

Deep Learning Based BCI Control of a Robotic Service Assistant Using Intelligent Goal Formulation

D. Kuhner^{a,1,*}, L.D.J. Fiederer^{b,c,1,*}, J. Aldinger^{a,1,*}, F. Burget^{a,*}, M. Völker^{a,b,*}, R.T. Schirrmeister^b, C. Do^a, J. Bödecker^a, B. Nebel^a, T. Ball^b, W. Burgard^a

^aDepartment of Computer Science, University of Freiburg, Germany

^bFaculty of Medicine, University of Freiburg, Germany

^cFaculty of Biology, University of Freiburg, Germany

Abstract

As autonomous service robots become more affordable and thus available for the general public, there is a growing need for user-friendly interfaces to control these systems. Control interfaces typically get more complicated with increasing complexity of the robotic tasks and the environment. Traditional control modalities as touch, speech or gesture commands are not necessarily suited for all users. While non-expert users can make the effort to familiarize themselves with a robotic system, paralyzed users may not be capable of controlling such systems even though they need robotic assistance most. In this paper, we present a novel framework, that allows these users to interact with a robotic service assistant in a closed-loop fashion, using only thoughts. The system is composed of several interacting components: non-invasive neuronal signal recording and co-adaptive deep learning which form the brain-computer interface (BCI), high-level task planning based on referring expressions, navigation and manipulation planning as well as environmental perception. We extensively evaluate the BCI in various tasks, determine the performance of the goal formulation user interface and investigate its intuitiveness in a user study. Furthermore, we demonstrate the applicability and robustness of the system in real world scenarios, considering fetch-and-carry tasks and tasks involving human-robot interaction. As our results show, the system is capable of adapting to frequent changes in the environment and reliably accomplishes given tasks within a reasonable amount of time. Combined with high-level planning using referring expressions and autonomous robotic systems, interesting new perspectives open up for non-invasive BCI-based human-robot interactions.

Keywords: EEG, Co-Adaptive Brain-Computer-Interface, Realtime Deep Learning, Autonomous Robotics, Referring Expression Generation, High-level Task Planning, Computer Vision

1. Highlights

- BCI-controlled autonomous robotic service assistant
- First online brain-computer-interface using deep learning
- Menu-driven language generation based on referring expression
- Modular ROS-based mobile robot interaction
- Experimental evaluation using a real robot

2. Introduction

Persons with impaired communication capabilities, such as severely paralyzed patients, rely on constant help of human care-takers. Robotic service assistants can re-establish some degree of autonomy for these patients, if they offer adequate interfaces and possess a sufficient level of intelligence. Generally, such systems require adaptive task- and motion-planning modules to determine appropriate task plans and motion trajectories for the robot to execute a task in the real world. Moreover, it requires a perception component to detect objects of interest or to avoid accidental collisions with obstacles. With increasing capabilities of autonomous systems intelligent control opportunities also become more important. Typical interfaces, such as haptic (buttons), audio (speech) or visual (gesture) interfaces, are well suited for healthy users. However, for persons with impaired communication skills these control opportunities are unreliable or impossible to use.

In this paper, we present and evaluate a novel framework, schematically depicted in Fig. 1, that allows closed-loop interaction between users with minimal communica-

*Equally contributing

Email addresses: kuhnerd@informatik.uni-freiburg.de (D. Kuhner), lukas.fiederer@uniklinik-freiburg.de (L.D.J. Fiederer), aldinger@informatik.uni-freiburg.de (J. Aldinger), burgetf@informatik.uni-freiburg.de (F. Burget), martin.voelker@uniklinik-freiburg.de (M. Völker), robin.schirrmeister@uniklinik-freiburg.de (R.T. Schirrmeister), do@informatik.uni-freiburg.de (C. Do), jboedeck@informatik.uni-freiburg.de (J. Bödecker), nebel@informatik.uni-freiburg.de (B. Nebel), tonio.ball@uniklinik-freiburg.de (T. Ball), burgard@informatik.uni-freiburg.de (W. Burgard)

¹Corresponding authors

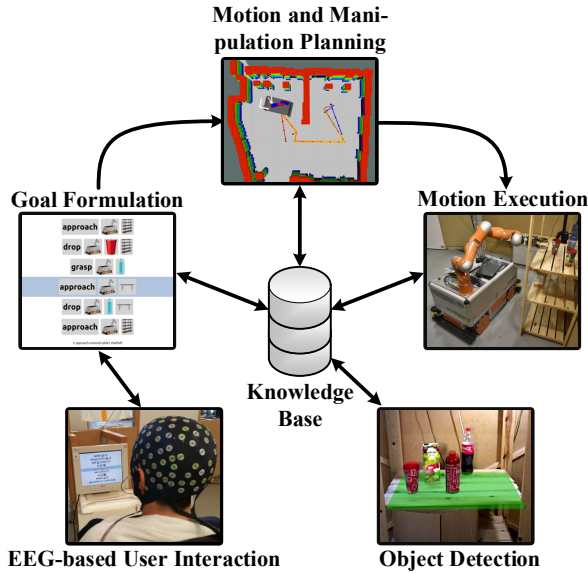


Figure 1: Our framework that unifies decoding of neuronal signals, high-level task planning based on referring expressions, low-level motion- and manipulation-planning, and scene perception with a centralized knowledge base at its core. Intuitive goal selection is provided through an adaptive graphical user interface.

tion capabilities and a robotic service assistant. To do so, we record neuronal activity elicited in the human brain, the common origin of all types of communication, with an electroencephalography (EEG) system. Furthermore, we employ a deep convolutional neural network (ConvNet) approach for online co-adaptive decoding of neuronal activity, in order to allow users to navigate through a graphical user interface (GUI) which is connected to a high-level task planner. It allows the intuitive selection of goals based on the generation of referring expressions that identify the objects to be manipulated. The set of feasible actions displayed in the GUI, depends in turn on the current state of the world, which is stored in a central knowledge base and continuously updated with information provided by the robot and a camera perception system. Once a task has been selected, it is decomposed into a sequence of atomic actions by the high-level planner. Subsequently, each action is resolved to a motion for the mobile manipulator using low-level motion-planning techniques. This approach minimizes the cognitive load required of the user, which is a crucial aspect in the design of a BCI. Furthermore, the intelligence and autonomy of the system make it possible to interface non-invasive BCIs, which currently have low throughput, with our robotic assistant composed of 11 degrees-of-freedom (DOF). In the following, we present the related work, describe the individual components shown in Fig. 1 and present a quantitative evaluation of the system regarding its performance and user-friendliness.

3. Related Work

The multi-disciplinary work presented in this paper relies on robotics, brain-computer interfaces and natural-language generation (NLG). This section outlines related work in these fields.

Robotic Assistants. Multiple previous studies have focused on robotic systems assisting people with disabilities. For example, Park *et al.* [1] implemented a system for autonomously feeding yogurt to a person. Chung *et al.* [2] focus on autonomous drinking which involves locating the drink, picking it up and bringing it to the person's mouth. Using a hybrid BCI and head movement control, Achic *et al.* [3] studies a setup with a moving wheelchair and an attached robotic arm. None of these systems use pure BCI control. In contrast, Wang *et al.* [4] employ a motor imagery BCI with three classes to achieve low-level control of a robotic arm. More relevant, Schröer *et al.* [5] propose a robotic system which receives a BCI command from a user and autonomously assists the user in drinking from a cup. However, this approach only considers a single object and a fixed-base manipulator. Grigorescu *et al.* [6] use steady-state visually evoked potentials to control the commercially available FRIEND III assistance robot. This work is perhaps closest to ours with respect to the number of possible commands (namely 5), the high-level control concept and the (semi-)autonomy of the assistance robot. In contrast to their work, we use active brain signals to control the graphical user interface and apply co-adaptive training and decoding. See the excellent review of Mladenović *et al.* [7] for details on co-adaptive BCIs. Additionally, we propose a specific design of the user interface to improve the human-robot interaction and show that our system is fully autonomous. Most recently, the work of Muelling *et al.* [8] presents a shared-control approach in the field of assistive robotics based on an invasive BCI. This is contrary to our approach which relies on a non-invasive BCI. Nonetheless, their approach could be combined with the goal formulation interface presented in this work.

Brain-Computer Interfaces. To ensure user acceptance, robust decoding of brain signals is required. Inspired by the successes of deep ConvNets in computer vision [9, 10] and speech recognition [11, 12], deep ConvNets have recently been applied more frequently to EEG brain-signal decoding and - related to this paper - to decode tasks in brain-computer interfaces. Lawhern *et al.* [13] use a deep ConvNet to decode P300 oddball signals, feedback error-related negativity and two movement-related tasks. In cross-participant evaluation (i. e., trained on some participants and evaluated on others), their ConvNet yields competitive accuracies compared to widely-used traditional brain-signal decoding algorithms. Tabar and Halici [14] combine a ConvNet and a convolutional stacked auto-encoder to decode motor imagery within-participant and improve accuracies compared to several non-ConvNet decoding algorithms. Schirrmeyer *et al.* [15] use a shallow and a deep

113 ConvNet to decode both motor imagery and motor execu-168
114 tion within-participant. Their approach results in com-169
115 parable or slightly better accuracies than the widely used-170
116 EEG motor-decoding algorithm *filter bank common spa-*171
117 *tial patterns* [16]. Bashivan *et al.* [17] estimate the mental-172
118 workload with a ConvNet trained on fourier-transformed-173
119 inputs. In addition to the above work on evaluating ConvNet-174
120 decoding accuracies, ConvNet visualization methods allow-175
121 us to get a sense of what brain-signal features the net-176
122 work is using [15, 17, 18, 19, 20, 21]. Taken together,177
123 these advances make deep ConvNets a viable alternative178
124 for brain-signal decoding in brain-computer interfaces. A179
125 first attempt at using shallow ConvNets for online BCI180
126 has recently been reported [22]. To the best of our knowl-181
127 edge, apart from our previous paper [23], there is no other182
128 work, which uses a deep ConvNet-based online control to183
129 implement an EEG-based brain-computer interface. 184

185
186 *Referring Expressions.* When humans communicate goals186
187 to other humans, they identify objects in the world by re-187
188 ferring expressions (e. g., *a red cup on the shelf*). The gen-188
189 eration of referring expressions has been subject to compu-189
190 tational linguistics research for years as one part of natural190
191 language generation (NLG) [24]. With recent advances in191
192 natural language processing, computer vision and the rise192
193 of neuronal networks, it is nowadays possible to identify193
194 objects in images by building referring expressions gen-194
195 erated from features [25]. Spatial references can be used195
196 to discriminate similar objects [26]. The NLG problem196
197 has been approached with planning techniques [27]. How-197
198 ever, such systems usually lack knowledge about the ac-198
199 tions that can be executed and the objects that can be ma-199
200 nipulated. To overcome this problem and to improve the200
201 human-robot interaction we propose a user interface that201
202 allows specifying actions in a domain-independent way and202
203 automatically adapts to changes in the environment. 203

204
205 *Task- and Manipulation-Planning.* In contrast to classical204
205 task planning, Task- and Manipulation-Planning (TAMP)205
206 algorithms also consider the motion capabilities of the robot206
207 to determine feasible task plans. There are various ap-207
208 proaches to solve the this problem. Common to most208
209 TAMP approaches is a hierarchical decomposition of the209
210 problem into task- and motion-planning layers. Due to210
211 the high dimensionality of the TAMP problem the decom-211
212 position can be understood as a way to guide the low-212
213 level planners based on the high-level plan solution and213
214 vice versa. For example, Kaelbling *et al.* [28, 29] propose214
215 an aggressively hierarchical planning method. Such a hi-215
216 erarchical decomposition allows handling problems with216
217 long horizons efficiently. De Silva *et al.* [30] show an ap-217
218 proach based on Hierarchical Task Networks (HTNs) to218
219 reason on abstract tasks and combine them with a geo-219
220 metric task planner which works in a discrete space of pre-220
221 computed grasp, drop and object positions. Recently, the221
222 work of Dantam *et al.* [31] introduce the probabilistically-222
223 complete Iteratively Deepened Task- and Motion-Planning223

(IDTMP) algorithm, which uses a constrained-based task
planner to create tentative task plans and sampling-based
motion planners for feasibility tests. Srivastava *et al.* [32]
focus on a planner-independent interface layer between
task- and motion-planners. Lozano-Pérez *et al.* [33] post-
pone the decision on motion plans to avoid expensive back-
tracking due to restrictions which might happen, if the low-
level planner is queried too early. Instead, they generate a
"skeleton" high-level plan and a set of constraints, which
need to be satisfied to achieve the goals of the high-level
planner. Dornhege *et al.* [34] integrate task- and motion-
planning by extending the TFD task planner [35] with
semantic attachments, i. e., modules which check the fea-
sibility of motion plans on demand to ensure that task
plans can be refined to motion plans. In this work, the
goal formulation interface outputs a task plan composed
of high-level actions. We assume that these actions can
be refined to motions of the robot, if the task plan is con-
sistent with the current world model. Thus, task- and
manipulation-planning is also considered in a hierarchical
way but we postpone the decision on the actual feasibil-
ity of motion plans to reduce the computational effort.
Nonetheless, due to the modular structure of our system
most of the principles applied in this field could be inte-
grated into our framework as well.

