

# MuMMER: Socially Intelligent Human-Robot Interaction in Public Spaces

Mary Ellen Foster,<sup>1</sup> Bart Craenen,<sup>1</sup> Amol Deshmukh,<sup>1</sup>  
 Oliver Lemon,<sup>2</sup> Emanuele Bastianelli,<sup>2</sup> Christian Dondrup,<sup>2</sup> Ioannis Papaioannou,<sup>2</sup> Andrea Vanzo,<sup>2</sup>  
 Jean-Marc Odobez,<sup>3</sup> Olivier Canévet,<sup>3</sup> Yuanzhouhan Cao,<sup>3</sup> Weipeng He,<sup>3</sup>  
 Angel Martínez-González,<sup>3</sup> Petr Motlicek,<sup>3</sup> Rémy Siegfried,<sup>3</sup>  
 Rachid Alami,<sup>4</sup> Kathleen Belhassein,<sup>4</sup> Guilhem Buisan,<sup>4</sup> Aurélie Clodic,<sup>4</sup> Amandine Mayima,<sup>4</sup>  
 Yoan Sallami,<sup>4</sup> Guillaume Sarthou,<sup>4</sup> Phani-Teja Singamaneni,<sup>4</sup> Jules Waldhart,<sup>4</sup>  
 Alexandre Mazel,<sup>5</sup> Maxime Caniot,<sup>5</sup>  
 Marketta Niemelä,<sup>6</sup> Päivi Heikkilä,<sup>6</sup> Hanna Lammi,<sup>6</sup> Antti Tammela<sup>6</sup>  
<sup>1</sup> University of Glasgow, Glasgow, UK <sup>2</sup> Heriot-Watt University, Edinburgh, UK  
<sup>3</sup> Idiap Research Institute, Martigny, Switzerland <sup>4</sup> LAAS-CNRS, Toulouse, France  
<sup>5</sup> SoftBank Robotics Europe, Paris, France <sup>6</sup> VTT Technical Research Centre of Finland, Tampere, Finland  
 MaryEllen.Foster@glasgow.ac.uk

## Abstract

In the EU-funded MuMMER project, we have developed a social robot designed to interact naturally and flexibly with users in public spaces such as a shopping mall. We present the latest version of the robot system developed during the project. This system encompasses audio-visual sensing, social signal processing, conversational interaction, perspective taking, geometric reasoning, and motion planning. It successfully combines all these components in an overarching framework using the Robot Operating System (ROS) and has been deployed to a shopping mall in Finland interacting with customers. In this paper, we describe the system components, their interplay, and the resulting robot behaviours and scenarios provided at the shopping mall.

## Introduction

In the EU-funded MuMMER project (<http://mummer-project.eu/>), we have developed a socially intelligent interactive robot designed to interact with the general public in open spaces, using SoftBank Robotics' Pepper humanoid robot as the primary platform (Foster et al. 2016). The MuMMER system provides an entertaining and engaging experience to enrich a human-robot interaction. Crucially, our robot exhibits behaviour that is *socially appropriate* and *engaging* by combining speech-based conversational interaction with non-verbal communication, and motion planning. To support this behaviour, we have developed and integrated new methods from audiovisual scene processing, social-signal processing, conversational AI, perspective taking, and geometric reasoning.

The primary MuMMER deployment location is Ideapark, a large shopping mall in Lempäälä, Finland. The MuMMER robot system has been taken to the shopping mall several times for short-term co-design activities with the mall customers and retailers (Heikkilä, Lammi, and Belhassein 2018; Heikkilä et al. 2019); the full robot system has been deployed for short periods in the mall in September 2018 (Figure 1), May 2019, and June 2019, and has been installed for a long-term, three-month deployment as of September 2019.



Figure 1: The MuMMER robot system interacting with a customer in the Ideapark shopping mall, September 2018.

The demo system supports a range of behaviours covering a variety of functional and entertainment tasks that are appropriate for a shopping-mall setting, including guidance to various locations within the mall, small-talk, and playing quiz games with customers. The activities during the deployment have included a number of data collection studies with real users: recording of customer interaction with the robot in guidance situations, sound localisation and automatic speech recognition in the noisy mall environment, and tests for AI-based conversation and localisation and navigation based on a partial 3D model of the mall and a complete semantic model.

In the remainder of this paper, we outline the technical contributions in each of the main MuMMER component areas: audiovisual sensing, social signal processing, conversational interaction, human-aware robot motion planning, knowledge representation and decision. At the end, we describe the details of the deployed robot system.

