

# Mapping Language to Vision in a Real-World Robotic Scenario

Karla Štepánová<sup>ID</sup>, Frederico Belmonte Klein, Angelo Cangelosi, and Michal Vavrečka

**Abstract**—Language has evolved over centuries and was gradually enriched and improved. The question, how people find assignment between meanings and referents, remains unanswered. There are many of computational models based on the statistical co-occurrence of meaning-reference pairs. Unfortunately, these mapping strategies show poor performance in an environment with a higher number of objects or noise. Therefore, we propose a more robust noise-resistant algorithm. We tested the performance of this novel algorithm with simulated and physical iCub robots. We developed a testing scenario consisting of objects with varying visual properties presented to the robot accompanied by utterances describing the given object. The results suggest that the proposed mapping procedure is robust, resistant against noise and shows better performance than one-step mapping for all levels of noise in the linguistic input, as well as slower performance degradation with increasing noise. Furthermore, the proposed procedure increases the clustering accuracy of both modalities.

**Index Terms**—Cognitive modeling, cross-situational learning, iCub robot, language acquisition, symbol grounding.

## I. INTRODUCTION

THE ESSENTIAL (and still not fully answered) question in language acquisition is how percepts are anchored in some arbitrary symbols, in other words, how words (symbols) get their meanings. This is the so-called symbol grounding problem [20]. For many years, cognitive modeling, neuroscience, psychology, and machine learning have jointly attempted to understand how humans can solve this “problem” [6]. The ability to learn language through perception and especially through visual grounding is not only important for understanding human cognition but is also applicable in many areas, such as verbal control of interactive robots [28], automatic sports commentators [15], car navigation systems, for the visually impaired, situated

Manuscript received June 27, 2017; revised December 6, 2017; accepted February 24, 2018. This work was supported in part by the European EU FP7 Research Project TRADR under Grant 609763, in part by TACR CAK under Grant TE01020197, in part by the CAPES Foundation, Ministry of Education of Brazil under Grant BEX 1084/13-5, in part by the CNPq Brazil under Grant 232590/2014-1, and in part by the U.K. EPSRC Project BABEL under Grant EP/J004561/1 and Grant EP/J00457X/1. (Corresponding author: Karla Štepánová.)

K. Štepánová and M. Vavrečka are with the Czech Institute of Informatics, Robotics, and Cybernetics, Czech Technical University in Prague, 16000 Prague, Czech Republic (e-mail: karla.stepanova@ciirc.cvut.cz).

F. B. Klein and A. Cangelosi are with the School of Computing, Electronics and Mathematics, Plymouth University, Plymouth PL4 8AA, U.K.

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TCDS.2018.2819359

speech understanding in computer games [19], automated generation of weather forecasts [18], tutoring children in foreign languages [22], etc.

Despite the extensive research in the area of language acquisition, the question how the word-to-meaning mapping is learned remains unanswered. There are basically two approaches for mapping language to other modalities. A classic approach for separate modality-dependent representations is advocated by Barsalou [1]; while an example of a system where language is mapped to an intermediate modal representation that can be derived by multiple modalities is [27]. Using the unsupervised approach, Li *et al.* [24] designed the DevLex model, consisting of two self-organizing networks that are bidirectionally connected. Gliozzi *et al.* [17] proposed an alternative with a multimodal representation layer: their unsupervised feature-based model was used to account for early category formation in young infants. This approach postulates the unsupervised role of linguistic labels that can affect categorization during the acquisition process, which has also been supported by Taniguchi *et al.* [47]. Vavrečka and Farkaš [51] recently introduced a multimodal architecture for the grounding of spatial words using a biologically inspired approach (separate “what” and “where” visual subsystems) in which the visual scenes (two objects in 2-D space in a spatial relation) are associated with their linguistic descriptions, thus leading to the integration of modalities.

However, a fully unsupervised architecture, which would be able to deal with language grounding [46], particularly language grounding in a case where sentences have variable structure and when there is more than one object in a scene, is not available. The current state-of-the-art on variable-length sentences is very restricted and deals only with static scenes [29]. Most of the recent models based on deep networks are oriented toward application in image-to-text [21] or video-to-text [11] mapping and do not take into account the psychological aspects of language acquisition (e.g., mutual exclusivity). Moreover, these systems are trained in a supervised manner without the advantage of transfer learning.

The difficulty of the task was described in a well-known experiment performed by Quine [35] who imagined the anthropologist meeting a native who pointed at the scene and said “gavagai.” When the anthropologist is stimulated in a situation by seeing a rabbit, he will suppose that the word represents the running rabbit in front of him, even though it could also mean “ground,” “sun,” “hello,” or whatever else. This problem is related to language relativity, as there are several

82 objects and their features that are described by words [33]. A  
 83 simplified version of this problem consists of a simple visual  
 84 scene and separate words that are grounded based on statistical  
 85 co-occurrence (cross-situational learning).

86 From a neuroscientific point of view, symbol grounding can  
 87 be viewed as a process of finding mappings between primary  
 88 unimodal visual and language brain areas. Where exactly the  
 89 integration is performed is still the subject of research, and  
 90 existing literature provides only incomplete accounts of the  
 91 cortical location of this convergence. For example, the study  
 92 by Büchel *et al.* [4] provides evidence for the involvement of  
 93 the left basal posterior temporal lobe (BA37) in the integration  
 94 of language and visual information. Other studies (e.g., [43])  
 95 propose that access to verbal meaning depends on anterior  
 96 and posterior heteromodal cortical systems within the tem-  
 97 poral lobe. The grounding of actions and motoric primitives  
 98 is associated with the activity in the dorsal stream and the  
 99 premotor cortex [8].

100 How language can be developed in an unsupervised man-  
 101 ner is also an important task in developmental robotics as most  
 102 language acquisition in humans is fully unsupervised. One of  
 103 the main long-term objectives of many teams worldwide is  
 104 building conversational robots, which will be able to partici-  
 105 pate in cooperative tasks mediated by a natural language. It has  
 106 been shown how robots can learn new symbols using already  
 107 grounded ones and their combination [5] and how to transfer  
 108 knowledge between agents [52]. Cangelosi *et al.* [5] presented  
 109 their research on language emergence and grounding in senso-  
 110 rimotor agents and robots. This model was further extended by  
 111 Tikhanoff *et al.* [49], who performed iCub simulation experi-  
 112 ments and focused on the integration of speech and action. The  
 113 grounding of higher-order concepts in action was also explored  
 114 by Stramandinoli *et al.* [44], who made use of recurrent neural  
 115 networks. Sugita and Tani [45] described an experiment deal-  
 116 ing with semantic compositionality: the capability of a robot to  
 117 use the compositional structure to generalize novel word com-  
 118 binations. Daoutis and Mavridis [9] summarized desiderata for  
 119 grounded semantic compositionality. Despite all the progress  
 120 in language grounding, however, the current state-of-the-art  
 121 on grounding variable-length sentences is very restricted and  
 122 deals only with static scenes [29], [38].

123 In this paper, we present this paper in the area of lan-  
 124 guage acquisition using a real-world robotic scenario. We  
 125 implemented a hierarchical cognitive architecture for lan-  
 126 guage acquisition that includes visual and language processing.  
 127 In particular, we chose to extend current models of cross-  
 128 situational learning by allowing vision-to-language mapping  
 129 in the case of a nonequal number of classes and by taking  
 130 into account situation-time dynamics.

131 We show that this can be accomplished more efficiently  
 132 by replacing one-shot mapping with sequential mapping and  
 133 adding inhibitory mechanisms to the connections. The best  
 134 mapped classes are gradually eliminated, and the clusterization  
 135 is adaptively changed. We see this paper as an extension of  
 136 the McMurray *et al.* [29] model, and we compare it with other  
 137 single-step mapping models. The mapping strategy presented  
 138 in this paper was shown to be very robust as it not only can find  
 139 the mapping in circumstances of very noisy real-world input  
 140 but also increases the clustering accuracy of both modalities.

Recently, we tested the proposed algorithm also on the task of  
 141 clustering body parts from simultaneous tactile and linguistic  
 142 input [59]. In that case, sequential mapping showed slower  
 143 degradation with increasing noise level in the linguistic input  
 144 and outperformed one-step mapping for all dataset sizes and  
 145 all levels of noise.  
 146

The rest of this paper is structured as follows. In Section II,  
 147 we compare different mapping algorithms used in cross-  
 148 situational learning. In particular, in Section II-B we provide  
 149 a mathematical formulation of the newly proposed sequential  
 150 mapping algorithm, and in Section II-C we describe the whole  
 151 cognitive architecture which incorporates unimodal processing  
 152 of vision and language and finding their association through  
 153 the mapping algorithm. The performance of the proposed  
 154 method on data from an iCub humanoid robot and from an  
 155 iCub simulator is evaluated in Section III. Finally, results are  
 156 discussed in Section IV with suggestions for future work.  
 157

## II. MATERIALS AND METHODS

In this section, we first present one-step and newly proposed  
 159 sequential mapping algorithms (Sections II-A and II-B). Then,  
 160 we describe in detail the whole cognitive architecture used to  
 161 process data from individual modalities (vision and language).  
 162 Finally, we describe the iCub robotic platform and the iCub  
 163 simulator in Section II-D and provide a description of the  
 164 evaluation in Section II-F.  
 165

### A. One-Step Mapping in Cross-Situational Learning

One possible way how to establish mapping between visual  
 167 concepts and linguistic elements is to use frequencies of refer-  
 168 ent and meaning co-occurrences; that is, the ones with the  
 169 highest co-occurrence are mapped together [40], [41], [54].  
 170 This method is usually called cross-situational learning and  
 171 supposes the availability of the ideal associative learner  
 172 who can keep track and store all co-occurrences in all tri-  
 173 als, internally memorizing and representing the word-object  
 174 co-occurrence matrix of the input. This allows the learner  
 175 to subsequently choose the most strongly associated refer-  
 176 ent [55], [56]. These models do not see the mapping as  
 177 dynamic competition but operate only with the static state.  
 178 Although some use likelihoods of different words and refer-  
 179 ents to perform Bayesian inference [16], [54], they do not take  
 180 into account how the similarity of two-word forms can affect  
 181 learning although it has been shown that the similarity affects  
 182 learning in children [36]. Another shortcoming of these strate-  
 183 gies is that they do not address how these similarities affect  
 184 learning in a dynamic competition.  
 185

The simplest one-step word-to-referent learning algorithm  
 186 simply accumulates word-referent pairs. This can be viewed  
 187 as Hebbian learning: the connection between a word and an  
 188 object is strengthened if the pair co-occurs in a trial. To extend  
 189 this basic idea, we can enable also forgetting by introducing  
 190 the parameter  $\eta$ , which can capture the memory decay. This  
 191 so-called dumb associative-learning model was implemented  
 192 by Yu and Smith [56]. Supposing that at each trial  $t$  we observe  
 193 an object  $\mathbf{o}_t^n$  and hear a corresponding word  $\mathbf{w}_t^n$  ( $N_t$  possible  
 194 associations), we can describe the update of the strength of the  
 195 association between word model  $L(i)$  and object model  $K(j)$   
 196

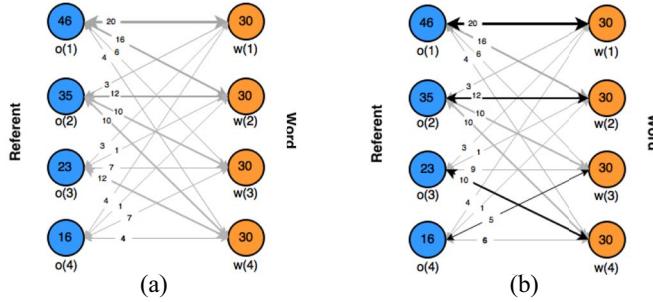


Fig. 1. One-step mapping. (a) In the first stage, weights between objects (referents) and words are changed using Hebbian learning. The connection between a word and an object is increased if the pair co-occurs in a trial [56]. (b) Then, word-to-referent mapping is found in one step: objects and words with the highest co-occurrence are mapped together. The number on each connection between word  $w(i)$  and object  $o(j)$  refers to the number of co-occurrences of  $o(j)$  and  $w(i)$ . In this example, we suppose that there are 30 occurrences of each word.

as follows:

$$A(i,j) = \sum_{t=1}^R \eta(t) \sum_{n=1}^{N_t} \delta(\mathbf{w}_t^n, i) \delta(\mathbf{o}_t^n, j) \quad (1)$$

where  $R$  is the number of trials,  $\delta$  is the Kronecker delta function (equal to 1 when both arguments are identical and 0 otherwise),  $\mathbf{w}_t^n$  and  $\mathbf{o}_t^n$  indicate the  $n$ th word-object association that the model attends to and attempts to learn in the trial  $t$ , and  $\eta(t)$  is the parameter controlling the gain of the strength of the association.

Now let us assume that word  $w(i)$  is modeled by model  $L_i$  in the language domain and object (referent)  $o(j)$  is modeled by the model  $K_{m(i)}$  in the visual domain. Our goal is to find the corresponding model  $K_{m(i)}$  from the visual subdomain for each model  $L_i$  from the language domain to assign them together. Indices  $m(i)$  are found as follows:

$$\forall i : m(i) = \operatorname{argmax}_i A(i,j) \quad (2)$$

where  $A$  is the co-occurrence matrix computed in (1) [element  $A(i,j)$  captures co-occurrence between word  $w(i)$  and object  $o(j)$ ].

In Fig. 1, one-step mapping is visualized as it was implemented in this paper.

Modifications of the basic model include Regier's [36] work. He proposed a mapping model, which stems from competition models and incorporates two-way associations between words and referents. This enables the model to capture selective attention to individual words and referents, as well as to provide a probability distribution over associated referents/words. Regier [36] also showed that for his model, learning of a novel word is most effective when memory interference is minimized.

### B. Proposed Sequential Mapping

Because we know that learning is not static but is a dynamic process, it seems reasonable to extend the basic idea of one-step cross-situational learning by incorporating dynamic competition mechanisms between words and referents in the

model. To capture the dynamic competition among models, we extend the basic one-step mapping algorithm with the sequential addition of inhibitory connections.

In this case, the process of finding word-referent associations resembles Hebbian learning with inhibitory connections. Once the word is associated with a corresponding object (referent), links from this referent to other words are inhibited. This idea also corresponds to the fact that children prefer mapping where an object has only one label to multiple labels, the so-called mutual exclusivity bias [25]. The inhibitory mechanisms and situation-time dynamics were already partially included in the model of cross-situational learning proposed by McMurray *et al.* [29].

Even though our model shares some similarities with the model proposed by McMurray *et al.* [29], our model stems from different computational mechanisms. The proposed sequential mapping is able to capture nondiscrete assignment to individual clusters, as well as dynamic competition mechanisms. The first mechanism is incorporated into the model by considering likelihoods that the observed data were generated by a given model instead of 1/0 assignment to models. In this way, similarities of individual meanings and referents, as well as the likelihood of their recognition in each trial, is taken into account. The second mechanism (dynamic competition) facilitates the sequential mapping as the best-mapped classes are gradually justified with inhibitory connections to other classes (i.e., after a reliable assignment between a language and a tactile model is found, inhibitory connections among this tactile model and all other language models are added). Thanks to this mechanism, the mutual exclusivity principle (the fact that children prefer mapping where an object has only one label to multiple labels [25]) is guaranteed. The assignment between visual models  $K_j$  and language models  $L_i$  is found using the following iterative procedure.

- Visual and language data are clustered separately and the corresponding posterior probabilities are found

$$p(L_i | \mathbf{w}_t^n) = \frac{p(\mathbf{w}_t^n | L_i)p(L_i)}{\sum_{i'} p(\mathbf{w}_t^n | L_{i'})p(L_{i'})}, \quad \forall i \in \{1, \dots, I\} \quad (267)$$

$$\forall t \in \{1, \dots, T\}, \quad \forall n \in \{1, \dots, N_t\} \quad (268)$$

$$p(K_j | \mathbf{o}_t^n) = \frac{p(\mathbf{o}_t^n | K_j)p(K_j)\mathbf{k}(j)}{\sum_{j'} p(\mathbf{o}_t^n | K_{j'})p(K_{j'})}, \quad \forall j \in \{1, \dots, J\} \quad (269)$$

$$\forall t \in \{1, \dots, T\}, \quad \forall n \in \{1, \dots, N_t\} \quad (270)$$

where  $I$  is the number of language models,  $J$  is the number of visual models,  $T$  is the number of trials, and  $N_t$  is the number of possible object-word associations in trial  $t$ .

