

# Communicating human intent to a robotic companion by multi-type gesture sentences

Petr Vanc<sup>1</sup>Jan Kristof Behrens<sup>1</sup>Karla Stepanova<sup>1</sup>Vaclav Hlavac<sup>1</sup>

**Abstract**—Human-Robot collaboration in home and industrial workspaces is on the rise. However, the communication between robots and humans is a bottleneck. Although people use a combination of different types of gestures to complement speech, only a few robotic systems utilize gestures for communication. In this paper, we propose a gesture pseudo-language and show how multiple types of gestures can be combined to express human intent to a robot (i.e., expressing both the desired action and its parameters - e.g., pointing to an object and showing that the object should be emptied into a bowl). The demonstrated gestures and the perceived tabletop scene (object poses detected by CosyPose) are processed in real-time) to extract the human’s intent. We utilize behavior trees to generate reactive robot behavior that handles various possible states of the world (e.g., a drawer has to be opened before an object is placed into it) and recovers from errors (e.g., when the scene changes). Furthermore, our system enables switching between direct teleoperation of the end-effector and high-level operation using the proposed gesture sentences. The system is evaluated on increasingly complex tasks using a real 7-DoF Franka Emika Panda manipulator. Controlling the robot via action gestures lowered the execution time by up to 60%, compared to direct teleoperation.

**Index Terms**—Human-robot collaboration, Intent recognition, Scene awareness, Gesture detection

## I. INTRODUCTION

There are many ways how to control a robot with hand gestures. Commonly used are teleoperation [1] and pointing (deictic) gestures [2]. Direct teleoperation enables users to control robots in real-time and is useful to operate the robots in places that are inaccessible or dangerous for humans (e.g., robotic surgery and safe and rescue operations in harsh environments). However, it requires the user’s full attention and becomes infeasible in situations with high communication delays. Space robotics is a prime example where communication delays make some autonomy necessary during teleoperation (ground-in-the-loop) [3].

In this paper, we propose a system that enables the user to communicate its high-level intent to the robot via gestures. Knowing the user’s intent enables the robot to execute many of the required actions autonomously. Both the robot and the human can request to change the control mode. In this way, the system can disambiguate the instructions or acquire further specifications.

Different modalities allow for communicating different concepts more efficiently and complement each other. While language might be good for establishing order, gestures

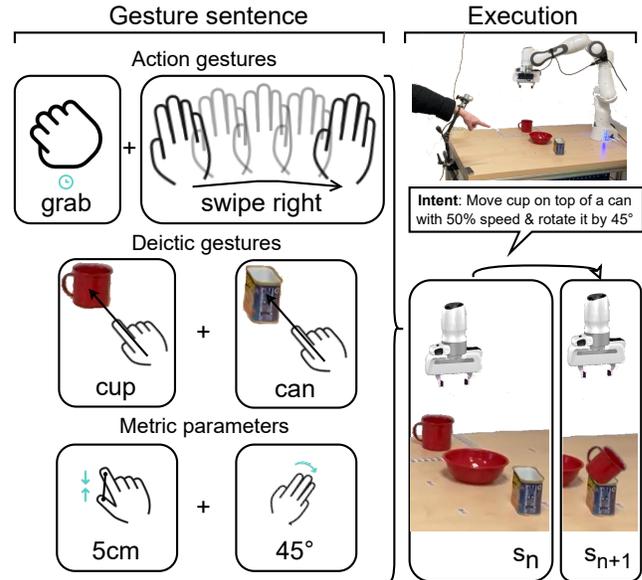


Fig. 1: Human operating a robot with gestures.

are more convenient for referring to objects and places, to express angle-related parameters or trajectory assignments. In general, gestures have a more limited vocabulary than speech. Therefore, one gesture might not be enough to express a complex user intent, and several gestures must be combined. In this paper, we specify gesture sentence that enables expression of a complex user intent to the robot via gesture pseudo-language.

In this work, we present a highly expressive gesture pseudo-language. Its purpose is to increase the efficiency and comfort of teleoperation of the robots. We show that in the application of tabletop manipulation tasks, the operator spends less time giving instructions compared to other systems. The main contributions of the paper are:

- Definition of a gesture sentence composed of different gesture types (i.e., action, deictic, and metric gestures).
- Extension of the gesture framework [4] to handle various types of gestures, disambiguate between different gesture meanings, and determine user intent based on the observed gesture sentence.
- Evaluation on the increasingly complex tasks in simulation and on the real setup.

The source code, accompanying video, and additional materials are available on the project webpage <http://imitrob.ciirc.cvut.cz/gestureSentence.html>.

<sup>1</sup>Czech Technical University in Prague, Czech Institute of Informatics, Robotics, and Cybernetics, petr.vanc@cvut.cz, jan.kristof.behrens@cvut.cz, karla.stepanova@cvut.cz, vaclav.hlavac@cvut.cz

## II. RELATED WORK

The current systems that utilize gestures to communicate with a human mainly focus on the teleoperation of the robot (e.g., [5]) or use gestures to command the robot without the option to connect the human gestures to the current state of the scene (e.g., [6], [7], [8]). Zhang et al. [5] uses similarly to us Leap motion sensor to detect hand movements and classify individual gestures. A basic set of gestures is used to teleoperate the robotic arm without connection to the given scene. In [6], wearable sensors were used to detect both dynamic and static gestures and users were provided with visual feedback. The Meguru system proposed in [7] enables users to combine several actions and connect gestures to higher-level actions. However, none of these systems enables users to specify flexibly objects to manipulate (e.g., by pointing) nor the parameters of the performed action (e.g., amount of degrees to rotate). Furthermore, these works do not enable flexible switching between direct teleoperation of the robot and gesture-based operation during execution.