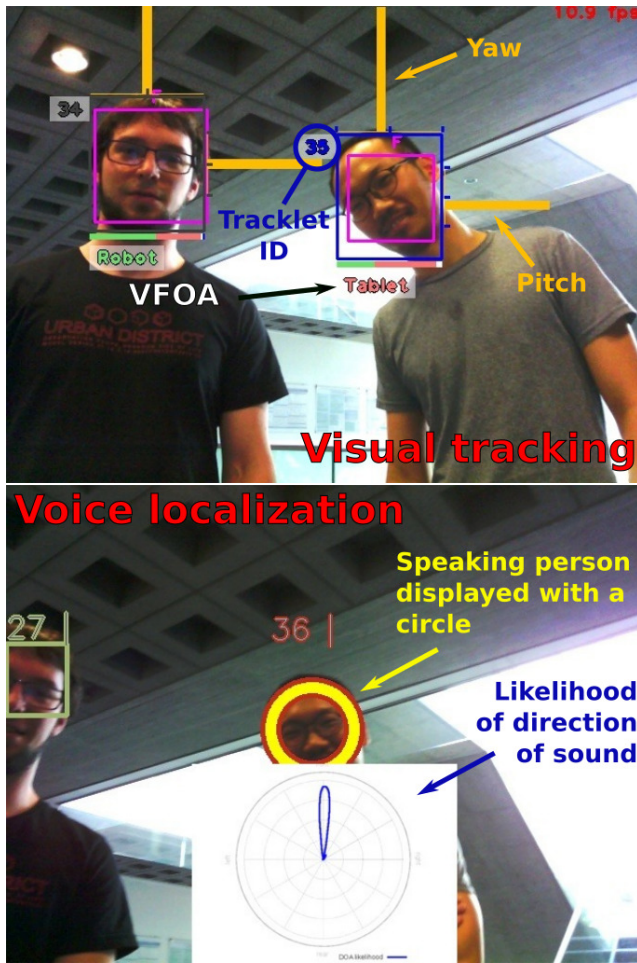


Figure 2: The perception systems tracks and re-identifies people leaving the field, and extracts other features: speaker turns, head pose, visual focus of attention and nods.

### Audio-visual sensing

For MuMMER, the main task of audio-visual perception is sensing people in general – that is, maintaining a representation of the persons around the robot, with a dedicated attention to people susceptible of interacting with it, or those who are (or have been) interacting with it. This requires several audio-visual algorithms to detect, track, re-identify people, and detect their non-verbal behaviors and activities, and also predict their position/behaviours even when they are not seen. At the same time, the representation of people needs to be defined and shared with other modules which are responsible for inferring other knowledge about people (for instance, to define a person’s goal in the interaction).

For visual tracking, we first detect the person with the convolutional pose machines (CPM) (Cao et al. 2017) which provide accurate locations of the body joints (nose, eyes, shoulders, etc.). The output of the algorithm is almost perfect when people are in the foreground of the image and up to 3 of 5 meters, depending on the resolution. This is our use-case definition of an entertainment robot in a shopping

mall. On top of the CPM, we use OpenHeadPose<sup>1</sup> (Cao, Canévet, and Odobez 2018), which makes use of the heatmap of the CPM to estimate the head pose of the person. Then, we perform head pose tracking (Khalidov and Odobez 2017) to maintain a consistent identity across adjacent frames (Figure 2). The head pose tracker is represented by a particle filter mainly based on color and face cues. As faces are tracked, we store OpenFace features (Amos, Ludwiczuk, and Satyanarayanan 2016) which are computed on an aligned face. When a new tracklet is created, the OpenFace features of the new tracklet are compared with the features previously accumulated, and the new tracklet is re-assigned the identity of the one which had the more votes.

For sound localization, we use a multi-task neural network (NN) which jointly performs speech/non-speech detection and sound source localisation (He, Motlicek, and Odobez 2018a; He, Motlicek, and Odobez 2018b) applied on top of the 4-channel microphone array (embedded on the robot). The NN uses as an input a 4-channel audio transformed into a frequency domain, and it outputs the likelihood values for the two tasks. Thanks to a semi-automated and synthetic data collection procedure taking advantage of the robotic platform as well as the use of a weak supervision learning approach, it is possible to quickly collect data to learn the models for a new sensor (He, Motlicek, and Odobez 2019). The fusion between the visual and audio parts is done by assigning the detected speech to the person who is standing in the given direction.

Finally, although a close range (up to 1.2m) gaze sensing module is available and can be applied for one selected person using a self-calibrated approach (Siegfried, Yu, and Odobez 2017), as a compromise between computation and robustness, we instead compute the visual focus of attention of each person based on the head pose (Sheikhi and Odobez 2015). The algorithm can reliably estimate the object the person is looking at (either the robot, the other persons, the targets, the shops, or the tablet embedded at the robot) which is a preliminary step to identify the addressee, and is also used in the context of perspective taking to determine whether the human has looked in the direction where the robot pointed (Sallami et al. 2019).

### Social signal processing

For Social Signal Processing, we focus on two primary tasks: fusing the provided audio-visual sensing data for social state estimation, and synthesising appropriate social signals for the robot to use when communicating with users. While detecting, tracking, (re)identifying users, as well as detecting their primary non-verbal behaviours and activities provide the basic signals, the multi-modal fusion of these signals allows for a more accurate and deeper understanding of the underlying social state, including gaining personality impressions from the user. The estimated social state is then made available to inform planning of the robot’s subsequent actions; who and how to converse with the users of the robot; and, how the robot is to move and behave (gestures) in the presence of the user(s).

