

Composing HARMONI: An Open-source Tool For Human And Robot Modular Open Interaction

Micol Spitale^{1*}, Chris Birmingham^{2*}, R. Michael Swan³, and Maja J Matarić⁴

Abstract—The research and development of socially interactive robots is a complex challenge because of the wide variety of capabilities needed for effective social human-robot interactions (HRI). Many of these capabilities, including perception, dialog, and control, have state of the art methods and solutions, but combining those into a comprehensive and seamless interaction is still an open challenge. We describe HARMONI, a multimodal, open-source tool for rapid social HRI development and deployment. HARMONI is centered around a ROS package for interaction development, including decision management and node orchestration. HARMONI systematically integrates with disparate functionalities needed to conduct a meaningful social human-robot interaction such as external cloud services, AI models, and modules for sensing, planning, and acting on a variety of platforms. HARMONI was applied to the QT robot platform and usability tests were conducted to evaluate the ease and speed of development and deployment. This paper describes the architecture and design of HARMONI and reports the results of a pilot study with novice users.

I. INTRODUCTION

The last decade has brought about major advances in computing for real-world systems. Computer vision [1], [2], [3] and natural language understanding [4], [5], [6], [7] capabilities have enabled the rise of semi-autonomous cars [8] and digital voice assistants [9]. However, socially interactive robots have not yet entered daily use [10], because of the less studied complexities of social interaction that require the integration of many capabilities across separate modalities [11], [12]. This complexity and the associated challenges grow in proportion to the diversity of sensor hardware and the improvements in robot hardware affordances provided by socially interactive robots, as well as the frequent mismatch between the computational power required for perception compared to what is typically available on robot platforms [13]. Prior work in developing socially interactive robots has adopted a diverse set of approaches to handling these challenges, including a wide variety of control architectures [14], [15], [16], [17] and a pragmatic focus on specific

interaction challenges rather than fully integrated interactive systems.

Although the diversity of technical approaches in social HRI research supports the varied needs of different interaction settings, there are benefits to the adoption of more standardized tools. Having a shared framework that can be deployed across different platforms allows for rapid replication of experiments in different contexts. Furthermore, a modular framework enables researchers to focus on individual interaction challenges while using state of the art methods from other fields. Finally, packaging and sharing functional modules allows researchers to readily create and pilot interactions.

This work presents HARMONI, an *open-source, robot-agnostic, modular, composable, and multimodal* tool designed to address many of the challenges outlined above. Section II provides a detailed explanation of HARMONI goals and requirements. HARMONI is built on top of the Robot Operating System (ROS) [18], leveraging the widespread ROS community adoption, capabilities, and support. At the core of HARMONI is a ROS package for standardizing orchestration and communication of the services supporting the various capabilities HARMONI integrates and for composing interactions within a standard decision framework. Even with ROS, integrating disparate systems remains challenging, so HARMONI packages these capabilities into a containerized docker environment [19] that can be rapidly deployed on most robot systems. This integration allows HARMONI to support services inside and outside of the ROS ecosystem, including cloud services like AWS Lex and Polly [20], AI models such as Facebook Wave2Letter [21], and local modules for display-based robot faces such as Cordial [22] and controlling gestures on underlying robot hardware such as QT [23] or Blossom [24]. The open-source implementation of HARMONI is found on Github¹. This paper contributes the following aspects of HARMONI:

- A description of the open-source tool for development and deployment of research into socially interactive robotics;
- Instructions and examples for using HARMONI in social HRI research studies;
- A brief pilot study demonstrating novice user impressions of HARMONI.

This paper is organized as follows. Section II explains the design goals driving the design choices for defining HARMONI. Section III describes the state of the art of

*Micol Spitale and Chris Birmingham contributed equally to this work. This work was supported by the National Science Foundation

¹Micol Spitale is a PhD Student at I3Lab, in the DEIB Department, Politecnico di Milano, Milan, Italy, and Visiting PhD Student at the Interaction Lab, in the CS Department, University of Southern California, Los Angeles, CA, The United States micol.spitale@polimi.it

²Chris Birmingham is a PhD Student at the Interaction Lab, in the CS Department, University of Southern California, Los Angeles, CA, The United States cbirming@usc.edu

³R. Michael Swan is an Masters Student with Interaction Lab, in the CS Department, University of Southern California, Los Angeles, CA, The United States rmswan@usc.edu

⁴Maja J Matarić is the PI of the Interaction Lab, Professor of Computer Science, Neuroscience, and Pediatrics at the University of Southern California, mataric@usc.edu

¹<https://github.com/interaction-lab/HARMONI>

the recent robot development tools in social HRI. Section IV presents the structure and use of HARMONI. Section V provides an example use case and a short usability study. Lastly, Section VI discusses future works of HARMONI.