In other works, such as [9], gestures are used to point to objects that should be manipulated. In general, gestures are used in current works as one-word commands and not as real sentences containing verbs, subjects and further parameters. This means that they cannot properly pass the human intent to the robotic system.

There are several works focused on human intent inference [10]. In [11], they used recursive Bayesian Inference, which can fuse multiple observations, formulated a mathematical model, and estimated human intent based on the human’s joystick inputs. The estimated intent then influenced the amount of shared autonomy between a human and a robot. Shared autonomy methods were also used in [12], where the robot actively reveals its chosen convention to speed up the novice introduction to the system. Their assist approach improves the success rate of a task based on observations from joystick input than without assistance. In these works, the intent is estimated based on joystick operation. The more natural way of communication, such as human gestures and the context of the current scene, is not taken into account. In our approach, we merge more input modalities (e.g., gestures and scene information).

In our previous work [4] we showed how to estimate human intent from a combination of the scene and gesture vector. However, this concept wasn’t tested in a real setup and only one type of gesture was used at a time with a simulated region of interest. In this paper, we propose a more general gesture pseudo-language that uses a sequence of gestures of various types (i.e., static, dynamic, deictic, teleoperation) to define human intent. The system is evaluated on a set of tasks with increasing complexity on a real robot.

## III. PROBLEM FORMULATION

*Problem definition:* Given an observation  $\mathbf{o}$ , compute human intent  $\mathbf{i}$  and based on it generate a sequence of robotic actions  $\mathbf{a}$  to fulfill this intent.

*Observation tuple* defined as  $\mathbf{o} = [\mathbf{h}, \mathbf{p}, \mathbf{s}]$  is composed of time-series of hand gesture features  $\mathbf{h}$  and deictic gesture features  $\mathbf{p}$ , both of them are computed from the raw hand data.  $\mathbf{s}$  represents the context of the situation and state of the system. It includes the configuration of the objects and the robot (robot end-effector, object positions, etc.).

The *human intent*  $\mathbf{i}$  represents the desired change of the world state (e.g., how objects on the scene should move, if the robot gripper should open, etc.). To be able to estimate the human intent from the observation  $\mathbf{o}$ , we have to know (or learn) also the mapping  $\mathcal{M}: \mathbf{i} = \mathcal{M}(\mathbf{o})$ . The intent is mapped to the robot actions by mapping  $\mathcal{A}: \mathbf{a} = \mathcal{A}(\mathbf{i})$ .

Individual gestures can signal different types of information (e.g., desired action, parameters of the action, or the target object). Our system fuses several gesture demonstrations into a gesture sentence with a subject, verb, objects, and further quantitative and qualitative parameters. Here we list the main types of gestures. More information on individual gestures can be found in the Sec. IV:

- 1) Action gestures (static and dynamic): individual gestures are mapped to various intended actions, such as pick up, open, or push. Based on the pre-trained set of gestures  $\mathbf{G}$ , hand observations  $\mathbf{h}$ , and classifier  $\mathcal{G}$ , our system outputs probabilities of individual gestures  $\mathbf{g}$ ,

$$\mathbf{g} = \mathcal{G}(\mathbf{h}, \mathbf{G}).$$

- 2) Deictic gesture. This gesture enables the user to select region of interest (objects) in the scene. The direction vector  $\mathbf{p}_{\text{line}}$  (direction of the pointing finger or hand palm direction) is extracted from the tracked hand configuration. Considering known object poses in the scene  $\mathbf{o}_{\text{poses}}$ , probabilities of individual objects to be the target object can be estimated based on  $\mathbf{t}_{\text{dists}}$ . Vector  $\mathbf{t}_{\text{dists}}$  represents the distance of each object from  $\mathcal{O}$  to the direction vector:

$$\mathbf{t}_{\text{dists}} = \mathcal{D}(\mathbf{p}_{\text{line}}, \mathbf{o}_{\text{poses}})$$

Individual distances are computed by finding the closest distance from the object  $s_i \in \mathcal{O}$  to the direction line:

$$d_i^2 = \frac{|(\mathbf{p}_2 - \mathbf{p}_1) \times (\mathbf{p}_1 - \mathbf{s})|^2}{|\mathbf{p}_2 - \mathbf{p}_1|^2},$$

where  $\mathbf{p}_1$  and  $\mathbf{p}_2$  are two random points from the 3D line and  $\mathbf{s}_i$  is a position of the object  $i$ , all vectors are 3D Cartesian vectors.