<sup>1</sup><https://gitlab.idiap.ch/software/openheadpose>

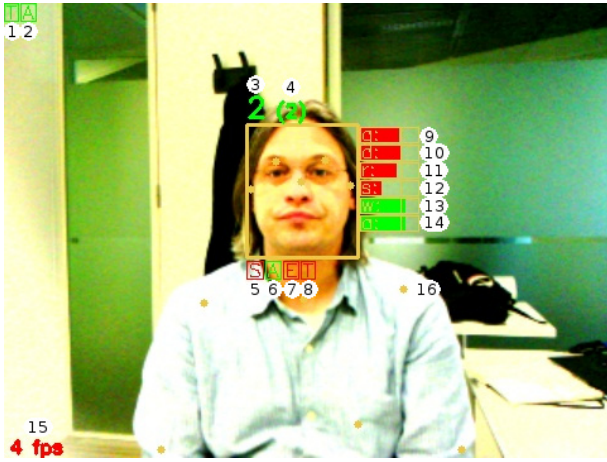


Figure 3: Social state estimator visualiser output, displaying all relevant information for fine-tuning.

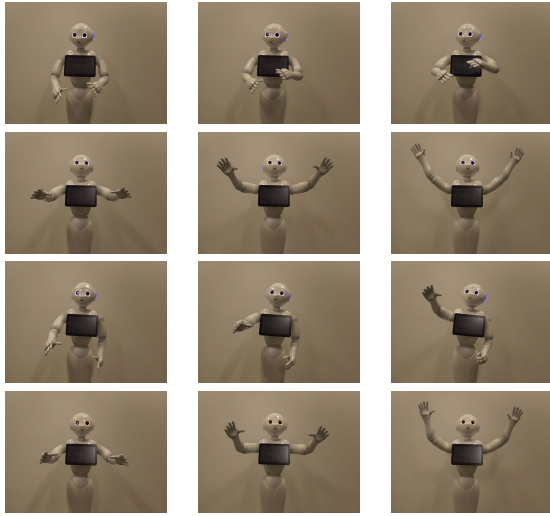


Figure 4: Gesture variations. Each row shows the same gesture using different parameters.

## Social state estimation

On the fusion side, the main function of the social state estimator is to determine which user the robot should initiate interaction with. We used the underlying assumption that the robot should initiate interaction with the user perceived to be the most willing to interact; which we took to be the user paying the most attention to the robot. We assume that the user paying the most attention to the robot is the user that is looking most directly at the robot, and who is most closely situated near the robot.

To this end, the social state estimator aggregates audio-visual sensing data about the head pose of users, whether the users are looking at the robot and/or the screen on the robot, and the distance between the users’ head and the robot. The head pose data of the users are used to calculate the (Euclidean) distance between the head pose of the users

and three centroids derived from clustering/classifying previously recorded lab and deployment data. This distance is then normalised to a value between zero and one, and used as a probability. The distances between the users and the robot are used as a penalty, and normalised between zero and one as a probability in such a way that users further away from the robot are penalised more than users closer to the robot.

This results in four probabilities, two taken directly from the audio-visual sensing data, and two derived from it. These four probabilities are then fused into one *attention* probability by calculating their (weighted) average. Choosing which user the robot should interact with is done by comparing the attention probability against a configurable minimum attention probability threshold, and then selecting the user with the highest attention probability.

To prevent immediately re-initiating an interaction with a user that the robot has just interacted with, the social state estimator also monitors the actions of the planning and dialogue components. The social state estimator then maintains a list of the users that the robot is, and has been, interacting with, and applies a penalty to their attention probability while they are interacting and for a short time afterwards.

The social state estimator is fully configurable by a set of parameters, with the initial parameter settings determined from extensive recorded lab data. The parameters were further fine-tuned during deployment to provide the most accurate and applicable social state estimates. To facilitate fine-tuning the parameters of social state estimator, a visualiser is provided to display all relevant features (Figure 3).

## Social signal generation

For the synthesis side, a repertoire of non-verbal social signals, including gestures and sounds, has been developed for the robot, available to be used in conjunction with moving and interacting with the users. The non-verbal behaviour of an embodied agent is at least as communicative as its verbal behaviour (Vinciarelli, Pantic, and Bourlard 2009), and in a noisy environment such as a mall it may even be more important, so understanding and controlling the robot's non-verbal signals is crucial. Examples of some robot gesture variants are shown in Figure 4.

In a series of perception experiments, we have examined how manipulating gesture parameters affect users' subjective responses to the robot as well as their perception of the robot's personality. These studies have found several clear relationships: for example, manipulating the amplitude and speed significantly affected users' perception of the Extraversion and Neuroticism of the robot, while the attributed personality also affected users' subjective reactions to the robot (Craenen et al. 2018b). In addition, it was found that while the majority of users preferred a robot that they perceived to have a similar personality to their own, a significant minority preferred a robot whose personality was perceived to be different than their own (Craenen et al. 2018a).