II. DESIGN GOALS

Open-source and *multimodal* support were the original driving requirements for HARMONI. We then surveyed socially interactive robotics literature for design principles [25], and identified these additional key design goals:

- Robot-agnostic
- Modular
- Composable

Open-source is an important feature because it provides wide access, allowing widespread adoption and community support. **Multimodality** is important for social HRI because humans produce multimodal signals and interact most naturally and effectively with embodied agents that also produce and understand multimodal signals [26]. **Robot-agnostic** is important for similar reasons; it allows a broad range of developers to make use of the tool’s capabilities and facilitates comparisons across different platforms. **Modularity** is important for extensibility to new plug-ins. **Composability** allows developers to bring together various different modules and capabilities to accomplish complex goals.

A. Robot-agnostic

Robotic platforms in HRI are constantly improving and being replaced so it is essential that our framework is adaptable to new platforms. The HARMONI API is designed to be independent of any particular robot platform. To make that possible, the hardware interface modules are implemented at the lowest layer of the API so that they can be readily replaced without creating dependency issues. Integrating the API with a new robot platform requires building modules around all hardware endpoints and satisfying the well-defined interface for each type of endpoint. While this design tends toward well abstracted code that offers high levels of flexibility, it can be more difficult to maintain and extend initially.

B. Modular

ROS [18] and microservices [27] provide a highly effective approach to modularity in distributed systems. Taking lessons from these, HARMONI aims to achieve similar success with its approach to modularity while following best practices from software engineering [28], [29]. We believe interface design is of particular importance in facilitating modularity in social HRI programming. ROS already has excellent interface design but is also very general, and such flexibility can be overwhelming for new developers and interaction designers. By taking a social HRI-centered approach, module interfaces can be appropriately constrained and therefore simplified. As is shown in Section IV through the simplified HARMONI interface, new researchers and developers can ramp up and learn faster with the benefit of clear module interconnection schemas.

C. Composable

Modules should be readily combined and controlled in order to compose or extend existing modules into more complex functional units. To accomplish this goal, the tool implicitly dictates that module control and error handling follow a hierarchy, as is illustrated in Section IV. This allows the highest level of command to supply very general instructions that are broken down as the chain of command is traversed downward.

Errors should be handled at the lowest level possible, leaving the top level to deal only with serious or unrecoverable errors as they are pushed up the chain. This has two benefits: it reduces coupling as modules typically have fewer contact points, and it decreases the amount of micromanagement code that programmers tend to include in high level modules.

III. BACKGROUND

Existing tools for developing social human-robot interaction can be organized based on their capabilities and intended use. In the following sections, we describe all social HRI programming tools we are aware of. We then provide a comparison of these tools along dimensions of interest described in Section II.

A. Robot-Specific Tools

Various platform-specific tools have been developed by robot manufacturers to support their robot products. For example, Choregraphe [30] is a visual programming tool built specifically for the Nao Robot [31] that allows different levels of programming: behavior control at the high level, fine tuning of complex joint or Cartesian motions, and programming in Python at the lowest level. Interaction Blocks is another Nao programming tool; it aids design exploration and prototyping of interactions with composable visual blocks for authoring social interactions [32]. HRI Framework [33] is a tool for human-robot collaboration.

B. Content Authoring Tools

Other robot programming tools support content authoring aimed at facilitating collaboration between developers and designers. Interaction Blocks, described above, is one such tool [34]. The Interaction Composer framework builds human-robot interaction with a visual blocks interface [35]. The tasks are split between programmers and developers: the former develop modules such as face recognition, while the latter design social interactions with the robot with behaviors and dialogue. RIZE [36] is a web based tool which provides GUI-based block programming similar to Interaction Composer. At the time of this writing RIZE was unavailable for testing. Another programming tool for both programmers and users is TiViPE, a textual robot command language [37], [38] that allows textual and visual programming for users with minimal understanding of programming.

Some tools do not use graphical interfaces for content generation, but still support authoring content in various markup languages. For example, *BONSAI* is a framework designed to support developers of home robots for social

interaction [39]. Another is the Robot Behavior Toolkit [40], which helps programmers to develop behaviors for socially interactive robots. These tools use a markup language for controlling behaviors at two levels, the cognitive level (i.e., what decision to make) and execution level (e.g., gaze direction, speech).

Unity² is a popular game and animation engine often used for developing simulated and extended reality (XR) applications. Unity has been used to build virtual social interactions [41] and socially assistive robot interactions [42]. Unity contains both visual and textual programming modalities in addition to a large library of tools/packages that support programmers and interaction designers. Commercial tools also exist for supporting the development of functionalities for socially interactive robots; for example, Arora³ is an enterprise platform for content creation and deployment on interactive personal robots.