- The most probable cluster to which the data point is assigned is selected (for each data point)

$$a_t^n = \operatorname{argmax}_i p(L_i | \mathbf{w}_t^n) \quad (277)$$

$$b_t^n = \operatorname{argmax}_j p(K_j | \mathbf{o}_t^n) \quad (278)$$

$$\forall t \in \{1, \dots, T\}, \forall n \in \{1, \dots, N_t\}. \quad (279)$$

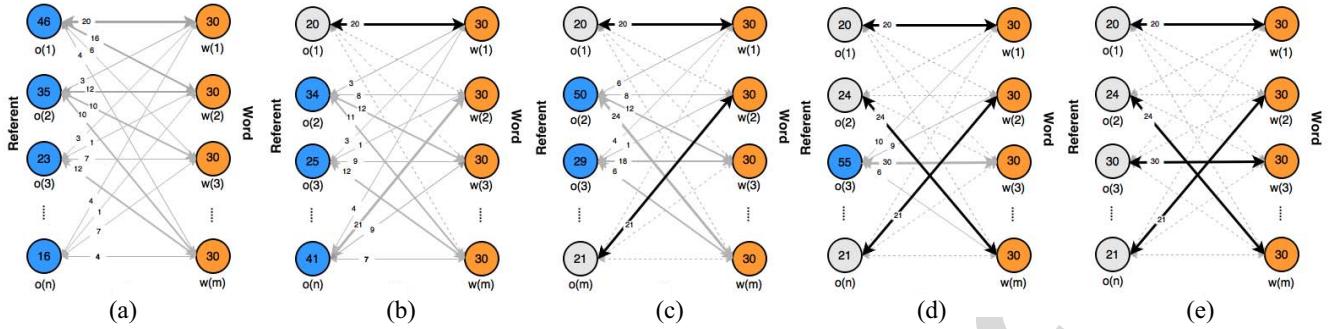


Fig. 2. Sequential mapping. The toy example of sequential mapping is shown to clarify the mechanism of finding the object-word assignment. In this example, we suppose that there are 30 occurrences of each word. The dotted line marks the inhibitory connection between object  $o(j)$  and word  $w(i)$  and the black line corresponds to the already found mapping. The number on each connection between word  $w(j)$  and object  $o(i)$  refers to the number of co-occurrences of  $o(j)$  and  $w(i)$ . Objects  $o(j)$  and words  $w(i)$  are assigned to corresponding models based on the given clustering mechanism.

280 3) Co-occurrence matrix  $A(i, j)$  is computed

$$281 \quad A(i, j) = \zeta(i, j) \sum_{t=1}^R \eta(t) \sum_{n=1}^{N_t} \delta(a_t^n, i) \delta(b_t^n, j) \quad (7)$$

282 where  $R$  is the number of trials,  $\zeta(i, j)$  is the matrix storing  
283 the strength of the connections between visual model  
284  $K_j$  and language model  $L_i$ , and  $\eta(t)$  is the parameter  
285 controlling the gain of the strength of association.

286 4) The best assignment is selected

$$287 \quad [im, m(im)] = \underset{i}{\operatorname{argmax}} \underset{j}{\operatorname{argmax}} A(i, j). \quad (8)$$

288 In this step, the visual model  $K_{m(im)}$  is assigned to the  
289 language model  $L_{im}$ .

290 5) Inhibition connections are added between the assigned  
291 visual model  $K_{m(im)}$  and all language models other than  
292  $L_{im}$  (mutual exclusivity)

$$293 \quad \zeta(i, m(im)) = \zeta(i, m(im))(1 - z_1), \quad \forall i \neq im \quad (9)$$

294 where  $z_i$  is the parameter capturing the strength of  
295 the inhibition (in our experiment, this is set to 1,  
296 which corresponds to the total inhibition of the given  
297 connection).

298 6) Inhibition is added to the assigned visual model  $K_{m(im)}$   
299 (a prior probability of the model is changed)

$$300 \quad k(m(im)) = k(m(im))(1 - z_2) \quad (10)$$

301 where  $z_2$  is the parameter capturing the inhibition of the  
302 assigned visual model (in our experiment, this parameter  
303 is set to 1, which corresponds to total inhibition of the  
304 given model).

305 7) The assigned points are deleted from the dataset (data  
306 points which belong to model  $K_{jm}$  and  $L_{im}$ )

$$307 \quad X = X \setminus \left\{ (\mathbf{o}_t^n, \mathbf{w}_t^n) \mid \underset{j}{\operatorname{argmax}} p(K_j | \mathbf{o}_t^n) == m(im) \right. \quad (11)$$

$$308 \quad \left. \wedge \underset{i}{\operatorname{argmax}} p(L_i | \mathbf{w}_j^n) == im \right\}. \quad (12)$$

309 8) If the dataset is not empty ( $X \neq \emptyset$ ) or  $\|\mathbf{k}\| > 0$  (some of  
310 the visual models are not totally inhibited), go to step 1,  
311 else stop.

312 The proposed algorithm where words are assigned to  
313 corresponding referents in a sequential manner is visualized  
314 in Fig. 2.

315 In the ideal case, unambiguous mapping between the two  
316 clusterizations will be found. In the real case (where the clus-  
317 terization in visual and language layers is not optimal), none  
318 or more than one model from the visual layer will be assigned  
319 to one cluster  $L_i$  in the language layer or vice versa.

### C. Specific Architecture

321 Our multimodal hierarchical architecture consists of  
322 multimodal and unimodal parts. The unimodal part has  
323 two layers performing separate processing of localist inputs:  
324 1) visual objects and 2) auditory word-forms. Both unimodal  
325 layers are subsequently mapped one to each other in the upper  
326 multimodal layer (see Fig. 3).

327 1) *Visual Layer:* Each data point (object  $\mathbf{o}_t^n$ ) can be consid-  
328 ered as a triplet of continuous-valued vectors for each visual  
329 feature:  $\mathbf{o}_t^n = (\mathbf{x}_{t,n}^{\text{size}}, \mathbf{x}_{t,n}^{\text{color}}, \mathbf{x}_{t,n}^{\text{shape}})$ . This enables us to write the  
330 visual dataset as  $X^{\text{vis}} = [X^{\text{size}} \ X^{\text{color}} \ X^{\text{shape}}]$  and process data for  
331 each visual feature separately. For processing visual data, the  
332 Gaussian mixture model (GMM) was used, which is a convex  
333 mixture of  $d$ -dimensional Gaussian densities  $l(\mathbf{x}^k | \theta_j^k)$ , where  
334  $k \in \{\text{size, color, shape}\}$ . In this case, each visual model  $K_j^k$  is  
335 described by a set of parameters  $\theta_j^k$ . The posterior probabilities  
336  $p(\theta_j^k | \mathbf{x}^k)$  are computed as follows:

$$337 \quad p(\theta_j^k | \mathbf{x}^k) = \sum_{j=1}^{J_k} r_j^k l(\mathbf{x}^k | \theta_j^k) \quad (13)$$

$$338 \quad l(\mathbf{x}^k | \theta_j^k) = \frac{1}{\sqrt{(2\pi)^d} \sqrt{|\mathbf{S}_j^k|}} \exp \left[ -\frac{1}{2} (\mathbf{x}^k - \mathbf{m}_j^k)^T (\mathbf{S}_j^k)^{-1} \right. \\ \left. \times (\mathbf{x}^k - \mathbf{m}_j^k) \right] \quad (14)$$

340 where  $k \in \{\text{size, color, shape}\}$ ,  $\mathbf{x}^k$  is a set of  $d$ -dimensional  
341 continuous-valued data vectors,  $r_j^k$  are the mixture weights,  
342  $J_k$  is the number of visual models for each visual feature  
343  $k$ , and parameters  $\theta_j^k$  are cluster centers  $\mathbf{m}_j^k$  and covariance  
344 matrices  $\mathbf{S}_j^k$ .

345 The Gaussian mixture is trained by the expectation-  
346 maximization algorithm [10]. An output of this layer for

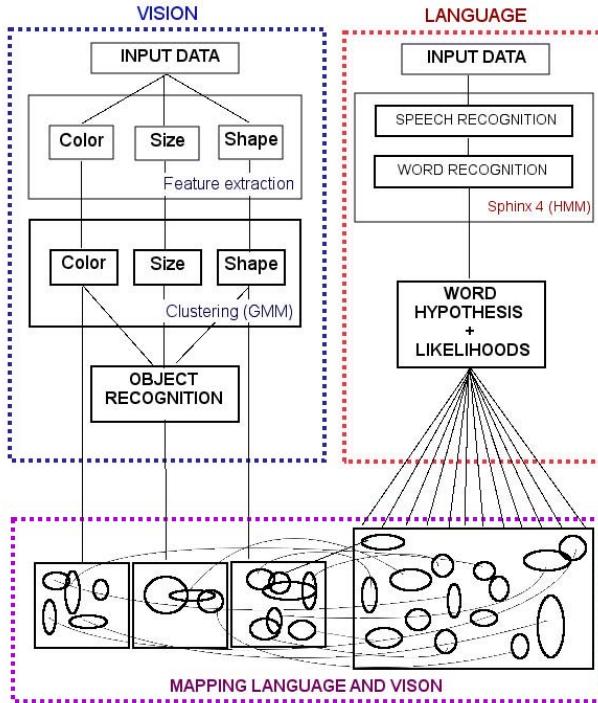


Fig. 3. Proposed multimodal architecture.

each data point  $\mathbf{x}_i^k$  is the vector  $\mathbf{y}_i^k$  of  $J_k$  output parameters describing the data point (the likelihood that the data point belongs to each individual cluster in a mixture). This corresponds to the fuzzy memberships (distributed representation). For a simpler evaluation, we used a localist representation (winner-takes-all), where only the cluster with the highest cluster membership probability is considered for further processing [see (5) and (6)]

$$M(K_j^k | O) = \begin{cases} 1 & \text{if } j = \operatorname{argmax}_j f(K_j^k | O) \\ 0 & \text{if } j \neq \operatorname{argmax}_j f(K_j^k | O) \end{cases} \quad (15)$$

where  $k \in \{\text{size, color, shape}\}$ ,  $j \in \{1, \dots, J_k\}$ .

2) *Language Layer*: The auditory word-forms are extracted from the language input which are sentences describing the image in the format <size> <color> <shape> (e.g., “small red triangle”). Afterward, individual word-forms are extracted from the sentences and compared to prelearned language models, and the log-scale score  $p(\mathbf{w}_t^n | L_i)$  of the audio matching the model is computed. Based on these data, posterior probability can be computed

$$p(L_i | \mathbf{w}_t^n) = \frac{p(\mathbf{w}_t^n | L_i)p(L_i)}{\sum_{i'} p(\mathbf{w}_t^n | L_{i'})p(L_{i'})}, \quad \forall i \in \{1, \dots, I\} \quad \forall t \in \{1, \dots, T\}, \quad \forall n \in \{1, \dots, N_t\} \quad (16)$$

where  $I$  is the number of language models,  $T$  is the number of trials (sentences) and  $N_t$  is the number of word-forms in the trial (sentence)  $t$ .

An output of this layer for each data point  $\mathbf{w}_t^n$  is the vector  $\mathbf{y}_i^k$  of  $I$  output parameters describing the data point (the likelihood that the data point belongs to each individual language model). This corresponds to the fuzzy memberships

---

**Algorithm 1** Sequential Mapping—Fixed Grammar

---

**Inputs:**

language clusters  $L_i$  ( $i \in 1 : I$ ), visual clusters  $K_j^k \sim N(\mathbf{m}_j^k, \mathbf{S}_j^k)$ ,  $j \in 1 : J_k$ , input data  $\mathbf{x}^k$  for each feature  $k \in \{\text{size, color, shape}\}$ , number of clusters  $J^k$  for each feature  $k$

**Output:**

mapping between all visual classes  $K_j^k$  and language classes  $L_i$

```

for  $k \in \{\text{size, color, shape}\}$  do
     $NCl \leftarrow J^k$ 
    while  $NCl > 0$  and  $\mathbf{X}^k$  is not empty do
        assign each data point from  $\mathbf{x}^k$  to visual and language cluster (Winner-takes-all, see Eq. (15))
        for  $j = 1 : NCl$  do
            for  $i = 1 : I$  do
                 $A_{ij} \leftarrow$  how many times was class  $i$  classified as  $j$ 
            end for
        end for
         $[im, jm] \leftarrow \operatorname{argmax}_i \operatorname{argmax}_j A_{ij}$ 
         $\mathbf{X}_{del}^k \leftarrow$  data points assigned to  $K_{jm}^k$  and  $L_{im}$ 
         $\Theta_{new}^k \leftarrow N(\mathbf{x}_{del}^k)$  learn Gaussian on the to be deleted data
         $\mathbf{X}^k \leftarrow \mathbf{X}^k \setminus \mathbf{X}_{del}^k$  delete data assigned to both  $K_{jm}^k$  and  $L_{im}$ 
         $NCl \leftarrow NCl - 1$ 
        relearn  $K^k \sim N(\mathbf{m}^k, \mathbf{S}^k)$  on new data  $\mathbf{x}^k$  with  $NCl$  clusters
    end while
end for
cluster visual data using new  $\theta_{new}^k$  parameters (cluster centers  $\mathbf{m}^k$ , covariance matrices  $\mathbf{S}^k$ ) and perform one-step mapping (Model 1)

```

---

(distributed representation). Linguistic and visual inputs are processed simultaneously.

3) *Mapping—Models 1 and 2*: After the visual and language data are clustered, the mapping between the two layers must be found. For each cluster  $L_i$  in the language layer, a corresponding cluster  $K_j^k$  in the visual layer (for each feature  $k \in \{\text{size, color, shape}\}$ ) is found. The mapping is found as follows. For each  $j$  and  $k$ , we find cluster  $L_{kmax_{jk}}$  from the language layer which will be assigned to cluster  $K_j^k$  from the visual layer. In this paper, we compare two different models to find indices  $kmax_{jk}$ . We compared one-step mapping (see Section II-A) and newly proposed sequential mapping (see Section II-B).

The exact algorithm used to find the mapping between visual and language models in a sequential manner is described in detail in Algorithm 1. Indices  $m(i)$  are found sequentially. In each step, the best-mapped data are excluded, and the rest of data are reclustered using GMMs. Then, one-step mapping is performed (see Algorithm 1). An extension of the algorithm for a variable-length sentence is described in the Appendix. Results for a variable length sentence using a fully artificial dataset with a controlled noise level are described in detail in [58].

#### D. iCub Robotic Platform and iCub Simulator

For the experiment, we used a simulated [48] and a physical stationary [30] iCub robot. The iCub [Fig. 1(c)] is an open-source humanoid robot the size of a three-and-a-half-year-old child, with the fully articulated hands and a head-and-eye system which makes him ideal for cognitive experiments. The iCub simulator has been designed to reproduce, as accurately as possible, the physics and dynamics of the robot and its environment [48]. The simulator and the actual robot have the same interface supporting YARP [31] which is a robot platform for interprocess communication and control of the physical and simulated robots in real-time.

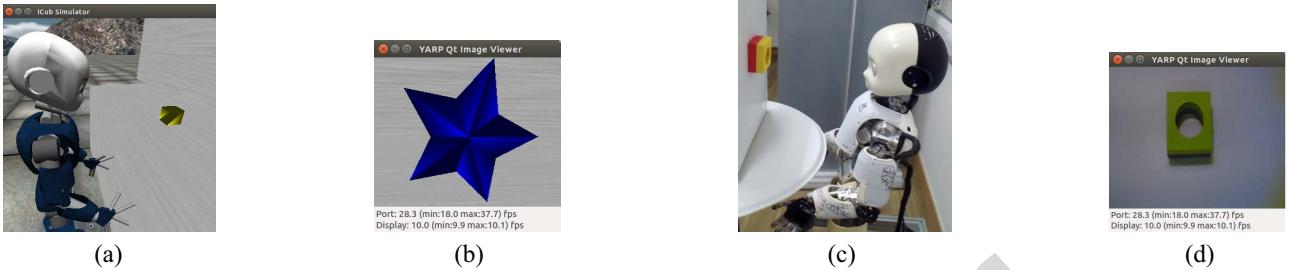


Fig. 4. Experiment design and corresponding input data. (a) iCub simulator. (b) Blender object. (c) Physical iCub. (d) Real object.

