and a robot. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Beijing, China, pp pp. 2620–2625, DOI 10.1109/IROS.2006.281941
20. Khambhaita H, Alami R (2020) Viewing robot navigation in human environment as a cooperative activity. In: Amato NM, Hager G, Thomas S, Torres-Torriti M (eds) *Robotics Research*, Springer International Publishing, pp pp. 285–300
21. Khambhaita H, Alami R (2020) Viewing robot navigation in human environment as a cooperative activity. In: *Robotics Research*, Springer, pp 285–300
22. Kruse T, Pandey AK, Alami R, Kirsch A (2013) Human-Aware Robot Navigation: A Survey. *Robotics and Autonomous Systems* vol. 61(12):pp. 1726–1743
23. Kulić D, Croft EA (2003) Estimating intent for human-robot interaction. In: *IEEE International Conference on Advanced Robotics*, pp pp. 810–815
24. Kulic D, Croft EA (2007) Affective state estimation for human-robot interaction. *IEEE Transactions on Robotics* vol. 23(5):pp. 991–1000
25. Lallement R, De Silva L, Alami R (2014) HATP: An HTN Planner for Robotics. In: *2nd ICAPS Workshop on Planning and Robotics*, Portsmouth, United States
26. Lemaignan S, Garcia F, Jacq A, Dillenbourg P (2016) From real-time attention assessment to “with-me-ness” in human-robot interaction. In: *11th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pp pp. 157–164
27. Lemaignan S, Warnier M, Sisbot EA, Clodic A, Alami R (2017) Artificial Cognition for Social Human-Robot Interaction: An Implementation. *Artificial Intelligence* vol. 247:pp. 45–69
28. Lemaignan S, Warnier M, Sisbot EA, Clodic A, Alami R (2017) Artificial cognition for social human-robot interaction: An implementation. *Artificial Intelligence* 247:45–69, DOI 10.1016/j.artint.2016.07.002, special Issue on AI and Robotics
29. Mayima A, Clodic A, Alami R (2019) Evaluation of the Quality of Interaction from the robot point of view in Human-Robot Interactions. In: *1st Edition of Quality of Interaction in Socially Assistive Robots (QISAR) Workshop*, The 11th International Conference on Social Robotics (ICSR 2019), Madrid, Spain, URL <https://hal.laas.fr/hal-02403081>
30. Mayima A, Clodic A, Alami R (2020) Toward a robot computing an online estimation of the quality of its interaction with its human partner. In: *2020 29th IEEE International Conference on Robot and Human Interactive Communication (RO-MAN)*, pp 291–298, DOI 10.1109/RO-MAN47096.2020.9223464
31. Michael J, Salice A (2017) The sense of commitment in human-robot interaction. *International Journal of Social Robotics* vol. 9(5):pp. 755–763, DOI 10.1007/s12369-016-0376-5
32. Michael J, Sebanz N, Knoblich G (2016) The sense of commitment: A minimal approach. *Frontiers in Psychology* vol. 6:1968, DOI 10.3389/fpsyg.2015.01968
33. Milliez G, Warnier M, Clodic A, Alami R (2014) A framework for endowing an interactive robot with reasoning capabilities about perspective-taking and belief management. In: *The 23rd IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*, Edinburgh, United Kingdom, pp pp. 1103–1109, DOI 10.1109/ROMAN.2014.6926399
34. Olsen DR, Goodrich MA (2003) Metrics for evaluating human-robot interaction. In: *PERMIS*, Gaithersburg, United States
35. Robinson JD (2012) Overall structural organization. In: Sidnell J, Stivers T (eds) *The Handbook of Conversation Analysis*, John Wiley & Sons, Ltd, pp pp. 257–280, DOI 10.1002/9781118325001.ch13
36. Sallami Y, Lemaignan S, Clodic A, Alami R (2019) Simulation-based physics reasoning for consistent scene estimation in an hri context. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp 1–8
37. Sanchez-Matilla R, Chatzilygeroudis K, Modas A, Duarte NF, Xompero A, Frossard P, Billard A, Cavallaro A (2020) Benchmark for human-to-robot handovers of unseen containers with unknown filling. *IEEE Robotics and Automation Letters* 5(2):1642–1649, DOI 10.1109/LRA.2020.2969200
38. Sarthou G, Alami R, Clodic A (2019) Semantic Spatial Representation: a unique representation of an environment based on an ontology for robotic applications. In: *SpLU-RoboNLP*, pp 50–60
39. Schegloff EA, Sacks H (1973) Opening up closings. *Semiotica* vol. 8(4):pp. 289–327
40. Sidner CL, Lee C (2003) Engagement rules for human-robot collaborative interactions. In: *IEEE International Conference on Systems, Man and Cybernetics (SMC)*, Washington DC, United States, pp pp. 3957–3962, DOI 10.1109/ICSMC.2003.1244506
41. Singamaneni PT, Alami R (2020) HATEB-2: Reactive Planning and Decision making in Human-Robot Co-navigation. In: *International Conference on Robot & Human Interactive Communication*, 2020, Online, Italy, DOI 10.1109/RO-MAN47096.2020.9223463
42. Steinfeld A, Fong T, Kaber D, Lewis M, Scholtz J, Schultz A, Goodrich M (2006) Common metrics for human-robot interaction. In: *Proceedings of the 1st ACM SIGCHI/SIGART Conference on Human-robot Interaction*, Salt Lake City, United States, pp pp. 33–40, DOI 10.1145/1121241.1121249
43. Tabrez A, Luebbbers MB, Hayes B (2020) A survey of mental modeling techniques in human-robot teaming. *Current Robotics Reports* DOI 10.1007/s43154-020-00019-0
44. Tanevska A, Rea G, Fand Sandini, Sciutti A (2017) Towards an Affective Cognitive Architecture for Human-Robot Interaction for the iCub Robot. In: *1st Workshop on “Behavior, Emotion and Representation: Building Blocks of Interaction”*, Bielefeld, Germany
45. Thomaz A, Hoffman G, Çakmak M (2016) Computational human-robot interaction. *Foundations and Trends in Robotics* vol. 4(2-3):pp. 105–223, DOI

10.1561/23000000049

46. Waldhart J, Clodic A, Alami R (2019) Reasoning on shared visual perspective to improve route directions. In: IEEE International Conference on Robot and Human Interactive Communication (RO-MAN), IEEE, pp 1–8