C. Dialog Tools

There has been rapid growth in the development of tools to support spoken dialogue. While these tools are not specifically designed for robotics, they are used in socially interactive robots that employ verbal interactions. One such tool is the Artificial Intelligence Markup Language (AIML), an XML-based markup language for authoring dialogue with natural language processing [43]. Another is ReTiCo, a Python-based dialogue framework to create spoken dialogue systems and simulations [44]. ReTiCo was expanded to support capabilities needed for use on robots and renamed the Robot-Ready Spoken Dialogue System (RR-SDS) [45]. Plato is a conversational AI platform in Python that supports conversational agents in general [46]. One tool that focuses on the multimodal aspects of interaction, including dialog, is the Platform for Situated Intelligence (PSI). Like RR-SDS, PSI focuses on real-time, synchronized signal processing for improved system response times and accuracy [47].

D. Tool Limitations and Comparisons

We briefly consider each of the tools described above in relation to the goals and requirements specified in Section II, as shown in Table I. A more extensive overview can be found in [36]. Most of the available tools do not have most or all of the selected features. The tools may be effective in particular use cases, but it is often difficult or impossible to make modifications or extensions, limiting the appeal and adoption by the robotics community. Notably, a few tools do include all of the features, with the following limitations. The Robot-Ready Spoken Dialog System effectively integrates the incremental spoken dialog of ReTiCo with robot control, but the tool’s focus is on supporting spoken capabilities. The Platform for Situated Intelligence includes a real-time synchronization system for processing multimodal sensory data. However, it currently lacks support for Linux systems and modules for actuating or controlling a robotic system.

²<https://unity.com/>

³semio.ai

TABLE I
A COMPARISON OF SOCIAL HRI PROGRAMMING TOOLS

Tool	Open Source	Robot Agnostic	Composable	Modular	Multimodal
Choregraphe			✓		✓
Interaction Blocks				✓	✓
Interaction Composer		✓	✓		✓
TiViPE				✓	✓
HRI Framework	✓		✓	✓	
AIML	✓	✓	✓		
BONSAI		✓	✓	✓	✓
Plato	✓	✓	✓		
Arora		✓	✓	✓	✓
ReTiCo	✓	✓	✓		
Robot Behavior Toolkit	✓	✓	✓		✓
Unity		✓	✓	✓	✓
rr-SDS	✓	✓	✓	✓	✓
Platform for Situated Intelligence	✓	✓	✓	✓	✓
RIZE	✓	✓	✓	✓	✓

In response to the available tools and their features, we developed a platform that enables socially interactive robotics developers to overcome the above limitations. The design goals for the tool are described in the next section.

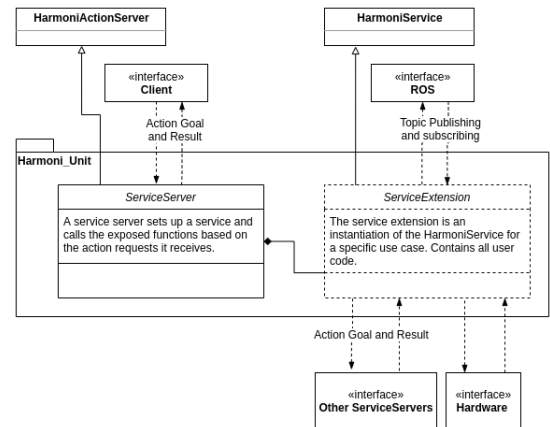


Fig. 1. The HARMONI Unit: a simplified UML diagram describing the structure of every service in HARMONI.

IV. HARMONI: TOOL STRUCTURE

HARMONI is a modular tool built around a simple design pattern called the HARMONI Unit. Every module in the tool follows the structured pattern and API of the HARMONI Unit. The central functionality needed to run HARMONI is supplied in a group of packages called the HARMONI Core. These packages manage and orchestrate all the other modules and functions of the tool. The rest of the tool is integrated into a modular architecture that supports basic interactions. These major components are described next.

A. The HARMONI Unit

The HARMONI Unit makes use of the ROS action server⁴ for handling state and control across the distributed modules that make up the tool. By using ROS to support communication, HARMONI falls under the larger umbrella of the ROS community, but it also ties it to weak guarantees on message passing and real-time computation. The HARMONI Unit design pattern, illustrated in the class diagram in Figure 1 describes and standardizes the communication needed for every unit to communicate with the rest of the tool, including other control units, informational topics, and stored files or parameters. Every unit exposes a predefined API, allowing other units to control or query it, without detailed knowledge of its internal structure. This simple API specifies functionality for a given unit to ‘do’ a certain action, to ‘request’ a response from an external service, or to ‘start’, ‘stop’, or ‘pause’ a running action. For example, a microphone node can be commanded to ‘start’, ‘stop’, or ‘pause’ listening while a text-to-speech node may receive a string ‘request’ to Amazon Polly and return audio and visemes. Similarly, a motor control node may be called to ‘do’ an arm wave. The difference between ‘do’ and ‘start’ is both in duration (‘start’ is intended for time-extended actions) and in the function arguments (‘do’ typically requires an argument describing what to do, whereas ‘start’ should be inherent to the unit).