We are currently integrating a finer-grained method of gesture control based on sentiment (Deshmukh, Foster, and Mazel 2019), as well as a set of affectively generated artificial sounds (Hastie et al. 2016), with the goal of further enhancing the robot’s expressiveness.



## Conversational interaction

The MuMMER system focuses on enabling an agent to combine a task-based dialogue system with chat-style open-domain social interaction, to fulfil the required tasks while at the same time being natural, entertaining, and engaging to interact with. The presented work is based on the “Alana” conversational framework, a finalist of the Amazon Alexa Challenge in both 2017 (Papaioannou et al. 2017) and 2018 (Curry et al. 2018). Alana was initially developed for the Amazon Echo as an open-domain social chatbot. For the needs of this project, Alana acts as the core module for every dialogue interaction with the user from every other module. This means that whenever a module requires to either verbally notify the user or get the user’s feedback, Alana will handle this task. In this way the conversation throughout the interaction will be more contextually relevant, and easier to maintain. Since the robot needs to engage in social dialogue as well as to complete tasks, Alana was enriched with so-called *task bots* to conversationally execute and monitor behaviours on a physical agent (Figure 5) (Papaioannou, Dondrup, and Lemon 2018).

In order to enable the functionality described above, a new Natural Language Understanding (NLU) module, HERMIT NLU, has been implemented and integrated into the Alana system, which is able to deal with social chit-chat but also extract the necessary information from commands to start tasks. HERMIT NLU (Vanzo, Bastianelli, and Lemon 2019) is thus used to decide if the *task bot* is triggered and to extract the required parameters for tasks, such as the name of the shop someone is looking for. While standard chatbots mostly rely on NLU that works on shallow semantic representations (e.g., intents + slots), task-based applications require richer characterisations. In line with (Dinarelli et al. 2009), we promote the idea that the user’s intent can be represented through the combination of existing theories, capturing different dimensions of the overall problem, namely *Dialogue Acts* and *Frame Semantics*. Existing approaches to NLU for dialogue systems are based on formal languages designed around the targeted domain. However, it has been widely demonstrated that the generalisation capability of statistical-based approaches is more robust towards lexical and domain variability (Bastianelli et al. 2016). We thus use a deep learning architecture based on a hierarchy of self-attention mechanisms and BiLSTM encoders followed by CRF tagging layers to perform multi-task learning over the aforementioned semantic dimensions (Rastogi, Gupta, and Hakkani-Tur 2018). The system effectively learns how to predict Dialogue Acts, Frames, and Frame Elements in a sequence labelling scheme, starting from a corpus of annotated sentences which we are currently developing.

After a task has been identified, executing it on a robot usually includes physical actions that require a finite amount of time to complete and are not instantaneous such as dialogue actions. While the robot is executing such an action, the user might want to continue the conversation, or give new instructions. In order to be able to support such a multi-threaded dialogue management of interleaving tasks with general chit-chat and other tasks, we build on the ideas presented in (Lemon et al. 2002). To this end, the execu-

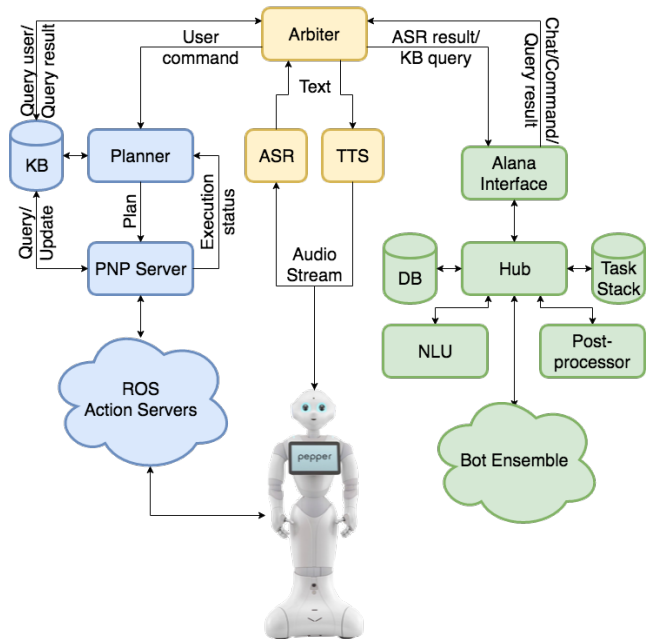


Figure 5: Architecture of the Dialogue system. The blue parts on the left represent the task management and execution system, the green parts on the right represent Alana as the dialogue system. The Bot Ensemble contains social chat bots and the *task bot* which is able to trigger tasks and handle communication between the task and the user. The yellow middle part is the task specific dialogue management system (Arbitrator), Text-To-Speech (TTS) and Automatic Speech Recognition (ASR).

tion system introduced in (Dondrup et al. 2017) has been extended to use so-called *recipes* that define dialogue and physical actions to execute in order to achieve the given goal (Papaioannou, Dondrup, and Lemon 2018). The execution framework described therein has been redesigned to support multi-threaded execution and an arbitration process has been put in place to manage the currently running tasks on the execution side and in the Alana system. This lets tasks be started, stopped, and paused at any time, with appropriate feedback to the user. If a task has been suspended by another action, it will be resumed after the new action finishes and any open questions will be re-raised to prompt the user.