4. Autonomous BCI-controlled Service Assistant

In this paper, we present an autonomous robotic ser-
vice assistant which uses a BCI and an intuitive goal for-
mulation framework to aid users in fetch-and-carry tasks.
Our system relies on multiple components which are de-
picted in Fig. 2. The communication between user and
robotic service assistant is established using an EEG-based
BCI. It decodes the brain signals using deep ConvNets and
is explained in Sec. 4.1. Sec. 4.2 describes the goal for-
mulation assistant that employs referring expressions and a
menu-driven user interface to allow an intuitive specifica-
tion of tasks. These are then processed by a high-level
task planner to break them into a set of executable sub-
tasks that are sent to the navigation- and manipulation-
planning algorithms. Furthermore, we apply a Monte-
Carlo localization approach to estimate the pose of the
robot in the environment. Based on these poses and the
map of the environment, the navigation module deter-
mines a collision-free trajectory, that allows the robot to
move between different locations. Additionally, roadmap-
based planning allows to execute manipulation tasks as
grasping and dropping objects. More details on motion
generation are available in Sec. 4.3. The goal formulation
interface and the TAMP algorithms depend on a percep-
tion module that dynamically detects relevant objects in
the environment. Autonomous drinking capabilities also
require that the framework is able to determine the user's
mouth location and the robust estimation of liquid level
heights to avoid spilling while serving a drink (Sec. 4.4).

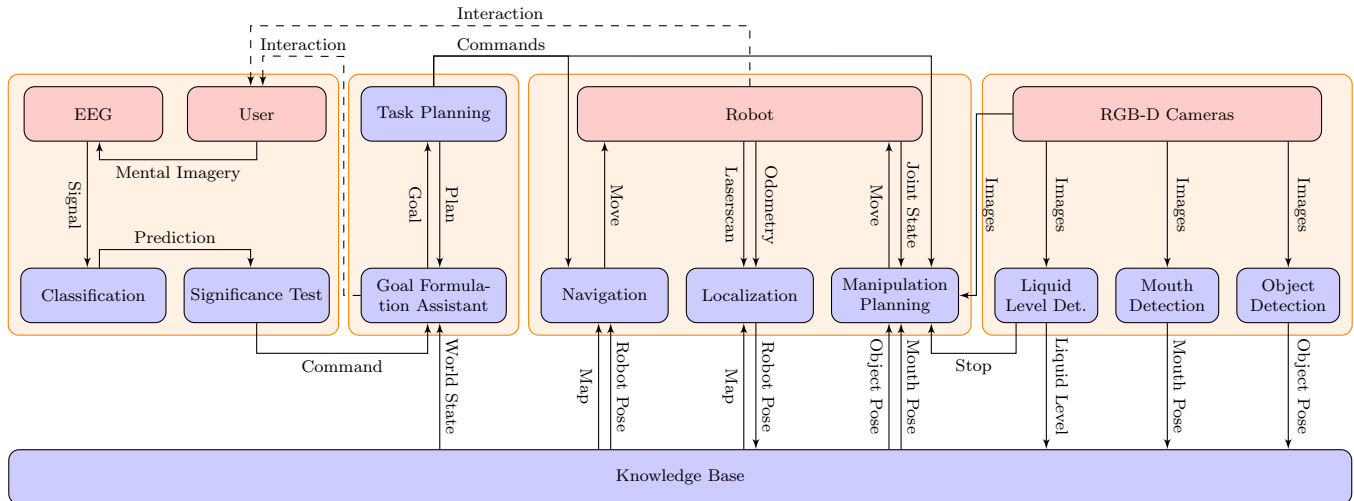


Figure 2: Detailed overview of our framework. It uses a brain-computer interface to decode the thoughts of the user. Thus, the user has control over a goal formulation assistant which is connected to a high-level planner. The commands send by the high-level planner are then processed by the low-level motion planners and executed on the robot. A perception system determines information on object poses, the user’s mouth position and liquid levels. Finally, a central knowledge base stores and provides data to establish a connection between all components.

Finally, to conduct the data communication a knowledge base connects all components (Sec. 4.5).

4.1. Online Decoding of Neuronal Signals

This section introduces the deep ConvNet and the strategies to train the network. Furthermore, we explain the online decoding pipeline to extract meaningful commands from EEG data, which are required to control the robotic assistant.

4.1.1. Deep Hybrid ConvNet Training

As reliable classification of brain signals related to directional commands cannot yet be achieved directly with non-invasive BCIs, we decode multiple surrogate mental tasks from EEG using a deep ConvNet approach [15]. This approach introduces a hybrid network, combining a deep ConvNet with a shallower ConvNet architecture. The deep part consists of four convolution-pooling blocks using exponential linear units (ELU) [36] and max-pooling, whereas the shallow part uses a single convolution-pooling block with squaring non-linearities and mean-pooling. Both parts use a final convolution with ELUs to produce the output features. These features are then concatenated and fed to a final classification layer. All details of the architecture are visualized in Fig. 3.

We train the subject-specific ConvNets on 40 Hz lowpass filtered EEG data to decode five mental tasks: sequential right-hand finger tapping, synchronous movement of all toes, object rotation, word generation and rest. These mental tasks evoke discernible brain patterns and are used as surrogate signals to control the GUI. The mental tasks map to the *select*, *go down*, *go back*, *go up*, and *rest* GUI actions, respectively. *Offline* training is conducted based on a cropped training strategy using shifted time windows,

which we call crops, within the trials as input data [15]. The crop size of ~ 2 s (522 samples @ 250 Hz) is given by the size of the ConvNet’s receptive field. Crops start ~ 1.5 s before trial onset and end with trial offset. This corresponds to the first output being predicted 500 ms after trial onset and the last output being predicted on trial offset. To speed up training, one super-crop consisting of 239 consecutive crops (760 samples) is processed in a single forward pass of the model. This results in 239 outputs for each forward pass. One training batch consists of 60 super-crops. We perform stochastic gradient descent using Adam [37] and a learning rate of 10^{-3} . To optimize the model we minimize the categorical cross entropy loss. For offline evaluation, we retain the last two runs as our test set, which corresponds to 20 min of data. The remaining data is split into training (80%) and validation (20%) sets. We train for 100 epochs in the first phase of the training and select the epoch’s model with the highest validation accuracy. Using the combined training and validation sets, retraining is performed until the validation loss reaches the training loss from the epoch selected in the previous phase. Throughout this paper, we report decoding accuracies on the test set. The initial model used for *online* co-adaptive training is trained on all available *offline* data by following the above mentioned scheme.

We perform *online* co-adaptive training with tenfold reduced learning rate, a super-crop size of 600 samples and a batch size of 45 super-crops. We keep all other parameters identical and train for five batches in all 2 s-breaks. During this time incoming EEG data is accumulated and processed once the newly trained ConvNet is available. Training initiates once ten trials have been accumulated in an experimental session. Only session specific data is used during the training. A session is defined as the time

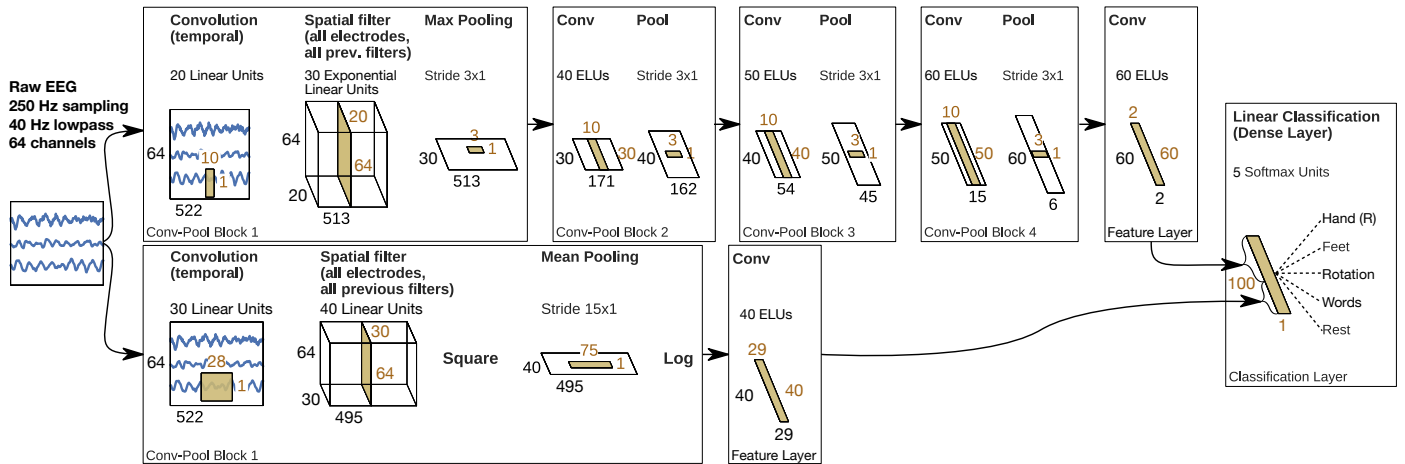


Figure 3: Hybrid deep convolutional neural network: Black numbers depict the dimensions of the input data to each layer. Orange numbers depict the dimensions of the kernels of each layer.

288 interval during which the participants continuously wear 322
 289 the EEG cap. As soon as the EEG cap is removed and 323
 290 reapplied a new session starts. 324

291 4.1.2. Participant Training 326

292 Based on our experience it is important to train the 327
 293 BCI decoder and participants in an environment that is 328
 294 as close as possible to the real application environment to 329
 295 avoid pronounced performance drops when transiting from 330
 296 training to application. Therefore, we designed a gradual 331
 297 training paradigm within the goal-formulation user inter- 332
 298 face (see Sec. 4.2) in which the displayed environment, 333
 299 timing and actions are identical to those of the real con- 334
 300 trol task. The training paradigm proceeds as follows. 335

301 *Offline Training.* We first train each participant *offline* 336
 302 using simulated feedback. Participants are aware of not 337
 303 being in control of the GUI. The mental tasks are cued us- 338
 304 ing modified versions of the BCI2000 [38] grayscale images 339
 305 that are presented for 0.5 s in the center of the display. To 340
 306 minimize eye movements the participants were instructed 341
 307 to look at a fixation circle, permanently displayed in the 342
 308 center of the GUI. After a random time interval of 1-7 s the 343
 309 fixation circle is switched to a disk for 0.2 s, which indicates 344
 310 the end of the mental task. At the same time the GUI ac- 345
 311 tion (*go up*, *go down*, *select*, *go back*, *rest*) corresponding 346
 312 to the cued mental task (cf. Sec. 4.1.1) is performed to 347
 313 update the GUI. The *rest* mental task is implicitly taking 348
 314 place for 2 s after every other task². To allow the partici- 349
 315 pant to blink and swallow, every 4th rest lasts 7 s. Fig. 4
 316 gives a graphical overview of the offline training paradigm. 350
 317 To keep training realistic we include a 20 % error rate, i. e. 351
 318 on average every fifth action is purposefully erroneous. We 352
 319 instruct the participants to count the error occurrences to 353
 320 assert their vigilance. This offline data is used to train the 354
 321 individual deep ConvNets as described in Sec. 4.1.1. 355

²We initially used 1 s intervals to maximize speed, but they were 356
 too short for proper mental task transition.

324 *Online Co-Adaptive Training.* After offline training, the
 325 participants transit to co-adaptive *online* training where
 326 the cued mental tasks are decoded by the ConvNets and
 327 performed in the GUI. The ConvNets were retrained after
 328 each trial during the 2 s break, as described in Sec. 4.1.1.
 329 The participants are conscious of being in control of the
 330 GUI and are instructed to count the errors they make. In
 331 doing so, the participants are aware of their performance,
 332 which potentially triggers learning processes and asserts
 333 their vigilance.