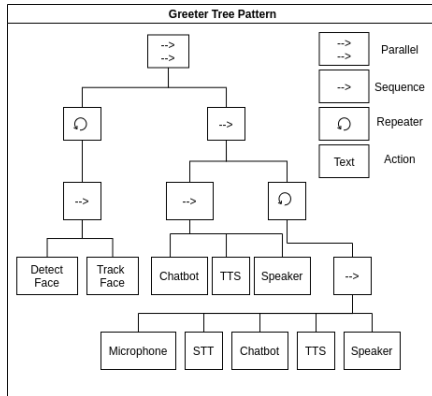


Fig. 2. A sample HARMONI behavior, structured as a behavior tree that starts with face detection, followed by a speech sequence, and then a dialog loop. Each leaf node is an API call to the specified HARMONI Unit. Note that some block definitions follow convention described in [48].

B. The HARMONI Core

The control and the orchestration of HARMONI are found in the HARMONI Core module, consisting of a set of ROS packages including the HARMONI Common Library, the Recorder, and the Pattern and Decision Manager. The HARMONI Common Library includes the specification of the HARMONI Unit and HARMONI Action Server, used by all modules to expose functionality to the rest of the tool. This functionality is controlled with a two-level controller. The lower-level controller is a ROS package called

⁴<http://wiki.ros.org/actionlib>

harmoni_pattern used to play any node that follows the HARMONI Unit convention. In this package, pattern nodes are scripted together in simple behavior trees to create higher-level functionality. For example, to create an interactive spoken dialog, a script might loop the sequence of [Microphone, Speech-To-Text, Chatbot, Text-To-Speech, Speaker], in which the output of each step is used as input for the next step. These patterns can be composed in numerous ways to generate various behaviors. A simple greeting interaction is shown in Figure 2.

Although behavior trees are capable of generating relatively sophisticated behaviors, *patterns* do not contain logic or the capability to make high-level decisions. Directing the flow of patterns and handling decision making is done by the *Decision Manager*. The Decision Manager operates as a higher-level controller that listens to the state of the system and receives feedback from the patterns in order to choose what pattern should be played next. The Decision Manager can be implemented as a decision tree, behavior tree, or custom script, depending on the length and needs of the interaction. The Decision Manager starts up nodes and handles any errors that arise. Lastly, in order to supply the Decision Manager with memory for long-term or repeated interactions, we supply the Recorder Package, which manages storage of information in a database that can be queried as the interaction progresses.

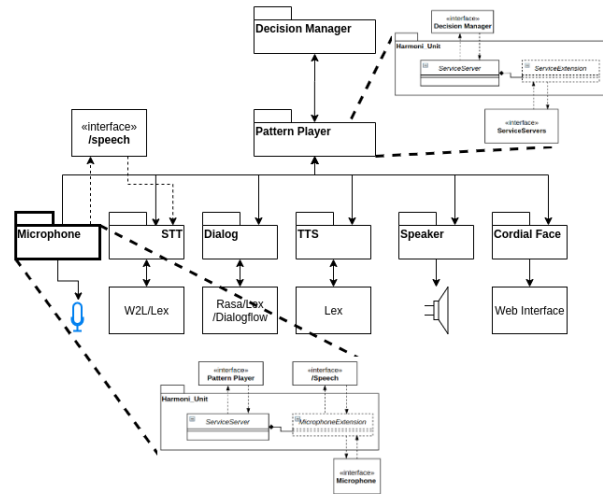


Fig. 3. A sample communication diagram between HARMONI Units for a dialog interaction, illustrating the hierarchical control and modular communication. Example HARMONI Unit UML diagrams highlight the internal structure.

C. System Integration and Architecture

The HARMONI Unit and Core, described above, make up the backbone of the tool. They give order and structure and enable communication between independent modules by constraining the supported design space. This structure enables the largest HARMONI contribution, the systematic integration of disparate functionalities needed to conduct a meaningful social human-robot interaction. Any existing

functionality, hardware, software, or cloud services, is simply added to HARMONI with a Python wrapper. To handle the multiple capabilities and different configurations that could be desirable for a social human-robot interaction, HARMONI adopts a standard hybrid sense-plan-act architecture. An example configuration is shown in Figure 3. As can be seen from this communication diagram, the Decision Manager directly controls the Pattern Player and monitors the status of each pattern as it plays. The Pattern Player can directly play each HARMONI Unit, serving as the interface and controller for the hardware or service it exposes. Units can communicate directly with each other when messages are large or need to stream continuously, as shown in the communication between the Microphone and the STT units in Figure 3. Otherwise, the result of an API call to one unit is piped to the next unit through the Pattern Player. In this architecture, new capabilities are readily added and only need to follow the simple HARMONI Unit convention for integration into a new or existing pattern.

