

NOIR: Neural Signal Operated Intelligent Robots for Everyday Activities

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** We present Neural Signal Operated Intelligent Robots (NOIR), a
2 general-purpose, intelligent brain-robot interface system that enables humans to
3 command robots to perform everyday activities through brain signals. Through
4 this interface, humans communicate their intended objects of interest and actions
5 to the robots using electroencephalography (EEG). Our novel system demon-
6 strates success in an expansive array of 20 challenging, everyday household ac-
7 tivities, including cooking, cleaning, personal care, and entertainment. The effec-
8 tiveness of the system is improved by its synergistic integration of robot learning
9 algorithms, allowing for NOIR to adapt to individual users and predict their inten-
10 tions. Our work enhances the way humans interact with robots, replacing tradi-
11 tional channels of interaction with direct, neural communication. Project website:
12 <https://sites.google.com/view/noir-corl2023>

13 **Keywords:** Brain-Robot Interface; Human-Robot Interaction

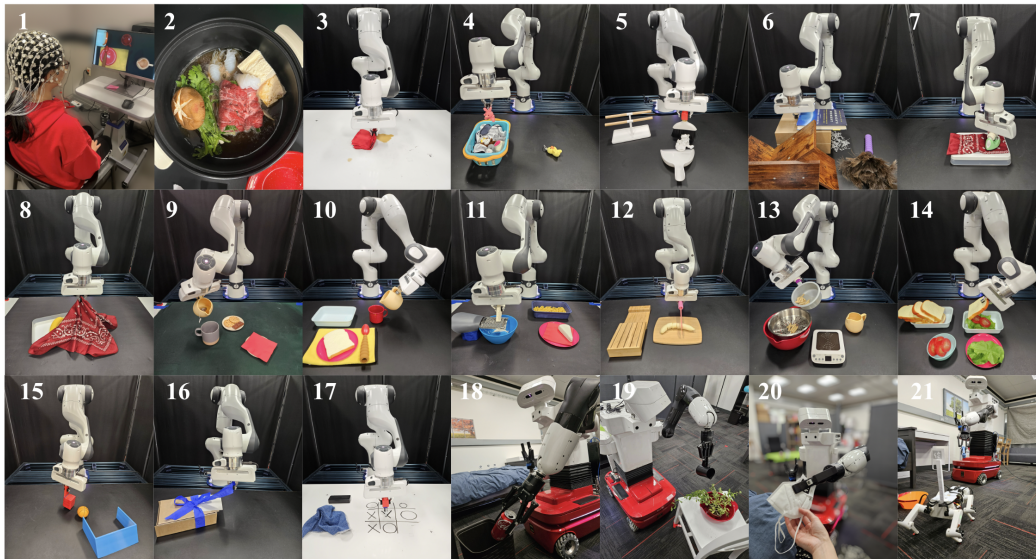


Figure 1: NOIR is a general-purpose brain-robot interface that allows humans to use their brain signals (1) to control robots to perform daily activities, such as making Suki-yaki (2), ironing clothes (7), playing Tic-Tac-Toe with friends (17), and petting a robot dog (21).

14 1 Introduction

15 Brain-robot interfaces (BRIs) are a pinnacle achievement in the realm of art, science, and engi-
16 neering. This aspiration, which features prominently in speculative fiction, innovative artwork, and
17 groundbreaking scientific studies, entails creating robotic systems that operate in perfect synergy

18 with humans. A critical component of such systems is their ability to communicate with humans.
19 In human-robot collaboration and robot learning, humans communicate their intents through ac-
20 tions [1], button presses [2, 3], gaze [4–7], facial expression [8], language [9, 10], etc [11, 12].
21 However, the prospect of direct communication through neural signals stands out to be the most
22 thrilling but challenging medium.

23 We present Neural Signal Operated Intelligent Robots (NOIR), a general-purpose, intelligent BRI
24 system with non-invasive electroencephalography (EEG). The primary principle of this system is
25 hierarchical shared autonomy, where humans define high-level goals while the robot actualizes the
26 goals through the execution of low-level motor commands. Taking advantage of the progress in
27 neuroscience, robotics, and machine learning, our system distinguishes itself by extending beyond
28 previous attempts to make the following contributions.

29 First, NOIR is *general-purpose* in its diversity of tasks and accessibility. We show that humans can
30 accomplish an expansive array of 20 daily everyday activities, in contrast to existing BRI systems
31 that are typically specialized at one or a few tasks or exist solely in simulation [13–22]. Additionally,
32 the system can be used by the general population, with a minimum amount of training.

33 Second, the “I” in NOIR means that our robots are *intelligent* and adaptive. The robots are equipped
34 with a library of diverse skills, allowing them to perform low-level actions without dense human su-
35 pervision. Human behavioral goals can naturally be communicated, interpreted, and executed by the
36 robots with *parameterized primitive skills*, such as `Pick(obj-A)` or `MoveTo(x, y)`. Additionally,
37 our robots are capable of learning human intended goals during their collaboration. We show that
38 by leveraging the recent progress in foundation models, we can make such a system more adaptive
39 with limited data. We show that this can significantly increase the efficiency of the system.

40 The key technical contributions of NOIR include a *modular* neural signal decoding pipeline for
41 human intentions. Decoding human intended goals (e.g., “pick up the mug from the handle”) from
42 neural signals is extremely challenging. We decompose human intention into three components:
43 *What* object to manipulate, *How* to interact with the object, and *Where* to interact, and show that
44 such signals can be decoded from different types of neural data. These decomposed signals naturally
45 correspond to parameterized robot skills and can be communicated effectively to the robots.

46 In 20 household activities involving tabletop or mobile manipulations, three human subjects suc-
47 cessfully used our system to accomplish these tasks with their brain signals. We demonstrate that
48 few-shot robot learning from humans can significantly improve the efficiency of our system. This
49 approach to building intelligent robotic systems, which utilizes human brain signals for collabora-
50 tion, holds immense potential for the development of critical assistive technologies for individuals
51 with or without disabilities and to improve the quality of their life.

52 **2 Brain-Robot Interface (BRI): Background**

53 Since Hans Berger’s discovery of EEG in 1924, several types of devices have been developed to
54 record human brain signals. We chose non-invasive, saline-based EEG due to its cost and acces-
55 sibility to the general population, signal-to-noise ratio, temporal and spatial resolutions, and types
56 of signals that can be decoded (see Appendix 2). EEG captures the spontaneous electrical activity
57 of the brain using electrodes placed on the scalp. EEG-based BRI has been applied to prosthetics,
58 wheelchairs, as well as navigation and manipulation robots. For comprehensive reviews, see [22–
59 25]. We utilize two types of EEG signals that are frequently employed in BRI, namely, steady-state
60 visually evoked potential (SSVEP) and motor imagery (MI).