334 *Online Training.* To evaluate the *uncued*, *online* perfor-
 335 mance of the BCI control, we stop cueing the mental tasks
 336 and let the participants select instructed goals in the GUI.
 337 The corresponding task plans are then executed by a simu-
 338 lated robot or – when available – the real mobile manipu-
 339 lator. To provide more control over the mobile manipu-
 340 lator and enhance the feeling of agency, participants have
 341 to confirm the execution of every planned action and can
 342 interrupt the chain of actions at any moment during their
 343 execution using the *go back* GUI action. BCI decoding
 344 accuracies for the label-less instructed tasks are assessed
 345 by manually rating each decoding based on the instructed
 346 task steps. Statistical significance of the decoding accu-
 347 racies are tested using a conventional permutation test with
 348 100 k random permutations of the labels (i. e., the p-value
 349 is the fraction of label permutations that would have led to
 350 better or equal accuracies than the accuracy of the original
 351 labels).

352 4.1.3. Online Decoding Pipeline

353 During *online* control of the GUI, the EEG data is low-
 354 pass filtered at 40 Hz, downsampled to 250 Hz and sent to
 355 a GPU server in blocks of 200 ms for decoding. During co-
 356 adaptive *online* training (cf. Sec. 4.1.2) the data is addi-
 tionally labeled (to identify mental tasks) before being sent
 to the GPU server for decoding, storing and subsequent
 training. On the GPU server 600 samples (one super-crop,

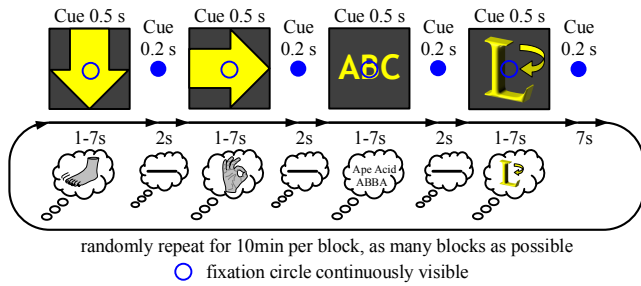


Figure 4: Offline training paradigm. Cue icons, modified from BCI2000 [38], indicate which mental task should be performed by the participant. The cue icons shown here have been modified for better visibility. In our experimental environment we use grayscale versions of the icons. The mental tasks are illustrated by lines in the smaller 'thought bubbles'. Each mental task maps to a GUI action: word generation \rightarrow *go up*, synchronous movement of all toes \rightarrow *go down*, sequential right-hand finger tapping \rightarrow *select*, object rotation \rightarrow *go back*, rest \rightarrow *rest*

2.4 s @ 250 Hz) are accumulated until the decoding process is initiated. Subsequently, a decoding step (forward pass of the ConvNet) is performed whenever 125 new samples (0.5 s @ 250 Hz) have accumulated. All predictions are sent back to the EEG-computer on which a growing ring buffer stores up to 14 predictions corresponding to 7 s of EEG data. Once the ring buffer contains two predictions (i. e., 1 s) our algorithm extracts the mental task with the largest mean prediction. A two-sample t-test is then used to determine if the predictions significantly deviate from 0.05. We define significance as $p < 0.2^3$. These two steps are repeated for all predictions until significance is reached. The ring buffer's size increases (max. 14 predictions) as long as the predictions are not significant. Once significance is reached the GUI action linked to the mental task is executed and the ring buffer is cleared.

4.2. Goal Formulation Planning

Our approach adopts domain-independent planning for high-level control of the robotic system. Whereas many automated planning approaches seek to find a sequence of actions to accomplish a predefined task, the intended goal in this paper is determined by the user. Specifying goals in the former case requires insight into the internal representation of objects in the planning domain. By using a dynamic knowledge base that contains the current world state and referring expressions that describe objects based on their type and attributes, we obstruct direct user access to the internal object representation. Furthermore, we are able to adapt the set of possible goals to changes in the environment. For this purpose, our automatic goal formulation assistant incrementally builds references to feasible goals in a menu-driven graphical user interface.

³Initially we defined significance as $p < 0.1$. Initial experiments however showed that the time required for accumulating evidence to push p from 0.2 to 0.1 was disproportionately large. We therefore define significance as $p < 0.2$ to speed-up the decoding at the cost of accuracy.



```
(: objects cup01 cup02 - cup
shelf01 shelf02 - shelf
omnirob - robot)
(: init (arm-empty omnirob)
(at omnirob shelf02)
(position cup01 shelf02)
(contains cup01 water))
```

Figure 5: Left: The red cup in the real world, referred to by *cup01*. Right: Exemplary PDDL problem description with objects and their initial state.

4.2.1. Domain-Independent Planning

Automated planning is used to transfer a system into a desired goal state by sequentially executing high-level actions. A planning task consists of a planning domain \mathcal{D} and a problem description Π . The former is a tuple $\mathcal{D} = \langle \mathcal{T}, \mathcal{C}_d, \mathcal{P}, \mathcal{O} \rangle$, where

- $\mathcal{T} = \langle T, \prec \rangle$ is the type system together with a partial ordering \prec that specifies the sub-type relations between types in T ,
- \mathcal{C}_d contains a set of domain constant symbols,
- \mathcal{P} is the set of predicate symbols, and
- \mathcal{O} corresponds to the set of planning operators and specifies their effects and preconditions.

The problem description $\Pi = \langle \mathcal{D}, \mathcal{C}_t, \mathcal{I} \rangle$ is defined as follows:

- \mathcal{D} is the domain description,
- \mathcal{C}_t are the additional task-dependent constant symbols, where $\mathcal{C}_d \cap \mathcal{C}_t = \emptyset$, and
- \mathcal{I} is the initial state.

We specify \mathcal{D} and Π using the Planning Domain Definition Language (PDDL) [39]. For example, in the service assistance domains that we use in our experiments, \mathcal{T} contains a type hierarchy, where *furniture* and *robot* are of super-type *base*, and *bottle* and *cup* are of super-type *vessel*. Furthermore, \mathcal{P} specifies attributes or relations between objects, e. g., *arm-empty* is an attribute indicating whether the robot's gripper is empty and *position* is a relation between objects of type *vessel* and *base*. Finally, \mathcal{O} defines actions as *grasp* and *move*. The problem description specifies the initial state \mathcal{I} including object instances, such as *cup01* of type *cup* and *shelf02* of type *shelf* as well as relations between them, e. g., the *position* of *cup01* is *shelf02*, as illustrated in Fig. 5.

4.2.2. Human and Machine Understandable References

A major challenge when trying to communicate goals to the user is the limited shared vocabulary between the user and the planning system, whose world is described by a PDDL planning task. The planner's most concise representation of the cup in Fig. 5 might be *cup01*, which is not

sufficiently clear for the user if there are multiple cups. To solve this problem, the goal generation and selection component uses a set of basic references shared between planner and user. These *shared references* can be combined to create *referring expressions* to objects or sets of objects in the world [40, 41]. Generally, a referring expression ϕ is a logical formula with a single free variable. We say that ϕ refers to an object o if $\mathcal{I} \models \phi(o)$, i. e., ϕ is valid in our PDDL domain theory. For example, we can refer to *cup01* by $\phi(x) \equiv \text{cup}(x) \wedge \text{contains}(x, \text{water})$. We restrict ourselves to references that are conjunctions of relations R_0, \dots, R_m and only allow existential quantifiers, i. e.,

$$\phi(x) = \exists x_1 \dots x_n R_1(x_{11}, \dots) \wedge \dots \wedge R_m(x_{m1}, \dots), \quad (1)$$

where each argument x_{ij} corresponds to one of the variables x_1, \dots, x_n . This is preferable for computational reasons and also allows us to incrementally refine references by adding constraints, e. g., adding $\text{contains}(x, \text{water})$ to $\text{cup}(x)$ restricts the set of all cups to the set of cups containing water. A reference $\phi(o)$ to an object o is *unique* if it refers to exactly one object:

$$\mathcal{I} \models \phi(o) \text{ and } \mathcal{I} \not\models \phi(o') \text{ for all } o \neq o'. \quad (2)$$

However, it is usually sufficient to create references to sets of objects, e. g., if the user wants a glass of water might not be necessary to refer to a specific glass as long as it contains water.

To reference objects in planning domains, we need to specify the components that are required to create *shared references*. We distinguish three fundamental reference types. **Individual references** describe objects that can be identified by their name, e. g., the *content* objects *water* or *apple-juice*, and the *omniRob* robot. Additionally, **type name references** are used to specify objects by their type. They allow referring to unspecific objects as a *shelf* or a *cup*. With **relational references** we can refer an object using a predicate in which the object occurs as an argument. In our scenario, most relational references are binary attribute relations whose first parameter is the object that is referred to, and the second parameter is an object in the domain of attribute values. In the example above, a cup can be described using its *content* by the binary relation $\text{contains}(x, \text{water})$.

The most natural way for the planner to represent a goal is a conjunction of predicates, e. g., $\text{cup}(x) \wedge \text{shelf}(y) \wedge \text{position}(x, y)$ to put a cup on a shelf. This, however, is a rather unnatural way to refer to goals for humans. We found that it is more natural to use the action that achieves the goal than the goal itself, e. g., $\text{action}(\text{put}, x, y) \wedge \text{cup}(x) \wedge \text{shelf}(y)$. Therefore, we include **action references**, a macro reference for all predicates in the action's effect, as additional building blocks to create references to objects in the world and allow the users to specify their goals.

4.2.3. Adaptive Graphical Goal Formulation Interface

In our aim for a flexible yet user-friendly control method to set the robot's goals, we use the references presented

in Sec. 4.2.2 to create a dynamic, menu-driven goal formulation user interface. We allow the user to incrementally refine references to the objects which occur as parameters of a desired goal. We distinguish three different levels of atomicity for the control signals of the GUI: a *step* is a directional command (i. e., *go up*, *go down*, *select*, *go back*) whereas a *stride* is the selection of one refinement option offered by the GUI. A *stride* is therefore a sequence of *steps* ending with either *go back* (to go back to the previous refinement level) or *select* (to further refine the reference to the object) which does not account for the *go up* and *go down* steps. Finally, a *parameter refinement* is the creation of a reference to one parameter. The goal selection procedure is depicted in Fig. 6. After the initial selection of a goal type, e. g., *drop* (Fig. 6.a), we have to determine objects for all parameters of the selected goal. We start by populating the action with the most specific reference that still matches all possible arguments, e. g., *omniRob*, *transportable(x)* and *base(y)*, assuming that *omniRob* corresponds to an individual reference and *transportable* and *base* are type-name references (Fig. 6.b). The current goal reference is displayed in the top row of the GUI. The user interface then provides choices to the user for further refinement of the argument. In our example, the first argument *omniRob* is the only object in the world that fits the parameter type *robot* which is why it does not have to be refined any further. Therefore, we start by offering choices for refining the second argument *transportable(x)* which yields the selections *bottle(x)*, *glass(x)*, *cup(x)* and *vase(x)*. This continues until the argument is either unique, it is impossible to further constrain the argument or any remaining option is acceptable for the user. In the example, we refine the first choice *bottle(x)* based on its *content* (Fig. 6.c) by adding a relation $\text{contains}(x, o)$ to the referring expression, where o is an object of type *content*. This procedure is repeated for all parameters of the goal, which will finally result in a single goal or set of goals (if the references are not unique) that are sent to the high-level planner.

Some features cannot be used to partition the remaining objects for one parameter (e. g., not all objects have the attribute *color*), in which case an entry for all other objects can be chosen. Additionally, we allow to skip the refinement of the current parameter and use an *arbitrary* object for it. Finally, we provide an entry to go *back* to the previous refinement stride.

In each stride the shown refined references form a partition, which corresponds to a set of references ϕ_i :

$$P = \left\{ \phi_1, \dots, \phi_n, \bigwedge_{i=1}^n \neg \phi_i \right\}. \quad (3)$$

The last term ensures, that the partition covers all objects of the previous reference.