V. USING HARMONI

A. Use Case

To validate HARMONI on a physical robot platform, we used a QT robot by LuxAI [23]. QT is a tabletop humanoid with static legs, 4 DOF arms, 2 DOF neck, and a screen face. It is equipped with a RGB-D camera, an array microphone, and a speaker. For this example, we designed an interaction intended for children using a socially assistive robot in a Story Retelling Activity. The activity involves implementing an interaction in which the robot, QT, tells a story to the user, and then asks the user to repeat back the story. If the user repeats the story correctly, the activity is completed successfully and it ends; otherwise, the robot prompts the user with questions about the missing parts of the story. For example, if the user forgets to mention the main characters of the story, the robot asks: “Who are the main characters of the story?”. This interaction involves dialogue; the robot must understand the user sufficiently and record the user’s answers to its questions in order to validate the correctness of the repeated story.

The interaction session with the robot was designed as a combination of the following behavior patterns: Greeting Pattern, Head Tracking Pattern, Retelling Pattern, and Goodbye Pattern, illustrated in Figure 4 and described here:

- **Greeter and head tracking.** To begin, the Face Detection and Tracking Face start. In parallel, the Recorder starts. As soon as a face has been detected, the robot introduces itself and starts the question-answer phase, where it asks the user some introductory questions, such as “What is your name?”. The full diagram for this pattern can be seen in Figure 2.
- **Retelling.** Next, the robot tells a story to the user, while tracking the user’s face to check if the face is always in the frame of the camera. This behavior pattern is defined as a loop of sequential behaviors and actions: dialogue (Chatbot), synthesize (Text-To-Speech), behave

(i.e., Gesture, Facial Expression, Speak), detect speech (Speech-To-Text). The loop breaks when the dialogue intent is fulfilled. The robot then asks the user to answer questions related to the parts of the story the user missed during the previous step. The robot keeps prompting the user until they successfully recall all of the missing parts of the story.

- **Goodbye.** Finally, the robot completes the conversation, and all the long running processes that were started at the beginning of the session are stopped (i.e., Tracking Face, Face Detection).

As the Decision Manager plays the simple behavior tree composed of the Greeter Pattern, Retelling Pattern, and Goodbye Pattern, it also monitors the state of each node and restarts nodes if they fail. It monitors the status of specified units (e.g., Detect Face, to determine if the user is still in the frame) and intervenes if necessary.

With existing HARMONI capabilities this interaction was simple to compose. This allowed the developer ample time to focus on the core of the relevant research, in this case the dialog and story analysis functionality implemented separately within an AWS Lex Chatbot.

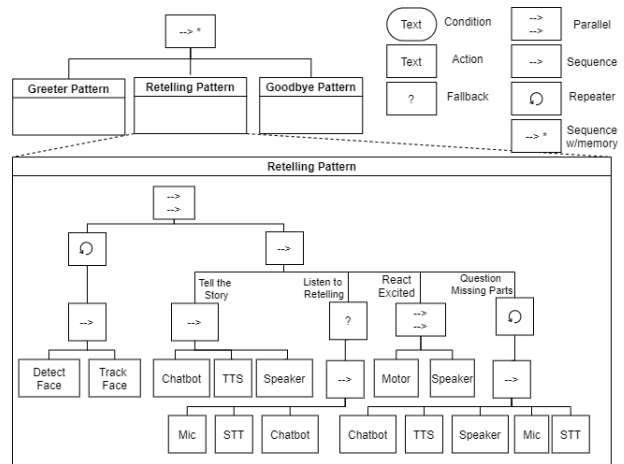


Fig. 4. Full decision behavior tree for the Story Retelling activity, including a call out of the Retelling Pattern behavior tree [49]

B. Pilot Study

1) **Method:** To evaluate whether HARMONI is simple and easy to use, we designed and conducted a pilot user study where 5 Computer Science PhD students that had no expertise in robotics attempted to setup and use HARMONI. The participants were not familiar with the HARMONI framework or APIs before the experiment. They were asked to perform three assignments in HARMONI:

- F1. Set up HARMONI on docker, which consists of the following steps: 1) clone the HARMONI repository, 2) setup cloud services with credentials we provided, 3) launch the complete system through docker-compose, 4) launch the HARMONI container and run *roscore*;
- F2. Run a test demonstration, where the user is required to setup their hardware (sensor and/or motors) endpoints

and launch the *demo* behavior pattern;

- F3. Modify an interaction in HARMONI, which consists of modifying the *demo* behavior pattern JSON we provided them and adding a new module (face of the robot displayed) to the interaction.

The sequence of tasks for each assignment can be found on the HARMONI Wiki page. At the end of each assignment, participants filled SUS [50] and NASA-TLX [51] questionnaires regarding their experience with HARMONI in terms of usability [52] and effort, respectively. The SUS [50] is a multiple-choice 10-item questionnaire which aims to measure the usability of a tool. The five response options follow the Likert scale from 1 to 5 (from Strongly agree to Strongly disagree). The NASA-TLX [51] is a 6-item questionnaire which evaluates the perceived workload to accomplish a task with a tool. The users score the perceived workload for each of the 6 categories (e.g., temporal demand, mental demand) scoring it from 1 to 20.