61 SSVEP is the brain’s exogenous response to periodic external visual stimulus [26], wherein the brain
62 generates periodic electrical activity at the same frequency as flickering visual stimulus. The appli-
63 cation of SSVEP in assistive robotics often involves the usage of flickering LED lights physically
64 affixed to different objects [27, 28]. Attending to an object (and its attached LED light) will increase
65 the EEG response at that stimulus frequency, allowing the object’s identity to be inferred. Inspired
66 by prior work [15], our system utilizes computer vision techniques to detect and segment objects,
67 attach virtual flickering masks to each object, and display them to the participants for selection.

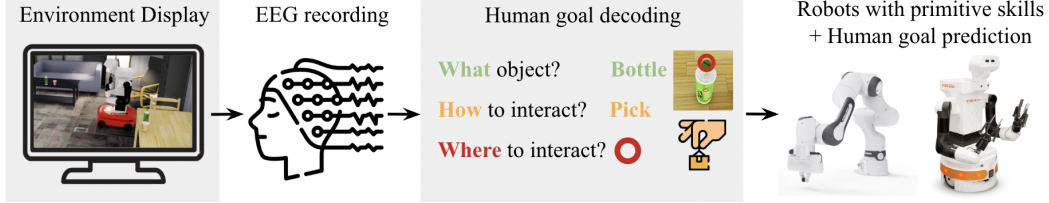


Figure 2: NOIR has two components, a modular pipeline for decoding goals from human brain signals, and a robotic system with a library of primitive skills. The robots possess the ability to learn to predict human intended goals hence reducing the human effort required for decoding.

68 Motor Imagery (MI) differs from SSVEP due to its endogenous nature, requiring individuals to
 69 mentally simulate specific actions, such as imagining oneself manipulating an object. The decoded
 70 signals can be used to indicate a human’s intended way of interacting with the object. This approach
 71 is widely used for rehabilitation, and for navigation tasks [29] in BRI systems. This approach often
 72 suffers from low decoding accuracy [22].

73 Much existing BRI research focuses on the fundamental problem of brain signal decoding, while
 74 several existing studies focus on how to make robots more intelligent and adaptive [13–17, 30]. In-
 75 spired by this line of work, we leverage few-shot policy learning algorithms to enable robots to learn
 76 human preferences and goals. This minimizes the necessity for extensive brain signal decoding,
 77 thereby streamlining the interaction process and enhancing overall efficiency.

78 Our study is grounded in substantial advancements in both the field of brain signal decoding and
 79 robot learning. Currently, many existing BRI systems target only one or a few specific tasks. To the
 80 best of our knowledge, no previous work has presented an intelligent, versatile system capable of
 81 successfully executing a wide range of complex tasks, as demonstrated in our study.

82 3 The NOIR System

83 The challenges we try to tackle are: 1) How do we build a general-purpose BRI system that works
 84 for a variety of tasks? 2) How do we decode relevant communication signals from human brains? 3)
 85 How do we make robots more intelligent and adaptive for more efficient collaboration? An overview
 86 of our system is shown in Fig. 2. Humans act as planning agents to perceive, plan, and communicate
 87 behavioral goals to the robot, while robots use pre-defined primitive skills to achieve these goals.

88 The overarching goal of building a general-purpose BRI system is achieved by synergistically inte-
 89 grating two designs together. First, we propose a novel *modular* brain decoding pipeline for human
 90 intentions, in which the human intended goal is decomposed into three components: what, how, and
 91 where (Sec. 3.1). Second, we equip the robots with a library of parameterized primitive skills to
 92 accomplish human-specified goals (Sec. 3.2). This design enables humans and robots to collaborate
 93 to accomplish a variety of challenging, long-horizon everyday tasks. At last, we show a key feature
 94 of NOIR to allow robots to act more efficiently and to be capable of adapting to individual users, we
 95 adopt few-shot imitation learning from humans (Sec. 3.3).

96 3.1 The brain: A modular decoding pipeline

97 We hypothesize that the key to building a general-purpose EEG decoding system is modularization.
 98 Decoding complete behavioral goals (e.g., in the form of natural language) is only feasible with ex-
 99 pensive devices like fMRI, and with many hours of training data for each individual [31]. As shown
 100 in Fig. 3, we decompose human intention into three components: (a) *What* object to manipulate; (b)
 101 *How* to interact with the object; (c) *Where* to interact. The decoding of specific user intents from
 102 EEG signals is challenging but can be done with steady-state visually evoked potential and motor
 103 imagery, as introduced in Sec. 2. For brevity, details of decoding algorithms are in Appendix 6.

104 **Selecting objects with steady-state visually evoked potential (SSVEP).** Upon showing the task
 105 set-up on a screen, we first infer the user’s intended object. We make objects on the screen flicker
 106 with different frequencies (Fig. 3a), which, when focused on by the user, evokes SSVEP [26]. By
 107 identifying which frequency is stronger in the EEG data, we may infer the frequency of the flick-

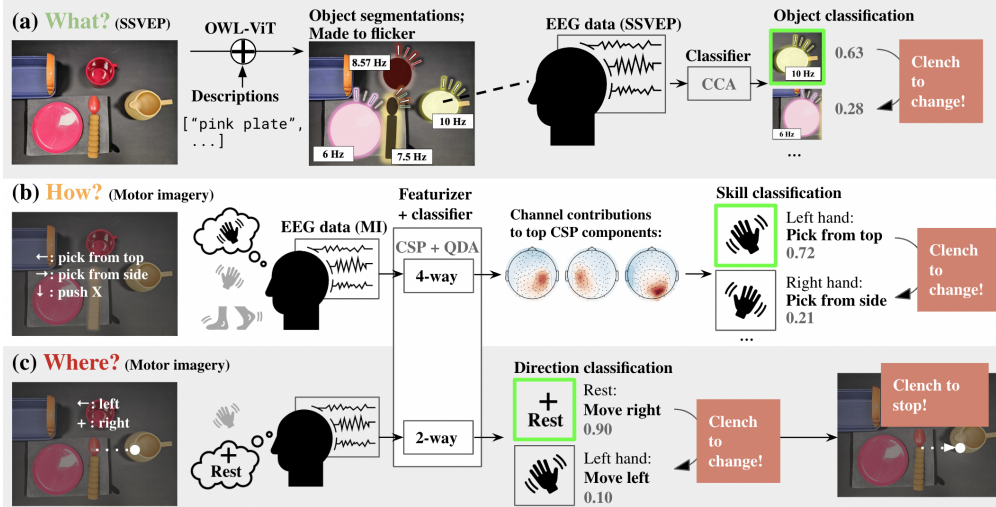


Figure 3: A modular pipeline for decoding human intended goals from EEG signals: (a) *What* object to manipulate, decoded from SSVEP signals using CCA classifiers; (b) *How* to interact with the object, decoded from MI signals using CSP+QDA algorithms; (c) *Where* to interact, decoded from MI signals. A safety mechanism that captures muscle tension from jaw clench is used to confirm or reject decoding results.