#### 409 E. Input Data Description and Preprocessing

410 The input for our model consisted of visual and language  
 411 data. The visual scene was composed of an object in the cen-  
 412 ter of the scene with a variable position. The visual features  
 413 (size, shape, and color) of the object also varied. We devel-  
 414 oped two separate datasets for training and testing purposes.  
 415 A real-world dataset has visual sensory data acquired from  
 416 the cameras of the physical iCub robot which observed sim-  
 417 ple objects placed on the white board in front of his eyes  
 418 at a distance of 1 m [see Fig. 4(c) and (d)] (210 instances,  
 419 three sizes, five colors, and seven shapes, in total 70 indi-  
 420 vidual objects; each seen three times in slightly different  
 421 placement). The simulated dataset is made in the iCub sim-  
 422 ulator [see Fig. 4(a) and (b)] as Blender-generated virtual  
 423 objects (432 instances, three sizes, six colors, and six shapes,  
 424 in total 108 individual objects; each seen four times in slightly  
 425 different placement).

426 The spoken language input consisted of sentences pro-  
 427 nounced by a non-native English speaker describing the image  
 428 in the format <size> <color> <shape> (e.g., small red trian-  
 429 gle) and was processed simultaneously with the visual input. In  
 430 the real-world dataset, the tutor was talking to a robot from an  
 431 approximate 2-m distance with natural background noise. The  
 432 linguistic input was captured using an external microphone.

433 1) *Speech Recognition*: CMU Sphinx (an open-source flex-  
 434 ible Markov model-based speech recognizer system) was used  
 435 for speech recognition [23]. Sphinx itself offers a large vocab-  
 436 ular, but we created our own task-specific smaller vocabulary  
 437 using the online IMtool that produces a dictionary based on a  
 438 CMU dictionary and matches its language model.

439 There is a probabilistic output from the CMU Sphinx. The  
 440 ten best hypotheses for a matching model with correspond-  
 441 ing scores were saved for each utterance (they are log-scale  
 442 scores of the audio matching the model). Because the scores  
 443 for the hypothesis of each word in the sentence were needed  
 444 for further evaluation, the words were pronounced with large  
 445 pauses, and the end of the sentence was marked by the word  
 446 “STOP.” An output of the language layer is a  $I$ -dimensional  
 447 continuous valued vector, where  $I$  is the number of language  
 448 clusters (corresponding to the number of possible utterances).

449 This vector contains ten nonzero values, and the rest is zero.  
 450 2) *Image Processing*: The image inputs are processed using  
 451 standard MATLAB functions. First, the image is morphologi-  
 452 cally opened with a disk-shaped structuring element (*imopen*)  
 453 to remove the noisy background of the image; then all grayish



Fig. 5. Image processing: original image, removal of the background, converting to BW image, and filling the holes.

pixels are removed, and the image is converted from the true 454  
 455 color RGB to the gray-scale intensity image by eliminating 456  
 457 the hue and saturation information while retaining the lumi- 458  
 459 nance (*rgb2gray*). Finally, the intensity image is converted to 460  
 461 a binary image using the threshold computed with Otsu’s [34] 462  
 463 method (*threshold*). An example of the preprocessed image is 464  
 465 shown in Fig. 5.

Afterward, the properties of the image regions are measured 466  
 467 using the function *regionprops*. Individual visual features 468  
 469 (shape, color, and size) are subsequently processed sepa- 470  
 471 rately. The following features were used: color (three features: 472  
 473 average RGB of the selected region), size (six features: 474  
 475 perimeter of an object, distance from the centroid to the left 476  
 477 corner of the bounding box, and width and length of the 478  
 479 bounding box), and shape (13 features: area, centroid, major 480  
 481 axis length, eccentricity, orientation, convexArea, FilledArea, 482  
 483 EulerNumber, EquivDiameter, solidity, extent, and perimeter). 484  
 To obtain the shape features, we automatically cropped and 485  
 486 resized the image to equalize the size of the objects.

Although the visual model is mainly mathematical and 487  
 488 implemented in a very “machine vision” sort of way, the 489  
 490 bases of its processing follow biological correlates of mammal 491  
 492 vision. More specifically, from the neuroanatomical point of 493  
 494 view, this corresponds to the processing of the visual input in 495  
 496 the separate higher visual centers in the brain, specifically to 497  
 498 the independent processing of the information about the posi- 499  
 500 tion and identification of an object in the ventral (what) and 501  
 502 dorsal (where) neural pathways, respectively [32]. Individual 503  
 504 object properties are identified in the separate visual centers 505  
 506 of the occipital lobe.

#### 484 F. Evaluation

In the case of the supervised GMM, growing when required 485  
 486 neural gas (GWR), and SOM algorithms, data were divided 487  
 488 into training and validation datasets in the ratio 70:30. For the 489  
 490 unsupervised GMM, hidden Markov model, and  $k$ -means algo- 491  
 492 rithms, we computed the accuracy differently. After the data 493  
 494

TABLE I  
COMPARISON OF CLUSTERIZATION AND CLASSIFICATION ACCURACY OF VISUAL DATA. THE MEAN AND STANDARD DEVIATION FROM 100 REPETITIONS ARE VISUALIZED

Accuracy [%]	Real data			Blender		
	Size	Color	Shape	Size	Color	Shape
GMM sup.	83.3 ± 0.0	99.0 ± 0.0	81.4 ± 0.0	98.6 ± 0.0	97.9 ± 0.0	93.1 ± 0.0
GMM unsup.	76.2 ± 6.8	76.1 ± 9.1	56.1 ± 6.2	74.2 ± 10.1	60.9 ± 9.0	64.3 ± 7.2
K-means	67.8 ± 6.2	81.2 ± 1.1	53.1 ± 4.2	66.3 ± 0.2	77.1 ± 10.7	72.8 ± 6.9
SOM	69.6 ± 5.6	78.9 ± 6.8	54.2 ± 4.1	66.1 ± 4.2	81.7 ± 7.6	59.3 ± 6.2
GWR	89.9 ± 2.1	99.5 ± 0.4	76.6 ± 1.4	88.9 ± 0.7	98.1 ± 0.9	94.2 ± 0.6

TABLE II  
COMPARISON OF ONE-STEP MAPPING AND SEQUENTIAL MAPPING FOR DATA FROM THE iCUB SIMULATOR (BLENDER) AND THE PHYSICAL iCUB (REAL DATA). THE MEAN AND STANDARD DEVIATION FROM 100 REPETITIONS ARE VISUALIZED

Accuracy [%]	Real data			Blender		
	Size	Color	Shape	Size	Color	Shape
Vision	76.2 ± 6.8	76.1 ± 9.1	56.1 ± 6.2	74.2 ± 10.1	60.9 ± 9.0	64.3 ± 7.6
Language	70.6 ± 0.0	82.4 ± 0.0	77.5 ± 0.0	98.1 ± 0.0	96.5 ± 0.0	98.1 ± 0.0
One-step mapping	54.1 ± 4.1	58.2 ± 10.3	52.2 ± 4.9	67.3 ± 8.2	56.2 ± 6.1	61.9 ± 3.2
Sequential mapping	74.2 ± 15.1	87.1 ± 10.2	72.9 ± 5.1	96.1 ± 31.2	95.2 ± 1.2	92.1 ± 0.9

490 is clustered, each cluster is assigned to the class that appears  
491 most frequently in the cluster, and then the accuracy of this  
492 assignment is measured by counting the number of correctly  
493 assigned data points (compared to the manual true labels)  
494 and dividing this number by the total number of data points.

495 The accuracy of the learned mapping is calculated in the fol-  
496 lowing manner. We cluster output activations from the visual  
497 layer and assign each data point to the most probable cluster.  
498 Then, we find indices  $m(i)$  for all clusters as defined in (2) for  
499 one-step mapping and (8) for sequential mapping. Based on  
500 this mapping, we can assign each data point to the language  
501 label. These language labels are subsequently compared to the  
502 ground truth. Accuracy is then computed as

$$503 \quad \text{acc} = \frac{\text{TP}}{N} \quad (17)$$

504 where TP (true positive) is the number of correctly assigned  
505 data points, and  $N$  is the number of all data points.

### 506 III. EXPERIMENTAL RESULTS

507 The first part of the results is dedicated to the performance  
508 of the model in the real-world scenario. The robot interacts  
509 with a human in a noisy condition that distorts speech input.

#### 510 A. Vision

511 In the first stage, we evaluated the Vision subpart of the  
512 model. Several algorithms were compared: namely the GMM  
513 algorithm, supervised GMM algorithm,  $k$ -means, SOM, and  
514 GWR algorithm [2], [26]. The SOM and GWR had 100 nodes.  
515 The results for the real-world dataset and the simulated dataset  
516 with Blender objects can be seen in Table I. Although the SOM  
517 and GWR algorithms are considered unsupervised algorithms,  
518 we adopted a technique for labeling inputs; thus, they should  
519 be compared with supervised algorithms. The number of clus-  
520 ters is also overestimated (the number of nodes corresponds  
521 to the number of clusters). It indicates that these algorithms

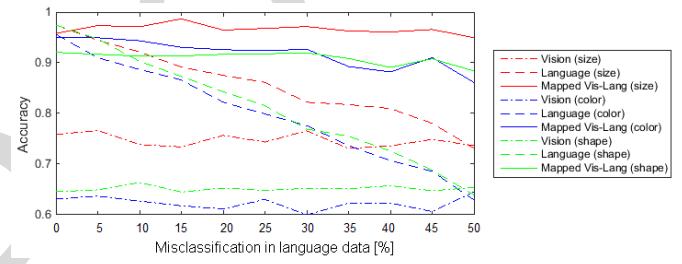


Fig. 6. Dependence of mapping accuracy on the misclassification in the language data for a fixed-length sentence (mean values over 50 repetitions are visualized). Different colors correspond to different visual features (red: size, blue: color, and green: shape). Visual data are generated in Blender and acquired through the iCub simulator, and language data are processed using Sphinx 4.

522 partly overfit the data so we divided the set into testing and validation data. 523

#### 524 B. Mapping

525 The performance of one-step mapping (vision and lan-  
526 guage are mapped in one step based on the frequency of  
527 co-occurrence) and sequential mapping (see Algorithm 1)  
528 is shown in Table II. We calculated the accuracy for  
529 real-world data from physical iCub and for the Blender  
530 objects placed in the iCub simulator. Language accuracy  
531 for Blender dataset is much higher compared to the real-  
532 world data as a tutor was speaking directly into the  
533 microphone.

534 Tolerance of the sequential mapping to misclassification  
535 in the language data is visualized in Fig. 6 for visual data  
536 from the iCub simulator in combination with the language  
537 data processed by Sphinx 4. The misclassification is added to  
538 the language data subsequently and evenly to all classes (the  
539 given proportion of the language inputs was randomly changed  
540 to random words). The misclassification (a synthetic error)  
541 was added artificially but can be interpreted either as white  
542 noise added to the input data or mistakes in labeling perceived

543 objects. We grouped them together into a misclassification  
 544 variable. The visual data are left intact so the only cause of  
 545 the observed variations in the accuracy is initialization. As can  
 546 be seen, the accuracy of sequential mapping remains very sta-  
 547 ble even though the accuracy of the language decreases and  
 548 outperforms the language and vision for almost all values of  
 549 the misclassification.

#### 550 IV. CONCLUSION

551 Current models of vision-to-language mapping often make  
 552 use of cross-situational learning while relying directly on sta-  
 553 tistical co-occurrence of meaning-referent pairs: the ones with  
 554 the highest co-occurrence are mapped together (e.g., [41]).  
 555 This approach shows inferior performance in cases where  
 556 the clustering of data is difficult (e.g., due to the high num-  
 557 ber of objects, high noise level or overlapping classes, etc.).  
 558 Therefore, we extended this basic model and introduced a  
 559 new more robust and noise-resistant mapping procedure. Our  
 560 approach incorporates situation-time dynamics and mutual  
 561 exclusivity and is able to deal with a nonequal number of  
 562 classes in individual subdomains.

563 Mathematical formulation of the newly proposed mapping  
 564 is provided (see Section II-B), and results for simulated and  
 565 real-world data from the iCub robot are compared to one-  
 566 step mapping (see Table II). It was shown that the method is  
 567 able to find mapping between language and vision, and the  
 568 method improves the accuracy of both individual subdomains  
 569 and shows very good resistance to noise or misclassification in  
 570 language (see Fig. 6). How to map in an unsupervised manner  
 571 several clusterings (e.g., for vision, action, and language) is an  
 572 important question not only in cognitive modeling but also in  
 573 general machine learning, where data acquired from different  
 574 sensors or in different situations can be independently clus-  
 575 tered and mapped to each other. A more detailed discussion  
 576 of the results follows.

577 The trivial one-step mapping can be imagined as basic  
 578 Hebbian learning. Our extension can be likened to Hebbian  
 579 learning with inhibitory connections. In recent years, sev-  
 580 eral approaches finding an alternative to the basic approach  
 581 appeared [29], [57]. Our model partially stems from the  
 582 McMurray *et al.* [29] approach, who showed that associative  
 583 learning can be sufficient for language acquisition and that the  
 584 main components of this type of learning are an online com-  
 585 petition of models and pruning incorrect associations which  
 586 enable gradual improvement of associations between models.  
 587

588 The mutual exclusivity principle is guaranteed in our  
 589 method thanks to the inhibitory connections which are grad-  
 590 ually created among models. Once the mapping between the  
 591 referent and the meaning is found, the connection from a given  
 592 meaning to other referents is inhibited. Dynamic competition  
 593 is addressed in the following way: when any association of  
 594 meaning and referent is found, other models compete again  
 595 for resources. First, well-mapped data are deleted, and then,  
 596 the resting data are reclustered. Furthermore, the likelihoods  
 597 can associate each data point to many separate models instead  
 598 of binary membership. In this way, our algorithm also emulates  
 599 how the similarity of two-word forms can affect learning as  
 it occurs in children's development [36]. The algorithm also

enables mapping together data in a case where we have an  
 600 uneven number of clusters in both subdomains. 601

602 However, the principle of mutual exclusivity is not suitable  
 603 for further stages of language acquisition, namely, the learning  
 604 of polysemic words. A polyseme is a word or phrase with  
 605 different but related senses (e.g., wood as a piece of tree or the  
 606 area with trees). The homonyms are a subset of the polysemes,  
 607 but the difference between the homonyms and the polysemes  
 608 is subtle and fuzzy. Homonyms represent a group of words  
 609 that share similar spelling (homographs) and the same sound  
 610 (homophones) but have different and unrelated meanings, for  
 611 example, the homograph bank standing for embankment or  
 612 place where money is kept. The learning of polysemic words  
 613 violates the principle of mutual exclusivity as the dynamic  
 614 competition does not allow to map one word to different visual  
 615 models. We are aware of this problem, and we would like to  
 616 extend it in the future iteration of the model. Similar to humans  
 617 who have to learn polysemic words as an exception, we will  
 618 incorporate this principle into the next version of our model. 619

In the next part, we analyze the ability of our architec-  
 620 ture to deal with ambiguous inputs. The mapping will find  
 621 reliable labeling for the visual input data (more generally for  
 622 data from any other modality) with a possibility of incorporat-  
 623 ing the fuzziness of this mapping. For some concepts, finding  
 624 an unambiguous mapping is very easy; for others, it is much  
 625 more difficult or impossible (such as abstract words, e.g., the  
 626 love has no dominant color, but the sky is usually blue). Since  
 627 the mapping is established only among the clusters where it  
 628 makes sense, dealing with a lot of redundant information is  
 629 avoided. A similar idea is used in classification algorithms  
 630 which use sparse matrices (e.g., [12] and [13]). We also ana-  
 631 lyzed the strong and weak points of the algorithms adopted  
 632 in our architecture. First, the ability of different algorithms  
 633 to classify unimodal visual data was compared. As expected,  
 634 for data which are well separated and mainly spherically  
 635 distributed (this is generally a case for simulated and artifi-  
 636 cially generated data), the  $k$ -means algorithm outperformed  
 637 the GMM algorithm. However, for nonspherical real data, the  
 638 GMM algorithm generally performed better (see Table I for  
 639 the comparison of the performance on simulated data placed  
 640 in an iCub simulator and data from real iCub robot cameras).  
 641 As the unsupervised algorithms are highly dependent on the  
 642 initialization, it can be seen that the standard deviation of the  
 643 data is quite high even though 20 repetitions were averaged.  
 644 We should conclude that the performance of the algorithms in  
 645 our tasks reflects their fundamental advances and limitations.  
 646 We are still missing the algorithm that is able to cope with  
 647 highly variable datasets in terms of their statistical properties.  
 648

The most significant differences between data acquired from  
 649 the real and the simulated iCub can be found in the visual sub-  
 650 domain. Affected by natural light and slightly altered points  
 651 of view, the shading of objects and their bevel differed. This  
 652 resulted in a different projection of objects to iCub cameras  
 653 and in the less accurate classification of shapes for the real-  
 654 world dataset. For colors, different lighting changed the true  
 655 color of the objects. Nevertheless, as can be seen in Table II,  
 656 this did not affect the ability of the algorithm to discrimi-  
 657 nate among individual colors. The real-world objects had a  
 658

658 variable size which resulted in a higher overlap of clusters for  
 659 individual sizes compared to the Blender dataset.