2) *Results*: For setting up HARMONI (F1), the results showed that most users considered the system easy to use ($Mdn=4$), and that the majority of the functions were well integrated in the system ($Mdn=4$). In addition, these non-expert users were able to use the system without any support ($Mdn=1$). However, most found the system complex ($Mdn=4$) and considered the system not so easy to learn ($Mdn=2$) for F1. Regarding assignments F2 and F3, users felt more confident in using the system for running the interaction ($Mdn=3$) than for setting it up ($Mdn=2$). Once they set up the system (F1), they considered the system less complex (F2: $Mdn=2$, F3: $Mdn=1$), but more difficult to use (F2: $Mdn=3$, F3: $Mdn=2$) than in F1 ($Mdn=4$). From the NASA-TLX questionnaire, users found the testing of functionalities mentally demanding in general (F2: $Mdn=14$, F3: $Mdn=14$), and for F1 in particular ($Mdn=16$). However they did not feel frustrated or annoyed during execution of testing of functionalities (F1: $Mdn=10$, F2: $Mdn=8$, F3: $Mdn=8$).

3) *Discussion*: Due to the small sample size, we can draw limited conclusions. Participants with little to no experience in robotics found HARMONI easy to use and, once set up, they found it simple. Setting up the open-source implementation on a new robot in HARMONI requires wrapping the sensor and actuator hardware according to the predefined API, which new users found moderately difficult the first time they used it, and reported that the required mental effort was high. After overcoming that hurdle, users found existing behaviors easy to map to the behavior API and new behaviors required low mental effort to define through the behavior definition JSON. If participants had the opportunity to familiarize themselves with HARMONI's APIs before the study testing, they might have found the setup process less challenging. We will include a familiarization phase for testing HARMONI in future studies. While this was a preliminary study with a small sample size, some of HARMONI's limitations were exposed, especially in relation to first-time use. That being said, once users became familiar with HARMONI, they were all successful in exploiting HARMONI's functions. We will incorporate user feedback

to improve initial use documentation in particular. Future studies will involve more participants to provide greater insights on the novice population.

VI. CONCLUSION & FUTURE WORK

We presented HARMONI, a robot interaction programming tool designed to be *open-source*, *robot-agnostic*, *modular*, *composable*, and *multimodal* to address many of the challenges inherent to social HRI. The tool is ROS-based and provides a well-documented interface. HARMONI aims to serve as a shared framework that can be deployed across different platforms and facilitate replication of social HRI experiments to improve the robustness of results. The modularity of the framework allows researchers to focus on HRI research challenges while benefiting from the state of the art from the rest of the field.

Our continuing work will aim to broaden the appeal of the tool to a wider robotics community. Within the current design pattern, development is ongoing to support a fuller range of interaction capabilities (including NLP and image processing) entirely offline. Additional support for behaviors that include locomotion and manipulation will also be added. Support for long running autonomy is also a high priority. Longer term work includes developing a visual programming interface for highest-level control and enabling multi-robot communication and support. A useful general-purpose open-source tool has the potential to radically increase the pace and quality of robotics research and development, as ROS has demonstrated. HARMONI's impact depends on usability and adoption. In addition to the planned improvements and expansions listed above, to reach its potential, the HARMONI project needs the adoption and support of the socially interactive robotics community.

ACKNOWLEDGMENT

This work was supported by the National Science Foundation Expeditions in Computing IIS-1139148.

REFERENCES

- [1] J. C. Russ, "Image processing," in *Computer-assisted microscopy*. Springer, 1990, pp. 33–69.
- [2] M. Sonka, V. Hlavac, and R. Boyle, *Image processing, analysis, and machine vision*. Cengage Learning, 2014.
- [3] A. Voulodimos, N. Doulamis, A. Doulamis, and E. Protopapadakis, "Deep learning for computer vision: A brief review," *Computational intelligence and neuroscience*, vol. 2018, 2018.
- [4] Y. Wilks, "Natural language understanding systems within the ai paradigm: A survey and some comparisons," *American Journal of Computational Linguistics*, 1976.
- [5] R. De Mori, "Spoken language understanding: a survey," in *2007 IEEE Workshop on Automatic Speech Recognition & Understanding (ASRU)*. IEEE, 2007, pp. 365–376.
- [6] X. Qiu, T. Sun, Y. Xu, Y. Shao, N. Dai, and X. Huang, "Pre-trained models for natural language processing: A survey," *arXiv preprint arXiv:2003.08271*, 2020.
- [7] D. W. Otter, J. R. Medina, and J. K. Kalita, "A survey of the usages of deep learning for natural language processing," *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [8] C. Badue, R. Guidolini, R. V. Carneiro, P. Azevedo, V. B. Cardoso, A. Forechi, L. Jesus, R. Berriel, T. M. Paixão, F. Mutz *et al.*, "Self-driving cars: A survey," *Expert Systems with Applications*, p. 113816, 2020.