- 3) Metric information. These gestures are optional and enhance the input gesture information (otherwise default values are used). They are continuous variables, e.g., distance between index finger and thumb is mapped to the parameter value such as distance of the push action or the angle when picking up an object (see examples in Fig. 2).
- 4) Direct teleoperation of the robot. Some action gestures might engage teleoperation mode. When active, the tracked hand center palm position and rotation around

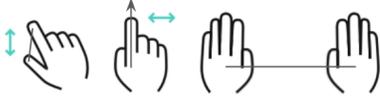


Fig. 2: Metric gesture example, they are auxiliary and optional. Pinch distance (left), point direction (center), hands distance (right)

$z$  axis is directly mapped onto the robot position and 7th DoF rotation, respectively.

Finally, gesture sentence  $S$  can be defined as a tuple of observed action gestures, objects specified by deictic gestures, and metric parameters specified by metric gestures:

$$\mathbf{S} = (\mathbf{action}, \mathit{object}(s), \mathit{metric\ parameters}). \quad (1)$$

While the **action** is always required, object specification is required only for specific actions and metric parameters are auxiliary (if not defined, default values are used).

Individual sentences corresponds to the user intent  $\mathbf{i}$  (as introduced in [4]), i.e. target action  $ta$ , target object  $to$ , target location, and other auxiliary metric parameters  $\mathbf{ap}$ :

$$\mathbf{S} \sim \mathbf{i} = (\underbrace{\underset{ta}{\operatorname{argmax}}(\mathbf{i}_{ta})}_{\mathcal{O}}, \underbrace{\underset{to}{\operatorname{argmin}}(\mathbf{to}_{\text{dists}})}_{\mathcal{O}}, \mathbf{ap}), \quad (2)$$

where  $\mathbf{i}_{ta} = \mathcal{M}_a(\mathbf{g})$  is *target action intent* vector probabilities generated from gestures vector  $\mathbf{g}$ ,  $\mathbf{to}_{\text{dists}} = \mathcal{D}_a(\mathbf{p}_{\text{line}}, \mathbf{o}_{\text{poses}})$  are *target object* vector probabilities.  $ta$  represents *target action intent*,  $to$  *target object*, and  $\mathbf{ap}$  set of auxiliary parameters.

An example of gesture sentence is:

$$S = (\mathbf{thumbsup}, [\mathit{mug}, \mathit{bowl}], [5\text{cm}]). \quad (3)$$

This corresponds to the user intent ("Move mug to the bowl by 50% speed"):

$$i = (\mathbf{move}, \mathit{mug}, [\mathit{bowl}, 50\%]) \quad (4)$$

The *robotic action sequence*  $\mathbf{a}$  is generated from estimated human intent  $\mathbf{i}$ . In this work, it is solved by behavior trees [13] (see example in Fig. 3). The whole structure can be seen in Fig. 5.

#### A. Complexity of the gesture sentence

The gesture sentences have variable complexity based on the number of variables that have to be specified for the given robotic action. The gesture sentence has only one required parameter ( $ta$ ). We define the complexity of the gesture sentence  $S = (ta, \mathbf{to}, \mathbf{ap})$  corresponding to robotic action  $a$  with  $N$  parameters  $\mathbf{p}$  ( $a(\mathbf{p})$ ) as:

$$S_c = N = \mathit{len}(\mathbf{to}) + \mathit{len}(\mathbf{ap})$$

. Simple robotic actions need to specify only the type of the action and have no additional parameters nor object dependency (e.g., move right based on grid world step). In this case, gesture sentence complexity  $S_c = 0$ . For actions

such as *Pick up a cup*, we must specify the action and the target object ( $S_c = 1$ ). Complex robotic actions can depend on two or more scene objects (e.g., swap the position of two objects).

Our system contains the following set of robotic actions with increasing number of parameters:

- 0-parameters: Cartesian moves, rotate, place
- 1-parameter: Pick up, Put, Pour, Open, Close
- 2-parameters: Put, Pour, Swap

Note that e.g. action *Put* can either have 1- parameter (if we are holding the target object then only action and target location has to be specified), or 2- parameters (if the gripper is empty we have to specify action, target object and target location).

## IV. MATERIALS AND METHODS

The problem described in Section 3 can be divided into several parts: user hand detection and bone reconstruction, gesture detection, user intent recognition, and robotic action generation and its execution by the robot. In our system, we don't consider hand detection and bone reconstruction. The overall system is visualized in Fig. 5.

Our system considers 1) Gesture recognition  $\mathcal{G}$ , 2) Intent classification  $\mathcal{M}$ , and 3) Robotics action generation  $\mathcal{A}$ . Parts of the system are introduced in more detail in the next sections and the last subsection connects all methods together (Sec. IV-F).

#### A. Real-time gesture recognition

We differentiate between static and dynamic gestures. Static gestures contain only a single time frame while typically having more features and dynamic have more time frames and have fewer features. Features are extracted from hand bone structure. We utilize the Leap Motion sensor [14] to obtain hand feature data. We utilize Gesture Toolbox [15] to manage gestures for a given session. The set of gestures may be adapted to current needs.

*Static gestures* feature set is hand-crafted, processes hand bone structure, and has 57 features: 1) distances between fingertip positions, alongside with the palm center position, and 2) joint angles between hand bones  $\alpha_1, \dots, \alpha_n$  constructed as angles between direction vectors. Direction vectors are made from each bone's start and end positions. The best-suitable method to classify static gestures has proven to be a probabilistic neural network model [16]. Other approaches, including convolutional neural network and deterministic (hand-picked threshold) approach, were tested while scoring worse (see Sec. VI-B).