The most important part of the refinement process is to compute possible successors that split the current partition. To make progress in selecting a goal, the successor reference needs to be strictly narrower than its parent. Ad-

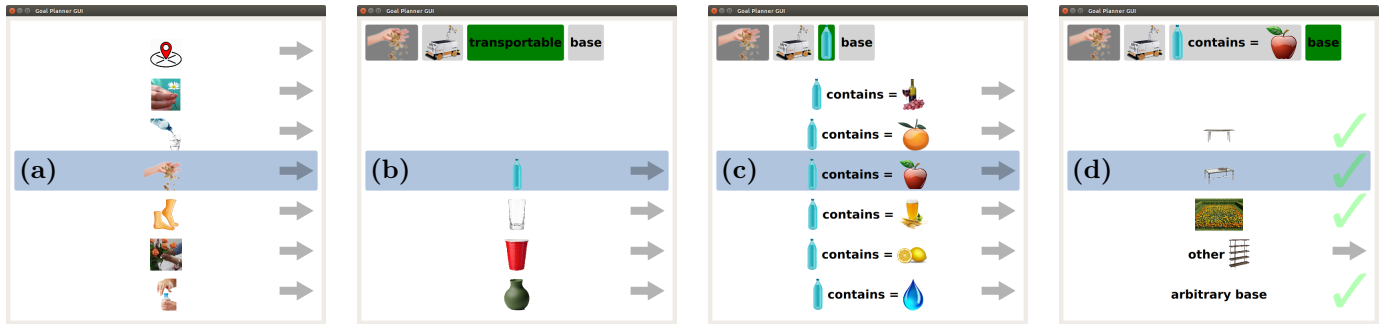


Figure 6: Graphical user interface of the goal formulation assistant. (a) Selection of the desired action. (b) Refinement of the first action parameter of type *transportable*. (c) Refinement of the argument based on *content*. (d) Refinement of the last action parameter of type *base*.

ditionally, forming a partition requires that the references in a candidate set are disjoint. The decision if a complete partition exists corresponds to the NP-complete Exact Cover problem. However, by applying a greedy search algorithm we can approximate the possibly incomplete successor candidate sets in a sufficient way. To ensure complete partitions we add a reference that covers all objects (the *other* entry in our menu) as depicted in Eq. (3). Finally, the references can be reused in the selection process and therefore computed once.

The decision on which successor reference to use for refining the current selection is based on maximizing the resulting partition's information content, which is similarly computed as in decision tree learning [42]. This strategy prefers to split the remaining objects in a way that reduces the total number of refinement strides. Moreover, the method allows to split the referable objects more equally, thus offering the user a meaningful choice at every stride. During the refinement process, we only offer choices that can result in an achievable goal, where goal reachability is efficiently approximated by *delete relaxation* [43]. For example, if all cups were out of reach of the robot, the choice $cup(x)$ would be removed from the selection above. This can result in a completely different selection being preferred, e.g., one that uses the *transportable's* color or position for distinction. If several objects satisfy the specified goal, the planner resolves this ambiguity by picking an arbitrary object among them.

4.3. Robot Motion Generation

For navigation planning of the mobile base, we apply the sampling-based planning framework B^2RRT^* [44]. Given a pair of terminal configurations, it performs a bidirectional search using uniform sampling in the configuration space until an initial sub-optimal solution path is found. This path is subsequently refined for the remaining planning time, adopting an informed sampling strategy, which yields a higher rate of convergence towards the optimal solution. Execution of paths is implemented via a closed-loop joint trajectory tracking algorithm using robot localization feedback.

In this work, we additionally adopt a probabilistic roadmap planner approach [45] to realize pick, place, pour and drink motions efficiently. Therefore, we sample poses in the task space which contains all possible end-effector poses. The poses are then connected by edges based on a user-defined radius. We apply an A^* -based graph search to find an optimal path between two nodes using the Euclidean distance as the cost and heuristic function. To perform robotic motions we need to map the end-effector to joint paths, which can be executed by the robot. We thus use a task space motion controller which uses the robot's Jacobian matrix to compute the joint velocities based on end-effector velocities. Additionally, collision checks ensures that there are no undesired contacts between the environment and the robot. Given a start- and end-pose of a manipulation task, the planner connects them to the existing graph and runs the mentioned search algorithm. For grasping objects, we randomly sample grasp poses around a given object and run the planner to determine a motion plan. Furthermore, we extract horizontal planes from the camera's point cloud and sample poses on these planes to find a suitable drop location for an object. Finally, special motions as drinking and pouring are defined by specifying a path in the cup's rim frame (the point which needs to be connected to the mouth during drinking) and the bottle's rim frame, respectively. Based on these paths the planner samples roadmaps that allow to find motion paths close to the given ones in order to react to small changes in the environment.

4.4. Perception

This section outlines the perception techniques applied in this work and explains how objects and mouth locations are determined and liquid levels are estimated.

Object Detection. In order to detect objects we employ the method of Pauwels *et al.* [46] that relies on dense motion and depth cues and applies sparse keypoint features to extract and track six-degrees-of-freedom object poses in the environment. The algorithm additionally requires models that describe the structure and texture of the detectable objects. It is able to track multiple objects in

597 realtime using a GPU-based solution. These poses and lo-646
598 cations (e. g., the shelf) are finally stored and continuously,647
599 updated in the knowledge base.648

600 *Pouring Liquids.* An important aspect of pouring liquids,649
601 is to be able to determine when to stop pouring. This prevents650
602 overflowing and spilling the liquid as well as opening651
603 up possibilities such as mixing drinks, or preparing meals652
604 where exact amounts of liquid are required. Our approach653
605 to detect the liquid level employs an Asus Xtion Pro cam-654
606 era, which determines depth based on active structured655
607 light. Using this type of sensor, liquids can be categorized656
608 as either opaque or transparent. Opaque liquids, such as657
609 milk or orange juice, *reflect* the infrared light and the ex-658
610 tracted liquid level represents the real liquid level (with659
611 some noise). In the case of transparent liquids, such as660
612 water and apple juice, the infrared light is *refracted* and661
613 the depth value is incorrect.

614 To detect the fill level of a transparent liquid, we base662
615 our approach on a feature further described by Hara *et al.* [47]
616 and Do *et al.* [48]. This feature is given as follows:663

$$h = \left(\frac{\sqrt{n_l^2 - 1 + \cos^2(\alpha)}}{\sqrt{n_l^2 - 1 + \cos^2(\alpha)} - \cos(\alpha)} \right) h_r. \quad (4)$$

617 Here h_r represents the raw depth measured liquid level and668
618 h the estimated liquid height. The index of refraction of669
619 the liquid is given by n_l and angle α is the incidence angle670
620 of infrared light from the camera projector with respect to671
621 the normal of the liquid surface. A Kalman filter is then672
622 used to track the liquid level and compensate for noise.673

623 Before pouring, we first detect the cup in the point674
624 cloud and determine a region within the cup boundaries675
625 where the liquid could be. During the pour, we extract676
626 the depth values for the liquid and estimate the real liquid677
627 height by either applying Eq. 4, in the case of transparent678
628 liquids, or using the extracted value directly, in the case of679
629 opaque liquids. The type of liquid and hence the index of680
630 refraction is given beforehand through the user's selection.681
631 The viewing angle α , can be determined from the depth682
632 data. Once it is detected that the liquid level has exceeded683
633 a user defined value, a stop signal is sent to terminate the684
634 pouring motion.685

635 *Face Detection.* We use a two-step approach to detect and687
636 localize the user's mouth. In the first step, we segment the688
637 image based on the output of a face detection algorithm
638 that uses Haar cascades [49, 50] in order to extract the
639 image region containing the user's mouth and eyes. Af-689
640 terwards, we detect the position of the mouth of the user,
641 considering only the obtained image patch. Regarding the
642 mouth orientation, we additionally consider the position
643 of the eyes in order to obtain a robust estimation of the
644 face orientation, hence compensating for slightly changing
645 angles of the head.690
691
692
693
694
695
696

4.5. Dynamic Knowledge Base

The knowledge base provides data storage and estab-
lishes the communication between all components. In our
work, it is initialized by a domain and problem descrip-
tion based on PDDL files. Once the knowledge base is ini-
tialized, it acts as a central database from which all par-
ticipating network nodes can retrieve information about
specific objects in the world as well as their attributes.
Dynamic behavior is achieved by an additional layer that
allows nodes to add, remove or update objects as well as
their attributes. Moreover, the knowledge base actively
spreads information about incoming changes as updates
on object attributes across the network. Based on this in-
formation each network node decides on its own whether
that information is relevant and which actions need to be
taken.

5. Implementation Details

In our system, we distribute the computation across
a network of seven computers that communicate among
each other via ROS. The decoding of neuronal signals has
four components. EEG measurements are performed using
Waveguard EEG caps with 64 electrodes and a *NeurOne*
amplifier in AC mode. Additionally, vertical and horizon-
tal electrooculograms (EOGs), electromyograms (EMGs)
of the four extremities, electrocardiogram (ECG), electro-
dermal activity (EDA) and respiration are recorded. The
additional data is used to control for ocular and muscular
artifacts, changes in heart beat frequency and skin conduc-
tance, and respiratory frequency, respectively. It is rou-
tinely recorded during EEG experiments in our lab. For
recording and online-preprocessing, we use BCI2000 [38]
and Matlab. We then transfer the data to a GPU server
where our deep ConvNet, implemented using Lasagne [51]
and Theano [52], classifies the data into five classes.

Furthermore, to find symbolic plans for the selected
goal we use the A*-configuration of the *Fast Downward*
planner [53]. The knowledge base is able to store ob-
jects with arbitrary attribute information. All changes
in the knowledge base automatically trigger updates in
the goal formulation GUI. Unexpected changes interrupt
the current motion trajectory execution. Finally, we use
SimTrack [46] for object pose detection and tracking and
OpenCV for face detection.

6. Experiments

We evaluated the proposed framework in multiple ex-
periments. Sec. 6.2 focuses on the BCI control of the whole
system. Afterwards, the results regarding the goal formu-
lation interface are presented in Sec. 6.3. We provide a
detailed performance analysis (Sec. 6.3.2) and a user sur-
vey that studies the friendliness and intuitiveness of the
goal formulation interface (Sec. 6.3.3) based on simulated

697 environments with up to 100 objects in scenarios as ex-
 698 emplarily depicted in Fig. 9. Furthermore, we conducted
 699 two experiments in the real-world environment which are
 700 explained in Sec. 6.1. The results in Sec. 6.4.1 show that
 701 the framework is capable of handling fetch-and-carry tasks
 702 even if there are undesired changes in the environment. Fi-
 703 nally, in Sec. 6.4.2 we discuss the results of combining all
 704 presented components into an autonomous robotic assist-
 705 tant that provides a drink to the user.

6.1. Real World Experimental Environment

707 We performed multiple experiments in the real world
 708 scenario depicted in Fig. 7. It contains two shelves and a
 709 table as potential locations for manipulation actions. The
 710 user sits in a wheelchair in front of a screen that displays
 711 the goal formulation GUI. The autonomous service assist-
 712 ant we use is an *omniRob* omni-directional mobile manip-
 713 ulator platform by KUKA Robotics. The robot is com-
 714 posed of 10 DOF, i. e., three DOF for the mobile base and
 715 seven DOF for the manipulator. Additionally, a Schunk
 716 *Dexterous Hand 2.0* with three fingers is attached to the
 717 manipulator’s flange and used to perform grasping and
 718 manipulation actions, thus adding another DOF for open-
 719 ing and closing the hand. The tasks we consider in our
 720 experiments require the robotic system to autonomously
 721 perform the following actions: drive from one location to
 722 another, pick up an object, drop an object on a shelf or ta-
 723 ble, pour liquids from bottles into cups and supply a user
 724 with a drink. Moreover, our experimental setup uses a per-
 725 ception system composed of five RGBD cameras. Three of
 726 them are statically mounted at the shelves and the table,
 727 in order to observe the scene and to report captured infor-
 728 mation as object locations and liquid levels to the knowl-
 729 edge base. The other two cameras are carried by the robot
 730 on-board. The first one is located at the mobile base and
 731 used to perform collision checks in manipulation planning,
 732 whereas the second camera is mounted at the robot’s end-
 733 effector and used for tasks involving physical human-robot
 734 interaction as serving a drink to a user. Demonstrations of
 735 our work can be found online: <http://www.informatik.uni-freiburg.de/~kuhnerd/neurobots/>.