108 ering visual stimulus, and hence the object that the user focuses on. We apply modern computer
 109 vision techniques to circumvent the problem of having to physically attach LED lights [27, 28].
 110 Specifically, we use the foundation model OWL-ViT [32] to detect and track objects, which takes in
 111 an image and object descriptions and outputs object segmentation masks. By overlaying each mask
 112 of different flickering frequencies (6Hz, 7.5Hz, 8.57Hz, and 10Hz [33, 34]), and having the user
 113 focus on the desired object for 10 seconds, we are able to identify the attended object.

114 We use only the signals from the visual cortex (Appendix 6) and preprocess the data with a notch
 115 filter. We then use Canonical Correlation Analysis (CCA) for classification [35]. We create a Canon-
 116 ical Reference Signal (CRS), which is a set of sin and cos waves, for each of our frequencies and
 117 their harmonics. We then use CCA to calculate the frequency whose CRS has the highest correlation
 118 with the EEG signal, and identify the object that was made to flicker at that frequency.

119 **Selecting skill and parameters with motor imagery (MI).** The user then chooses a skill and its
 120 parameters. We frame this as a k -way ($k \leq 4$) MI classification problem, where we aim to decode
 121 which of the k pre-decided actions the user was imagining. Unlike SSVEP, a small amount of
 122 calibration data (10-min) is required due to the distinct nature of each user’s MI signals. The four
 123 classes are: Left Hand, Right Hand, Legs, and Rest; the class names describe the body parts that
 124 users imagine using to execute some skills (e.g. pushing a pedal with feet). Upon being presented
 125 with the list of k skill options, we record a 5-second EEG signal, and use a classifier trained on the
 126 calibration data. The user then guides a cursor on the screen to the appropriate location for executing
 127 the skill. To move the cursor along the x axis, the user is prompted to imagine moving their Left
 128 hand for leftward cursor movement. We record another five seconds of data and utilize a 2-way
 129 classifier. This process is repeated for x , y , and z axes.

130 For decoding, we use only EEG channels around the brain areas related to motor imagery (Appendix
 131 6). The data is band-pass-filtered between 8Hz and 30Hz to include μ -band and β -band frequency
 132 ranges correlated with MI activity [36]. The classification algorithm is based on the common spatial
 133 pattern (CSP) [37–40] algorithm and quadratic discriminant analysis (QDA). Due to its simplicity,
 134 CSP+QDA is explainable and amenable to small training datasets. Contour maps of electrode con-
 135 tributions to the top few CSP-space principal components are shown in the middle row of Fig. 3.
 136 There are distinct concentrations around the right and left motor areas, as well as the visual cortex
 137 (which correlates with the Rest class).

138 **Confirming or interrupting with muscle tension.** Safety is critical in BRI due to noisy decoding.
139 We follow a common practice and collect electrical signals generated from facial muscle tension
140 (Electromyography, or EMG). This signal appears when users frown or clench their jaws, indicating
141 a negative response. This signal is strong with near-perfect decoding accuracy, and thus we use it
142 to confirm or reject object, skill, or parameter selections. With a pre-determined threshold value
143 obtained through the calibration stage, we can reliably detect muscle tension from 500-ms windows.

144 **3.2 The robot: Parameterized primitive skills**

145 Our robots must be able to solve a diverse set of manipulation tasks under the guidance of humans,
146 which can be achieved by equipping them with a set of parameterized primitive skills. The benefits
147 of using these skills are that they can be combined and reused across tasks. Moreover, these skills
148 are intuitive to humans. Since skill-augmented robots have shown promising results in solving long-
149 horizon tasks, we follow recent works in robotics with parameterized skills [41–52], and augment
150 the action space of our robots with a set of primitive skills and their parameters. Neither the human
151 nor the agent requires knowledge of the underlying control mechanism for these skills, thus the skills
152 can be implemented in any method as long as they are robust and adaptive to various tasks.

153 We use two robots in our experiment: A Franka Emika Panda arm for tabletop manipulation tasks,
154 and a PAL Tiago robot for mobile manipulation tasks (see Appendix for hardware details). Skills for
155 the Franka robot use the operational space pose controller (OSC) [53] from the Deoxys API [54].
156 For example, *Reaching* skill trajectories are generated by numerical 3D trajectory interpolation
157 conditioned on the current robot end-effector 6D pose and target pose. Then OSC controls the
158 robot to reach the waypoints along the trajectory orderly. The Tiago robot’s navigation skill is
159 implemented using the ROS MoveBase package, while all other skills are implemented using MoveIt
160 motion planning framework [55]. A complete list of skills for both robots is in Appendix 3. Later,
161 we will show that humans and robots can work together using these skills to solve all the tasks.

162 **3.3 Leveraging robot learning for efficient BRI**

163 The modular decoding pipeline and the primitive skill library lay the foundation for NOIR. How-
164 ever, the efficiency of such a system can be further improved. During the collaboration, the robots
165 should learn the user’s object, skill, and parameter selection preferences, hence in future trials, the
166 robot can predict users’ intended goals and be more autonomous, hence reducing the effort required
167 for decoding. Learning and generalization are required since the location, pose, arrangement, and
168 instance of the objects could differ from trial to trial. Meanwhile, the learning algorithms should be
169 sample-efficient since human data is expensive to collect.

170 **Retrieval-based few-shot object and skill selection.** In NOIR, human effort can be reduced if the
171 robot intelligently learns to propose appropriate object-skill selections for a given state in the task.
172 Inspired by retrieval-based imitation learning [56–58], our proposed method learns a latent state
173 representation from observed states. Given a new state observation, it finds the most similar state in
174 the latent space and the corresponding action. Our method is shown in Fig. 4. During task execu-
175 tion, we record data points that consist of images and the object-skill pairs selected by the human.
176 The images are first encoded by a pre-trained R3M model [59] to extract useful features for robot
177 manipulation tasks, and are then passed through several trainable, fully-connected layers. These
178 layers are trained using contrastive learning with a triplet loss[60] that encourages the images with
179 the same object-skill label to be embedded closer in the latent space. The learned image embeddings
180 and object-skill labels are stored in the memory. During test time, the model retrieves the nearest
181 data point in the latent space and suggests the object-action pair associated with that data point to
182 the human. Details of the algorithm can be found in Appendix 7.1.