*Dynamic gestures* are evaluated as hand trajectories over time. The used sample configuration is composed only of the palm center position (Cartesian) in each time step:

$$\mathbf{h} = \left[ \begin{array}{c|c|c} x_1 & x_2 & x_N \\ y_1 & y_2 & y_N \\ z_1 & z_2 & z_N \end{array} \right], \quad (5)$$

Where  $N$  is a number of trajectory points obtained from recording with a frequency around 90Hz, we reduce the

frequency of received frames to  $f = 20\text{Hz}$  without any significant loss [17]. We set the detection frequency to 10Hz, each detection takes the last time window with adjustable length (here set to 1s).

Dynamic gestures are learnt from a set of demonstrations and represented as probabilistic motion primitives [18]. As a classifier, we use the Dynamic Time Warping method (DTW) [19], which can compare two motions regardless of speed and deformations. The comparison parameter is distance  $\delta$ . Every gesture has its saved representative motion  $(\mathbf{h}_1, \dots, \mathbf{h}_D)$ , where  $D$  is a number of dynamic gestures in the current set, and  $\mathbf{h}$  is defined in 5. Then the most probable gesture is obtained with  $\text{argmin}((\delta)_1^D) = \text{argmin}(\text{DTW}((\mathbf{h}_1^D, \mathbf{h}_s))$ , where  $\mathbf{h}_s$  is sample motion, which we want to classify.

### B. Deictic gesture

The deictic gesture is a special type of static gestures with a special methodology for detecting intended target point or object (see Sec. 1 how the closest object to the direction line is found). In the case, that we want to select a given position on the table instead of a target object (e.g., when placing an object to position), the deictic approach finds the contact of pointing 3D line with the workspace table.

The solution (the closest object) may be computed in each frame, but its information is used only when the user points to an object or place. This gesture together with action gestures enables us to detect the desired human intent. The accuracy and usability of this gesture is evaluated in a separate experiment in Sec. VI-A.

### C. Direct robot teleoperation

Direct teleoperation, when enabled, maps the chosen point in the user’s hand (e.g., the center of hand palm) to the robot end-effector in real-time. To achieve smooth control, we update new values to the robot with a frequency of at least 10Hz. We use position control wrapped around velocity control to avoid shaking when changing the robot’s course while the robot is moving.

The mapping feature vector comprises a palm Cartesian position ( $p = [x, y, z]$ ) and hand rotation around the  $z$ -axis  $\theta$  to accommodate the robot’s 7th degree of freedom. This is important when we want to grasp an object in the wrong configuration. We tested mapping the full rotation of the human hand, but the teleoperation experience and accuracy were worse because of the user concentration on that many angles.

Direct teleoperation enables a demonstration of unknown lower-level motion primitive to teach the robot a new skill. Later it can be mapped to higher-level commands (gestures). Furthermore, it enables to adjust position of the robot when the object detections are not perfect.

### D. Intended action estimation

The intended target action  $\mathbf{i}_{ta}$  is estimated from the hand gestures  $\mathbf{g}$  and context information such as scene  $\mathbf{s}$  same as in [4]. This mapping task ( $\mathcal{M}$ ) is treated as a classification problem, a shallow probabilistic neural network is sufficient

for having few scene features on input, having two fully connected hidden layers. The output is a categorical distribution [20] commonly used for getting a discrete output. The automatic differentiation variational inference (ADVI) [21] fits the network weights using a Kullback-Leibler divergence loss function. We are using PyMC [16] framework.

### E. Action sequence generation and execution

The intent  $\mathbf{i}$  tuple (*target action, target object, (optionally) auxiliary parameters*) is mapped to a set of robotics actions  $\mathbf{a}$ . We utilize a behavior tree to encode this high-level policy which enables us to handle more complex tasks, for example, grasping an object that is not directly visible. The high-level policy of behavior tree recalculates all tasks for every executed action, so it can handle situations when the previous action fails. The behavior tree is embedded as generating precondition actions to fulfil the desired intent without the need to complete it manually (see Fig. 3 for an example of behavior tree used in our experiments).

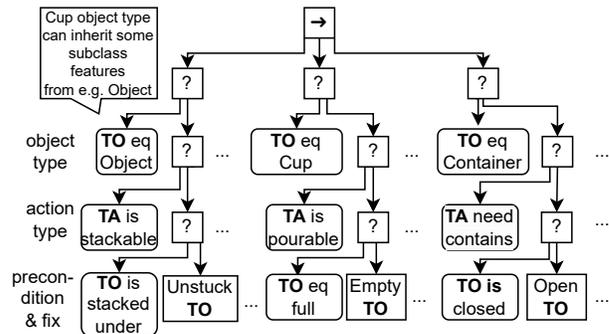


Fig. 3: Behavior tree structure example is solved as completing object state preconditions. For example, picking up an object stacked up with other objects requires to first unstacking all objects above it.

### F. Human-Robot interaction

The interaction between a human and a robot uses simple gesture sentences, where consecutive gestures compose a gesture sentence (see Eq. 1) and define the user intent (see Eq. 2). The user is expressing his intent through gestures while getting visual feedback (e.g. detected gesture, selected object, etc.) both from the graphical user interface and directly from the robot (see the accompanying video). The diagram of human-robot interaction is shown in Fig. 5. The user expresses each piece of information (i.e. action intent, object intent, auxiliary parameters) within a short interaction which we call *an episode*.