- [9] A. S. Tulshan and S. N. Dhage, "Survey on virtual assistant: Google assistant, siri, cortana, alexa," in *International Symposium on Signal Processing and Intelligent Recognition Systems*. Springer, 2018, pp. 190–201.
- [10] T. B. Sheridan, "Human–robot interaction: status and challenges," *Human factors*, vol. 58, no. 4, pp. 525–532, 2016.
- [11] K. Dautenhahn, "Socially intelligent robots: dimensions of human–robot interaction," *Philosophical transactions of the royal society B: Biological sciences*, vol. 362, no. 1480, pp. 679–704, 2007.
- [12] T. Belpaeme, P. Baxter, J. De Greeff, J. Kennedy, R. Read, R. Looije, M. Neerinx, I. Baroni, and M. C. Zelati, "Child-robot interaction: Perspectives and challenges," in *International Conference on Social Robotics*. Springer, 2013, pp. 452–459.
- [13] P. Tsarouchi, S. Makris, and G. Chryssolouris, "Human–robot interaction review and challenges on task planning and programming," *International Journal of Computer Integrated Manufacturing*, vol. 29, no. 8, pp. 916–931, 2016.
- [14] T. Fong, I. Nourbakhsh, and K. Dautenhahn, "A survey of socially interactive robots," *Robotics and autonomous systems*, vol. 42, no. 3–4, pp. 143–166, 2003.
- [15] R. Mead, E. Wade, P. Johnson, A. S. Clair, S. Chen, and M. J. Mataric, "An architecture for rehabilitation task practice in socially assistive human-robot interaction," in *19th International Symposium in Robot and Human Interactive Communication*. IEEE, 2010, pp. 404–409.
- [16] B. R. Duffy, M. Dragone, and G. M. O'Hare, "Social robot architecture: A framework for explicit social interaction," in *Android Science: Towards Social Mechanisms, CogSci 2005 Workshop, Stresa, Italy*, 2005, pp. 3–4.
- [17] D. Feil-Seifer and M. J. Mataric, "B 3 ia: A control architecture for autonomous robot-assisted behavior intervention for children with autism spectrum disorders," in *RO-MAN 2008-The 17th IEEE International Symposium on Robot and Human Interactive Communication*. IEEE, 2008, pp. 328–333.
- [18] M. Quigley, K. Conley, B. Gerkey, J. Faust, T. Foote, J. Leibs, R. Wheeler, and A. Y. Ng, "Ros: an open-source robot operating system," in *ICRA workshop on open source software*, vol. 3, no. 3.2. Kobe, Japan, 2009, p. 5.
- [19] D. Merkel, "Docker: lightweight linux containers for consistent development and deployment," *Linux journal*, vol. 2014, no. 239, p. 2, 2014.
- [20] S. Gulabani, *Amazon Web Services Bootcamp: Develop a scalable, reliable, and highly available cloud environment with AWS*. Packt Publishing Ltd, 2018.
- [21] R. Collobert, C. Puhersch, and G. Synnaeve, "Wav2letter: an end-to-end convnet-based speech recognition system," *arXiv preprint arXiv:1609.03193*, 2016.
- [22] E. Short and M. Mataric, "Towards autonomous moderation of an assembly game," in *Workshop on Groups in Human-Robot Interaction at RO-MAN*, 2016.
- [23] "Qtrobot: Humanoid social robot for research and teaching." [Online]. Available: <http://luxai.com/qtrobot-for-research/>
- [24] M. Suguitan and G. Hoffman, "Blossom: A handcrafted open-source robot," *ACM Transactions on Human-Robot Interaction (THRI)*, vol. 8, no. 1, pp. 1–27, 2019.
- [25] N. Lazzeri, D. Mazzei, L. Cominelli, A. Cisternino, and D. E. De Rossi, "Designing the mind of a social robot," *Applied Sciences*, vol. 8, no. 2, p. 302, 2018.
- [26] E. Deng, B. Mutlu, M. J. Mataric *et al.*, "Embodiment in socially interactive robots," *Foundations and Trends® in Robotics*, vol. 7, no. 4, pp. 251–356, 2019.
- [27] N. Dragoni, S. Giallorenzo, A. L. Lafuente, M. Mazzara, F. Montesi, R. Mustafin, and L. Safina, "Microservices: yesterday, today, and tomorrow," in *Present and ulterior software engineering*. Springer, 2017, pp. 195–216.
- [28] S. R. Schach, *Object-oriented and classical software engineering*. McGraw-Hill Companies,, 2011.
- [29] R. C. Martin, *Clean code: a handbook of agile software craftsmanship*. Pearson Education, 2009.
- [30] E. Pot, J. Monceaux, R. Gelin, and B. Maisonnier, "Choregraphe: a graphical tool for humanoid robot programming," in *RO-MAN 2009-The 18th IEEE International Symposium on Robot and Human Interactive Communication*. IEEE, 2009, pp. 46–51.