660 Afterward, we focused on the mapping between vision and  
 661 language and compared two different approaches: one-step  
 662 mapping to sequential mapping which in a stepwise man-  
 663 ner finds the best-mapped clusters while constantly relearning  
 664 clusterization of visual data. As can be seen in the results in  
 665 Table II, the novel sequential mapping led to an improvement  
 666 in effectiveness compared to the method which maps vision  
 667 to language in one step. This can be seen in more accurate  
 668 mapping which leads to a better estimation of the clustered  
 669 data labels and consequently to lower classification error for  
 670 all of the evaluated datasets.

671 The accuracy of multimodal mapping outperforms vision,  
 672 language or both. This is an important finding, because  
 673 sequential mapping improves not only the accuracy of visual  
 674 clustering but can also fix mistakes in language recognition,  
 675 which provides the labels. Furthermore, this result suggests  
 676 that we are able not only to find mapping between more clus-  
 677 terings, but we can also improve the clusterization accuracy  
 678 by combining individual classifiers. This is not easily seen  
 679 in the presented dataset as there is high recognition accu-  
 680 racy of Sphinx software (especially in the case of sentences  
 681 recorded for the simulated dataset). Therefore, we also tested  
 682 whether the misclassification (or noise) in the language data  
 683 affects the correct mapping between vision and language. The  
 684 result can be seen in Fig. 6. Although a synthetic error (mis-  
 685 classification) is added to the language data, the ability to  
 686 find the mapping remains nearly intact. The mapping accu-  
 687 racy decreases only very slightly and remains at around 90%  
 688 because the accuracy of language recognition drops from the  
 689 original approximately 95% to approximately 70% (depend-  
 690 ing on the specific visual feature). In the future, we would  
 691 like to explore the effect of noise in individual modalities on  
 692 classification accuracy in more detail by adding a controlled  
 693 amount of noise directly to the input. The error propagation  
 694 from individual input modalities to the multimodal layer was  
 695 explored, for example, in [37]. They showed that the visual  
 696 context can steer speech recognition and vice versa.

697 We also analyzed how the complexity of environment  
 698 affects the accuracy of our architecture. We suppose that  
 699 the performance of one-step mapping will decrease with the  
 700 increasing complexity of the task (more clusters, higher over-  
 701 lap, and higher dimensionality) as it is more difficult to find  
 702 reliable clustering of the data. This hypothesis is supported by  
 703 our preliminary results for clustering body parts from simu-  
 704 taneous tactile and linguistic input [59] and by the results  
 705 presented in this paper in Table II. The quality of one-step  
 706 mapping correlates with the quality of the visual data clus-  
 707 tering. For Blender data, the worst performance was achieved  
 708 for mapping words to the visual feature shape ( $52 \pm 5\%$ )  
 709 and for physical iCub for the feature color ( $62 \pm 3\%$ ). The  
 710 feature shape has the highest number of clusters (10), and  
 711 the feature color has the second highest number of clusters  
 712 (9) and the highest overlap of the clusters for the physical  
 713 iCub. Real-world tasks are much more complex, and we can  
 714 expect tens of different object shapes. In that case, the clus-  
 715 tering performance is crucial, and one-step mapping would

not be able to provide reliable mapping. It can be seen from  
 716 the results that the proposed mapping which enables grad-  
 717 ual re-estimation of the model parameters and works in a  
 718 dynamic fashion achieves much higher accuracy even for the  
 719 cases where one-step mapping fails. We suppose that the  
 720 mapping accuracy of the proposed method decreases more  
 721 slowly with the decreasing accuracy of clustering of individual  
 722 modalities. Unfortunately, this factor was not studied in our  
 723 restricted scenario, but the preliminary results on mapping tac-  
 724 tile and linguistic input [59] support this hypothesis. We plan  
 725 to investigate this phenomenon more deeply in future research.  
 726

The language dataset differs considerably from natural lan-  
 727 guage. However, the dataset reflects some characteristics from  
 728 the findings of Werker and McLeod [53] as infant-directed  
 729 words are usually kept short with large pauses between words.  
 730 Moreover, Brent and Siskind [3] showed that frequency of  
 731 exposure to a word in isolation predicts better whether that  
 732 word will be learned than the total frequency of exposure to  
 733 that word. In addition, Snow [42] found that mother's speech  
 734 to two-year-old is much simpler and less redundant than their  
 735 speech to ten-year-old. This indicates that young children have  
 736 available a sample of speech which is simpler, more redundant  
 737 and less confusing than normal adult speech.  
 738

The proposed algorithm was tested on language to vision  
 739 mapping and on language to tactile mapping [59], and it  
 740 can be easily extended to language to any other modalities  
 741 mapping. The mapping between multiple modalities and  
 742 words was researched in [47]. Fazly *et al.* [14] proposed a  
 743 probabilistic model of cross-situational learning where they  
 744 considered sentences that contain objects and their motion.  
 745 Monaghan *et al.* [33] studied differences between cross-  
 746 situational learning of nouns and verbs on human participants  
 747 as an extension of Tomasello and Akhtar's [50] work and  
 748 Schwartz and Terrell's [39] work. Monaghan *et al.* [33] con-  
 749 sidered learning of verbs as difficult as learning of nouns  
 750 when presented in a syntactic context. The authors noticed  
 751 that nouns are learned quicker, but verbs and nouns can be  
 752 acquired simultaneously.  
 753

An important direction for extensions of the proposed model  
 754 deals with the grounding of language related not just to static  
 755 sequences but also to events which require a temporal axis. For  
 756 example, Crick and Scassellati [7] studied the interconnection  
 757 between verbal narratives and episodes of intentional relative  
 758 motion with the goal that the robot can learn from observing  
 759 a game the rules of the game, relationships between players  
 760 and their goals and intentions and then participate in the play.  
 761

The main goal of this paper is to analyze mapping between  
 762 modalities. Thus, we kept processing of the individual modalities  
 763 deliberately simple for better understanding of cross-  
 764 situational learning. We performed further experiments with an  
 765 artificially generated dataset where the visual features, as well  
 766 as noise in the linguistic and the visual domain, were under  
 767 full control. These experiments are described in detail in [58].  
 768

## APPENDIX

769  
 Here we describe an extension of the proposed algorithm  
 770 presented in Section II-B for a case where we have sen-  
 771 tences with variable structure. This means that we cannot  
 772

**Algorithm 2** Sequential Mapping—Variable Sentence

**Inputs:**  
language clusters  $L_i$  ( $i \in 1 : I$ ), visual clusters  $K_j^k \sim N(\mathbf{m}^k, \mathbf{S}^k)$ ,  
 $j \in 1 : J^k$  for each feature  $k$ , visual input  $\mathbf{x}_t = \{\mathbf{x}_t^1, \dots, \mathbf{x}_t^K\}$  and  
corresponding language data  $W_t = \{\mathbf{w}_t^1, \dots, \mathbf{w}_t^{N_t}\}$  for each trial  $t$   
and number of clusters  $J^k$

**Output:**  
mapping between all visual classes  $K_j^k$  and language classes  $L_i$

---

```

while  $\sum J^k > 0$  and  $\mathbf{X}$  is not empty do
     $l_t^n \leftarrow$  assign each word  $\mathbf{w}_t^n$  from each sentence  $t$  to a language cluster (Winner-takes-all, see Eq. (15),  $l_t^n = \text{argmax}_i(P(w_t^n | L_i))$ )
    for  $k \in \{\text{size, color, orientation, texture, shape}\}$  do
         $v_t^k \leftarrow$  assign each data point  $\mathbf{x}_t^k$  to a visual cluster (Winner-takes-all, see Eq. (15),  $v_t^k = \text{argmax}_j(P(\mathbf{x}_t^k | K_j^k))$ )
        for  $j = 1 : J^k$  do
            for  $i = 1 : I$  do
                 $T_{ij}^k \leftarrow$  how many times did visual class  $i$  co-occur with language class  $j$  ( $T_{ij}^k = \sum_{t: l_t^n=j} \sum_n \delta(l_t^n, i)$ , where  $\delta$  is Kronecker delta)
            end for
        end for
        end for
         $[km, im, jm] \leftarrow \text{argmax}_k \text{argmax}_i \text{argmax}_j T_{ij}^k$  (the visual cluster  $K_{jm}^{km}$  is mapped to the language cluster  $L_{im}$ )
         $\mathbf{X}_{del}^{km} \leftarrow$  data points assigned to  $K_{jm}^{km}$  and  $L_{im}$ 
         $\theta_{new,jk}^k \leftarrow N(\mathbf{x}_{del}^k)$  learn Gaussian on the to be deleted data
         $\mathbf{X}^{km} \leftarrow \mathbf{X}^{km} \setminus \mathbf{X}_{del}^{km}$  delete data assigned to  $K_{jm}^{km}$  and  $L_{im}$ 
         $J^{km} \leftarrow J^{km} - 1$ 
        relearn  $K_j^{km} \sim N(\mathbf{m}^k, \mathbf{S}^k)$  on new data  $\mathbf{x}^{km}$  with  $J^{km}$  clusters
    end while
cluster visual data using new  $\theta_{new}^k$  parameters (cluster centers  $\mathbf{m}^k$ , covariance matrices  $\mathbf{S}^k$ ) and perform one-step mapping (Model 1)

```

---

773 directly associate words from the sentence with individual  
774 visual features, and therefore, associations to all visual features  
775 must be taken into account. The whole algorithm is described  
776 in Algorithm 2 in the form of a pseudocode.

777 Algorithm 2 can be further extended when we incorpo-  
778 rate the language model and corresponding probabilities of  
779 sequences of individual visual features (see [58]).

780

## REFERENCES

- 781 [1] L. W. Barsalou, "Grounded cognition," *Annu. Rev. Psychol.*, vol. 59,  
782 pp. 617–645, Jan. 2008.
- 783 [2] F. B. Klein. (2016). *GWR and GNG Classifier—File Exchange—MATLAB Central*. [Online]. Available:  
784 <http://uk.mathworks.com/matlabcentral/fileexchange/57798-gwr-and-gng-classifier>
- 785 [3] M. R. Brent and J. M. Siskind, "The role of exposure to isolated words in  
786 early vocabulary development," *Cognition*, vol. 81, no. 2, pp. B33–B44,  
787 2001.
- 788 [4] C. Büchel, C. Price, and K. Friston, "A multimodal language region in  
789 the ventral visual pathway," *Nature*, vol. 394, no. 6690, pp. 274–277,  
790 1998.
- 791 [5] A. Angelosi, A. Greco, and S. Harnad, "From robotic toil to symbolic  
792 theft: Grounding transfer from entry-level to higher-level categories,"  
793 *Connection Sci.*, vol. 12, no. 2, pp. 143–162, 2000.
- 794 [6] S. Coradeschi, A. Loutfi, and B. Wrede, "A short review of symbol  
795 grounding in robotic and intelligent systems," *Künstliche Intelligenz*,  
796 vol. 27, no. 2, pp. 129–136, 2013, doi: [10.1007/s13218-013-0247-2](https://doi.org/10.1007/s13218-013-0247-2).
- 797 [7] C. Crick and B. Scassellati, "Controlling a robot with intention derived  
798 from motion," *Topics Cogn. Sci.*, vol. 2, no. 1, pp. 114–126, 2010.
- 799 [8] J. C. Culham and K. F. Valyear, "Human parietal cortex in action,"  
800 *Current Opin. Neurobiol.*, vol. 16, no. 2, pp. 205–212, 2006.
- 801 [9] M. Daoutis and N. Mavridis, "Towards a model for grounding semantic  
802 composition," in *Proc. Artif. Intell. Simulat. Behav. (AISB)*, 2014.
- 803 [10] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood  
804 from incomplete data via the EM algorithm," *J. Roy. Stat. Soc. B (Methodol.)*,  
805 vol. 39, no. 1, pp. 1–38, 1977.
- 806 [11] J. Donahue *et al.*, "Long-term recurrent convolutional networks for  
807 visual recognition and description," in *Proc. IEEE Conf. Comput. Vis.  
808 Pattern Recognit.*, Boston, MA, USA, 2015, pp. 2625–2634.
- 809 [12] E. Elhamifar and R. Vidal, "Robust classification using structured sparse  
810 representation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.  
811 (CVPR)*, Colorado Springs, CO, USA, 2011, pp. 1873–1879.
- 812 [13] E. Elhamifar and R. Vidal, "Sparse subspace clustering: Algorithm, the-  
813 ory, and applications," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35,  
814 no. 11, pp. 2765–2781, Nov. 2013.
- 815 [14] A. Fazly, A. Alishahi, and S. Stevenson, "A probabilistic computational  
816 model of cross-situational word learning," *Cogn. Sci.*, vol. 34, no. 6,  
817 pp. 1017–1063, 2010.
- 818 [15] M. Fleischman and D. Roy, "Grounded language modeling for automatic  
819 speech recognition of sports video," in *Proc. ACL*, Columbus, OH, USA,  
820 2008, pp. 121–129.
- 821 [16] M. C. Frank, N. D. Goodman, and J. B. Tenenbaum, "Using speakers'  
822 referential intentions to model early cross-situational word learning,"  
823 *Psychol. Sci.*, vol. 20, no. 5, pp. 578–585, 2009.
- 824 [17] V. Gligozzi, J. Mayor, J. F. Hu, and K. Plunkett, "Labels as features  
825 (not names) for infant categorization: A neurocomputational approach,"  
826 *Cogn. Sci.*, vol. 33, no. 4, pp. 709–738, 2009.
- 827 [18] E. Goldberg, N. Driedger, and R. I. Kittredge, "Using natural-language  
828 processing to produce weather forecasts," *IEEE Expert*, vol. 9, no. 2,  
829 pp. 45–53, Apr. 1994.
- 830 [19] P. Gorniak and D. Roy, "Speaking with your sidekick: Understanding  
831 situated speech in computer role playing games," in *Proc. AIIDE*,  
832 Marina Del Rey, CA, USA, 2005, pp. 57–62.
- 833 [20] S. Harnad, "The symbol grounding problem," *Phys. D Nonlin.  
834 Phenomena*, vol. 42, nos. 1–3, pp. 335–346, 1990.
- 835 [21] A. Karpathy and L. Fei-Fei, "Deep visual-semantic alignments for gen-  
836 erating image descriptions," in *Proc. IEEE Conf. Comput. Vis. Pattern  
837 Recognit.*, Boston, MA, USA, 2015, pp. 3128–3137.
- 838 [22] J. Kennedy, P. Baxter, and T. Belpaeme, "Comparing robot embodiments  
839 in a guided discovery learning interaction with children," *Int. J. Soc.  
840 Robot.*, vol. 7, no. 2, pp. 293–308, 2015.
- 841 [23] P. Lamere *et al.*, "The CMU SPHINX-4 speech recognition system," in  
842 *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, vol. 1,  
843 Hong Kong, 2003, pp. 2–5.
- 844 [24] P. Li, I. Farkas, and B. MacWhinney, "Early lexical development in  
845 a self-organizing neural network," *Neural Netw.*, vol. 17, nos. 8–9,  
846 pp. 1345–1362, 2004.
- 847 [25] E. M. Markman, "Constraints children place on word meanings," *Cogn.  
848 Sci.*, vol. 14, no. 1, pp. 57–77, 1990.
- 849 [26] S. Marsland, J. Shapiro, and U. Nehmzow, "A self-organising  
850 network that grows when required," *Neural Netw.*, vol. 15, nos. 8–9,  
851 pp. 1041–1058, 2002.
- 852 [27] N. Mavridis, "Grounded situation models for situated conversational  
853 assistants," Ph.D. dissertation, Massachusetts Inst. Technol., Cambridge,  
854 MA, USA, 2007.
- 855 [28] N. Mavridis, "A review of verbal and non-verbal human–robot  
856 interactive communication," *J. Robot. Auton. Syst.*, vol. 63, pp. 22–35,  
857 Jan. 2015.
- 858 [29] B. McMurray, J. S. Horst, and L. K. Samuelson, "Word learning emerges  
859 from the interaction of online referent selection and slow associative  
860 learning," *Psychol. Rev.*, vol. 119, no. 4, pp. 831–877, 2012.
- 861 [30] G. Metta *et al.*, "The iCub humanoid robot: An open platform for  
862 research in embodied cognition," in *Proc. ACM 8th Workshop Perform.  
863 Metrics Intell. Syst.*, Gaithersburg, MD, USA, 2008, pp. 50–56.
- 864 [31] G. Metta, P. Fitzpatrick, and L. Natale, "YARP: Yet another robot  
865 platform," *Int. J. Adv. Robot. Syst.*, vol. 3, no. 1, pp. 43–48, 2006.
- 866 [32] M. Mishkin, L. G. Ungerleider, and K. A. Macko, "Object vision and  
867 spatial vision: Two cortical pathways," *Trends Neurosci.*, vol. 6, no. 10,  
868 pp. 414–417, 1983.
- 869 [33] P. Monaghan, K. Mattock, R. A. I. Davies, and A. C. Smith,  
870 "Gavagai is as Gavagai does: Learning nouns and verbs from cross-  
871 situational statistics," *Cogn. Sci.*, vol. 39, no. 5, pp. 1099–1112, 2015,  
872 doi: [10.1111/cogs.12186](https://doi.org/10.1111/cogs.12186).
- 873 [34] N. Otsu, "A threshold selection method from gray-level histograms,"  
874 *IEEE Trans. Syst., Man, Cybern.*, vol. SMC-9, no. 1, pp. 62–66,  
875 Jan. 1979.
- 876 [35] W. V. Quine, "On the reasons for indeterminacy of translation," *J.  
877 Philos.*, vol. 67, no. 6, pp. 178–183, 1970.
- 878 [36] T. Regier, "The emergence of words: Attentional learning in form and  
879 meaning," *Cogn. Sci.*, vol. 29, no. 6, pp. 819–865, 2005.
- 880 [37] D. Roy and N. Mukherjee, "Towards situated speech understanding:  
881 Visual context priming of language models," *Comput. Speech Lang.*,  
882 vol. 19, no. 2, pp. 227–248, 2005.
- 883 [38] D. K. Roy, "Learning visually grounded words and syntax for a scene  
884 description task," *Comput. Speech Lang.*, vol. 16, no. 3, pp. 353–385,  
885 2002.
- 886