6.2. Online Decoding of Neuronal Signals

738 We evaluated the BCI control setup with four healthy
 739 participants (P1-4, all right-handed, three females, aged
 740 26.75 ± 5.9). In total, 133 runs have been recorded (90
 741 with the real robot) where the participants selected vari-
 742 ous instructed goals and executed the corresponding task
 743 plans in the goal formulation GUI. For 43 runs, we used
 744 simulated feedback from the GUI in order to generate a
 745 larger amount of data for the evaluation. In this case, we
 746 simulated action executions by simply applying the corre-
 747 sponding effects to the knowledge base. Finally, 38 runs
 748 were discarded because of technical issues with the online
 749 decoding setup.

750 The performance of the BCI decoding during the remain-
 751 ing 95 runs was assessed using video recordings of

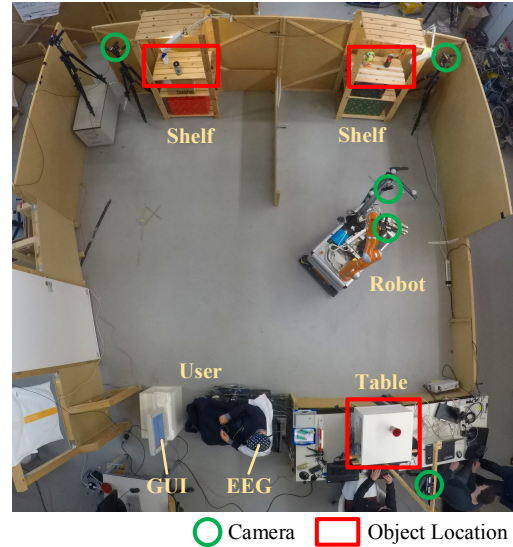


Figure 7: Physical experimental environment: Two shelves and a table could be considered by the robot for performing manipulation actions. Five RGBD sensors observed the environment. A human operator selected a goal using EEG control and the high-level planner GUI.

752 interactions with the GUI. We rated GUI actions as correct
 753 if they correspond to the instructed path and incorrect
 754 otherwise. Actions which were necessary to remediate
 755 a previous error were interpreted as correct if the correc-
 756 tion was intentionally clear. Finally, we rated *rest* actions
 757 as correct during the (simulated) robot executions and ig-
 758 nored them otherwise. For evaluation, five metrics have
 759 been extracted from the video recordings: (i) the accu-
 760 racy of the control, (ii) the time it took the participants
 761 to define a high-level plan, (iii) the number of steps used
 762 to define a high-level plan, (iv) the path optimality, i. e.,
 763 the ratio of the minimally possible number of steps to the
 764 number of steps used (e.g. 1 is a perfect path, while 2
 765 indicates that the actual path was twice longer than the
 766 optimal path), and (v) the average time per step. We
 767 summarized the results in Table 1. In total, a 76.95% cor-
 768 rect BCI control was achieved, which required 9 s per step.
 769 Defining a plan using the GUI took on average 123 s and
 770 required the user to perform on average 13.53 steps in the
 771 GUI of the high-level planner. The path formed by these
 772 steps was on average 1.64 times longer than the optimal
 773 path, mainly because of decoding errors which had to be
 774 corrected by the participants, requiring additional steps.
 775 The decoding accuracy of every participant is significantly
 776 above chance ($p < 10^{-6}$).

The participant-averaged EEG data used to train the hybrid ConvNets and the decoding results of the train/test transfer are visualized in Fig. 8. In Fig. 8(a) we show the signal-to-noise ratio (SNR) of all five classes \mathcal{C} of the labeled datasets. We define the SNR for a given frequency f , time t and channel c as

$$\text{SNR}_{f,t,c} = \frac{\text{IQR}(\{\text{median}(\mathcal{M}_i)\})}{\text{median}(\{\text{IQR}(\mathcal{M}_i)\})} \quad i \in \mathcal{C}, \quad (5)$$

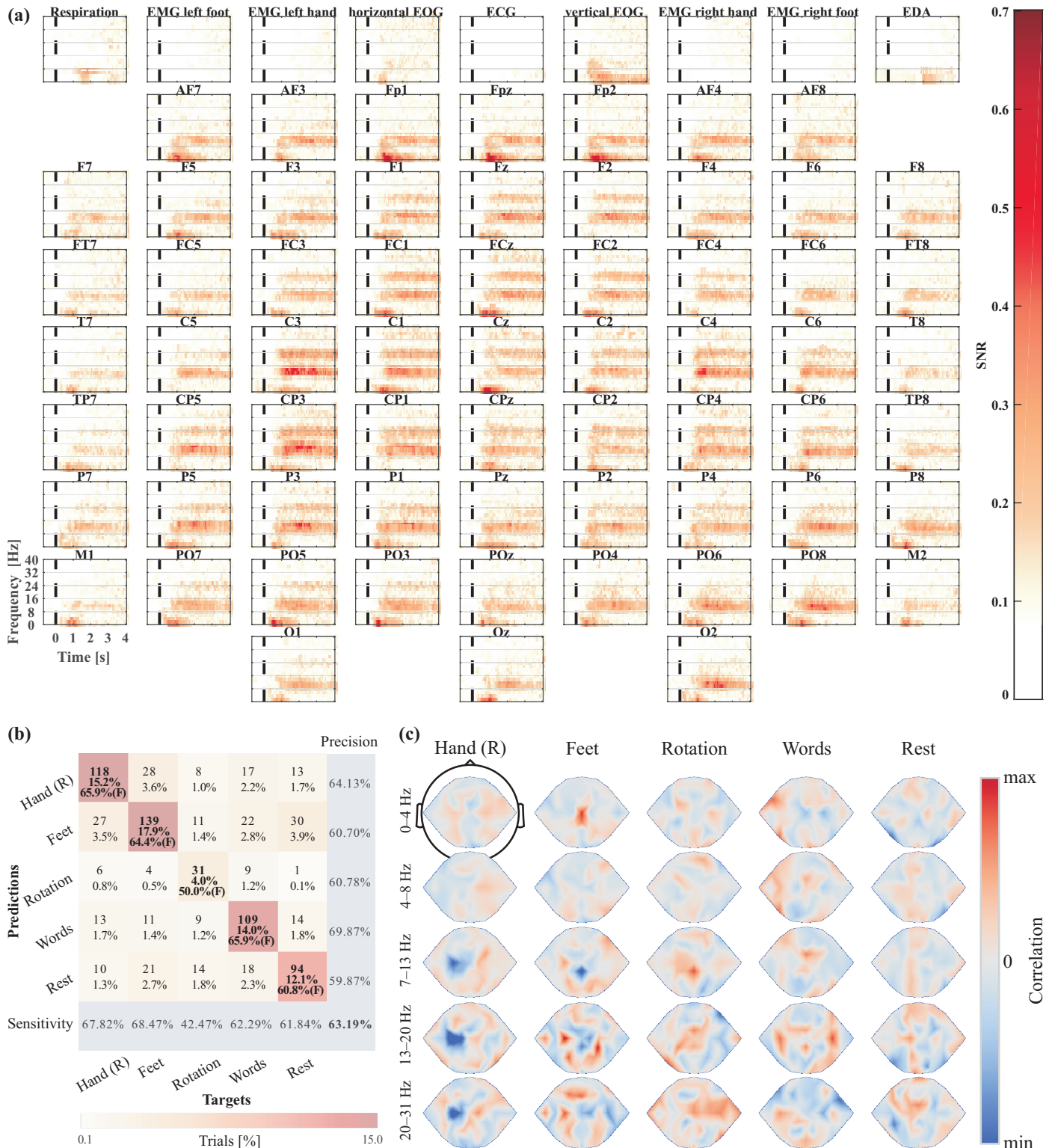


Figure 8: Offline EEG data, offline decoding results and learned features. **(a)** Participant-averaged SNR of the first 4 s of data used to train the hybrid ConvNet. The dashed line indicates the time at which the participants were instructed to start a mental task. Highest SNR can be observed in the alpha (7-14 Hz) and lower beta (16-26 Hz) bands. These frequency bands are robust markers of task related mental activity. Note that the non-EEG channels (top row) were withheld from the ConvNets at any time and are displayed as negative control. The position of most channels was adjusted to achieve a compact layout. **(b)** Confusion matrix of decoding accuracies for the offline train/test transfer pooled over all subjects. Numbers indicate the amount of trials in each cell. Percentages indicate the amount of trials in each cell relative to the total number of trials. (F) indicates F1 score. Dark/light colors indicate that a large/small portion of the targets were predicted for a given class, respectively. **(c)** Topographically plausible input-perturbation network-prediction correlation maps in the delta (0-4 Hz), theta (4-8 Hz), alpha (7-13 Hz), low beta (13-20 Hz) and high beta (20-31 Hz) frequency bands averaged over all participants. The colormap is scaled individually for every frequency band. For details on the visualization technique we refer the reader to [15].

Table 1: Aggregated mean \pm std results for 95 BCI control runs (Exp. 6.2), * p-value $< 10^{-6}$

	Runs #	Accuracy [%]*	Time [s]	Steps #	Path Optimality	Time/Step [s]	
	P1	18	84.1 \pm 6.1	125 \pm 84	12.9 \pm 7.7	1.6 \pm 0.6	9 \pm 2
	P2	14	76.8 \pm 14.1	90 \pm 32	10.1 \pm 2.8	1.1 \pm 0.2	9 \pm 3
	P3	28	78.8 \pm 9.5	173 \pm 144	16.9 \pm 11.5	2.1 \pm 1.3	10 \pm 4
	P4	35	68.1 \pm 16.3	103 \pm 69	14 \pm 7.6	1.7 \pm 0.74	7 \pm 2
	95		76.9 \pm 9.1	123 \pm 36	13.5 \pm 2.8	1.6 \pm 0.4	9 \pm 1

777 where \mathcal{M}_i corresponds to the set of values at position⁸²²
778 (f, t, c) of the i -th task, with $|\mathcal{M}_i|$ being the number of⁸²³
779 repetitions. $\text{median}(\cdot)$ and $\text{IQR}(\cdot)$ is the median and in-⁸²⁴
780 terquartile range (IQR), respectively. The upper part de-⁸²⁵
781 scribes the variance of the class medians, i. e., a large⁸²⁶
782 variance means more distinguishable class clusters and a⁸²⁷
783 higher SNR. The denominator corresponds to the variance⁸²⁸
784 of values in each class, i. e., a lower variance of values re-⁸²⁹
785 sults in a higher SNR.⁸³⁰

786 In all non-peripheral EEG electrodes a clear and sus-⁸³¹
787 tained increase in SNR is visible in the alpha (\sim 8-14 Hz)⁸³²
788 and beta (\sim 14-30 Hz) frequency bands, starting around⁸³³
789 500 ms after the cue. These frequency bands are robust⁸³⁴
790 markers of brain activity. The partial absence of the in-⁸³⁵
791 creased beta band SNR in peripheral channels further sup-⁸³⁶
792 ports the neuronal origin of the signal [54]. An increased⁸³⁷
793 SNR is also visible in both EOG channels which could⁸³⁸
794 indicate a contamination of the EEG data by ocular arti-⁸³⁹
795 facts. The slight increase in SNR in the horizontal EOG⁸⁴⁰
796 channel in the delta (\sim 0-4 Hz) and theta (\sim 4-8 Hz) bands⁸⁴¹
797 0.5-1 s after the cue is most probably due to residual *neu-*⁸⁴²
798 *ronal* activity recorded by the EOG. Support for this as-⁸⁴³
799 sumption is based on the fact that this increase is stronger⁸⁴⁴
800 in most EEG electrodes, suggesting a generator located⁸⁴⁵
801 some distance from the eyes, i. e., in the brain. The sus-⁸⁴⁶
802 tained increase in SNR in the delta band visible in the⁸⁴⁷
803 vertical EOG is likely due to unconscious eye movements.⁸⁴⁸
804 As this increase in the delta band SNR is only visible in⁸⁴⁹
805 the three front-most EEG electrodes and weaker than the⁸⁵⁰
806 increased SNR of unambiguous neuronal origin described⁸⁵¹
807 above, we are confident that the hybrid ConvNets will not⁸⁵²
808 have learned to use this activity to differentiate between⁸⁵³
809 the mental tasks. The visualizations shown in Fig. 8(c)⁸⁵⁴
810 support this idea as no correlations are visible for frontal⁸⁵⁵
811 EEG electrodes in the delta band. The increased SNR in⁸⁵⁶
812 the lower frequencies of the respiration and EDA chan-⁸⁵⁷
813 nels is probably related to task engagement. A crosstalk⁸⁵⁸
814 between these signals and the EEG is unlikely and not sup-⁸⁵⁹
815 ported by the SNR analysis. The extremely low SNR in⁸⁶⁰
816 all EMG channels shows that the participants performed⁸⁶¹
817 pure imagery, without activating their limb muscles. In⁸⁶²
818 summary, the SNR analysis revealed that the offline train-⁸⁶³
819 ing data contains informative neuronal activity which the⁸⁶⁴
820 hybrid ConvNets should have been able to learn from.⁸⁶⁵