## Route guidance supervision

One of the core tasks for the MuMMER robot in the mall is the guiding of users to specific locations in the mall, by pointing at places and explaining the route to the wanted location. The task is triggered when a human asks for a location. A supervision system, based on Jason (Bordini, Hübner, and Wooldridge 2007), a BDI agent-oriented framework handles the execution through Jason reactive plans. Throughout the task, the robot supervises the execution and, depending what goes wrong, the robot has multiple possible responses. For example, it is able to handle nominal scenar-

ios of route guidance while being able to take into account contingencies such as the human lack of visibility of the direction, his/her ability or not to take the stairs, his/her understanding of the message, etc. Finally, if at some point, the human is not perceived during a certain time, the robot ends the task, assuming that the human has left.

### Route computing and route verbalization

The entire description of the route, from the search for the best route to get to the final destination to the verbalization of this route, is based on the SSR (Semantic Spatial Representation) (Sarhou, Alami, and Clodic 2019b). This representation is used to describe the topology of an indoor environment as well as semantic information (type of stores or items sold by stores) in a single ontology. This ontology is managed by Ontogenius (Sarhou, Alami, and Clodic 2019a), a lightweight open-source ROS-compatible package which stores semantic knowledge, reasons with it, and shares that information to all the other system components.

### Geometric reasoning

Geometric reasoning uses Underworlds (Lemaignan et al. 2018; Sallami et al. 2019), a lightweight framework for *cascading spatio-temporal situation assessment* in robotics. It represents the environment as real-time distributed data structures, containing scene graph (for representation of 3D geometries). Underworlds supports *cascading* representations: the environment is viewed as a set of *worlds* that can each have different spatial granularities, and may inherit from each other. It also provides a set of high-level client libraries and tools to introspect and manipulate the environment models. Based on a 3D model of the mall (Figure 6), it maintains what the robot knows about the scene as well as alternative world states. These states represent the estimation of the human’s beliefs about the scene. It also provides the symbolic relations among entities with stamped predicates (e.g.  $[isInsideArea(person, area)]$  or  $[isSpeakingTo(X, Y)]$  when  $X$  speaks and looks at  $Y$  (given by perception)).

### Motion planning

The navigation of the robot is implemented using the ROS navigation stack, with navfn as the global planner and a Timed Elastic Band (TEB) (Rösmann, Hoffmann, and Bertram 2017) planner as the local planner. For MuMMER, the local planner was modified in order to accommodate humans into planning inspired from (Khambhaita and Alami 2017), resulting a new planner called and this new planner called Social TEB (S-TEB). This algorithm is able to plan and execute trajectories while ensuring satisfaction of robot kinematics constraints, avoiding static and moving non-human obstacles and planning navigation solutions respecting social constraints with humans perceived. The planner ensures the safety of humans by re-planning a local plan at each control loop.

### SVP planner

Although the target robot location in the mall is in a large square, several elements of the environment can block the

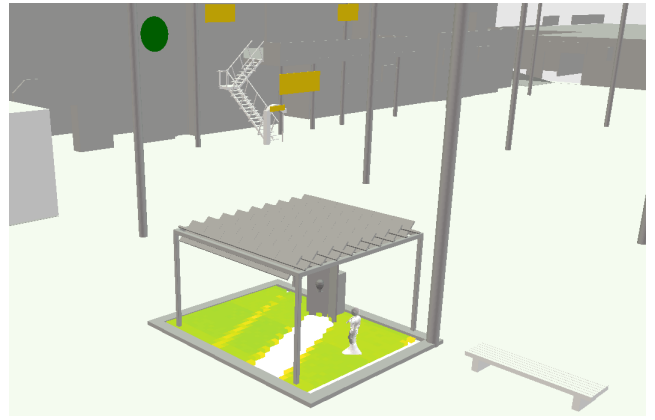


Figure 6: Visualization of the visibility grids of a landmark on the 3D model of the central square of the Ideapark shopping center.

visibility of important landmarks for the proper understanding of the route to take. The purpose of the SVP (Shared Visual Perspective) planner (Waldhart, Clodic, and Alami 2019) is therefore to try to find a position where the human will have to go in order to observe an element of the environment such as a passage, a staircase or a store. To do this, a visibility grid is computed for each possible landmark, as shown on figure 6. Having determine a good position for the human, the planer also allows to determine the good position for the robot so as to have a human-robot-landmark conformation allowing both to point the landmark and to look at the human.