AQ5

- 888 [39] R. G. Schwartz and B. Y. Terrell, "The role of input frequency in lexical  
889 acquisition," *J. Child Lang.*, vol. 10, no. 1, pp. 57–64, 1983.  
 890 [40] J. M. Siskind, "A computational study of cross-situational techniques  
891 for learning word-to-meaning mappings," *Cognition*, vol. 61, nos. 1–2,  
892 pp. 39–91, 1996.  
 893 [41] K. Smith, A. D. M. Smith, R. A. Blythe, and P. Vogt, *Cross-Situational  
894 Learning: A Mathematical Approach* (LNCS 4211). Heidelberg,  
895 Germany: Springer, 2006, pp. 31–44.  
 896 [42] C. E. Snow, "Mothers' speech to children learning language," *Child  
897 Dev.*, vol. 43, no. 2, pp. 549–565, 1972.  
 898 [43] G. Spitsyna, J. E. Warren, S. K. Scott, F. E. Turkheimer, and R. J. Wise,  
899 "Converging language streams in the human temporal lobe," *J. Neurosci.*,  
900 vol. 26, no. 28, pp. 7328–7336, 2006.  
 901 [44] F. Stramandinoli, D. Marocco, and A. Cangelosi, "The grounding of  
902 higher order concepts in action and language: A cognitive robotics  
903 model," *Neural Netw.*, vol. 32, pp. 165–173, Aug. 2012.  
 904 [45] Y. Sugita and J. Tani, "Learning semantic combinatoriality from the  
905 interaction between linguistic and behavioral processes," *Adapt. Behav.*,  
906 vol. 13, no. 1, pp. 33–52, 2005.  
 907 [46] M. Taddeo and L. Floridi, "Solving the symbol grounding problem: A  
908 critical review of fifteen years of research," *J. Exp. Theor. Artif. Intell.*,  
909 vol. 17, no. 4, pp. 419–445, 2005.  
 910 [47] A. Taniguchi, T. Taniguchi, and A. Cangelosi, "Multiple categorization  
911 by iCub: Learning relationships between multiple modalities and  
912 words," in *Proc. IROS Workshop Mach. Learn. Methods High Level  
913 Cogn. Capabilities Robot.*, 2016.  
 914 [48] V. Tikhanoff *et al.*, "An open-source simulator for cognitive robotics  
915 research: The prototype of the iCub humanoid robot simulator," in *Proc.  
916 ACM 8th Workshop Perform. Metrics Intell. Syst.*, Gaithersburg, MD,  
917 USA, 2008, pp. 57–61.  
 918 [49] V. Tikhanoff, A. Cangelosi, and G. Metta, "Integration of speech and  
919 action in humanoid robots: iCub simulation experiments," *IEEE Trans.  
920 Auton. Mental Develop.*, vol. 3, no. 1, pp. 17–29, Mar. 2011.  
 921 [50] M. Tomasello and N. Akhtar, "Two-year-olds use pragmatic cues to  
922 differentiate reference to objects and actions," *Cogn. Dev.*, vol. 10, no. 2,  
923 pp. 201–224, 1995.  
 924 [51] M. Vavrečka and I. Farkaš, "A multimodal connectionist architecture  
925 for unsupervised grounding of spatial language," *Cogn. Comput.*, vol. 6,  
926 no. 1, pp. 101–112, 2014.  
 927 [52] P. Vogt, "Language evolution and robotics: Issues on symbol grounding,"  
928 in *Artificial Cognition Systems*. Hershey, PA, USA: IGI Glob., 2006,  
929 p. 176.  
 930 [53] J. F. Werker and P. J. McLeod, "Infant preference for both male and  
931 female infant-directed talk: A developmental study of attentional and  
932 affective responsiveness," *Can. J. Psychol. Revue Can. De Psychol.*,  
933 vol. 43, no. 2, pp. 230–246, 1989.  
 934 [54] F. Xu and J. B. Tenenbaum, "Word learning as Bayesian inference,"  
935 *Psychol. Rev.*, vol. 114, no. 2, pp. 245–272, 2007.  
 936 [55] C. Yu and L. B. Smith, "Rapid word learning under uncertainty via  
937 cross-situational statistics," *Psychol. Sci.*, vol. 18, no. 5, pp. 414–420,  
938 2007.  
 939 [56] C. Yu and L. B. Smith, "Modeling cross-situational word-referent  
940 learning: Prior questions," *Psychol. Rev.*, vol. 119, no. 1, pp. 21–39,  
941 2012.  
 942 [57] D. Yurovsky, C. Yu, and L. B. Smith, "Competitive processes in cross-  
943 situational word learning," *Cogn. Sci.*, vol. 37, no. 5, pp. 891–921, 2013.  
 944 [58] K. Štepánová, "Hierarchical probabilistic model of language acquisi-  
945 tion," Ph.D. dissertation, Dept. Cybern., Czech Tech. Univ. Prague,  
946 Prague, Czech Republic, 2016.  
 947 [59] K. Štepánová *et al.*, "Where is my forearm? Clustering body parts from  
948 simultaneous tactile and linguistic input," in *Proc. Cogn. Artif. Life  
949 (KUZ XVII)*, 2017.



**Karla Štepánová** received the master's degree in condensed matter physics from Charles University, Prague, Czech Republic, and the Ph.D. degree in artificial intelligence and biocybernetics from Czech Technical University in Prague, Prague, in 2017.

She is a Researcher with the Czech Institute of Informatics, Robotics, and Cybernetics, Prague. She was a Visiting Researcher with Plymouth University, Plymouth, U.K., in 2016. Her current research interests include probabilistic models of cognition, unsupervised learning, language acquisition, and multimodal integration.



**Frederico Belmonte Klein** was born in Porto Alegre, Brazil, in 1982. He received the Electrical Engineering degree and the medical degree from the Federal University of Rio Grande do Sul, Porto Alegre, in 2006 and 2014, respectively. He is currently pursuing the Ph.D. degree in cognitive robotics with the University of Plymouth, Plymouth, U.K.



**Angelo Cangelosi** received the degree in psychology and cognitive science from the University of Rome La Sapienza, Rome, Italy, and the University of Genoa, Genoa, Italy.

He is a Professor of artificial intelligence and cognition and the Director of the Centre for Robotics and Neural Systems with Plymouth University, Plymouth, U.K. and was a Visiting Scholar with the University of California at San Diego, San Diego, CA, USA, and the University of Southampton, Southampton, U.K. His current research interests include language grounding and embodiment in humanoid robots, developmental robotics, human–robot interaction, and on the application of neuromorphic systems for robot learning.



**Michal Vavrečka** received the Ph.D. degree in general psychology from the Faculty of Social Studies, in 2008.

He is with the Czech Institution for Informatics Robotics and Cybernetics, Czech Technical University in Prague, Prague, Czech Republic. He develops cognitive architecture for symbol grounding, language acquisition, and knowledge representation.

## AUTHOR QUERIES

### AUTHOR PLEASE ANSWER ALL QUERIES

**PLEASE NOTE:** We cannot accept new source files as corrections for your paper. If possible, please annotate the PDF proof we have sent you with your corrections and upload it via the Author Gateway. Alternatively, you may send us your corrections in list format. You may also upload revised graphics via the Author Gateway.

AQ1: Please confirm/give details of funding source.

AQ2: Please provide the postal code for “Plymouth University Plymouth, U.K.”

AQ3: Please provide the page range for References [9], [47], and [59].

AQ4: Please provide the department name for Reference [27].

AQ5: Please confirm if the location and publisher information for Reference [41] is correct as set.

AQ6: Please specify the type of the degree and year attained by the author “A. Cangelosi.”

AQ7: Please provide the organization name and location for the Ph.D. degree obtained by “M. Vavrečka” biography.

# Mapping Language to Vision in a Real-World Robotic Scenario

Karla Štepánová<sup>ID</sup>, Frederico Belmonte Klein, Angelo Cangelosi, and Michal Vavrečka

**Abstract**—Language has evolved over centuries and was gradually enriched and improved. The question, how people find assignment between meanings and referents, remains unanswered. There are many of computational models based on the statistical co-occurrence of meaning-reference pairs. Unfortunately, these mapping strategies show poor performance in an environment with a higher number of objects or noise. Therefore, we propose a more robust noise-resistant algorithm. We tested the performance of this novel algorithm with simulated and physical iCub robots. We developed a testing scenario consisting of objects with varying visual properties presented to the robot accompanied by utterances describing the given object. The results suggest that the proposed mapping procedure is robust, resistant against noise and shows better performance than one-step mapping for all levels of noise in the linguistic input, as well as slower performance degradation with increasing noise. Furthermore, the proposed procedure increases the clustering accuracy of both modalities.

**Index Terms**—Cognitive modeling, cross-situational learning, iCub robot, language acquisition, symbol grounding.

speech understanding in computer games [19], automated generation of weather forecasts [18], tutoring children in foreign languages [22], etc.

Despite the extensive research in the area of language acquisition, the question how the word-to-meaning mapping is learned remains unanswered. There are basically two approaches for mapping language to other modalities. A classic approach for separate modality-dependent representations is advocated by Barsalou [1]; while an example of a system where language is mapped to an intermediate modal representation that can be derived by multiple modalities is [27]. Using the unsupervised approach, Li *et al.* [24] designed the DevLex model, consisting of two self-organizing networks that are bidirectionally connected. Gliozzi *et al.* [17] proposed an alternative with a multimodal representation layer: their unsupervised feature-based model was used to account for early category formation in young infants. This approach postulates the unsupervised role of linguistic labels that can affect categorization during the acquisition process, which has also been supported by Taniguchi *et al.* [47]. Vavrečka and Farkaš [51] recently introduced a multimodal architecture for the grounding of spatial words using a biologically inspired approach (separate “what” and “where” visual subsystems) in which the visual scenes (two objects in 2-D space in a spatial relation) are associated with their linguistic descriptions, thus leading to the integration of modalities.

However, a fully unsupervised architecture, which would be able to deal with language grounding [46], particularly language grounding in a case where sentences have variable structure and when there is more than one object in a scene, is not available. The current state-of-the-art on variable-length sentences is very restricted and deals only with static scenes [29]. Most of the recent models based on deep networks are oriented toward application in image-to-text [21] or video-to-text [11] mapping and do not take into account the psychological aspects of language acquisition (e.g., mutual exclusivity). Moreover, these systems are trained in a supervised manner without the advantage of transfer learning.

The difficulty of the task was described in a well-known experiment performed by Quine [35] who imagined the anthropologist meeting a native who pointed at the scene and said “gavagai.” When the anthropologist is stimulated in a situation by seeing a rabbit, he will suppose that the word represents the running rabbit in front of him, even though it could also mean “ground,” “sun,” “hello,” or whatever else. This problem is related to language relativity, as there are several

21

## I. INTRODUCTION

THE ESSENTIAL (and still not fully answered) question in language acquisition is how percepts are anchored in some arbitrary symbols, in other words, how words (symbols) get their meanings. This is the so-called symbol grounding problem [20]. For many years, cognitive modeling, neuroscience, psychology, and machine learning have jointly attempted to understand how humans can solve this “problem” [6]. The ability to learn language through perception and especially through visual grounding is not only important for understanding human cognition but is also applicable in many areas, such as verbal control of interactive robots [28], automatic sports commentators [15], car navigation systems, for the visually impaired, situated

Manuscript received June 27, 2017; revised December 6, 2017; accepted February 24, 2018. This work was supported in part by the European EU FP7 Research Project TRADR under Grant 609763, in part by TACR CAK under Grant TE01020197, in part by the CAPES Foundation, Ministry of Education of Brazil under Grant BEX 1084/13-5, in part by the CNPq Brazil under Grant 232590/2014-1, and in part by the U.K. EPSRC Project BABEL under Grant EP/J004561/1 and Grant EP/J00457X/1. (*Corresponding author:* Karla Štepánová.)

K. Štepánová and M. Vavrečka are with the Czech Institute of Informatics, Robotics, and Cybernetics, Czech Technical University in Prague, 16000 Prague, Czech Republic (e-mail: karla.stepanova@ciirc.cvut.cz).

F. B. Klein and A. Cangelosi are with the School of Computing, Electronics and Mathematics, Plymouth University, Plymouth PL4 8AA, U.K.

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TCDS.2018.2819359

AQ1

AQ2

82 objects and their features that are described by words [33]. A  
 83 simplified version of this problem consists of a simple visual  
 84 scene and separate words that are grounded based on statistical  
 85 co-occurrence (cross-situational learning).

86 From a neuroscientific point of view, symbol grounding can  
 87 be viewed as a process of finding mappings between primary  
 88 unimodal visual and language brain areas. Where exactly the  
 89 integration is performed is still the subject of research, and  
 90 existing literature provides only incomplete accounts of the  
 91 cortical location of this convergence. For example, the study  
 92 by Büchel *et al.* [4] provides evidence for the involvement of  
 93 the left basal posterior temporal lobe (BA37) in the integration  
 94 of language and visual information. Other studies (e.g., [43])  
 95 propose that access to verbal meaning depends on anterior  
 96 and posterior heteromodal cortical systems within the tem-  
 97 poral lobe. The grounding of actions and motoric primitives  
 98 is associated with the activity in the dorsal stream and the  
 99 premotor cortex [8].

100 How language can be developed in an unsupervised man-  
 101 ner is also an important task in developmental robotics as most  
 102 language acquisition in humans is fully unsupervised. One of  
 103 the main long-term objectives of many teams worldwide is  
 104 building conversational robots, which will be able to partici-  
 105 pate in cooperative tasks mediated by a natural language. It has  
 106 been shown how robots can learn new symbols using already  
 107 grounded ones and their combination [5] and how to transfer  
 108 knowledge between agents [52]. Cangelosi *et al.* [5] presented  
 109 their research on language emergence and grounding in senso-  
 110 rimotor agents and robots. This model was further extended by  
 111 Tikhanoff *et al.* [49], who performed iCub simulation experi-  
 112 ments and focused on the integration of speech and action. The  
 113 grounding of higher-order concepts in action was also explored  
 114 by Stramandinoli *et al.* [44], who made use of recurrent neural  
 115 networks. Sugita and Tani [45] described an experiment deal-  
 116 ing with semantic compositionality: the capability of a robot to  
 117 use the compositional structure to generalize novel word com-  
 118 binations. Daoutis and Mavridis [9] summarized desiderata for  
 119 grounded semantic compositionality. Despite all the progress  
 120 in language grounding, however, the current state-of-the-art  
 121 on grounding variable-length sentences is very restricted and  
 122 deals only with static scenes [29], [38].

123 In this paper, we present this paper in the area of lan-  
 124 guage acquisition using a real-world robotic scenario. We  
 125 implemented a hierarchical cognitive architecture for lan-  
 126 guage acquisition that includes visual and language processing.  
 127 In particular, we chose to extend current models of cross-  
 128 situational learning by allowing vision-to-language mapping  
 129 in the case of a nonequal number of classes and by taking  
 130 into account situation-time dynamics.