821 Indeed, the decoding accuracies (mean 63.0%, P1 70.7%⁸⁶⁶

P2 49.2%, P3 73.1%, P4 58.8%) resulting from the test dataset after initial training of the ConvNets are well above the theoretical chance level of 20%. These are visualized in Fig. 8(b) in the form of a pooled confusion matrix. Right hand and feet motor imagery were most often confused with each other, mental rotation was evenly confused with all other classes and word generation and rest were most often confused with feet motor imagery. The co-adaptive online training which took place between the initial training of the ConvNets and the online evaluation increased the decoding accuracy from 63.0% to 76.9%, which is a clear indication for the efficacy of our approach. It should further be noted that the increase in accuracy occurred from an *offline, cued* evaluation to an *online, uncued* evaluation, which is quite remarkable. It has to be mentioned however that the online accuracy is a subjective measure as the intentions of the participants had to be inferred from the instructions (cf. Sec. 4.1.2). The offline accuracy was fully objective because of the presented cues. Nevertheless, the online evaluation decoding accuracy leaves room for improvements. Preliminary offline steps have been undertaken using the data collected during the offline and online co-adaptive training to detect decoding errors directly from the neuronal signals [20]. This first attempt already yielded mean error detections of 69.33%. The detection accuracy could potentially be increased by including error sensitive peripheral measures as EDA, respiration and ECG into the decoding. Access to the high-gamma (\sim 60-90 Hz) band frequency range could further increase the decoding accuracy of both mental tasks [15] and error signals [55]. Once transferred to an online experiment one could use this error detection to undo the error, generate a new decoding and retrain the decoder. Lastly, detection of robotic errors could also be achieved from the ongoing EEG [56, 57, 58, 21] and used as both emergency stop and teaching signals.

To further support the neural origin of the BCI control signals, Fig. 8(c) shows physiologically plausible input-perturbation network-prediction correlation results (see [15] for methods). Specifically, predictions for right hand and feet motor imagery classes were negatively correlated with input-perturbations (see [15]) in the alpha and beta bands at EEG electrodes located directly above their motor and somatosensory cortex representations. This means that increasing the power in the given frequency bands at the

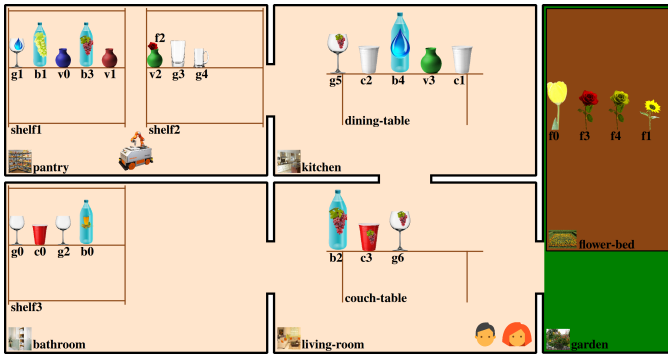


Figure 9: An exemplary scenario as used in our user experiments with four rooms, a garden, two humans, a robot and multiple objects.

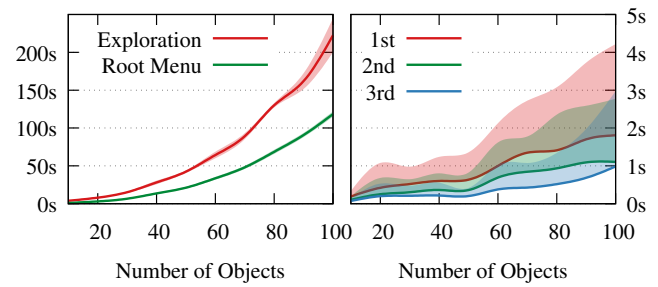


Figure 10: Evaluation of the computation time for different numbers of objects in the environment averaged over random actions. *Left*: The plot shows the mean and standard deviation of building the menu structure at the beginning and includes initial exploration and root menu creation. *Right*: Refinements of a goal can be done efficiently. It shows the mean and positive standard deviation times of the first three refinements.

867 specific electrodes resulted in reduced predictions. By
 868 symmetry, reduced power resulted into increased predictions.
 869 These correlations fit well with the neuronal basis of the event
 870 related desynchronisation of the alpha and beta bands during
 871 motor imagery [59]. A positive correlation is also apparent
 872 above the foot motor and somatosensory cortex for feet motor
 873 imagery in the delta band. This positive correlation probably
 874 reflects the feet motor potential [60]. For mental rotation,
 875 word generation and rest the input-perturbation network-
 876 prediction correlation results are less easily interpretable,
 877 mostly due to the lack of extensive electrophysiological
 878 reports. A positive correlation is visible for the mental
 879 rotation above the medial parietal cortex in the alpha band
 880 which could reflect the involvement of cortical representations
 881 of space. Similarly, positive correlations are visible bilaterally
 882 above the lateral central cortex and temporal cortex in the
 883 low beta band during word generation. They could reflect the
 884 involvement of speech and auditory brain areas. Further
 885 investigations will be needed to delineate these effects.

887 6.3. Goal Formulation Interface

888 In this section we present a performance experiment to
 889 evaluate the runtime required by the GUI and the results of a
 890 preliminary user study, which examines the user-friendliness
 891 and intuitiveness of the system. Moreover, we discuss how
 892 humans use references to objects.

893 6.3.1. Scenario Setup

894 We created a virtual scenario with five rooms as depicted
 895 in Fig. 9: a kitchen with a dining table, a living room with
 896 a couch table, a pantry with two shelves, a bathroom with
 897 one shelf and a garden containing a flowerbed. Bottles, cups,
 898 glasses and vases are distributed among the furniture. There
 899 are three types of flowers (e.g., *rose*), seven drinking con-
 900 tents (e.g., *red-wine*), five colors (e.g., *red*) for cups and
 901 vases and three for flowers and finally, four glass shapes
 902 (e.g., *balloon*). Flowers can be put into vases but may also
 903 be placed directly on furniture. The *omniRob* robot has the
 904 ability to move between the rooms and serve the two persons
 905 (*me* and *friend*). Finally, the

available actions are: *arrange* a flower in a vase, *pick* a
 flower out of a vase, *grasp* and *drop* an object, *give* an
 object to a human, *pour* a liquid from one vessel to another,
drink to assist a human with drinking a drink, *move* the
 robot between rooms and *approach* a furniture or human
 for further interaction.

922 6.3.2. Performance

923 In this experiment we evaluated the performance of the
 924 goal formulation interface. We used a scenario generator
 925 which randomly creates instances of the planning problem.
 926 To assess the performance, we measured the time required
 927 to start the user interface and select parameters of random
 928 actions. The experiment was repeated 100 times and averaged
 929 to retrieve reliable results. The performance of our Python
 930 implementation was determined using an Intel i7-3770K (3.5
 931 GHz) processor and 16 GB of memory. Fig. 10 illustrates the
 932 run times needed for several operations as a function of the
 933 number of objects present in the world. The most time-con-
 934 suming component is given by the reference exploration, where
 935 initial partitions are chosen (Fig. 10 left, red). Another
 936 computationally expensive operation is the root menu genera-
 937 tion, which determines potentially reachable goals for all
 938 actions based on delete relaxation (Fig. 10 left, green). In
 939 contrast, the reference refinements for the current parameter
 940 of an action requires in average less than 2s even for scenar-
 941 ios containing numerous objects (Fig. 10 right). However,
 942 this assertion only holds as long as the world and thus the
 943 references do not change. Considering dynamic environments,
 944 changes of the world are frequently triggered by, e.g., ac-
 945 tions taken by the robotic service assistant. For example,
 946 when the robot has grasped a cup, the system should no longer
 947 refer to the cup as *the cup on the table*. Instead, the refer-
 948 ence must be rebuilt given the updated environment state
 949 yielding *the cup at the robots gripper*. For simplicity, our
 950 approach rebuilds all object references when an environment
 951 change has been detected. In the future, only obsolete
 952 references should be recomputed in order to scale well

944 on larger scenarios.

945 6.3.3. User Study

946 *Participants.* A total of 20 participants (3 female, 17 male,
947 25 – 45 years) took part in the user study and gave their
948 consent for the anonymized processing of the collected
949 data. The participants were students in computer science
950 and administrative employees of the university. They used
951 our system the first time and were not familiar with it.

952 *Data Collection and Measures.* The participants had to
953 use our system to accomplish tasks in five simulated scen-
954 arios, which were generated beforehand to get compara-
955 ble results. The five scenarios with increasing complexity
956 were: (S1) Move the robot to the garden, (S2) Drink beer
957 using a beer mug, (S3) Arrange a red flower in a red vase,
958 (S4) Place a red rose on the couch table, and (S5) Give
959 a red wine glass with red wine to your friend. After in-
960 troducing the user interface by explaining the individual
961 components of the system, the participants had to accom-
962 plish the five tasks using the GUI. Since there were no
963 time constraints and sub-optimal strategies were allowed,
964 all users managed to reach the requested goal states. We
965 counted the number of *steps* the participants required to
966 finish the predefined tasks successfully, where a step is ei-
967 ther a refinement of an attribute or the selection of the
968 *back* entry in the menu.

969 For each scenario the participants had to rate if the dis-
970 played control opportunities offered by the user interface
971 *comply* to their expectations in a questionnaire, where the
972 *compliance* levels ranged from 1 (unreasonable) to 5 (fully
973 comply). Moreover, we asked the participants to rate the
974 overall *intuitiveness* of the GUI in the range of 1 (not in-
975 tuitive) to 5 (excellent). We then asked whether the par-
976 ticipants prefer to describe objects using references or via
977 internal names (e.g., *v2*). Additionally, we evaluated the
978 *subjective quality* of object references ranging from 1 (not
979 prefer at all) to 5 (highly prefer). We proposed four refer-
980 ences to objects depicted in Fig. 9 and let the users rate
981 how well each of those references describes the correspond-
982 ing object. Moreover, subjects were asked to generate refer-
983 ences to these objects in natural language themselves
984 in the way they would tell a friend to find an object. In
985 particular, we considered the green vase with the red rose
986 located in the pantry (*v2*) and the glass, filled with red
987 wine (*g6*), located on the couch table in the living room.
988 The proposed references ranged from under-determined to
989 over-determined descriptions, e.g., *the green vase* vs. *the*
990 *green vase located in the right shelf in the pantry which*
991 *contains a red rose.*

992 *Result.* Fig. 11 shows the quantitative result of the user
993 study. We counted the number of steps performed by each
994 of the participants to achieve the predefined tasks success-
995 fully. The figure shows box plots for each scenario. Ad-
996 ditionally, the plot contains the optimal number of steps
997 which are required to successfully achieve the goal.

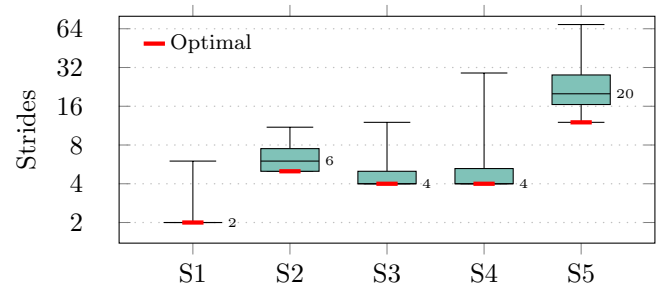


Figure 11: The box plots illustrate the number of steps required by our participants to achieve a given goal in five different scenarios S1-S5 (optimal number of steps indicated in red, numbers denote the median)

Most of the participants were able to find a near-optimal strategy to solve the task. The outliers in the first four scenarios are mainly caused by the user exploring the possibilities of the user interface. The increased number of steps in the last scenario can be traced back to the following reasons. First, the scenario required two actions to be able to achieve the task: fill a balloon shaped glass with red wine and give this glass to the friend. Only a few users were able to determine this fact at the beginning. Therefore, the participants had to correct their decisions which results in a higher number of steps in the fifth scenario. Second, the pour action as defined in our scenarios required to specify three parameters: the vessel to pour from, the vessel to pour to and the liquid that is poured. Our system usually refers to the first vessel by its content, so the redundant refinement of the liquid as last parameter is not intuitive to the users. Finally, we split a partition based on its information content to reduce the number of refinements. This strategy can lead to unexpected refinements of object attributes since the user might prefer these in a different order.