We define the *episode* as the time window when a human hand appears in the hand detection area and ends when it is no longer visible or the time exceeds the limit (set to  $t = 3\text{s}$ ). In one episode, we obtain a list of triggered gestures. Static and dynamic gestures are set up as two independent threads, where we get the recognition data from each. To be signaled as detected, the gesture needs to exceed the given threshold, i.e., the amount of time when the gesture exceeds

the threshold limit must be greater than  $0.3s$ . An example of gesture detection provided as optional feedback to the user is in Fig. 4.

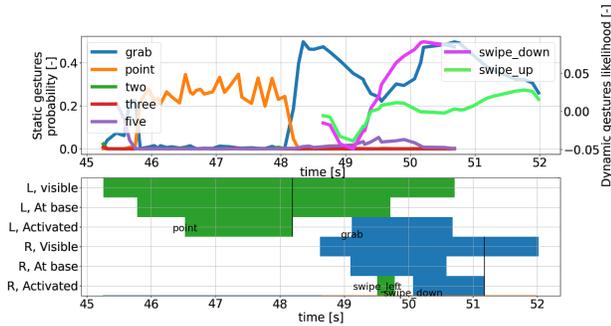


Fig. 4: The gesture detection output of a single episode. The upper plot shows the likelihood of static (left legend) and dynamic (right legend) gestures. The bottom plot shows the visibility of the hands and the activation of the gestures. First was detected the static gesture point followed by grab and dynamic gesture swipe down.

The interaction has the following steps (see Fig. 5): 1) The user signals the intended action by making an action gesture, which is recognized by gesture classifier  $\mathcal{G}$ . The recognizer is running the whole time when the hand is visible. 2) The action intent is estimated. Although multiple actions might be detected, in these experiments, the last detected gesture is used for determining the target action. 3) Based on the intended action, the user is asked to provide additional information about the performed action and add the target object and auxiliary parameters (e.g. target location, distance, etc.). 4) When the sentence information is completed (or the user does not want to provide additional information), the intent tuple is filled. 5) The behavior tree of actions may be generated, and a sequence of actions is executed, leading the robot to fulfil the intended task.

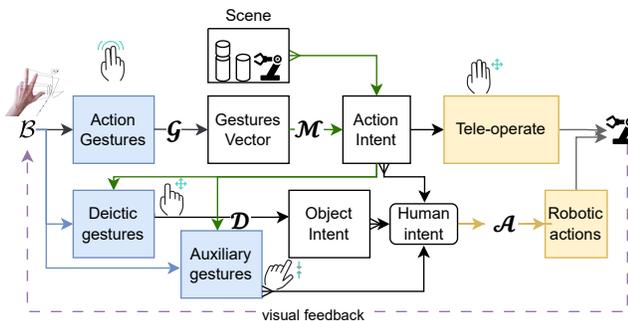


Fig. 5: System diagram.  $\mathcal{G}$  represents gesture classifier,  $\mathcal{M}$  is mapping from gestures to human intent,  $\mathcal{A}$  is robotic action generation using Behaviour tree,  $\mathcal{D}$  is Deictic gesture execution. Blue blocks represent user inputs, and yellow blocks are robot output modes.

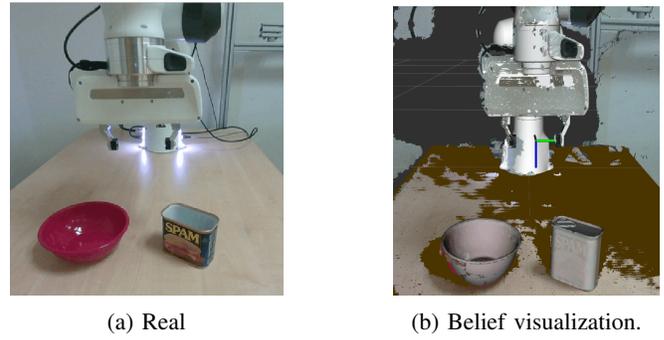


Fig. 6: Cosy Pose [24] object detection view using RViz visualization, which displays two detected objects. In the visualization, the meshes of detected objects (in grey) are displayed over the colored point cloud.

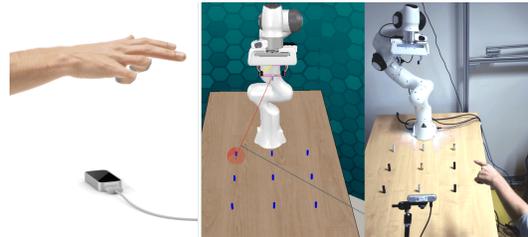


Fig. 7: Experimental setup. Leap Motion Controller (left), Simulation scene in Coppelia Sim (center), Real scene (right).

## V. EXPERIMENTAL SETUP

### A. Robot environment

Our system is tested in simulated and real environments (see Fig. 6 and Fig. 7). For the simulator, we use Coppelia Sim [22] with PyRep extension tool [23]. The scene contains a Franka Emika Panda with 7 DOF and several manipulation objects. An Intel Realsense D455 RGBD camera opposing the robot is used for object detection. A LeapMotion sensor is attached next to the camera on the table.