131 We show that this can be accomplished more efficiently  
 132 by replacing one-shot mapping with sequential mapping and  
 133 adding inhibitory mechanisms to the connections. The best  
 134 mapped classes are gradually eliminated, and the clusterization  
 135 is adaptively changed. We see this paper as an extension of  
 136 the McMurray *et al.* [29] model, and we compare it with other  
 137 single-step mapping models. The mapping strategy presented  
 138 in this paper was shown to be very robust as it not only can find  
 139 the mapping in circumstances of very noisy real-world input  
 140 but also increases the clustering accuracy of both modalities.

Recently, we tested the proposed algorithm also on the task of  
 141 clustering body parts from simultaneous tactile and linguistic  
 142 input [59]. In that case, sequential mapping showed slower  
 143 degradation with increasing noise level in the linguistic input  
 144 and outperformed one-step mapping for all dataset sizes and  
 145 all levels of noise.  
 146

The rest of this paper is structured as follows. In Section II,  
 147 we compare different mapping algorithms used in cross-  
 148 situational learning. In particular, in Section II-B we provide  
 149 a mathematical formulation of the newly proposed sequential  
 150 mapping algorithm, and in Section II-C we describe the whole  
 151 cognitive architecture which incorporates unimodal processing  
 152 of vision and language and finding their association through  
 153 the mapping algorithm. The performance of the proposed  
 154 method on data from an iCub humanoid robot and from an  
 155 iCub simulator is evaluated in Section III. Finally, results are  
 156 discussed in Section IV with suggestions for future work.  
 157

## II. MATERIALS AND METHODS

In this section, we first present one-step and newly proposed  
 159 sequential mapping algorithms (Sections II-A and II-B). Then,  
 160 we describe in detail the whole cognitive architecture used to  
 161 process data from individual modalities (vision and language).  
 162 Finally, we describe the iCub robotic platform and the iCub  
 163 simulator in Section II-D and provide a description of the  
 164 evaluation in Section II-F.  
 165

### A. One-Step Mapping in Cross-Situational Learning

One possible way how to establish mapping between visual  
 167 concepts and linguistic elements is to use frequencies of refer-  
 168 ent and meaning co-occurrences; that is, the ones with the  
 169 highest co-occurrence are mapped together [40], [41], [54].  
 170 This method is usually called cross-situational learning and  
 171 supposes the availability of the ideal associative learner  
 172 who can keep track and store all co-occurrences in all tri-  
 173 als, internally memorizing and representing the word-object  
 174 co-occurrence matrix of the input. This allows the learner  
 175 to subsequently choose the most strongly associated refer-  
 176 ent [55], [56]. These models do not see the mapping as  
 177 dynamic competition but operate only with the static state.  
 178 Although some use likelihoods of different words and refer-  
 179 ents to perform Bayesian inference [16], [54], they do not take  
 180 into account how the similarity of two-word forms can affect  
 181 learning although it has been shown that the similarity affects  
 182 learning in children [36]. Another shortcoming of these strate-  
 183 gies is that they do not address how these similarities affect  
 184 learning in a dynamic competition.  
 185

The simplest one-step word-to-referent learning algorithm  
 186 simply accumulates word-referent pairs. This can be viewed  
 187 as Hebbian learning: the connection between a word and an  
 188 object is strengthened if the pair co-occurs in a trial. To extend  
 189 this basic idea, we can enable also forgetting by introducing  
 190 the parameter  $\eta$ , which can capture the memory decay. This  
 191 so-called dumb associative-learning model was implemented  
 192 by Yu and Smith [56]. Supposing that at each trial  $t$  we observe  
 193 an object  $\mathbf{o}_t^n$  and hear a corresponding word  $\mathbf{w}_t^n$  ( $N_t$  possible  
 194 associations), we can describe the update of the strength of the  
 195 association between word model  $L(i)$  and object model  $K(j)$   
 196

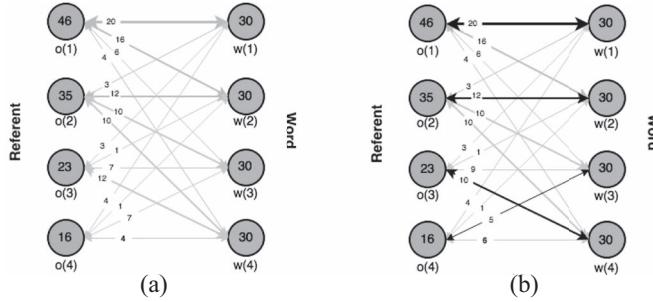


Fig. 1. One-step mapping. (a) In the first stage, weights between objects (referents) and words are changed using Hebbian learning. The connection between a word and an object is increased if the pair co-occurs in a trial [56]. (b) Then, word-to-referent mapping is found in one step: objects and words with the highest co-occurrence are mapped together. The number on each connection between word  $w(i)$  and object  $o(j)$  refers to the number of co-occurrences of  $o(j)$  and  $w(i)$ . In this example, we suppose that there are 30 occurrences of each word.

as follows:

$$A(i,j) = \sum_{t=1}^R \eta(t) \sum_{n=1}^{N_t} \delta(\mathbf{w}_t^n, i) \delta(\mathbf{o}_t^n, j) \quad (1)$$

where  $R$  is the number of trials,  $\delta$  is the Kronecker delta function (equal to 1 when both arguments are identical and 0 otherwise),  $\mathbf{w}_t^n$  and  $\mathbf{o}_t^n$  indicate the  $n$ th word-object association that the model attends to and attempts to learn in the trial  $t$ , and  $\eta(t)$  is the parameter controlling the gain of the strength of the association.

Now let us assume that word  $w(i)$  is modeled by model  $L_i$  in the language domain and object (referent)  $o(j)$  is modeled by the model  $K_{m(i)}$  in the visual domain. Our goal is to find the corresponding model  $K_{m(i)}$  from the visual subdomain for each model  $L_i$  from the language domain to assign them together. Indices  $m(i)$  are found as follows:

$$\forall i : m(i) = \operatorname{argmax}_i A(i,j) \quad (2)$$

where  $A$  is the co-occurrence matrix computed in (1) [element  $A(i,j)$  captures co-occurrence between word  $w(i)$  and object  $o(j)$ ].

In Fig. 1, one-step mapping is visualized as it was implemented in this paper.

Modifications of the basic model include Regier's [36] work. He proposed a mapping model, which stems from competition models and incorporates two-way associations between words and referents. This enables the model to capture selective attention to individual words and referents, as well as to provide a probability distribution over associated referents/words. Regier [36] also showed that for his model, learning of a novel word is most effective when memory interference is minimized.

### B. Proposed Sequential Mapping

Because we know that learning is not static but is a dynamic process, it seems reasonable to extend the basic idea of one-step cross-situational learning by incorporating dynamic competition mechanisms between words and referents in the

model. To capture the dynamic competition among models, we extend the basic one-step mapping algorithm with the sequential addition of inhibitory connections.

In this case, the process of finding word-referent associations resembles Hebbian learning with inhibitory connections. Once the word is associated with a corresponding object (referent), links from this referent to other words are inhibited. This idea also corresponds to the fact that children prefer mapping where an object has only one label to multiple labels, the so-called mutual exclusivity bias [25]. The inhibitory mechanisms and situation-time dynamics were already partially included in the model of cross-situational learning proposed by McMurray *et al.* [29].

Even though our model shares some similarities with the model proposed by McMurray *et al.* [29], our model stems from different computational mechanisms. The proposed sequential mapping is able to capture nondiscrete assignment to individual clusters, as well as dynamic competition mechanisms. The first mechanism is incorporated into the model by considering likelihoods that the observed data were generated by a given model instead of 1/0 assignment to models. In this way, similarities of individual meanings and referents, as well as the likelihood of their recognition in each trial, is taken into account. The second mechanism (dynamic competition) facilitates the sequential mapping as the best-mapped classes are gradually justified with inhibitory connections to other classes (i.e., after a reliable assignment between a language and a tactile model is found, inhibitory connections among this tactile model and all other language models are added). Thanks to this mechanism, the mutual exclusivity principle (the fact that children prefer mapping where an object has only one label to multiple labels [25]) is guaranteed. The assignment between visual models  $K_j$  and language models  $L_i$  is found using the following iterative procedure.

- 1) Visual and language data are clustered separately and the corresponding posterior probabilities are found

$$p(L_i | \mathbf{w}_t^n) = \frac{p(\mathbf{w}_t^n | L_i)p(L_i)}{\sum_{i'} p(\mathbf{w}_t^n | L_{i'})p(L_{i'})}, \quad \forall i \in \{1, \dots, I\} \quad (267)$$

$$\forall t \in \{1, \dots, T\}, \quad \forall n \in \{1, \dots, N_t\} \quad (268)$$

$$p(K_j | \mathbf{o}_t^n) = \frac{p(\mathbf{o}_t^n | K_j)p(K_j)\mathbf{k}(j)}{\sum_{j'} p(\mathbf{o}_t^n | K_{j'})p(K_{j'})}, \quad \forall j \in \{1, \dots, J\} \quad (269)$$

$$\forall t \in \{1, \dots, T\}, \quad \forall n \in \{1, \dots, N_t\} \quad (270)$$

where  $I$  is the number of language models,  $J$  is the number of visual models,  $T$  is the number of trials, and  $N_t$  is the number of possible object-word associations in trial  $t$ .

- 2) The most probable cluster to which the data point is assigned is selected (for each data point)

$$a_t^n = \operatorname{argmax}_i p(L_i | \mathbf{w}_t^n) \quad (277)$$

$$b_t^n = \operatorname{argmax}_j p(K_j | \mathbf{o}_t^n) \quad (278)$$

$$\forall t \in \{1, \dots, T\}, \forall n \in \{1, \dots, N_t\}. \quad (279)$$

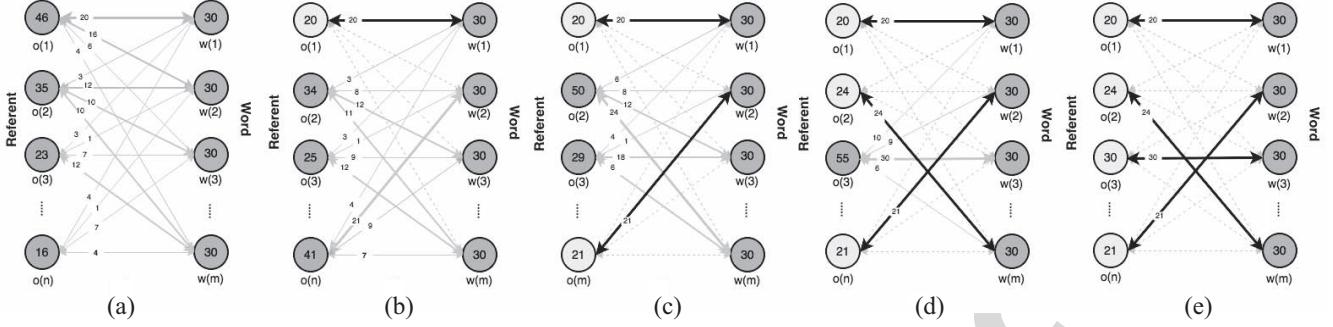


Fig. 2. Sequential mapping. The toy example of sequential mapping is shown to clarify the mechanism of finding the object-word assignment. In this example, we suppose that there are 30 occurrences of each word. The dotted line marks the inhibitory connection between object  $o(j)$  and word  $w(i)$  and the black line corresponds to the already found mapping. The number on each connection between word  $w(j)$  and object  $o(i)$  refers to the number of co-occurrences of  $o(j)$  and  $w(i)$ . Objects  $o(j)$  and words  $w(i)$  are assigned to corresponding models based on the given clustering mechanism.

280 3) Co-occurrence matrix  $A(i, j)$  is computed

$$281 \quad A(i, j) = \zeta(i, j) \sum_{t=1}^R \eta(t) \sum_{n=1}^{N_t} \delta(a_t^n, i) \delta(b_t^n, j) \quad (7)$$

282 where  $R$  is the number of trials,  $\zeta(i, j)$  is the matrix storing  
283 the strength of the connections between visual model  
284  $K_j$  and language model  $L_i$ , and  $\eta(t)$  is the parameter  
285 controlling the gain of the strength of association.

286 4) The best assignment is selected

$$287 \quad [im, m(im)] = \underset{i}{\operatorname{argmax}} \underset{j}{\operatorname{argmax}} A(i, j). \quad (8)$$

288 In this step, the visual model  $K_{m(im)}$  is assigned to the  
289 language model  $L_{im}$ .

290 5) Inhibition connections are added between the assigned  
291 visual model  $K_{m(im)}$  and all language models other than  
292  $L_{im}$  (mutual exclusivity)

$$293 \quad \zeta(i, m(im)) = \zeta(i, m(im))(1 - z_1), \quad \forall i \neq im \quad (9)$$

294 where  $z_i$  is the parameter capturing the strength of  
295 the inhibition (in our experiment, this is set to 1,  
296 which corresponds to the total inhibition of the given  
297 connection).

298 6) Inhibition is added to the assigned visual model  $K_{m(im)}$   
299 (a prior probability of the model is changed)

$$300 \quad k(m(im)) = k(m(im))(1 - z_2) \quad (10)$$

301 where  $z_2$  is the parameter capturing the inhibition of the  
302 assigned visual model (in our experiment, this parameter  
303 is set to 1, which corresponds to total inhibition of the  
304 given model).

305 7) The assigned points are deleted from the dataset (data  
306 points which belong to model  $K_{jm}$  and  $L_{im}$ )

$$307 \quad X = X \setminus \left\{ (\mathbf{o}_t^n, \mathbf{w}_t^n) \mid \underset{j}{\operatorname{argmax}} p(K_j | \mathbf{o}_t^n) == m(im) \right. \quad (11)$$

$$308 \quad \left. \wedge \underset{i}{\operatorname{argmax}} p(L_i | \mathbf{w}_j^n) == im \right\}. \quad (12)$$

309 8) If the dataset is not empty ( $X \neq \emptyset$ ) or  $\|\mathbf{k}\| > 0$  (some of  
310 the visual models are not totally inhibited), go to step 1,  
311 else stop.

312 The proposed algorithm where words are assigned to  
313 corresponding referents in a sequential manner is visualized  
314 in Fig. 2.

315 In the ideal case, unambiguous mapping between the two  
316 clusterizations will be found. In the real case (where the clus-  
317 terization in visual and language layers is not optimal), none  
318 or more than one model from the visual layer will be assigned  
319 to one cluster  $L_i$  in the language layer or vice versa.

### C. Specific Architecture

321 Our multimodal hierarchical architecture consists of  
322 multimodal and unimodal parts. The unimodal part has  
323 two layers performing separate processing of localist inputs:  
324 1) visual objects and 2) auditory word-forms. Both unimodal  
325 layers are subsequently mapped one to each other in the upper  
326 multimodal layer (see Fig. 3).

327 1) *Visual Layer:* Each data point (object  $\mathbf{o}_t^n$ ) can be consid-  
328 ered as a triplet of continuous-valued vectors for each visual  
329 feature:  $\mathbf{o}_t^n = (\mathbf{x}_{t,n}^{\text{size}}, \mathbf{x}_{t,n}^{\text{color}}, \mathbf{x}_{t,n}^{\text{shape}})$ . This enables us to write the  
330 visual dataset as  $X^{\text{vis}} = [X^{\text{size}} \ X^{\text{color}} \ X^{\text{shape}}]$  and process data for  
331 each visual feature separately. For processing visual data, the  
332 Gaussian mixture model (GMM) was used, which is a convex  
333 mixture of  $d$ -dimensional Gaussian densities  $l(\mathbf{x}^k | \theta_j^k)$ , where  
334  $k \in \{\text{size, color, shape}\}$ . In this case, each visual model  $K_j^k$  is  
335 described by a set of parameters  $\theta_j^k$ . The posterior probabilities  
336  $p(\theta_j^k | \mathbf{x}^k)$  are computed as follows:

$$337 \quad p(\theta_j^k | \mathbf{x}^k) = \sum_{j=1}^{J_k} r_j^k l(\mathbf{x}^k | \theta_j^k) \quad (13)$$

$$338 \quad l(\mathbf{x}^k | \theta_j^k) = \frac{1}{\sqrt{(2\pi)^d} \sqrt{|\mathbf{S}_j^k|}} \exp \left[ -\frac{1}{2} (\mathbf{x}^k - \mathbf{m}_j^k)^T (\mathbf{S}_j^k)^{-1} \right. \\ \left. \times (\mathbf{x}^k - \mathbf{m}_j^k) \right] \quad (14)$$

340 where  $k \in \{\text{size, color, shape}\}$ ,  $\mathbf{x}^k$  is a set of  $d$ -dimensional  
341 continuous-valued data vectors,  $r_j^k$  are the mixture weights,  
342  $J_k$  is the number of visual models for each visual feature  
343  $k$ , and parameters  $\theta_j^k$  are cluster centers  $\mathbf{m}_j^k$  and covariance  
344 matrices  $\mathbf{S}_j^k$ .