183 **One-shot skill parameter learning.** Parameter selection requires a lot of human effort as it needs
184 precise cursor manipulation through MI. To reduce human effort, we propose a learning algorithm
185 for predicting parameters given an object-skill pair as an initial point for cursor control. Assuming
186 that the user has once successfully pinpointed the precise key point to pick a mug’s handle, does
187 this parameter need to be specified again in the future? Recent advancement in foundation models
188 such as DINOv2 [61] allows us to find corresponding semantic key points, eliminating the need

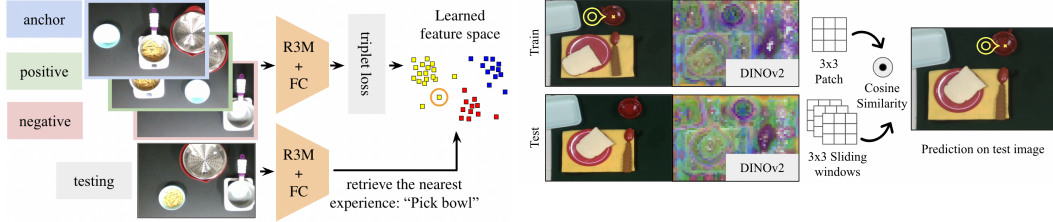


Figure 4: Left: Retrieval-based few-shot object and skill selection model. The model learns a latent representation for observations. Given a new observation, it finds the most relevant experience in the memory and selects the corresponding skill and object. Right: One-shot skill parameter learning algorithm, which finds a semantically corresponding point in the test image given a reference point in the training image. The feature visualization shows 3 of the 768 DINOv2 tokens used.

189 for parameter re-specification. Compared to previous works, our algorithm is one-shot [62–66] and
 190 predicts specific 2D points instead of semantic segments [67, 68]. As shown in Fig. 4, given a
 191 training image (360×240) and parameter choice (x, y) , we predict the semantically corresponding
 192 point in the test images, in which positions, orientations, instances of the target object, and contexts
 193 may vary. We utilize a pre-trained DINOv2 model to obtain semantic features [61]. We input both
 194 train and test images into the model and generate 768 patch tokens, each as a pixel-wise feature
 195 map of dimension 75×100 . We then extract a 3×3 patch centered around the provided training
 196 parameter and search for a matching feature in the test image, using cosine similarity as the distance
 197 metric. Details of this algorithm can be found in Appendix 7.2.

198 4 Experiments

199 **Tasks.** NOIR can greatly benefit those who require assistance with everyday activities. We select
 200 tasks from the BEHAVIOR benchmark [69] and Activities of Daily Living [70] to capture actual
 201 human needs. The tasks are shown in Fig. 1, and consist of 16 tabletop tasks and four mobile
 202 manipulation tasks. The tasks encompass various categories, including eight meal preparation tasks,
 203 six cleaning tasks, three personal care tasks, and three entertainment tasks. For systematic evaluation
 204 of task success, we provide formal definitions of these activities in the BDDL language format [69,
 205 71], which specifies the initial and goal conditions of a task using first-order logic. Task definitions
 206 and figures can be found in Appendix 4.

207 **Procedure.** The human study conducted has received approval from Institutional Review Board.
 208 Three healthy human participants (2 male, 1 female) performed all 15 Franka tasks. Suki-yaki, four
 209 Tiago tasks, and learning tasks are performed by one user. We use the EGI NetStation EEG system,
 210 which is completely non-invasive, making almost everyone an ideal subject. Before experiments,
 211 users are familiarized with task definitions and system interfaces. During the experiment, users stay
 212 in an isolated room, remain stationary, watch the robot on a screen, and solely rely on brain signals
 213 to communicate with the robots (more details about the procedure can be found in Appendix 5).

214 5 Results

215 We seek to provide answers to the following questions through extensive evaluation: 1) Is NOIR
 216 truly general-purpose, in that it allows all of our human subjects to accomplish the diverse set of
 217 everyday tasks we have proposed? 2) Does our decoding pipeline provide accurate decoding results?
 218 3) Does our proposed robot learning and intention prediction algorithm improve NOIR’s efficiency?

219 **System performance.** Table 1 summarizes the performance based on two metrics: the number of
 220 attempts until success and the time to complete the task in successful trials. When the participant
 221 reached an unrecoverable state in task execution, we reset the environment and the participant re-
 222 attempted the task from the beginning. Task horizons (number of primitive skills executed) are
 223 included as a reference. Although these tasks are long-horizon and challenging, NOIR shows very
 224 encouraging results: on average, tasks can be completed with only 1.83 attempts. The reason for
 225 task failures is human errors in skill and parameter selection, i.e., the users pick the wrong skills or
 226 parameters, which leads to non-recoverable states and needs manual resets. Decoding errors or robot

Task	WipeSpill	CollectToy	SweepTrash	CleanBook	IronCloth	OpenBasket	PourTea	SetTable	GrateCheese	CutBanana
Task horizon	4.33	7.67	5.67	7.00	4.67	5.33	4.00	8.33	7.00	5.33
# Attempts	1.00	1.33	2.33	3.33	2.33	1.67	1.67	5.67	1.33	1.67
Time (min)	14.74	25.24	20.59	27.73	16.95	15.90	13.53	20.91	24.98	17.68
Human time (%)	79.02	83.97	82.34	80.00	79.56	82.03	83.15	81.15	81.79	81.21
Task	CookPasta	Sandwich	Hockey	OpenGift	TicTacToe	Sukiyaki	TrashDisposal	CovidCare	WaterPlant	PetDog
Task horizon	8.33	9.00	5.00	7.00	14.33	13.00	8.00	8.00	4.00	6.00
# Attempts	1.67	1.67	1.33	2.67	2.00	1.00	1.00	1.00	1.00	1.00
Time (min)	30.06	27.87	15.83	23.57	43.08	43.45	7.25	8.80	3.00	4.58
Human time (%)	83.26	82.71	82.00	79.90	80.54	84.85	55.32	62.29	87.41	87.53

Table 1: NOIR system performance. Task horizon is the average number of primitive skills executed. # attempts indicate the average number of attempts until the first success (1 means success on the first attempt). Time indicates the task completion time in successful trials. Human time is the percentage of the total time spent by human users, this includes decision-making time and decoding time. With only a few attempts, all users can accomplish these challenging tasks.