Our real-world setup (see Fig. 1) takes advantage of object pose detector CosyPose [24] (see Fig. 6), this allows us to get object poses in real-time. We use the YCB dataset [25] which contains various commonly used objects in the household. This allows us to create various scene configurations.

### B. Gesture classification

Our hand tracking device includes Leap Motion Controller [14], which provides us with hand bone structure in real-time. We utilize our Gesture Classification framework [15] to manage gesture sets, handle gesture recognition and visualization of each gesture probability, as well as a graphical interface and feedback for the user.

Static and dynamic gestures used in experiments include: grab, pinch, point, two, three, four five, thumbs up, swipe up, swipe down, swipe left, and swipe right gestures. The *no gesture* is used to represent no movement in case of dynamic gestures and it is excluded from the list.

### C. Description of experimental scenarios

We conducted a set of typical manipulation experiments of increasing complexity of the gesture sentence (see Sec. III-A for more details,  $x$  signs the object in the gripper):

- $S_c = 0$ ,  $S = (\mathbf{ta}, -, -)$ : (*rotate, x, -*), (*place, x, -*), (*move, -, -*)
- $S_c = 1$ ,  $S = (\mathbf{ta}, \mathbf{to}, -)$ : (*pick, can*), (*open, drawer*), (*close, drawer*)
- $S_c = 1$ ,  $S = (\mathbf{ta}, -, \mathbf{loc})$ : (*pour, x, bowl*), (*put, x, bowl*)
- $S_c = 2$ ,  $S = (\mathbf{ta}, \mathbf{to}, \mathbf{loc})$ : (*pour, spam, bowl*), (*swap, can, bowl*), (*put, spam, drawer*)
- $S_c = 3$ ,  $S = (\mathbf{ta}, \mathbf{to}, \mathbf{loc}, \mathbf{par})$ : (*rotate, spam, ang = 180^\circ*), (*pour, can, bowl, ang = 60^\circ*)
- **Multistep tasks**: stack 3 objects, tidy up 3 objects, pour 2 objects into the bowl
- **Infeasible tasks**: move object to occupied bowl, put object to closed drawer

We are interested in the total time of task completion, the success rate for each scenario, and the user interaction time. Execution of all these tasks is visible in the accompanying videos. For experiments with the users, we selected a subset of these tasks: 1) Put an object into the bowl, 2) Swap objects, and 3) Place the object in an object-occupied bowl, demonstrating our system's ability to first move the object out automatically. The following three assistance modes are evaluated: Tasks are tested based on three modes of assistance: 1) direct teleoperation, 2) low-level control via action gestures, and 3) high-level control via action gestures.

## VI. EXPERIMENTAL RESULTS

We show results for the following experiments:

- 1) Deictic gestures - 1) time to complete, 2) accuracy
- 2) Static and dynamic gestures accuracy
- 3) Scenario completion

All the scenes and tasks are set so that the robot is able to complete tasks using objects in the scene.

### A. Deictic gestures evaluation

A rough grid was created on a working table consisting of 9 small objects (see Fig. 7). The grid distance between two objects is set to  $d = 0.2m$ . The overall accuracy in detecting the intended (pointed) target object depends on the Leap Motion hand detection accuracy and accuracy of each transformation (e.g., the sensor to base, or hand to sensor).

Firstly, we tracked the hand-pointing finger bone positions, we observed high variances across measurements as the requirement is that the pointing finger is parallel to palm direction and requires an experienced user. Therefore we used values from hand palm.

Accuracies for detecting the target object in the grid for two users are listed in Tab. I. Feedback from the simulator (see Fig. 7) shows to the user which object is in aim. The better way than observing values on display may be operating the robot over the object that is currently a target. This increases overall system ergonomics and immersion with the

system. More detailed results with accuracies for individual positions are shown in Fig.8.

TABLE I: Total accuracies of Deictic gesture evaluated on 9 objects grid on 5 users. Firstly, users were given a blind test with no feedback. Secondly, they observed feedback on the screen. Finally, after training the accuracy was evaluated again with no feedback.

Accuracy [%]	Group members	Others
naive, no feedback	59.3	40.7
trained, no feedback	68.5	66.7
with feedback	100	100

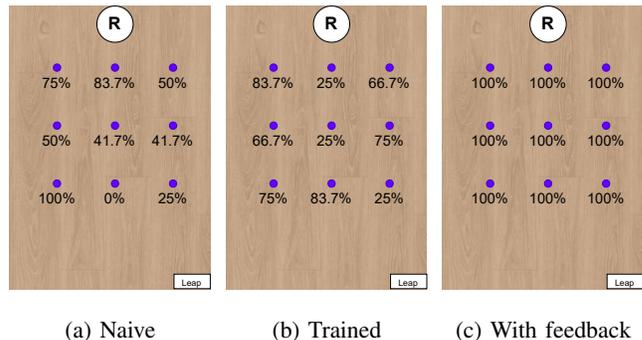


Fig. 8: Deictic gesture evaluation results. R represents the robot base, Leap represents the Hand sensor device and blue circles are target objects. The percentage represents the accuracy of both tested users for each object.