## Deployment

A long-term deployment (three months from September 2019) will allow the study of the customer behaviours around a helpful and entertaining robot over an extended period of time. This section gives details on how the hardware and software that are being employed in the final deployment, as well as the scenarios that are supported.

### Setup

The fully integrated MuMMER system consists of several hardware components to allow the computation to be performed on the appropriate platforms.

The robot we are using is an updated, custom version of the Pepper platform, which is equipped with an Intel D435 camera and a NVIDIA Jetson TX2 in addition to the traditional sensors that are found on the previous versions of the robot. We use the Robot Operating System (ROS) to enable the communication between the processing nodes. All the streams (audio, video, robot states) are sent to a remote laptop which performs all the computation. The laptop has a NVIDIA RTX 2080 graphics card (for the deep learning part) and 12 CPU cores. The perception algorithms process the Intel images at a resolution of  $320 \times 180$  for the detection and tracking parts, and at a resolution of  $640 \times 360$  for the re-identification part, which enables fast tracking and a good

re-identification quality with OpenFace. The 4 microphone streams are processed at a frequency of 16000 Hz, and the full perception system delivers the output at 10 fps.

To transcribe the user's speech signal we use the Google Automatic Speech Recognition (ASR) API<sup>2</sup> which receives an enhanced audio signal from a delay and sum beamformer based on the location of the speaking person determined by the audio-visual sensing. A dedicated ROS node streams the audio to the ASR which in real-time returns an incrementally updated string transcribing the utterance. Using silence to mark the end of speech, this transcription is enhanced using the context of the sentence to provide for a more coherent result. Finally, the text output of the Google ASR is sent to the Alana framework to perform the dialogue task, through the arbitration module as explained above.

The system is deployed in two languages, English and Finnish, though due to the vast linguistic differences between the two languages, the two versions have been kept separated, and the whole interaction can either be in one or the other. Due to the complexity of the NLU module in the Finnish version of the system, the user's utterance is being translated into English using Google Translate API<sup>3</sup>. The result of this translation is sent to the Alana conversational framework and goes through the NLU pipeline described in detail in (Curry et al. 2018). In the English version of the system, Alana then returns the reply to be verbalised. Due to the relatively poor performance of Google translate when it comes to translating English into Finnish (as remarked upon by our Finnish partners), the Finnish version of Alana has a much reduced set of bots in its ensemble (see Figure 5). These bots mainly return answers based on templates that have been translated into Finnish beforehand.

## Scenarios deployed

As a proof of concept, a real-time autonomous system has been built to integrate all the components described in the sections above. The following types of interactions can be triggered by the user:

**Chat** The staple of the interaction is social dialogue (Curry et al. 2018). During all other modes of interaction, the user can always default to simply chat to the robot irrespective of whether it is currently executing a goal/task (e.g. the user requires guidance to a specific shop) or not. For example the user might approach the robot and start discussing various topics. At specific points throughout this conversation, the system might explain its capabilities to the user in order to recover from a conversational stalemate or to simply make them aware of the fact that it can also be helpful in finding your way around the mall (see below).

**Quiz** In this scenario, multiple choice questions are asked by the robot, and the human replies by stating the number of the answer they think is correct.

**Route Description (Dialogue only)** When the human asks how to get to a specific shop, the robot gives him/her the

route description. In this most "simple" form, the system uses only verbal interaction. This means that especially the route description is merely presented as a string of synthesised text.

**Route Guidance (Dialogue + Pointing)** In this version, the robot guides the human to specific locations in the mall, by pointing at places and explaining the route to the wanted location. To do so, the robot first computes positions so that the human will be able to see what the robot is pointing for him/her. Then, the robot navigates to its position (this part is optional), expecting that the human will join it once it stopped and checks the human's visibility. Then, the robot explains to the human how to reach the destination. According to a human-human guidance study (Belhassein et al. 2017), the robot points, first at the location direction and then points at the access point (a corridor, stairs or an escalator) to go through to reach the location. While pointing, the robot verbalizes the route. Finally, the robot checks that the human knows how to reach the goal and leaves open the possibility to repeat if needed. All along the task, the robot supervises the task and adapts accordingly.

All these modes of interaction can be interleaved at the user's discretion. This means, for example, that during the quiz the user could revert to social dialogue. If they do so, the system might occasionally try to bring them back to the quiz by re-raising the last question. The same holds true if the person chooses to abandon a route guidance task before it was finished.

## Conclusions

The MuMMER project has built a fully autonomous entertainment robot to perform HRI scenarios in a shopping mall, in which the main goal is to have entertainment interaction (quiz, chat), as well as route guidance. The system is real-time, by leveraging the heavy deep learning computation on a remote laptop, the ASR on the Google platform, and the Alana conversational AI system on a remote server. This system enables a natural interaction with the participants; it has been tested and was tested in real conditions for several short sessions, and as of September 2019 is fully deployed for a three-month long-term user study.

Further work for large scale deployment could include some software optimizations to run more components on the robot itself, and to reduce the lag which sometime exists between the human speech and the robot reply.