227 execution errors are avoided thanks to our safety mechanism with confirmation and interruption.
 228 Although our primitive skill library is limited, human users find novel usage of these skills to solve
 229 tasks creatively. Hence we observe emerging capabilities such as extrinsic dexterity. For example,
 230 in task CleanBook (Fig. 1.6), Franka’s Pick skill is not designed to grasp a book from the table, but
 231 users learn to push the book towards the edge of the table and grasp it from the side. In CutBanana
 232 (Fig. 1.12), users utilize Push skill to cut. The average task completion time is 20.29 minutes.
 233 Note that the time humans spent on decision-making and decoding is relatively long (80% of total
 234 time), partially due to the safety mechanism. Later, we will show that our proposed robot learning
 235 algorithms can address this issue effectively.

236 **Decoding accuracy.** A key to our system’s success is the accuracy in decoding brain signals. Ta-
 237 ble 2 summarizes the decoding accuracy of different stages. We find that CCA on SSVEP produces a
 238 high accuracy of 81.2%, meaning that object selection is mostly accurate. As for CSP + QDA on MI
 239 for parameter selection, the 2-way classification model performs at 73.9% accuracy, which is con-
 240 sistent with current literature [36]. The 4-way skill-selection classification models perform at about
 241 42.2% accuracy. Though this may not seem high, it is competitive considering inconsistencies at-
 242 tributed to long task duration (hence the discrepancy between calibration and task-time accuracies).
 243 Our calibration time is only 10 minutes, which is significantly shorter compared to the duration of
 244 typical MI calibration and training sessions by several orders of magnitude [21]. More calibration
 245 provides more data for training more robust classifiers, and allows human users to practice more
 246 which typically yields stronger brain signals. Overall, the decoding accuracy is satisfactory, and
 247 with the safety mechanism, there has been no instance of task failure caused by incorrect decoding.

248 **Object and skill selection results.** We then answer the third question: Does our proposed robot
 249 learning algorithm improve NOIR’s efficiency? First, we evaluate object and skill selection learn-
 250 ing. We collect a dataset offline with 15 training samples for each object-skill pair in MakePasta
 251 task. Given an image, a prediction is considered correct if both the object and the skill are pre-
 252 dicted correctly. Results are shown in Table 3. While a simple image classification model using
 253 ResNet [72] achieves an average accuracy of 0.31, our method with a pre-trained ResNet backbone
 254 achieves significantly higher accuracy at 0.73, highlighting the importance of contrastive learning
 255 and retrieval-based learning. Using R3M as the feature extractor further improves the performance
 256 to 0.94. The generalization ability of the algorithm is tested on the same MakePasta task. For
 257 instance-level generalization, 20 different types of pasta are used; for context generalization, we
 258 randomly select and place 20 task-irrelevant objects in the background. Results are shown in Table
 259 3. In all variations, our model achieves accuracy over 93%, meaning that the human can skip the
 260 skill and object selection 93% of the time, significantly reducing their time and effort. We further
 261 test our algorithm during actual task execution (Fig. 5). A human user completes the task with and
 262 without object-skill prediction two times each. With object and skill learning, the average time re-
 263 quired for each object-skill selection is reduced by 60% from 45.7 to 18.1 seconds. More details
 264 about the experiments and visualization of learned representation can be found in Appendix 7.1.

265 **One-shot parameter learning results.** First, using our pre-collected dataset (see Appendix 7.2),
 266 we compare our algorithm against multiple baselines. The MSE values of the predictions are shown
 267 in Table 4. *Random sample* shows the average error when randomly predicting points in the 2D
 268 space. *Sample on objects* randomly predicts a point on objects and not on the background; the ob-

Decoding Stage	Signal	Technique	Calibration Acc.	Task-Time Acc.
Object selection (What?)	SSVEP	CCA (4-way)	-	0.812
Skill selection (How?)	MI	CSP + QDA (4-way)	0.580	0.422
Parameter selection (Where?)	MI	CSP + QDA (2-way)	0.882	0.739
Confirmation / interruption	EMG	Thresholding (2-way)	1.0	1.0

Table 2: Decoding accuracy at different stages of the experiment.

269 ject masks here are detected with the Segment Anything Model (SAM) [73]. For *Pixel similarity*,
 270 we employ the cosine similarity and sliding window techniques used in our algorithm, but on raw
 271 images without using DINOv2 features. All of the baselines are drastically outperformed by our
 272 algorithm. Second, our one-shot learning method demonstrates robust generalization capability, as
 273 tested on the respective dataset; table 4 presents the results. The low prediction error means that
 274 users spend much less effort in controlling the cursor to move to the desired position. We fur-
 275 ther demonstrate the effectiveness of the parameter learning algorithm in actual task execution for
 276 SetTable, quantified in terms of saved human effort in controlling the cursor movement (Fig. 5).
 277 Without learning, the cursor starts at the chosen object or the center of the screen. The predicted
 278 result is used as the starting location for cursor control which led to a considerable decrease in cursor
 279 movement, with the mean distance reduced by 41%. These findings highlight the potential of pa-
 280 rameter learning in improving efficiency and reducing human effort. More results and visualizations
 281 can be found in Appendix 7.2.

Method	Acc.↑	Generalization	Acc.↑
Random	0.12±0.02	Position	0.95±0.04
Classification (ResNet)	0.31±0.11	Pose	0.94±0.04
Ours (ResNet)	0.73±0.09	Instance	0.93±0.02
Ours (R3M)	0.94±0.04	Context	0.98±0.02

Table 3: Object-skill learning results. Our method is highly accurate and robust.

Method	MSE↓	Generalization	MSE↓
Random sample	175.8±29.7	Position	5.6±6.0
Sample on objects	137.2±55.7	Orientation	12.0±11.7
Pixel similarity	45.9±50.1	Instance	16.4±22.2
Ours	15.8±23.8	Context	26.8±62.5

Table 4: One-shot parameter learning results. Our method is highly accurate and generalizes well.

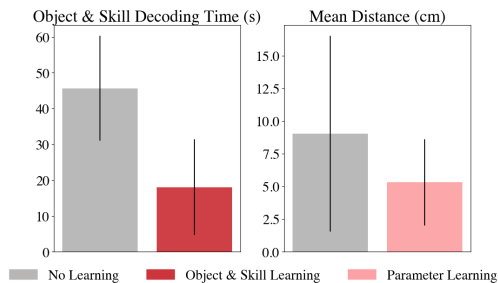


Figure 5: Left: Object and skill selection learning reduces the decoding time by 60%. Right: Parameter learning decreases cursor movement distance by 41%.

284 6 Conclusion, Limitations, and Ethical Concerns

285 In this work, we presented a general-purpose, intelligent BRI system that allows human users to
 286 control a robot to accomplish a diverse, challenging set of real-world activities using brain signals.
 287 NOIR enables human intention prediction through few-shot learning, thereby facilitating a more
 288 efficient collaborative interaction. NOIR holds a significant potential to augment human capabilities
 289 and enable critical assistive technology for individuals who require everyday support.