Fig. 12 shows the results on how well the choices offered by the high-level planning GUI actually comply with the expectations of the users. A large percentage of them comply with the refinements provided by the GUI in the scenarios S1 to S4. Due to the previously mentioned problems however, S5 has been rated worse. A short training period of the users to get familiar with the interface might help to improve the compliance in S5. Overall, 80% of the participants rated the GUI as intuitive, i.e., according to the aforementioned metric they rated the intuitiveness with at least 3 (acceptable). In particular, 85% of the participants preferred referring to objects by incremental referencing over internal names (e.g., *green vase on the couch table* vs. *v1*).

In the last user experiment, we evaluated the *subjective quality* of object references. According to our results, preferred references highly depend on whether the spatial context of the agents in the world is considered or not. One group of users only preferred references that uniquely identify the objects independent from the location of the agents. This group preferred references such as *the vase*

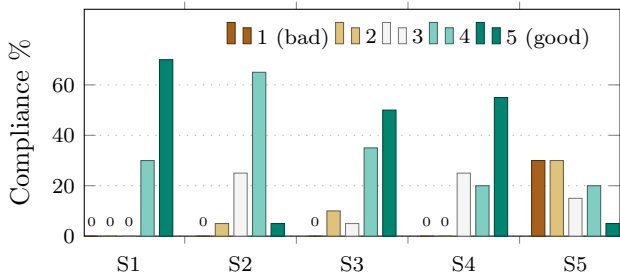


Figure 12: Compliance of the offered choices with the users' expectation for five tasks in different scenarios. The participants had to select compliance levels from 1 (unreasonable) to 5 (fully comply).

containing a rose or occasionally also the vase in the right shelf for $v2$ and the red wine glass on the couch table for $v6$. Another group preferred under-determined references as they considered the spatial context of the agents. This group preferred references such as the green vase for $v2$ and the red wine glass for $v6$. Interestingly, the capability of users to impersonate the acting agent has also a strong influence on the references preferred by the second group. For referring to $v2$, some users of the second group additionally specified the room or the content of the vase assuming that the assisting agent is also located in the living room and therefore requires a more detailed object description, while they preferred under-specified references for objects on the couch table. Detailed over-specified references were refused by all participants, but more firmly by the second group. Summarizing, our evaluation revealed that incrementally building object references is suitable to describe objects precisely. Room for improvement was identified in updating object references that change during plan execution and in the consideration of temporal and spatial context.

6.4. Robotic Service Assistant

We performed two experiments to evaluate the system in the real world using a mobile robot. The first one explores how the system reacts to unexpected changes. In the second experiment, we present the results of the whole system involving all components.

6.4.1. Fetch and Carry Task with Disturbances

In a dynamic world, unexpected changes such as adding or removing objects can occur at all times. With this experiment, we examine how our system adapts to disturbances, i.e. unexpected changes in the environment.

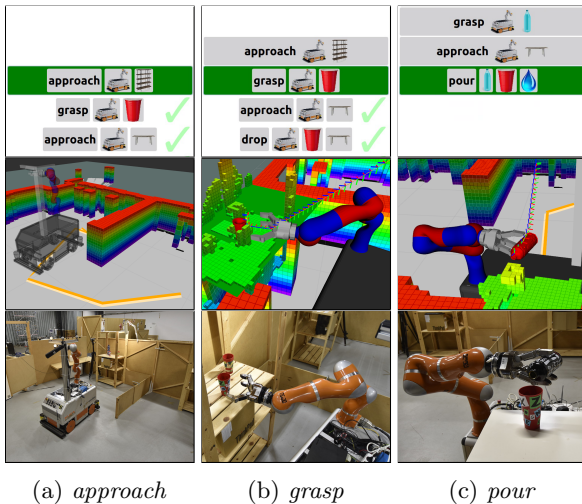
We performed the experiments in a way that unexpected world changes may occur at any time through actions taken by another unknown agent. In practice, this agent could refer to a human taking actions that directly affect the execution of the current high-level plan. Therefore, we initially placed a cup on one of the shelves and queried the goal formulation assistant to generate a sequence of actions leading to the goal state *cup on table*, i.e., *approach(shelf with cup)*, *grasp(cup)*, *approach(table)*

drop(cup). Once the robot arrived at the corresponding shelf in the execution phase of the plan, a human agent took the cup while the robot was about to grasp it and transferred it to the other shelf. In order to obtain quantitative results on the performance of our framework in such a scenario, we ran this experiment 10 times with different initial cup placements and evaluated its ability to generate the goal state in the real world despite the external disturbance introduced by the human agent. For all runs, our perception system correctly updated the information on the cup in the knowledge base, in turn triggering a re-planning step. The updated action sequence always contained two additional actions, namely moving to the shelf where the human agent dropped the cup and grasping it again. In total, 59 out of 60 (98.33%) scheduled actions were successfully executed and thus 90% of the runs succeeded in generating the goal state. Only one run failed in the action execution phase due to the inability of the low-level motion planning algorithm to generate a solution path for the mobile base within the prescribed planning time. On average, our system required an overall time of 258.7 ± 28.21 s for achieving the goal state.

6.4.2. Drinking Task

The last experiment evaluates the direct interaction between user and robot. Therefore, we implemented an autonomous robotic drinking assistant. Our approach enabled the robot to fill a cup with a liquid, move the robot to the user and finally provide the drink to the user by execution of the corresponding drinking motion in front of the user's mouth. Fig. 13 shows examples for the actions *move*, *grasp* and *pour*. The first row contains the task-plan visualizations of the goal formulation GUI, which are displayed after a goal has been selected. Additionally, the second row depicts the planning environment as used by the navigation and manipulation planners to generate collision-free motions. The corresponding view of the real world is shown in the last line.

Table 2 shows the averaged results for the experiment. Again, the user is one of the authors. Here, only 3.75% of the 160 scheduled actions had to be repeated in order to complete the task successfully. In one run, plan recovery was not possible leading to abortion of the task. Thus, our system achieved in total a success rate of 90% for the drinking task. Planning and execution required on average 545.56 ± 67.38 s. For the evaluation of the liquid level detection approach, we specified a desired fill level and executed 10 runs of the pour action. The resulting mean error and standard deviation is 6.9 ± 8.9 mm. In some instances the bottle obstructed the camera view, resulting in poor liquid level detection and a higher error. We have begun investigation possible improvements for the monitoring of liquid levels by additionally considering the brain activity of an observer [58, 21]. Our latest results show that events where the liquid spills over the cup can be detected with an accuracy of 78.2 ± 8.4 % (mean over 5 subjects) using deep ConvNets and EEG [21]. This feedback can be used



(a) *approach* (b) *grasp* (c) *pour*

Figure 13: Snapshots of our experiments for the actions *approach*, *grasp* and *pour*. The first line shows the corresponding step in the high-level planner user interface. The results of the motion and manipulation planning is depicted in the second row. Finally, the third row shows the robot system, which executes the actions.

Actions	# Executions (# Scheduled)	Success Exec. [%]	Runtime [s]	
			Mean	Std
Grasp	34 (30)	91.0	40.42	10.31
Drop	30 (30)	97.0	37.59	4.83
Approach	80 (80)	100.0	20.91	7.68
Pour	10 (10)	100.0	62.90	7.19
Drink	13 (10)	77.0	57.10	8.20
Total	167 (160)	95.86	32.46	15.51

Table 2: Aggregated results for 10 runs (Exp. 6.4.2)

to inform the liquid level detection and pouring procedure of the failure and trigger an adaptation of the algorithms to prevent future spills. To completely prevent errors, detection prior to a spill event will have to be achieved in future work.

7. Conclusions

In this paper, we presented a novel framework that allows users to control a mobile robotic service assistant by thought. This is particularly interesting for severely paralyzed patients who constantly rely on human caretakers as some independence is thereby restored. Our system performs complex tasks in dynamic real-world environments, including fetch-and-carry tasks and close-range human-robot interactions. Our experiments revealed that the five-class-BCI has an uncued online decoding accuracy of 76.9%, which enables users to specify robotic tasks using intelligent goal formulation. Furthermore, a user study substantiates that participants perceive the goal formulation interface as user-friendly and intuitive. Finally, we conducted experiments in which the proposed autonomous robotic service assistant successfully provides drinks to

humans. By combining techniques from brain signal decoding, natural language generation, task planning, robot motion generation, and computer vision we overcome the curse of dimensionality typically encountered in robotic BCI control schemes. This opens up new perspectives for human-robot interaction scenarios.

8. Acknowledgements

This research was supported by the German Research Foundation (DFG, grant number EXC1086), the Federal Ministry of Education and Research (BMBF, grant number Motor-BIC 13GW0053D) and grant BMI-Bot (grant number ROB-7) by the Baden-Württemberg Stiftung.

9. References

- [1] D. Park, Y. K. Kim, Z. M. Erickson, C. C. Kemp, Towards assistive feeding with a general-purpose mobile manipulator, arXiv:1605.07996.
- [2] C. S. Chung, H. Wang, R. A. Cooper, Autonomous function of wheelchair-mounted robotic manipulators to perform daily activities, in: 2013 IEEE 13th International Conference on Rehabilitation Robotics (ICORR), 2013, pp. 1–6. doi:10.1109/ICORR.2013.6650378.
- [3] F. Achic, J. Montero, C. Penalzoa, F. Cuellar, Hybrid bci system to operate an electric wheelchair and a robotic arm for navigation and manipulation tasks, in: Workshop on Advanced Robotics and its Social Impacts (ARSO), IEEE, 2016, pp. 249–254.
- [4] C. Wang, B. Xia, J. Li, W. Yang, D. Xiao, A. C. Velez, H. Yang, Motor imagery bci-based robot arm system, in: Natural Computation (ICNC), 2011 Seventh International Conference on, Vol. 1, IEEE, 2011, pp. 181–184.
- [5] S. Schröer, I. Killmann, B. Frank, M. Völker, L. D. J. Fiederer, T. Ball, W. Burgard, An autonomous robotic assistant for drinking, in: International Conference on Robotics and Automation (ICRA), IEEE, 2015, pp. 6482–6487.
- [6] S. M. Grigorescu, T. Lüth, C. Fragkopoulos, M. Cyriacks, A. Gräser, A BCI-controlled robotic assistant for quadriplegic people in domestic and professional life, *Robotica* 30 (03) (2012) 419–431. doi:10.1017/S0263574711000737.
- [7] J. Mladenović, J. Mattout, F. Lotte, A generic framework for adaptive EEG-based BCI training and operation, arXiv:1707.07935 [cs, q-bio]ArXiv: 1707.07935. URL <http://arxiv.org/abs/1707.07935>
- [8] K. Muelling, A. Venkatraman, J.-S. Valois, J. E. Downey, J. Weiss, S. Javdani, M. Hebert, A. B. Schwartz, J. L. Collinger, J. A. Bagnell, Autonomy infused teleoperation with application to brain computer interface controlled manipulation, *Autonomous Robots* (2017) 1–22doi:10.1007/s10514-017-9622-4.
- [9] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: Advances in neural information processing systems, 2012, pp. 1097–1105.
- [10] K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Image Recognition, arXiv:1512.03385.
- [11] T. N. Sainath, B. Kingsbury, G. Saon, H. Soltau, A.-r. Mohamed, G. Dahl, B. Ramabhadran, Deep Convolutional Neural Networks for Large-scale Speech Tasks, *Neural Networks* 64 (2015) 39–48. doi:10.1016/j.neunet.2014.08.005.
- [12] T. Sercu, C. Puhersch, B. Kingsbury, Y. LeCun, Very deep multilingual convolutional neural networks for LVCSR, in: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2016, pp. 4955–4959. doi:10.1109/ICASSP.2016.7472620.