## Acknowledgements

This research has been partially funded by the European Union's Horizon 2020 research and innovation program under grant agreement no. 688147 (MuMMER, <http://mummer-project.eu/>).

<sup>2</sup><https://cloud.google.com/speech-to-text/>

<sup>3</sup><https://cloud.google.com/translate/>

## References

- [Amos, Ludwiczuk, and Satyanarayanan 2016] Amos, B.; Ludwiczuk, B.; and Satyanarayanan, M. 2016. Openface: A general-purpose face recognition library with mobile applications. Technical report, CMU-CS-16-118, CMU School of Computer Science.
- [Bastianelli et al. 2016] Bastianelli, E.; Croce, D.; Vanzo, A.; Basili, R.; and Nardi, D. 2016. A discriminative approach to grounded spoken language understanding in interactive robotics. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI'16*, 2747–2753.
- [Belhassein et al. 2017] Belhassein, K.; Clodic, A.; Cochet, H.; Niemelä, M.; Heikkilä, P.; Lammi, H.; and Tammela, A. 2017. Human-Human Guidance Study. Technical Report 17596, LAAS.
- [Bordini, Hübner, and Wooldridge 2007] Bordini, R. H.; Hübner, J. F.; and Wooldridge, M. 2007. *Programming Multi-Agent Systems in AgentSpeak Using Jason (Wiley Series in Agent Technology)*. USA: John Wiley & Sons, Inc.
- [Cao et al. 2017] Cao, Z.; Simon, T.; Wei, S.-E.; and Sheikh, Y. 2017. Realtime multi-person 2d pose estimation using part affinity fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 7291–7299.
- [Cao, Canévet, and Odobez 2018] Cao, Y.; Canévet, O.; and Odobez, J.-M. 2018. Leveraging convolutional pose machines for fast and accurate head pose estimation. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*.
- [Craenen et al. 2018a] Craenen, B.; Deshmukh, A.; Foster, M. E.; and Vinciarelli, A. 2018a. Do we really like robots that match our personality? the case of big-five traits, god-speed scores and robotic gestures. In *Proceedings of the 27th IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*.
- [Craenen et al. 2018b] Craenen, B. G.; Deshmukh, A.; Foster, M. E.; and Vinciarelli, A. 2018b. Shaping gestures to shape personalities: The relationship between gesture parameters, attributed personality traits, and Godspeed scores. In *Proceedings of the 27th IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*, 699–704.
- [Curry et al. 2018] Curry, A. C.; Papaioannou, I.; Suglia, A.; Agarwal, S.; Shalymov, I.; Xu, X.; Dušek, O.; Eshghi, A.; Konstas, I.; Rieser, V.; et al. 2018. Alana v2: Entertaining and informative open-domain social dialogue using ontologies and entity linking. *Alexa Prize Proceedings*.
- [Deshmukh, Foster, and Mazel 2019] Deshmukh, A.; Foster, M. E.; and Mazel, A. 2019. Contextual non-verbal behaviour generation for humanoid robot using text sentiment. In *Proceedings of the 28th IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*.
- [Dinarelli et al. 2009] Dinarelli, M.; Quarteroni, S.; Tonelli, S.; Moschitti, A.; and Riccardi, G. 2009. Annotating spoken dialogs: From speech segments to dialog acts and frame semantics. In *Proceedings of SRSI 2009, the 2nd Workshop on Semantic Representation of Spoken Language*, 34–41.
- [Dondrup et al. 2017] Dondrup, C.; Papaioannou, I.; Novikova, J.; and Lemon, O. 2017. Introducing a ROS based planning and execution framework for human-robot interaction. In *Proceedings of the 1st ACM SIGCHI International Workshop on Investigating Social Interactions with Artificial Agents, ISIAA 2017*, 27–28.
- [Foster et al. 2016] Foster, M. E.; Alami, R.; Gestranian, O.; Lemon, O.; Niemelä, M.; Odobez, J.-M.; and Pandey, A. K. 2016. The MuMMER project: Engaging human-robot interaction in real-world public spaces. In *Social Robotics*, 753–763. Cham: Springer International Publishing.
- [Hastie et al. 2016] Hastie, H.; Dente, P.; Küster, D.; and Kappas, A. 2016. Sound emblems for affective multimodal output of a robotic tutor: A perception study. In *Proceedings of the 18th ACM International Conference on Multimodal Interaction*, 256–260.
- [He, Motlicek, and Odobez 2018a] He, W.; Motlicek, P.; and Odobez, J. 2018a. Deep neural networks for multiple speaker detection and localization. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 74–79.
- [He, Motlicek, and Odobez 2018b] He, W.; Motlicek, P.; and Odobez, J.-M. 2018b. Joint localization and classification of multiple sound sources using a multi-task neural network. In *Proceedings of Interspeech 2018*, 312–316.
- [He, Motlicek, and Odobez 2019] He, W.; Motlicek, P.; and Odobez, J. 2019. Adaptation of multiple sound source localization neural networks with weak supervision and domain-adversarial training. In *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 770–774.
- [Heikkilä et al. 2019] Heikkilä, P.; Niemelä, M.; Belhassein, K.; Sarthou, G.; Tammela, A.; Clodic, A.; and Alami, R. 2019. Should a robot guide like a human? a qualitative four-phase study of a shopping mall robot. In *International Conference on Social Robotics (ICSR)*.
- [Heikkilä, Lammi, and Belhassein 2018] Heikkilä, P.; Lammi, H.; and Belhassein, K. 2018. Where can I find a pharmacy? - human-driven design of a service robot's guidance behaviour. In *Proceedings of PubRob 2018*.
- [Khalidov and Odobez 2017] Khalidov, V., and Odobez, J.-M. 2017. Real-time multiple head tracking using texture and colour cues. *Idiap-RR Idiap-RR-02-2017*, Idiap.
- [Khambhaita and Alami 2017] Khambhaita, H., and Alami, R. 2017. Viewing Robot Navigation in Human Environment as a Cooperative Activity. In *International Symposium on Robotics Research (ISSR 2017)*, 18p.
- [Lemaignan et al. 2018] Lemaignan, S.; Sallami, Y.; Wallbridge, C.; Clodic, A.; Belpaeme, T.; and Alami, R. 2018. UNDERWORLDS: Cascading Situation Assessment for Robots. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems*.
- [Lemon et al. 2002] Lemon, O.; Gruenstein, A.; Battle, A.; and Peters, S. 2002. Multi-tasking and collaborative activi-