290 NOIR represents a pioneering effort in the field, unveiling potential opportunities while simultane-
 291 ously raising questions about its limitations and potential ethical risks which we address in Appendix
 292 1. The decoding speed, as it currently stands, restricts tasks to those devoid of time-sensitive inter-
 293 actions. However, advancements in the field of neural signal decoding hold promise for alleviating
 294 this concern. Furthermore, the compilation of a comprehensive library of primitive skills presents a
 295 long-term challenge in robotics, necessitating additional exploration and development. Nonetheless,
 296 we maintain that once a robust set of skills is successfully established, human users will indeed be
 297 capable of applying these existing skills to complete new tasks.

298 References

299 [1] A. Hussein, M. M. Gaber, E. Elyan, and C. Jayne. Imitation learning: A survey of learning
 300 methods. *ACM Computing Surveys (CSUR)*, 50(2):1–35, 2017.

- 301 [2] R. Zhang, D. Bansal, Y. Hao, A. Hiranaka, J. Gao, C. Wang, R. Martín-Martín, L. Fei-Fei,
302 and J. Wu. A dual representation framework for robot learning with human guidance. In
303 *Conference on Robot Learning*, pages 738–750. PMLR, 2023.
- 304 [3] L. Guan, M. Verma, S. S. Guo, R. Zhang, and S. Kambhampati. Widening the pipeline in
305 human-guided reinforcement learning with explanation and context-aware data augmentation.
306 *Advances in Neural Information Processing Systems*, 34:21885–21897, 2021.
- 307 [4] H. Admoni and B. Scassellati. Social eye gaze in human-robot interaction: a review. *Journal*
308 *of Human-Robot Interaction*, 6(1):25–63, 2017.
- 309 [5] A. Saran, R. Zhang, E. S. Short, and S. Niekum. Efficiently guiding imitation learning agents
310 with human gaze. In *Proceedings of the 20th International Conference on Autonomous Agents*
311 *and MultiAgent Systems*, pages 1109–1117, 2021.
- 312 [6] R. Zhang, Z. Liu, L. Zhang, J. A. Whritner, K. S. Muller, M. M. Hayhoe, and D. H. Ballard.
313 Agil: Learning attention from human for visuomotor tasks. In *Proceedings of the european*
314 *conference on computer vision (eccv)*, pages 663–679, 2018.
- 315 [7] R. Zhang, A. Saran, B. Liu, Y. Zhu, S. Guo, S. Niekum, D. Ballard, and M. Hayhoe. Human
316 gaze assisted artificial intelligence: A review. In *IJCAI: Proceedings of the Conference*, volume
317 2020, page 4951. NIH Public Access, 2020.
- 318 [8] Y. Cui, Q. Zhang, B. Knox, A. Allievi, P. Stone, and S. Niekum. The empathic framework
319 for task learning from implicit human feedback. In *Conference on Robot Learning*, pages
320 604–626. PMLR, 2021.
- 321 [9] A. Brohan, Y. Chebotar, C. Finn, K. Hausman, A. Herzog, D. Ho, J. Ibarz, A. Irpan, E. Jang,
322 R. Julian, et al. Do as i can, not as i say: Grounding language in robotic affordances. In
323 *Conference on Robot Learning*, pages 287–318. PMLR, 2023.
- 324 [10] W. Huang, C. Wang, R. Zhang, Y. Li, J. Wu, and L. Fei-Fei. Voxposer: Composable 3d value
325 maps for robotic manipulation with language models. In *7th Annual Conference on Robot*
326 *Learning*, 2023.
- 327 [11] R. Zhang, F. Torabi, L. Guan, D. H. Ballard, and P. Stone. Leveraging human guidance for
328 deep reinforcement learning tasks. *arXiv preprint arXiv:1909.09906*, 2019.
- 329 [12] R. Zhang, F. Torabi, G. Warnell, and P. Stone. Recent advances in leveraging human guidance
330 for sequential decision-making tasks. *Autonomous Agents and Multi-Agent Systems*, 35(2):31,
331 2021.
- 332 [13] I. Akinola, Z. Wang, J. Shi, X. He, P. Lapborisuth, J. Xu, D. Watkins-Valls, P. Sajda, and
333 P. Allen. Accelerated robot learning via human brain signals. In *2020 IEEE international*
334 *conference on robotics and automation (ICRA)*, pages 3799–3805. IEEE, 2020.
- 335 [14] Z. Wang, J. Shi, I. Akinola, and P. Allen. Maximizing bci human feedback using active learn-
336 ing. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*,
337 pages 10945–10951. IEEE, 2020.
- 338 [15] I. Akinola, B. Chen, J. Koss, A. Patankar, J. Varley, and P. Allen. Task level hierarchical
339 system for bci-enabled shared autonomy. In *2017 IEEE-RAS 17th International Conference*
340 *on Humanoid Robotics (Humanoids)*, pages 219–225. IEEE, 2017.
- 341 [16] A. F. Salazar-Gomez, J. DelPreto, S. Gil, F. H. Guenther, and D. Rus. Correcting robot mis-
342 takes in real time using eeg signals. In *2017 IEEE international conference on robotics and*
343 *automation (ICRA)*, pages 6570–6577. IEEE, 2017.