345 The Gaussian mixture is trained by the expectation-  
346 maximization algorithm [10]. An output of this layer for

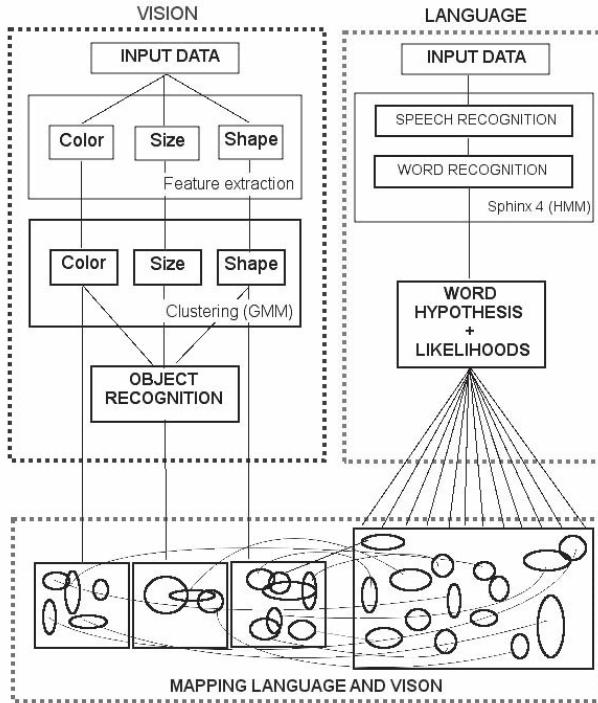


Fig. 3. Proposed multimodal architecture.

each data point  $\mathbf{x}_i^k$  is the vector  $\mathbf{y}_i^k$  of  $J_k$  output parameters describing the data point (the likelihood that the data point belongs to each individual cluster in a mixture). This corresponds to the fuzzy memberships (distributed representation). For a simpler evaluation, we used a localist representation (winner-takes-all), where only the cluster with the highest cluster membership probability is considered for further processing [see (5) and (6)]

$$M(K_j^k | O) = \begin{cases} 1 & \text{if } j = \operatorname{argmax}_j f(K_j^k | O) \\ 0 & \text{if } j \neq \operatorname{argmax}_j f(K_j^k | O) \end{cases} \quad (15)$$

where  $k \in \{\text{size, color, shape}\}$ ,  $j \in \{1, \dots, J_k\}$ .

2) *Language Layer*: The auditory word-forms are extracted from the language input which are sentences describing the image in the format <size> <color> <shape> (e.g., “small red triangle”). Afterward, individual word-forms are extracted from the sentences and compared to prelearned language models, and the log-scale score  $p(\mathbf{w}_t^n | L_i)$  of the audio matching the model is computed. Based on these data, posterior probability can be computed

$$p(L_i | \mathbf{w}_t^n) = \frac{p(\mathbf{w}_t^n | L_i)p(L_i)}{\sum_{i'} p(\mathbf{w}_t^n | L_{i'})p(L_{i'})}, \quad \forall i \in \{1, \dots, I\} \quad \forall t \in \{1, \dots, T\}, \quad \forall n \in \{1, \dots, N_t\} \quad (16)$$

where  $I$  is the number of language models,  $T$  is the number of trials (sentences) and  $N_t$  is the number of word-forms in the trial (sentence)  $t$ .

An output of this layer for each data point  $\mathbf{w}_t^n$  is the vector  $\mathbf{y}_i^k$  of  $I$  output parameters describing the data point (the likelihood that the data point belongs to each individual language model). This corresponds to the fuzzy memberships

---

**Algorithm 1** Sequential Mapping—Fixed Grammar

---

**Inputs:**  
language clusters  $L_i$  ( $i \in 1 : I$ ), visual clusters  $K_j^k \sim N(\mathbf{m}_j^k, \mathbf{S}_j^k)$ ,  
 $j \in 1 : J_k$ , input data  $\mathbf{x}^k$  for each feature  $k \in \{\text{size, color, shape}\}$ ,  
number of clusters  $J^k$  for each feature  $k$

**Output:**  
mapping between all visual classes  $K_j^k$  and language classes  $L_i$

```

for  $k \in \{\text{size, color, shape}\}$  do
     $NCl \leftarrow J^k$ 
    while  $NCl > 0$  and  $\mathbf{X}^k$  is not empty do
        assign each data point from  $\mathbf{x}^k$  to visual and language cluster (Winner-takes-all, see Eq. (15))
        for  $j = 1 : NCl$  do
            for  $i = 1 : I$  do
                 $A_{ij} \leftarrow$  how many times was class  $i$  classified as  $j$ 
            end for
        end for
         $[im, jm] \leftarrow \operatorname{argmax}_i \operatorname{argmax}_j A_{ij}$ 
         $\mathbf{X}_{del}^k \leftarrow$  data points assigned to  $K_{jm}^k$  and  $L_{im}$ 
         $\Theta_{new}^k \leftarrow N(\mathbf{x}_{del}^k)$  learn Gaussian on the to be deleted data
         $\mathbf{X}^k \leftarrow \mathbf{X}^k \setminus \mathbf{X}_{del}^k$  delete data assigned to both  $K_{jm}^k$  and  $L_{im}$ 
         $NCl \leftarrow NCl - 1$ 
        relearn  $K^k \sim N(\mathbf{m}^k, \mathbf{S}^k)$  on new data  $\mathbf{x}^k$  with  $NCl$  clusters
    end while
end for
cluster visual data using new  $\theta_{new}^k$  parameters (cluster centers  $\mathbf{m}^k$ , covariance matrices  $\mathbf{S}^k$ ) and perform one-step mapping (Model 1)

```

---

(distributed representation). Linguistic and visual inputs are processed simultaneously.

3) *Mapping—Models 1 and 2*: After the visual and language data are clustered, the mapping between the two layers must be found. For each cluster  $L_i$  in the language layer, a corresponding cluster  $K_j^k$  in the visual layer (for each feature  $k \in \{\text{size, color, shape}\}$ ) is found. The mapping is found as follows. For each  $j$  and  $k$ , we find cluster  $L_{kmax_{jk}}$  from the language layer which will be assigned to cluster  $K_j^k$  from the visual layer. In this paper, we compare two different models to find indices  $kmax_{jk}$ . We compared one-step mapping (see Section II-A) and newly proposed sequential mapping (see Section II-B).

The exact algorithm used to find the mapping between visual and language models in a sequential manner is described in detail in Algorithm 1. Indices  $m(i)$  are found sequentially. In each step, the best-mapped data are excluded, and the rest of data are reclustered using GMMs. Then, one-step mapping is performed (see Algorithm 1). An extension of the algorithm for a variable-length sentence is described in the Appendix. Results for a variable length sentence using a fully artificial dataset with a controlled noise level are described in detail in [58].

#### D. iCub Robotic Platform and iCub Simulator

For the experiment, we used a simulated [48] and a physical stationary [30] iCub robot. The iCub [Fig. 1(c)] is an open-source humanoid robot the size of a three-and-a-half-year-old child, with the fully articulated hands and a head-and-eye system which makes him ideal for cognitive experiments. The iCub simulator has been designed to reproduce, as accurately as possible, the physics and dynamics of the robot and its environment [48]. The simulator and the actual robot have the same interface supporting YARP [31] which is a robot platform for interprocess communication and control of the physical and simulated robots in real-time.

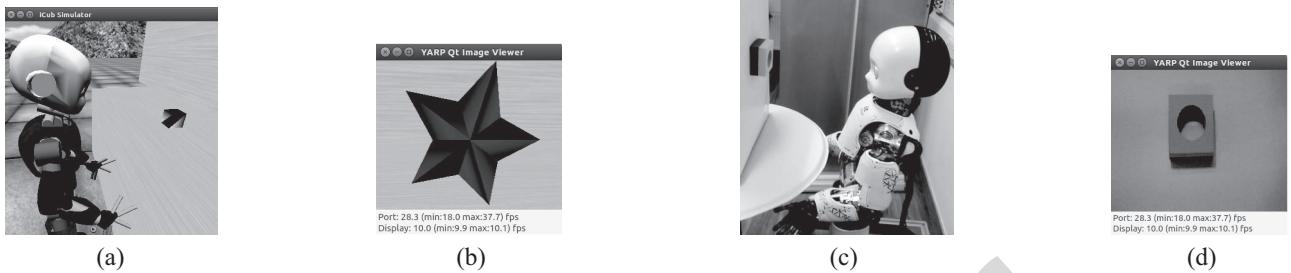


Fig. 4. Experiment design and corresponding input data. (a) iCub simulator. (b) Blender object. (c) Physical iCub. (d) Real object.

#### 409 E. Input Data Description and Preprocessing

410 The input for our model consisted of visual and language  
 411 data. The visual scene was composed of an object in the cen-  
 412 ter of the scene with a variable position. The visual features  
 413 (size, shape, and color) of the object also varied. We devel-  
 414 oped two separate datasets for training and testing purposes.  
 415 A real-world dataset has visual sensory data acquired from  
 416 the cameras of the physical iCub robot which observed sim-  
 417 ple objects placed on the white board in front of his eyes  
 418 at a distance of 1 m [see Fig. 4(c) and (d)] (210 instances,  
 419 three sizes, five colors, and seven shapes, in total 70 indi-  
 420 vidual objects; each seen three times in slightly different  
 421 placement). The simulated dataset is made in the iCub sim-  
 422 ulator [see Fig. 4(a) and (b)] as Blender-generated virtual  
 423 objects (432 instances, three sizes, six colors, and six shapes,  
 424 in total 108 individual objects; each seen four times in slightly  
 425 different placement).

426 The spoken language input consisted of sentences pro-  
 427 nounced by a non-native English speaker describing the image  
 428 in the format <size> <color> <shape> (e.g., small red trian-  
 429 gle) and was processed simultaneously with the visual input. In  
 430 the real-world dataset, the tutor was talking to a robot from an  
 431 approximate 2-m distance with natural background noise. The  
 432 linguistic input was captured using an external microphone.

433 1) *Speech Recognition*: CMU Sphinx (an open-source flex-  
 434 ible Markov model-based speech recognizer system) was used  
 435 for speech recognition [23]. Sphinx itself offers a large vocab-  
 436 ular, but we created our own task-specific smaller vocabulary  
 437 using the online IMtool that produces a dictionary based on a  
 438 CMU dictionary and matches its language model.

439 There is a probabilistic output from the CMU Sphinx. The  
 440 ten best hypotheses for a matching model with correspond-  
 441 ing scores were saved for each utterance (they are log-scale  
 442 scores of the audio matching the model). Because the scores  
 443 for the hypothesis of each word in the sentence were needed  
 444 for further evaluation, the words were pronounced with large  
 445 pauses, and the end of the sentence was marked by the word  
 446 “STOP.” An output of the language layer is a  $I$ -dimensional  
 447 continuous valued vector, where  $I$  is the number of language  
 448 clusters (corresponding to the number of possible utterances).

449 This vector contains ten nonzero values, and the rest is zero.  
 450 2) *Image Processing*: The image inputs are processed using  
 451 standard MATLAB functions. First, the image is morphologi-  
 452 cally opened with a disk-shaped structuring element (*imopen*)  
 453 to remove the noisy background of the image; then all grayish



Fig. 5. Image processing: original image, removal of the background, converting to BW image, and filling the holes.

454 pixels are removed, and the image is converted from the true 455 color RGB to the gray-scale intensity image by eliminating 456 the hue and saturation information while retaining the lumi- 457 nance (*rgb2gray*). Finally, the intensity image is converted to 458 a binary image using the threshold computed with Otsu’s [34] 459 method (*threshold*). An example of the preprocessed image is 460 shown in Fig. 5.

461 Afterward, the properties of the image regions are mea- 462 sured using the function *regionprops*. Individual visual features 463 (shape, color, and size) are subsequently processed sepa- 464 rately. The following features were used: color (three features: 465 average RGB of the selected region), size (six features: 466 perimeter of an object, distance from the centroid to the left 467 corner of the bounding box, and width and length of the 468 bounding box), and shape (13 features: area, centroid, major 469 axis length, eccentricity, orientation, convexArea, FilledArea, 470 EulerNumber, EquivDiameter, solidity, extent, and perimeter). 471 To obtain the shape features, we automatically cropped and 472 resized the image to equalize the size of the objects. 473

474 Although the visual model is mainly mathematical and 475 implemented in a very “machine vision” sort of way, the 476 bases of its processing follow biological correlates of mammal 477 vision. More specifically, from the neuroanatomical point of 478 view, this corresponds to the processing of the visual input in 479 the separate higher visual centers in the brain, specifically to 480 the independent processing of the information about the posi- 481 tion and identification of an object in the ventral (what) and 482 dorsal (where) neural pathways, respectively [32]. Individual 483 object properties are identified in the separate visual centers

#### 484 F. Evaluation

485 In the case of the supervised GMM, growing when required 486 neural gas (GWR), and SOM algorithms, data were divided 487 into training and validation datasets in the ratio 70:30. For the 488 unsupervised GMM, hidden Markov model, and  $k$ -means algo- 489 rithms, we computed the accuracy differently. After the data 490

TABLE I  
COMPARISON OF CLUSTERIZATION AND CLASSIFICATION ACCURACY OF VISUAL DATA. THE MEAN AND STANDARD DEVIATION FROM 100 REPETITIONS ARE VISUALIZED

Accuracy [%]	Real data			Blender		
	Size	Color	Shape	Size	Color	Shape
GMM sup.	83.3 ± 0.0	99.0 ± 0.0	81.4 ± 0.0	98.6 ± 0.0	97.9 ± 0.0	93.1 ± 0.0
GMM unsup.	76.2 ± 6.8	76.1 ± 9.1	56.1 ± 6.2	74.2 ± 10.1	60.9 ± 9.0	64.3 ± 7.2
K-means	67.8 ± 6.2	81.2 ± 1.1	53.1 ± 4.2	66.3 ± 0.2	77.1 ± 10.7	72.8 ± 6.9
SOM	69.6 ± 5.6	78.9 ± 6.8	54.2 ± 4.1	66.1 ± 4.2	81.7 ± 7.6	59.3 ± 6.2
GWR	89.9 ± 2.1	99.5 ± 0.4	76.6 ± 1.4	88.9 ± 0.7	98.1 ± 0.9	94.2 ± 0.6

TABLE II  
COMPARISON OF ONE-STEP MAPPING AND SEQUENTIAL MAPPING FOR DATA FROM THE iCUB SIMULATOR (BLENDER) AND THE PHYSICAL iCUB (REAL DATA). THE MEAN AND STANDARD DEVIATION FROM 100 REPETITIONS ARE VISUALIZED

Accuracy [%]	Real data			Blender		
	Size	Color	Shape	Size	Color	Shape
Vision	76.2 ± 6.8	76.1 ± 9.1	56.1 ± 6.2	74.2 ± 10.1	60.9 ± 9.0	64.3 ± 7.6
Language	70.6 ± 0.0	82.4 ± 0.0	77.5 ± 0.0	98.1 ± 0.0	96.5 ± 0.0	98.1 ± 0.0
One-step mapping	54.1 ± 4.1	58.2 ± 10.3	52.2 ± 4.9	67.3 ± 8.2	56.2 ± 6.1	61.9 ± 3.2
Sequential mapping	74.2 ± 15.1	87.1 ± 10.2	72.9 ± 5.1	96.1 ± 31.2	95.2 ± 1.2	92.1 ± 0.9

490 is clustered, each cluster is assigned to the class that appears  
491 most frequently in the cluster, and then the accuracy of this  
492 assignment is measured by counting the number of correctly  
493 assigned data points (compared to the manual true labels)  
494 and dividing this number by the total number of data points.

495 The accuracy of the learned mapping is calculated in the fol-  
496 lowing manner. We cluster output activations from the visual  
497 layer and assign each data point to the most probable cluster.  
498 Then, we find indices  $m(i)$  for all clusters as defined in (2) for  
499 one-step mapping and (8) for sequential mapping. Based on  
500 this mapping, we can assign each data point to the language  
501 label. These language labels are subsequently compared to the  
502 ground truth. Accuracy is then computed as

$$503 \quad \text{acc} = \frac{\text{TP}}{N} \quad (17)$$

504 where TP (true positive) is the number of correctly assigned  
505 data points, and  $N$  is the number of all data points.