- ties in dialogue systems. In *Proceedings of the 3rd SIGdial Workshop on Discourse and Dialogue - Volume 2*, 113–124.
- [Papaioannou et al. 2017] Papaioannou, I.; Cercas Curry, A.; Part, J.; Shalyminov, I.; Xinnuo, X.; Yu, Y.; Dusek, O.; Rieser, V.; and Lemon, O. 2017. Alana: Social dialogue using an ensemble model and a ranker trained on user feedback. In *2017 Alexa Prize Proceedings*.
- [Papaioannou, Dondrup, and Lemon 2018] Papaioannou, I.; Dondrup, C.; and Lemon, O. 2018. Human-robot interaction requires more than slot filling-multi-threaded dialogue for collaborative tasks and social conversation. In *FAIM/ISCA Workshop on Artificial Intelligence for Multimodal Human Robot Interaction*, 61–64.
- [Rastogi, Gupta, and Hakkani-Tur 2018] Rastogi, A.; Gupta, R.; and Hakkani-Tur, D. 2018. Multi-task learning for joint language understanding and dialogue state tracking. In *Proceedings of the 19th Annual SIGdial Meeting on Discourse and Dialogue*, 376–384.
- [Rösmann, Hoffmann, and Bertram 2017] Rösmann, C.; Hoffmann, F.; and Bertram, T. 2017. Integrated online trajectory planning and optimization in distinctive topologies. *Robotics and Autonomous Systems* 88:142–153.
- [Sallami et al. 2019] Sallami, Y.; Lemaignan, S.; Clodic, A.; and 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 2019)*. To appear.
- [Sarhou, Alami, and Clodic 2019a] Sarhou, G.; Alami, R.; and Clodic, A. 2019a. Ontologenius : A long-term semantic memory for robotic agents. In *RO-MAN 2019*.
- [Sarhou, Alami, and Clodic 2019b] Sarhou, G.; Alami, R.; and Clodic, A. 2019b. Semantic Spatial Representation: a unique representation of an environment based on an ontology for robotic applications. In *SpLU-RoboNLP 2019*, 50 – 60.
- [Sheikhi and Odobez 2015] Sheikhi, S., and Odobez, J. 2015. Combining dynamic head pose and gaze mapping with the robot conversational state or attention recognition in human-robot interactions. *Pattern Recognition Letters* 66:81–90.
- [Siegfried, Yu, and Odobez 2017] Siegfried, R.; Yu, Y.; and Odobez, J.-M. 2017. Towards the use of social interaction conventions as prior for gaze model adaptation. In *19th ACM International Conference on Multimodal Interaction (ICMI)*.
- [Vanzo, Bastianelli, and Lemon 2019] Vanzo, A.; Bastianelli, E.; and Lemon, O. 2019. Hierarchical multi-task natural language understanding for cross-domain conversational ai: HERMIT NLU. In *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*, to appear. Stockholm, Sweden: Association for Computational Linguistics.
- [Vinciarelli, Pantic, and Bourlard 2009] Vinciarelli, A.; Pantic, M.; and Bourlard, H. 2009. Social signal processing: Survey of an emerging domain. *Image and Vision Computing* 27(12):1743–1759.
- [Waldhart, Clodic, and Alami 2019] Waldhart, J.; Clodic, A.; and Alami, R. 2019. Reasoning on Shared Visual Perspective to Improve Route Directions. In *2019 28th IEEE International Conference on Robot & Human Interactive Communication*.