- [13] V. J. Lawhern, A. J. Solon, N. R. Waytowich, S. M. Gordon, C. P. Hung, B. J. Lance, Eegnet: A compact convolutional neural network for eeg-based brain-computer interfaces, *CoRR* abs/1611.08024. 1294
- [14] Y. R. Tabar, U. Halici, A novel deep learning approach for classification of EEG motor imagery signals, *Journal of Neural Engineering* 14 (1). doi:10.1088/1741-2560/14/1/016003. 1297
- [15] R. T. Schirrmeister, J. T. Springenberg, L. D. J. Fiederer, M. Glasstetter, K. Eggensperger, M. Tangermann, F. Hutter, W. Burgard, T. Ball, Deep learning with convolutional neural networks for EEG decoding and visualization, *Human Brain Mapping* 38 (11) (2017) 5391–5420. doi:10.1002/hbm.23730. 1302
- [16] K. K. Ang, Z. Y. Chin, H. Zhang, C. Guan, Filter Bank Common Spatial Pattern (FBCSP) in Brain-Computer Interface, in: *IEEE International Joint Conference on Neural Networks (IJCNN)*, 2008, pp. 2390–2397. doi:10.1109/IJCNN.2008.4634130. 1307
- [17] P. Bashivan, I. Rish, M. Yeasin, N. Codella, Learning Representations from EEG with Deep Recurrent-Convolutional Neural Networks, *ArXiv e-prints*. 1310
- [18] S. Stober, Learning Discriminative Features from Electroencephalography Recordings by Encoding Similarity Constraints in: *Bernstein Conference 2016*, 2016. doi:10.12751/nncn.bc2016.0223. 1314
- [19] K. G. Hartmann, R. T. Schirrmeister, T. Ball, Hierarchical representation of spectral features in deep convolutional networks trained for EEG decoding, in: *2018 6th International Conference on Brain-Computer Interface (BCI)*, 2018, pp. 1–6. doi:10.1109/IWW-BCI.2018.8311493. 1319
- [20] M. Völker, R. T. Schirrmeister, L. D. J. Fiederer, W. Burgard, T. Ball, Deep transfer learning for error decoding from non-invasive EEG, in: *2018 6th International Conference on Brain-Computer Interface (BCI)*, 2018, pp. 1–6. doi:10.1109/IWW-BCI.2018.8311491. 1324
- [21] J. Behncke, R. T. Schirrmeister, W. Burgard, T. Ball, The signature of robot action success in EEG signals of a human observer: Decoding and visualization using deep convolutional neural networks, in: *2018 6th International Conference on Brain-Computer Interface (BCI)*, 2018, pp. 1–6. doi:10.1109/IWW-BCI.2018.8311531. 1330
- [22] P. Ortega, C. Colas, A. Faisal, Convolutional neural network work, personalised, closed-loop Brain-Computer Interfaces for multi-way control mode switching in real-time, *bioRxiv* (2018) 256701. doi:10.1101/256701. 1334
- [23] F. Burget, L. D. J. Fiederer, D. Kuhner, M. Völker, J. Aldinger, R. T. Schirrmeister, C. Do, J. Boedecker, B. Nebel, T. Ball, W. Burgard, Acting thoughts: Towards a mobile robotic service assistant for users with limited communication skills, in: *2017 European Conference on Mobile Robotics (ECMR)*, 2017, pp. 1339–1340. doi:10.1109/ECMR.2017.8098658. 1340
- [24] E. Krahmer, K. Van Deemter, Computational generation of referring expressions: A survey, *Computational Linguistics* 38 (1) (2012) 173–218. 1343
- [25] M. Shridhar, D. Hsu, Grounding spatio-semantic referring expressions for human-robot interaction, *CoRR* abs/1707.05720. 1345
- [26] L. Yu, P. Poirson, S. Yang, A. C. Berg, T. L. Berg, Modeling context in referring expressions, *CoRR*. 1347
- [27] A. Koller, R. P. Petrick, Experiences with planning for natural language generation, *Computational Intelligence* 27. 1349
- [28] L. P. Kaelbling, T. Lozano-Pérez, Hierarchical task and motion planning in the now, in: *IEEE International Conference on Robotics and Automation (ICRA)*, IEEE, 2011, pp. 1470–1477. 1352
- [29] L. P. Kaelbling, T. Lozano-Pérez, Integrated task and motion planning in belief space, *I. J. Robotics Res.* 32 (9-10) (2013) 1194–1227. doi:10.1177/0278364913484072. 1355
- [30] L. De Silva, A. K. Pandey, M. Gharbi, R. Alami, Towards combining HTN planning and geometric task planning, *CoRR* abs/1307.1482. 1358
- [31] N. T. Dantam, Z. K. Kingston, S. Chaudhuri, L. E. Kavrakı, Incremental task and motion planning: A constraint-based approach, in: *Proceedings of Robotics: Science and Systems*, Ann Arbor, Michigan, 2016, pp. 1–6. doi:10.15607/RSS.2016.XII.002. 1360
- [32] S. Srivastava, E. Fang, L. Riano, R. Chitnis, S. Russell, P. Abbeel, Combined task and motion planning through an extensible planner-independent interface layer, in: *IEEE International Conference on Robotics and Automation (ICRA)*, IEEE, 2014, pp. 639–646. 1362
- [33] T. Lozano-Pérez, L. P. Kaelbling, A constraint-based method for solving sequential manipulation planning problems, in: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, IEEE, 2014, pp. 3684–3691. 1364
- [34] C. Dornhege, P. Eyerich, T. Keller, S. Trüg, M. Brenner, B. Nebel, Semantic Attachments for Domain-independent Planning Systems, in: *Proceedings of the 19th International Conference on Automated Planning and Scheduling (ICAPS)*, 2009. 1366
- [35] P. Eyerich, R. Mattmüller, G. Röger, Using the Context-enhanced Additive Heuristic for Temporal and Numeric Planning, in: *Proceedings of the 19th International Conference on Automated Planning and Scheduling (ICAPS 2009)*, 2009, pp. 130–137. 1368
- [36] D.-A. Clevert, T. Unterthiner, S. Hochreiter, Fast and Accurate Deep Network Learning by Exponential Linear Units (ELUs), in: *ArXiv e-prints*, Vol. 1511, 2016, p. arXiv:1511.07289. 1370
- [37] D. P. Kingma, J. Ba, Adam: A Method for Stochastic Optimization, in: *International Conference on Learning Representations (ICLR)*, 2015. 1372
- [38] G. Schalk, D. McFarland, T. Hinterberger, N. Birbaumer, J. Wolpaw, BCI2000: A General-Purpose Brain-Computer Interface (BCI) System, *IEEE Transactions on Biomedical Engineering* 51 (6) (2004) 1034–1043. doi:10.1109/TBME.2004.827072. 1374
- [39] D. McDermott, PDDL - The Planning Domain Definition Language, *Tech. rep.*, AIPS-98 Planning Competition Committee (1998). 1376
- [40] R. Dale, E. Reiter, Computational interpretations of the greece maxims in the generation of referring expressions, *Cognitive science* 19 (2) (1995) 233–263. 1378
- [41] M. Göbelbecker, Assisting with Goal Formulation for Domain Independent Planning, in: *KI 2015: Advances in Artificial Intelligence*, Springer, 2015, pp. 87–99. 1380
- [42] J. Quinlan, Induction of decision trees, *Machine Learning* 1 (1986) 81–106. 1382
- [43] B. Bonet, G. Loerincs, H. Geffner, A Robust and Fast Action Selection Mechanism for Planning, in: *Proceedings of the 14th National Conference on Artificial Intelligence and 9th Innovative Applications of Artificial Intelligence Conference (AAAI 1997/ IAAI 1997)*, 1997, pp. 714–719. 1384
- [44] F. Burget, M. Bennewitz, W. Burgard, BI²RRT*: An efficient sampling-based path planning framework for task-constrained mobile manipulation, in: *International Conference on Intelligent Robots and Systems (IROS)*, IEEE/RSJ, Daejeon, Korea, 2016, pp. 3714–3721. doi:10.1109/IROS.2016.7759547. 1386
- [45] L. E. Kavrakı, J.-C. Latombe, *Probabilistic Roadmaps for Robot Path Planning*, John Wiley, 1998, pp. 33–53. 1388
- [46] K. Pauwels, D. Kragic, Simtrack: A simulation-based framework for scalable real-time object pose detection and tracking, in: *International Conference on Intelligent Robots and Systems (IROS)*, IEEE/RSJ, 2015, pp. 1300–1307. 1390
- [47] Y. Hara, F. Honda, T. Tsubouchi, A. Ohya, Detection of liquids in cups based on the refraction of light with a depth camera using triangulation, in: *Intelligent Robots and Systems (IROS 2014)*, 2014 IEEE/RSJ International Conference on, IEEE, 2014, pp. 5049–5055. 1392
- [48] C. Do, T. Schubert, W. Burgard, A probabilistic approach to liquid level detection in cups using an RGB-D camera, in: *International Conference on Intelligent Robots and Systems (IROS)*, IEEE/RSJ, 2016, pp. 2075–2080. 1394
- [49] P. Viola, M. Jones, Rapid object detection using a boosted cascade of simple features, in: *Computer Vision and Pattern Recognition*, 2001. *CVPR 2001. Proceedings of the 2001 IEEE Computer Society Conference on*, Vol. 1, IEEE, 2001, pp. I–I. 1396

- 1362 [50] R. Lienhart, J. Maydt, An extended set of haar-like features for
1363 rapid object detection, in: Image Processing. 2002. Proceedings.
1364 2002 International Conference on, Vol. 1, IEEE, 2002, pp. I-I.
- 1365 [51] S. Dieleman, J. Schlter, C. Raffel, E. Olson, S. K. Snderby,
1366 D. Nouri, et al., Lasagne: First release. (Aug. 2015). doi:10.
1367 5281/zenodo.27878.
- 1368 [52] Theano Development Team, Theano: A Python framework for
1369 fast computation of mathematical expressions, arXiv e-prints
1370 abs/1605.02688.
- 1371 [53] M. Helmert, The Fast Downward Planning System, Journal of
1372 Artificial Intelligence Research 26 (JAIR 2006) (2006) 191–246.
- 1373 [54] I. I. Goncharova, D. J. McFarland, T. M. Vaughan, J. R. Wol-
1374 paw, EMG contamination of EEG: spectral and topographical
1375 characteristics, Clinical Neurophysiology 114 (9) (2003) 1580–
1376 1593. doi:10.1016/S1388-2457(03)00093-2.
- 1377 [55] M. Völker, L. D. J. Fiederer, S. Berberich, J. Hammer,
1378 J. Behncke, P. Kršek, M. Tomášek, P. Marusič, P. C. Reinacher,
1379 V. A. Coenen, M. Helias, A. Schulze-Bonhage, W. Burgard,
1380 T. Ball, The dynamics of error processing in the human brain as
1381 reflected by high-gamma activity in noninvasive and intracranial
1382 EEG, NeuroImage doi:10.1016/j.neuroimage.2018.01.059.
- 1383 [56] I. Iturrate, R. Chavarriaga, L. Montesano, J. Minguez, J. d. R.
1384 Millán, Teaching brain-machine interfaces as an alternative
1385 paradigm to neuroprosthetics control, Scientific Reports 5. doi:
1386 10.1038/srep13893.
- 1387 [57] A. F. Salazar-Gomez, J. DelPreto, S. Gil, F. H. Guenther,
1388 D. Rus, Correcting robot mistakes in real time using EEG sig-
1389 nals, in: 2017 IEEE International Conference on Robotics and
1390 Automation (ICRA), 2017, pp. 6570–6577. doi:10.1109/ICRA.
1391 2017.7989777.
- 1392 [58] D. Welke, J. Behncke, M. Hader, R. T. Schirrmeister,
1393 A. Schönau, B. Eßmann, O. Müller, W. Burgard, T. Ball,
1394 Brain Responses During Robot-Error Observation, Kognitive
1395 Systeme, 2017 - 1 doi:10.17185/dupublico/44533.
- 1396 [59] G. Pfurtscheller, F. H. Lopes da Silva, Event-related
1397 EEG/MEG synchronization and desynchronization: basic prin-
1398 ciples, Clinical Neurophysiology 110 (11) (1999) 1842–1857.
1399 doi:10.1016/S1388-2457(99)00141-8.
- 1400 [60] L. Gilden, H. G. Vaughan, L. D. Costa, Summated human EEG
1401 potentials with voluntary movement, Electroencephalography
1402 and Clinical Neurophysiology 20 (5) (1966) 433–438. doi:10.
1403 1016/0013-4694(66)90100-3.