- [31] "Nao the humanoid and programmable robot: Softbank robotics." [Online]. Available: <https://www.softbankrobotics.com/emea/en/nao>
- [32] A. Sauppé and B. Mutlu, "Design patterns for exploring and prototyping human-robot interactions," in *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, 2014, pp. 1439–1448.
- [33] V. Rajendran, P. Carreno-Medrano, W. Fisher, A. Werner, and D. Kulić, "A framework for human-robot interaction user studies," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 6215–6222.
- [34] D. Glas, S. Satake, T. Kanda, and N. Hagita, "An interaction design framework for social robots," in *Robotics: Science and Systems*, vol. 7, 2012, p. 89.
- [35] D. F. Glas, T. Kanda, and H. Ishiguro, "Human-robot interaction design using interaction composer eight years of lessons learned," in *2016 11th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, 2016, pp. 303–310.
- [36] E. Coronado, F. Mastrogiovanni, B. Indurkha, and G. Venture, "Visual programming environments for end-user development of intelligent and social robots, a systematic review," *Journal of Computer Languages*, vol. 58, p. 100970, 2020.
- [37] T. Lourens and E. Barakova, "User-friendly robot environment for creation of social scenarios," in *International Work-Conference on the Interplay between Natural and Artificial Computation*. Springer, 2011, pp. 212–221.
- [38] E. I. Barakova, J. C. Gillesen, B. E. Huskens, and T. Lourens, "End-user programming architecture facilitates the uptake of robots in social therapies," *Robotics and Autonomous Systems*, vol. 61, no. 7, pp. 704–713, 2013.
- [39] F. Siepman, "Behavior coordination for reusable system design in interactive robotics," *pub.uni-bielefeld.de*, 2013.
- [40] C.-M. Huang and B. Mutlu, "Robot behavior toolkit: generating effective social behaviors for robots," in *2012 7th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, 2012, pp. 25–32.
- [41] A. Hartholt, D. Traum, S. C. Marsella, A. Shapiro, G. Stratou, A. Leuski, L.-P. Morency, and J. Gratch, "All together now," in *International Workshop on Intelligent Virtual Agents*. Springer, 2013, pp. 368–381.
- [42] T. Groechel, Z. Shi, R. Pakkar, and M. J. Mataric, "Using socially expressive mixed reality arms for enhancing low-expressivity robots," in *2019 28th IEEE International Conference on Robot and Human Interactive Communication (RO-MAN)*. IEEE, 2019, pp. 1–8.
- [43] R. Wallace, "The elements of aiml style. alice ai foundation, inc.," 2016.
- [44] T. Michael, "Retico: An incremental framework for spoken dialogue systems," in *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, 2020, pp. 49–52.
- [45] C. Kennington, D. Moro, L. Marchand, J. Carns, and D. McNeill, "rrsds: Towards a robot-ready spoken dialogue system," in *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, 2020, pp. 132–135.
- [46] A. Papangelis, M. Namazifar, C. Khatri, Y.-C. Wang, P. Molino, and G. Tur, "Plato dialogue system: A flexible conversational ai research platform," *arXiv preprint arXiv:2001.06463*, 2020.
- [47] D. Bohus, S. Andrist, and M. Jalobeanu, "Rapid development of multimodal interactive systems: A demonstration of platform for situated intelligence," in *Proceedings of the 19th ACM International Conference on Multimodal Interaction*, ser. ICMI '17. New York, NY, USA: Association for Computing Machinery, 2017, p. 493–494. [Online]. Available: <https://doi.org/10.1145/3136755.3143021>
- [48] "Behavior manager," Apr 2018. [Online]. Available: <https://opsive.com/support/documentation/behavior-designer/behavior-manager/>
- [49] M. Colledanchise and P. Ögren, "Behavior trees in robotics and ai: an introduction," *arXiv preprint arXiv:1709.00084*, 2017.
- [50] J. Brooke, "Sus: a 'quick and dirty' usability," *Usability evaluation in industry*, p. 189, 1996.
- [51] S. G. Hart, "Nasa-task load index (nasa-tlx); 20 years later," in *Proceedings of the human factors and ergonomics society annual meeting*, vol. 50, no. 9. Sage publications Sage CA: Los Angeles, CA, 2006, pp. 904–908.
- [52] S. C. Peres, T. Pham, and R. Phillips, "Validation of the system usability scale (sus) sus in the wild," in *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, vol. 57, no. 1. SAGE Publications Sage CA: Los Angeles, CA, 2013, pp. 192–196.