### B. Gesture classification evaluation

The gesture set needs to be chosen correctly, any two gestures cannot overlay between each other, then the classifier can distinguish between each gesture and correctly classify them. For our sample gesture set, we used 8 static gestures (grab, pinch, point, two, three, four, five, thumbs up), see Tab. II for results.

Dynamic gesture set consists of 5 gesture including the *no gesture*. The gestures represent directional swipes and are tuned for right-hand usage, the training dataset consists of  $\sim 300$  demonstrations, and the benchmark (see Tab. II) is made on the training dataset because no learning is happening.

TABLE II: Gesture classification benchmarks. Accuracy represents balanced accuracy  $BA = \frac{1}{2} \cdot (\frac{TP}{TP+FN} + \frac{TN}{TN+FP})$ , where T/F is True/False and P/N is Positive/Negative (more in [link]) and is evaluated on the test dataset. The static category has around 12000 samples and 4000 test samples.

Type	Method	Accuracy
Static	Deterministic	84.3%
	Probabilistic NN	99.1%
Dynamic	Euler	72.0%
	DTW	87.3%

### C. Scenarios evaluation

Users were introduced to gestures by which they can command the robot. After users got familiar with execution, they were able to perform low-level tasks (e.g., pick up, pour) with no problems. Most of the errors happened thanks to the object misdetections: e.g., CosyPose didn't recognize objects on the scene, added ghost objects on the scene randomly, or was crashing. Exemplar execution of the scenarios listed in Sec. V-C can be seen in the accompanying video and on the project webpage <http://imitrob.ciirc.cvut.cz/gestureSentence.html>.

The three selected scenarios (see Sec. V-C) were evaluated more thoroughly. Parameters like time completion, success rate, and percentage of human intervention were noted. The completion time of each scenario when using Teleoperation, Low-level actions and High-level actions is visualized in Fig. 9. You can see that teleoperating the robot directly needed the highest execution time and the user had to be also properly trained. The time needed to fulfil the task was decreasing during individual trials.

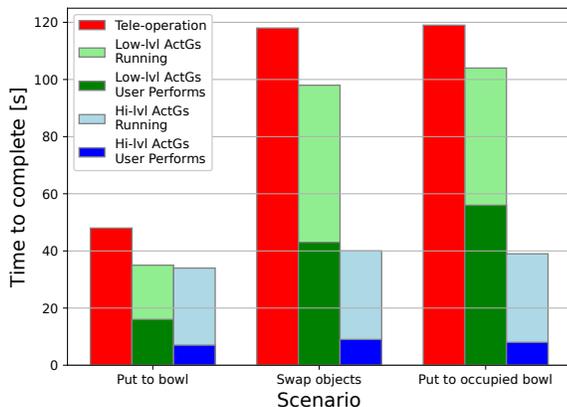


Fig. 9: Task completion by time per scenario. Only successful attempts were included. The system reduced user time needed by around 50% than total running time. The unfeasible placement task was the hardest to command the robot. Therefore it had the highest failure rate, it can be seen that even higher time commanding the robot.

## VII. CONCLUSION AND DISCUSSION

We proposed a system that can handle various gesture expressions while handling different gesture types (Action gestures, Deictic, etc.) and combining its information to express and execute the user intent.

Deictic test on users gave us information that visual feedback is a crucial component of such a system.

From the scenario test, we can see that having preprogrammed general, but high-level tasks results in task completion at the lowest time when gestures are performed correctly. Teleoperation is the slowest approach, but it can handle more fine-grained tasks. In the future, it may serve us as teaching

the robot new skill. When the skill is remembered, the low-level action gesture may be mapped to this skill. Then the action intent doesn't have to be manually programmed. Controlling the robot via action gestures lowered the execution task by up to 60%, than using a direct teleoperation method.

## VIII. ACKNOWLEDGMENT

This work was supported by the European Regional Development Fund under project Robotics for Industry 4.0 (reg. no. CZ.02.1.01/0.0/0.0/15\_003/0000470), MPO TRIO project num. FV40319, and by the Czech Science Foundation (project no. GA21-31000S). P.V. by CTU Student Grant Agency (reg. no. SGS23/138/OHK3-027/23).