- 344 [17] L. Schiatti, J. Tessadori, N. Deshpande, G. Barresi, L. C. King, and L. S. Mattos. Human in the
345 loop of robot learning: Eeg-based reward signal for target identification and reaching task. In
346 *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4473–4480.
347 IEEE, 2018.
- 348 [18] M. Aljalal, R. Djemal, and S. Ibrahim. Robot navigation using a brain computer interface
349 based on motor imagery. *Journal of Medical and Biological Engineering*, 39:508–522, 2019.
- 350 [19] Y. Xu, C. Ding, X. Shu, K. Gui, Y. Bezsudnova, X. Sheng, and D. Zhang. Shared control
351 of a robotic arm using non-invasive brain–computer interface and computer vision guidance.
352 *Robotics and Autonomous Systems*, 115:121–129, 2019.
- 353 [20] X. Chen, B. Zhao, Y. Wang, S. Xu, and X. Gao. Control of a 7-dof robotic arm system with an
354 ssvpep-based bci. *International journal of neural systems*, 28(08):1850018, 2018.
- 355 [21] J. Meng, S. Zhang, A. Bekyo, J. Olsoe, B. Baxter, and B. He. Noninvasive electroencephalo-
356 gram based control of a robotic arm for reach and grasp tasks. *Scientific Reports*, 6(1):38565,
357 2016.
- 358 [22] M. Aljalal, S. Ibrahim, R. Djemal, and W. Ko. Comprehensive review on brain-controlled
359 mobile robots and robotic arms based on electroencephalography signals. *Intelligent Service
360 Robotics*, 13:539–563, 2020.
- 361 [23] L. F. Nicolas-Alonso and J. Gomez-Gil. Brain computer interfaces, a review. *sensors*, 12(2):
362 1211–1279, 2012.
- 363 [24] L. Bi, X.-A. Fan, and Y. Liu. Eeg-based brain-controlled mobile robots: a survey. *IEEE
364 transactions on human-machine systems*, 43(2):161–176, 2013.
- 365 [25] N. M. Krishnan, M. Mariappan, K. Muthukaruppan, M. H. A. Hijazi, and W. W. Kitt. Elec-
366 troencephalography (eeg) based control in assistive mobile robots: A review. In *IOP Confer-
367 ence Series: Materials Science and Engineering*, volume 121, page 012017. IOP Publishing,
368 2016.
- 369 [26] E. D. Adrian and B. H. Matthews. The berger rhythm: potential changes from the occipital
370 lobes in man. *Brain*, 57(4):355–385, 1934.
- 371 [27] C. J. Perera, I. Naotunna, C. Sadaruwan, R. A. R. C. Gopura, and T. D. Lalitharatne. Ssvpep
372 based bmi for a meal assistance robot. In *2016 IEEE International Conference on Systems,
373 Man, and Cybernetics (SMC)*, pages 002295–002300. IEEE, 2016.
- 374 [28] J. Ha, S. Park, C.-H. Im, and L. Kim. A hybrid brain–computer interface for real-life meal-
375 assist robot control. *Sensors*, 21(13):4578, 2021.
- 376 [29] J. Zhang and M. Wang. A survey on robots controlled by motor imagery brain-computer
377 interfaces. *Cognitive Robotics*, 1:12–24, 2021.
- 378 [30] X. Mao, M. Li, W. Li, L. Niu, B. Xian, M. Zeng, and G. Chen. Progress in eeg-based brain
379 robot interaction systems. *Computational intelligence and neuroscience*, 2017, 2017.
- 380 [31] J. Tang, A. LeBel, S. Jain, and A. G. Huth. Semantic reconstruction of continuous language
381 from non-invasive brain recordings. *Nature Neuroscience*, pages 1–9, 2023.
- 382 [32] M. Minderer, A. Gritsenko, A. Stone, M. Neumann, D. Weissenborn, A. Dosovitskiy, A. Ma-
383 hendran, A. Arnab, M. Dehghani, Z. Shen, X. Wang, X. Zhai, T. Kipf, and N. Houlsby. Simple
384 open-vocabulary object detection with vision transformers, 2022.
- 385 [33] D. Zhu, J. Bieger, G. G. Molina, and R. M. Aarts. A survey of stimulation methods used in
386 ssvpep-based bcis. *Computational intelligence and neuroscience*, 2010:1–12, 2010.

- 387 [34] R. Kuś, A. Duszyk, P. Milanowski, M. Łabecki, M. Bierzyńska, Z. Radzikowska, M. Michal-
388 ska, J. Żygierewicz, P. Suffczyński, and P. J. Durka. On the quantification of ssvep frequency
389 responses in human eeg in realistic bci conditions. *PLoS one*, 8(10):e77536, 2013.
- 390 [35] L. Shao, L. Zhang, A. N. Belkacem, Y. Zhang, X. Chen, J. Li, and H. Liu. Eeg-controlled
391 wall-crawling cleaning robot using ssvep-based brain-computer interface, 2020.
- 392 [36] N. Padfield, J. Zabalza, H. Zhao, V. Masero, and J. Ren. Eeg-based brain-computer interfaces
393 using motor-imagery: Techniques and challenges. *Sensors*, 19(6):1423, 2019.
- 394 [37] S.-L. Wu, C.-W. Wu, N. R. Pal, C.-Y. Chen, S.-A. Chen, and C.-T. Lin. Common spatial pattern
395 and linear discriminant analysis for motor imagery classification. In *2013 IEEE Symposium on*
396 *Computational Intelligence, Cognitive Algorithms, Mind, and Brain (CCMB)*, pages 146–151.
397 IEEE, 2013.
- 398 [38] S. Kumar, A. Sharma, and T. Tsunoda. An improved discriminative filter bank selection ap-
399 proach for motor imagery eeg signal classification using mutual information. *BMC bioinforma-*
400 *tics*, 18:125–137, 2017.
- 401 [39] Y. Zhang, Y. Wang, J. Jin, and X. Wang. Sparse bayesian learning for obtaining sparsity of eeg
402 frequency bands based feature vectors in motor imagery classification. *International journal*
403 *of neural systems*, 27(02):1650032, 2017.
- 404 [40] B. Yang, H. Li, Q. Wang, and Y. Zhang. Subject-based feature extraction by using fisher
405 wpd-csp in brain-computer interfaces. *Computer methods and programs in biomedicine*, 129:
406 21–28, 2016.
- 407 [41] R. Chitnis, T. Silver, J. B. Tenenbaum, T. Lozano-Perez, and L. P. Kaelbling. Learning neuro-
408 symbolic relational transition models for bilevel planning. In *2022 IEEE/RSJ International*
409 *Conference on Intelligent Robots and Systems (IROS)*, pages 4166–4173. IEEE, 2022.
- 410 [42] S. Nasiriany, H. Liu, and Y. Zhu. Augmenting reinforcement learning with behavior primitives
411 for diverse manipulation tasks. In *2022 International Conference on Robotics and Automation*
412 *(ICRA)*, pages 7477–7484. IEEE, 2022.
- 413 [43] Y. Zhu, J. Tremblay, S. Birchfield, and Y. Zhu. Hierarchical planning for long-horizon manip-
414 ulation with geometric and symbolic scene graphs. In *2021 IEEE International Conference on*
415 *Robotics and Automation (ICRA)*, pages 6541–6548. IEEE, 2021.
- 416 [44] M. Shridhar, L. Manuelli, and D. Fox. Cliport: What and where pathways for robotic manipu-
417 lation. In *Conference on Robot Learning*, pages 894–906. PMLR, 2022.
- 418 [45] M. Shridhar, L. Manuelli, and D. Fox. Perceiver-actor: A multi-task transformer for robotic
419 manipulation. *arXiv preprint arXiv:2209.05451*, 2022.
- 420 [46] W. Liu, C. Paxton, T. Hermans, and D. Fox. Structformer: Learning spatial structure for
421 language-guided semantic rearrangement of novel objects. In *2022 International Conference*
422 *on Robotics and Automation (ICRA)*, pages 6322–6329. IEEE, 2022.
- 423 [47] D. Xu, A. Mandlekar, R. Martín-Martín, Y. Zhu, S. Savarese, and L. Fei-Fei. Deep afford-
424 ance foresight: Planning through what can be done in the future. In *2021 IEEE International*
425 *Conference on Robotics and Automation (ICRA)*, pages 6206–6213. IEEE, 2021.
- 426 [48] C. Wang, D. Xu, and L. Fei-Fei. Generalizable task planning through representation pretrain-
427 ing. *IEEE Robotics and Automation Letters*, 7(3):8299–8306, 2022.
- 428 [49] S. Cheng and D. Xu. Guided skill learning and abstraction for long-horizon manipulation.
429 *arXiv preprint arXiv:2210.12631*, 2022.