## REFERENCES

- [1] A. D. Dragan and S. S. Srinivasa, *Formalizing assistive teleoperation*. MIT Press, July, 2012, vol. 376.
- [2] A. Kobsa, J. Allgayer, C. Reddig, N. Reithinger, D. Schmauks, K. Harbusch, and W. Wahlster, "Combining deictic gestures and natural language for referent identification," in *Proceedings of the 11th Conference on Computational Linguistics*, ser. COLING '86. USA: Association for Computational Linguistics, 1986, p. 356–361. [Online]. Available: <https://doi.org/10.3115/991365.991471>
- [3] I. A. Nesnas, L. M. Fesq, and R. A. Volpe, "Autonomy for space robots: Past, present, and future," *Current Robotics Reports*, vol. 2, no. 3, p. 251–263, Sep 2021.
- [4] P. Vanc, J. K. Behrens, and K. Stepanova, "Context-aware robot control using gesture episodes," in *2018 IEEE/RSJ ICRA*, 2023.
- [5] W. Zhang, H. Cheng, L. Zhao, L. Hao, M. Tao, and C. Xiang, "A gesture-based teleoperation system for compliant robot motion," *Applied Sciences*, vol. 9, no. 24, p. 5290, 2019.
- [6] P. Neto, M. Simão, N. Mendes, and M. Safeea, "Gesture-based human-robot interaction for human assistance in manufacturing," *The International Journal of Advanced Manufacturing Technology*, vol. 101, pp. 119–135, 2019.
- [7] C. Nuzzi, S. Pasinetti, R. Pagani, S. Ghidini, M. Beschi, G. Coffetti, and G. Sansoni, "Meguru: a gesture-based robot program builder for meta-collaborative workstations," *Robotics and Computer-Integrated Manufacturing*, vol. 68, p. 102085, 2021.
- [8] O. Mazhar, B. Navarro, S. Ramdani, R. Passama, and A. Cherubini, "A real-time human-robot interaction framework with robust background invariant hand gesture detection," *Robotics and Computer-Integrated Manufacturing*, vol. 60, pp. 34–48, 2019.
- [9] J. K. Behrens, K. Stepanova, R. Lange, and R. Skoviera, "Specifying dual-arm robot planning problems through natural language and demonstration," *IEEE Robotics and Automation Letters*, vol. 4, no. 3, pp. 2622–2629, 2019.
- [10] D. P. Losey, C. G. McDonald, E. Battaglia, and M. K. O'Malley, "A review of intent detection, arbitration, and communication aspects of shared control for physical human-robot interaction," *Applied Mechanics Reviews*, vol. 70, no. 1, 2018.
- [11] S. Jain and B. Argall, "Recursive bayesian human intent recognition in shared-control robotics," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2018, pp. 3905–3912.
- [12] A. Jonnavittula and D. P. Losey, "Communicating robot conventions through shared autonomy," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 7423–7429.
- [13] M. Colledanchise and P. Ögren, *Behavior Trees in Robotics and AI*. CRC Press, jul 2018. [Online]. Available: <https://doi.org/10.1201%2F9780429489105>
- [14] F. Weichert, D. Bachmann, B. Rudak, and D. Fisseler, "Analysis of the accuracy and robustness of the leap motion controller," *Sensors*, vol. 13, no. 55, p. 6380–6393, May 2013.
- [15] P. Vanc, "Gesture teleoperation toolbox v.0.1," <https://github.com/imitrob/teleopgesturetoolbox/>, 2022, [Online; accessed 2-March-2023].
- [16] D. Emaasit, "Pymc-learn: Practical probabilistic machine learning in python," no. arXiv:1811.00542, Oct 2018, arXiv:1811.00542 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/1811.00542>

- [17] K. Forbes and E. Fiume, "An efficient search algorithm for motion data using weighted pca," in *Proceedings of the 2005 ACM SIGGRAPH/Eurographics symposium on Computer animation - SCA '05*. Los Angeles, California: ACM Press, 2005, p. 67. [Online]. Available: <http://portal.acm.org/citation.cfm?doid=1073368.1073377>
- [18] A. Paraschos, C. Daniel, J. R. Peters, and G. Neumann, "Probabilistic movement primitives," in *Advances in Neural Information Processing Systems*, vol. 26. Curran Associates, Inc., 2013. [Online]. Available: <https://proceedings.neurips.cc/paper/2013/hash/e53a0a2978c28872a4505bdb51db06dc-Abstract.html>
- [19] *Dynamic Time Warping*. Berlin, Heidelberg: Springer, 2007, p. 69–84. [Online]. Available: [https://doi.org/10.1007/978-3-540-74048-3\\_4](https://doi.org/10.1007/978-3-540-74048-3_4)
- [20] K. P. Murphy, *Machine learning : a probabilistic perspective*. Cambridge, Mass. [u.a.]: MIT Press, 2013. [Online]. Available: [https://www.amazon.com/Machine-Learning-Probabilistic-Perspective-Computation/dp/0262018020/ref=sr\\_1\\_2?ie=UTF8&qid=1336857747&sr=8-2](https://www.amazon.com/Machine-Learning-Probabilistic-Perspective-Computation/dp/0262018020/ref=sr_1_2?ie=UTF8&qid=1336857747&sr=8-2)
- [21] A. Kucukelbir, D. Tran, R. Ranganath, A. Gelman, and D. M. Blei, "Automatic differentiation variational inference," *Journal of machine learning research*, 2017.
- [22] E. Rohmer, S. P. N. Singh, and M. Freese, "V-rep: A versatile and scalable robot simulation framework," in *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems*, Nov 2013, p. 1321–1326.
- [23] S. James, M. Freese, and A. J. Davison, "Pyrep: Bringing v-rep to deep robot learning," no. arXiv:1906.11176, Jun 2019, arXiv:1906.11176 [cs]. [Online]. Available: <http://arxiv.org/abs/1906.11176>
- [24] Y. Labbé, J. Carpentier, M. Aubry, and J. Sivic, "Cosypose: Consistent multi-view multi-object 6d pose estimation," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVII 16*. Springer, 2020, pp. 574–591.
- [25] B. Calli, A. Walsman, A. Singh, S. Srinivasa, P. Abbeel, and A. M. Dollar, "Benchmarking in manipulation research: Using the yale-cmu-berkeley object and model set," *IEEE Robotics & Automation Magazine*, vol. 22, no. 3, pp. 36–52, 2015.