- 430 [50] C. Agia, T. Migimatsu, J. Wu, and J. Bohg. Taps: Task-agnostic policy sequencing. *arXiv preprint arXiv:2210.12250*, 2022.
431
- 432 [51] C. Li, R. Zhang, J. Wong, C. Gokmen, S. Srivastava, R. Martín-Martín, C. Wang, G. Levine,
433 M. Lingelbach, J. Sun, et al. Behavior-1k: A benchmark for embodied ai with 1,000 everyday
434 activities and realistic simulation. In *6th Annual Conference on Robot Learning*, 2022.
- 435 [52] A. Hiranaka, M. Hwang, S. Lee, C. Wang, L. Fei-Fei, J. Wu, and R. Zhang. Primitive skill-
436 based robot learning from human evaluative feedback. *arXiv preprint arXiv:2307.15801*, 2023.
- 437 [53] O. Khatib. A unified approach for motion and force control of robot manipulators: The opera-
438 tional space formulation. *IEEE Journal on Robotics and Automation*, 3(1):43–53, 1987.
- 439 [54] Y. Zhu, A. Joshi, P. Stone, and Y. Zhu. Viola: Object-centric imitation learning for vision-based
440 robot manipulation. In *Conference on Robot Learning*, pages 1199–1210. PMLR, 2023.
- 441 [55] D. Coleman, I. Sucas, S. Chitta, and N. Correll. Reducing the barrier to entry of complex
442 robotic software: a moveit! case study, 2014.
- 443 [56] E. Mansimov and K. Cho. Simple nearest neighbor policy method for continuous control tasks,
444 2018. URL <https://openreview.net/forum?id=ByL48G-AW>.
- 445 [57] S. Nasiriany, T. Gao, A. Mandlekar, and Y. Zhu. Learning and retrieval from prior data for
446 skill-based imitation learning. In *6th Annual Conference on Robot Learning*.
- 447 [58] M. Du, S. Nair, D. Sadigh, and C. Finn. Behavior retrieval: Few-shot imitation learning by
448 querying unlabeled datasets, 2023.
- 449 [59] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta. R3m: A universal visual represen-
450 tation for robot manipulation. In *6th Annual Conference on Robot Learning*, 2021.
- 451 [60] K. Q. Weinberger and L. K. Saul. Distance metric learning for large margin nearest neighbor
452 classification. *Journal of machine learning research*, 10(2), 2009.
- 453 [61] M. Oquab, T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haz-
454 iza, F. Massa, A. El-Nouby, M. Assran, N. Ballas, W. Galuba, R. Howes, P.-Y. Huang, S.-W. Li,
455 I. Misra, M. Rabbat, V. Sharma, G. Synnaeve, H. Xu, H. Jegou, J. Mairal, P. Labatut, A. Joulin,
456 and P. Bojanowski. DINOv2: Learning robust visual features without supervision, 2023.
- 457 [62] T. Luddecke and F. Worgotter. Learning to segment affordances. In *2017 IEEE International
458 Conference on Computer Vision Workshops (ICCVW)*, pages 769–776. IEEE, 2017.
- 459 [63] A. Nguyen, D. Kanoulas, D. G. Caldwell, and N. G. Tsagarakis. Object-based affordances
460 detection with convolutional neural networks and dense conditional random fields. In *2017
461 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5908–
462 5915. IEEE, 2017.
- 463 [64] A. Nguyen, D. Kanoulas, D. G. Caldwell, and N. G. Tsagarakis. Translating videos to com-
464 mands for robotic manipulation with deep recurrent neural networks. In *IEEE International
465 Conference on Robotics and Automation (ICRA)*, pages 3782–3788. IEEE, 2018.
- 466 [65] T. Mar, V. Tikhonoff, G. Metta, and L. Natale. Self-supervised learning of grasp dependent
467 tool affordances on the icub humanoid robot. In *IEEE International Conference on Robotics
468 and Automation (ICRA)*, pages 3200–3206. IEEE, 2015.
- 469 [66] T. Mar, V. Tikhonoff, G. Metta, and L. Natale. Self-supervised learning of tool affordances
470 from 3d tool representation through parallel som mapping. In *IEEE International Conference
471 on Robotics and Automation (ICRA)*, pages 894–901. IEEE, 2017.

- 472 [67] H. Luo, W. Zhai, J. Zhang, Y. Cao, and D. Tao. One-shot affordance detection. In *Proceedings*
473 *of the 30th International Joint Conference on Artificial Intelligence (IJCAI)*, 2021.
- 474 [68] D. Hadjivelichkov, S. Zwane, M. P. Deisenroth, L. Agapito, and D. Kanoulas. One-shot transfer
475 of affordance regions? affcorr! In *6th Conference on Robot Learning (CoRL)*, 2022.
- 476 [69] C. Li, R. Zhang, J. Wong, C. Gokmen, S. Srivastava, R. Martín-Martín, C. Wang, G. Levine,
477 M. Lingelbach, J. Sun, et al. Behavior-1k: A benchmark for embodied ai with 1,000 everyday
478 activities and realistic simulation. In *Conference on Robot Learning*, pages 80–93. PMLR,
479 2023.
- 480 [70] S. Katz. Assessing self-maintenance: activities of daily living, mobility, and instrumental
481 activities of daily living. *Journal of the American Geriatrics Society*, 1983.
- 482 [71] S. Srivastava, C. Li, M. Lingelbach, R. Martín-Martín, F. Xia, K. E. Vainio, Z. Lian, C. Gok-
483 men, S. Buch, K. Liu, et al. Behavior: Benchmark for everyday household activities in virtual,
484 interactive, and ecological environments. In *Conference on Robot Learning*, pages 477–490.
485 PMLR, 2022.
- 486 [72] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Pro-*
487 *ceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778,
488 2016.
- 489 [73] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead,
490 A. C. Berg, W.-Y. Lo, P. Dollár, and R. Girshick. Segment anything. *arXiv:2304.02643*, 2023.