### 506 III. EXPERIMENTAL RESULTS

507 The first part of the results is dedicated to the performance  
508 of the model in the real-world scenario. The robot interacts  
509 with a human in a noisy condition that distorts speech input.

#### 510 A. Vision

511 In the first stage, we evaluated the Vision subpart of the  
512 model. Several algorithms were compared: namely the GMM  
513 algorithm, supervised GMM algorithm,  $k$ -means, SOM, and  
514 GWR algorithm [2], [26]. The SOM and GWR had 100 nodes.  
515 The results for the real-world dataset and the simulated dataset  
516 with Blender objects can be seen in Table I. Although the SOM  
517 and GWR algorithms are considered unsupervised algorithms,  
518 we adopted a technique for labeling inputs; thus, they should  
519 be compared with supervised algorithms. The number of clus-  
520 ters is also overestimated (the number of nodes corresponds  
521 to the number of clusters). It indicates that these algorithms

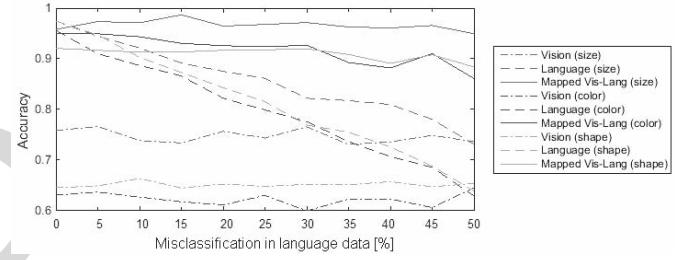


Fig. 6. Dependence of mapping accuracy on the misclassification in the language data for a fixed-length sentence (mean values over 50 repetitions are visualized). Different colors correspond to different visual features (red: size, blue: color, and green: shape). Visual data are generated in Blender and acquired through the iCub simulator, and language data are processed using Sphinx 4.

522 partly overfit the data so we divided the set into testing and validation data. 523

#### 524 B. Mapping

525 The performance of one-step mapping (vision and lan-  
526 guage are mapped in one step based on the frequency of  
527 co-occurrence) and sequential mapping (see Algorithm 1)  
528 is shown in Table II. We calculated the accuracy for  
529 real-world data from physical iCub and for the Blender  
530 objects placed in the iCub simulator. Language accuracy  
531 for Blender dataset is much higher compared to the real-  
532 world data as a tutor was speaking directly into the  
533 microphone.

534 Tolerance of the sequential mapping to misclassification  
535 in the language data is visualized in Fig. 6 for visual data  
536 from the iCub simulator in combination with the language  
537 data processed by Sphinx 4. The misclassification is added to  
538 the language data subsequently and evenly to all classes (the  
539 given proportion of the language inputs was randomly changed  
540 to random words). The misclassification (a synthetic error)  
541 was added artificially but can be interpreted either as white  
542 noise added to the input data or mistakes in labeling perceived

543 objects. We grouped them together into a misclassification  
 544 variable. The visual data are left intact so the only cause of  
 545 the observed variations in the accuracy is initialization. As can  
 546 be seen, the accuracy of sequential mapping remains very sta-  
 547 ble even though the accuracy of the language decreases and  
 548 outperforms the language and vision for almost all values of  
 549 the misclassification.

#### 550 IV. CONCLUSION

551 Current models of vision-to-language mapping often make  
 552 use of cross-situational learning while relying directly on sta-  
 553 tistical co-occurrence of meaning-referent pairs: the ones with  
 554 the highest co-occurrence are mapped together (e.g., [41]).  
 555 This approach shows inferior performance in cases where  
 556 the clustering of data is difficult (e.g., due to the high num-  
 557 ber of objects, high noise level or overlapping classes, etc.).  
 558 Therefore, we extended this basic model and introduced a  
 559 new more robust and noise-resistant mapping procedure. Our  
 560 approach incorporates situation-time dynamics and mutual  
 561 exclusivity and is able to deal with a nonequal number of  
 562 classes in individual subdomains.

563 Mathematical formulation of the newly proposed mapping  
 564 is provided (see Section II-B), and results for simulated and  
 565 real-world data from the iCub robot are compared to one-  
 566 step mapping (see Table II). It was shown that the method is  
 567 able to find mapping between language and vision, and the  
 568 method improves the accuracy of both individual subdomains  
 569 and shows very good resistance to noise or misclassification in  
 570 language (see Fig. 6). How to map in an unsupervised manner  
 571 several clusterings (e.g., for vision, action, and language) is an  
 572 important question not only in cognitive modeling but also in  
 573 general machine learning, where data acquired from different  
 574 sensors or in different situations can be independently clus-  
 575 tered and mapped to each other. A more detailed discussion  
 576 of the results follows.

577 The trivial one-step mapping can be imagined as basic  
 578 Hebbian learning. Our extension can be likened to Hebbian  
 579 learning with inhibitory connections. In recent years, sev-  
 580 eral approaches finding an alternative to the basic approach  
 581 appeared [29], [57]. Our model partially stems from the  
 582 McMurray *et al.* [29] approach, who showed that associative  
 583 learning can be sufficient for language acquisition and that the  
 584 main components of this type of learning are an online com-  
 585 petition of models and pruning incorrect associations which  
 586 enable gradual improvement of associations between models.

587 The mutual exclusivity principle is guaranteed in our  
 588 method thanks to the inhibitory connections which are grad-  
 589 ually created among models. Once the mapping between the  
 590 referent and the meaning is found, the connection from a given  
 591 meaning to other referents is inhibited. Dynamic competition  
 592 is addressed in the following way: when any association of  
 593 meaning and referent is found, other models compete again  
 594 for resources. First, well-mapped data are deleted, and then,  
 595 the resting data are reclustered. Furthermore, the likelihoods  
 596 can associate each data point to many separate models instead  
 597 of binary membership. In this way, our algorithm also emulates  
 598 how the similarity of two-word forms can affect learning as  
 599 it occurs in children's development [36]. The algorithm also

enables mapping together data in a case where we have an  
 600 uneven number of clusters in both subdomains.  
 601

602 However, the principle of mutual exclusivity is not suitable  
 603 for further stages of language acquisition, namely, the learning  
 604 of polysemic words. A polyseme is a word or phrase with  
 605 different but related senses (e.g., wood as a piece of tree or the  
 606 area with trees). The homonyms are a subset of the polysemes,  
 607 but the difference between the homonyms and the polysemes  
 608 is subtle and fuzzy. Homonyms represent a group of words  
 609 that share similar spelling (homographs) and the same sound  
 610 (homophones) but have different and unrelated meanings, for  
 611 example, the homograph bank standing for embankment or  
 612 place where money is kept. The learning of polysemic words  
 613 violates the principle of mutual exclusivity as the dynamic  
 614 competition does not allow to map one word to different visual  
 615 models. We are aware of this problem, and we would like to  
 616 extend it in the future iteration of the model. Similar to humans  
 617 who have to learn polysemic words as an exception, we will  
 618 incorporate this principle into the next version of our model.  
 619

620 In the next part, we analyze the ability of our architec-  
 621 ture to deal with ambiguous inputs. The mapping will find  
 622 reliable labeling for the visual input data (more generally for  
 623 data from any other modality) with a possibility of incorporat-  
 624 ing the fuzziness of this mapping. For some concepts, finding  
 625 an unambiguous mapping is very easy; for others, it is much  
 626 more difficult or impossible (such as abstract words, e.g., the  
 627 love has no dominant color, but the sky is usually blue). Since  
 628 the mapping is established only among the clusters where it  
 629 makes sense, dealing with a lot of redundant information is  
 630 avoided. A similar idea is used in classification algorithms  
 631 which use sparse matrices (e.g., [12] and [13]). We also ana-  
 632 lyzed the strong and weak points of the algorithms adopted  
 633 in our architecture. First, the ability of different algorithms  
 634 to classify unimodal visual data was compared. As expected,  
 635 for data which are well separated and mainly spherically  
 636 distributed (this is generally a case for simulated and artifi-  
 637 cially generated data), the  $k$ -means algorithm outperformed  
 638 the GMM algorithm. However, for nonspherical real data, the  
 639 GMM algorithm generally performed better (see Table I for  
 640 the comparison of the performance on simulated data placed  
 641 in an iCub simulator and data from real iCub robot cameras).  
 642 As the unsupervised algorithms are highly dependent on the  
 643 initialization, it can be seen that the standard deviation of the  
 644 data is quite high even though 20 repetitions were averaged.  
 645 We should conclude that the performance of the algorithms in  
 646 our tasks reflects their fundamental advances and limitations.  
 647 We are still missing the algorithm that is able to cope with  
 648 highly variable datasets in terms of their statistical properties.

649 The most significant differences between data acquired from  
 650 the real and the simulated iCub can be found in the visual sub-  
 651 domain. Affected by natural light and slightly altered points  
 652 of view, the shading of objects and their bevel differed. This  
 653 resulted in a different projection of objects to iCub cameras  
 654 and in the less accurate classification of shapes for the real-  
 655 world dataset. For colors, different lighting changed the true  
 656 color of the objects. Nevertheless, as can be seen in Table II,  
 657 this did not affect the ability of the algorithm to discrimi-  
 658 nate among individual colors. The real-world objects had a  
 659

658 variable size which resulted in a higher overlap of clusters for  
 659 individual sizes compared to the Blender dataset.

660 Afterward, we focused on the mapping between vision and  
 661 language and compared two different approaches: one-step  
 662 mapping to sequential mapping which in a stepwise man-  
 663 ner finds the best-mapped clusters while constantly relearning  
 664 clusterization of visual data. As can be seen in the results in  
 665 Table II, the novel sequential mapping led to an improvement  
 666 in effectiveness compared to the method which maps vision  
 667 to language in one step. This can be seen in more accurate  
 668 mapping which leads to a better estimation of the clustered  
 669 data labels and consequently to lower classification error for  
 670 all of the evaluated datasets.

671 The accuracy of multimodal mapping outperforms vision,  
 672 language or both. This is an important finding, because  
 673 sequential mapping improves not only the accuracy of visual  
 674 clustering but can also fix mistakes in language recognition,  
 675 which provides the labels. Furthermore, this result suggests  
 676 that we are able not only to find mapping between more clus-  
 677 terings, but we can also improve the clusterization accuracy  
 678 by combining individual classifiers. This is not easily seen  
 679 in the presented dataset as there is high recognition accu-  
 680 racy of Sphinx software (especially in the case of sentences  
 681 recorded for the simulated dataset). Therefore, we also tested  
 682 whether the misclassification (or noise) in the language data  
 683 affects the correct mapping between vision and language. The  
 684 result can be seen in Fig. 6. Although a synthetic error (mis-  
 685 classification) is added to the language data, the ability to  
 686 find the mapping remains nearly intact. The mapping accu-  
 687 racy decreases only very slightly and remains at around 90%  
 688 because the accuracy of language recognition drops from the  
 689 original approximately 95% to approximately 70% (depend-  
 690 ing on the specific visual feature). In the future, we would  
 691 like to explore the effect of noise in individual modalities on  
 692 classification accuracy in more detail by adding a controlled  
 693 amount of noise directly to the input. The error propagation  
 694 from individual input modalities to the multimodal layer was  
 695 explored, for example, in [37]. They showed that the visual  
 696 context can steer speech recognition and vice versa.

697 We also analyzed how the complexity of environment  
 698 affects the accuracy of our architecture. We suppose that  
 699 the performance of one-step mapping will decrease with the  
 700 increasing complexity of the task (more clusters, higher over-  
 701 lap, and higher dimensionality) as it is more difficult to find  
 702 reliable clustering of the data. This hypothesis is supported by  
 703 our preliminary results for clustering body parts from simu-  
 704 taneous tactile and linguistic input [59] and by the results  
 705 presented in this paper in Table II. The quality of one-step  
 706 mapping correlates with the quality of the visual data clus-  
 707 tering. For Blender data, the worst performance was achieved  
 708 for mapping words to the visual feature shape ( $52 \pm 5\%$ )  
 709 and for physical iCub for the feature color ( $62 \pm 3\%$ ). The  
 710 feature shape has the highest number of clusters (10), and  
 711 the feature color has the second highest number of clusters  
 712 (9) and the highest overlap of the clusters for the physical  
 713 iCub. Real-world tasks are much more complex, and we can  
 714 expect tens of different object shapes. In that case, the clus-  
 715 tering performance is crucial, and one-step mapping would

not be able to provide reliable mapping. It can be seen from  
 716 the results that the proposed mapping which enables grad-  
 717 ual re-estimation of the model parameters and works in a  
 718 dynamic fashion achieves much higher accuracy even for the  
 719 cases where one-step mapping fails. We suppose that the  
 720 mapping accuracy of the proposed method decreases more  
 721 slowly with the decreasing accuracy of clustering of individual  
 722 modalities. Unfortunately, this factor was not studied in our  
 723 restricted scenario, but the preliminary results on mapping tac-  
 724 tile and linguistic input [59] support this hypothesis. We plan  
 725 to investigate this phenomenon more deeply in future research.  
 726

The language dataset differs considerably from natural lan-  
 727 guage. However, the dataset reflects some characteristics from  
 728 the findings of Werker and McLeod [53] as infant-directed  
 729 words are usually kept short with large pauses between words.  
 730 Moreover, Brent and Siskind [3] showed that frequency of  
 731 exposure to a word in isolation predicts better whether that  
 732 word will be learned than the total frequency of exposure to  
 733 that word. In addition, Snow [42] found that mother's speech  
 734 to two-year-old is much simpler and less redundant than their  
 735 speech to ten-year-old. This indicates that young children have  
 736 available a sample of speech which is simpler, more redundant  
 737 and less confusing than normal adult speech.  
 738

The proposed algorithm was tested on language to vision  
 739 mapping and on language to tactile mapping [59], and it  
 740 can be easily extended to language to any other modalities  
 741 mapping. The mapping between multiple modalities and  
 742 words was researched in [47]. Fazly *et al.* [14] proposed a  
 743 probabilistic model of cross-situational learning where they  
 744 considered sentences that contain objects and their motion.  
 745 Monaghan *et al.* [33] studied differences between cross-  
 746 situational learning of nouns and verbs on human participants  
 747 as an extension of Tomasello and Akhtar's [50] work and  
 748 Schwartz and Terrell's [39] work. Monaghan *et al.* [33] con-  
 749 sidered learning of verbs as difficult as learning of nouns  
 750 when presented in a syntactic context. The authors noticed  
 751 that nouns are learned quicker, but verbs and nouns can be  
 752 acquired simultaneously.  
 753

An important direction for extensions of the proposed model  
 754 deals with the grounding of language related not just to static  
 755 sequences but also to events which require a temporal axis. For  
 756 example, Crick and Scassellati [7] studied the interconnection  
 757 between verbal narratives and episodes of intentional relative  
 758 motion with the goal that the robot can learn from observing  
 759 a game the rules of the game, relationships between players  
 760 and their goals and intentions and then participate in the play.  
 761

The main goal of this paper is to analyze mapping between  
 762 modalities. Thus, we kept processing of the individual modalities  
 763 deliberately simple for better understanding of cross-  
 764 situational learning. We performed further experiments with an  
 765 artificially generated dataset where the visual features, as well  
 766 as noise in the linguistic and the visual domain, were under  
 767 full control. These experiments are described in detail in [58].  
 768

## APPENDIX

769  
 Here we describe an extension of the proposed algorithm  
 770 presented in Section II-B for a case where we have sen-  
 771 tences with variable structure. This means that we cannot  
 772

**Algorithm 2** Sequential Mapping—Variable Sentence

**Inputs:**  
language clusters  $L_i$  ( $i \in 1 : I$ ), visual clusters  $K_j^k \sim N(\mathbf{m}^k, \mathbf{S}^k)$ ,  
 $j \in 1 : J^k$  for each feature  $k$ , visual input  $\mathbf{x}_t = \{\mathbf{x}_t^1, \dots, \mathbf{x}_t^K\}$  and  
corresponding language data  $W_t = \{\mathbf{w}_t^1, \dots, \mathbf{w}_t^{N_t}\}$  for each trial  $t$   
and number of clusters  $J^k$

**Output:**  
mapping between all visual classes  $K_j^k$  and language classes  $L_i$

---

```

while  $\sum J^k > 0$  and  $\mathbf{X}$  is not empty do
     $\eta_t^n \leftarrow$  assign each word  $\mathbf{w}_t^n$  from each sentence  $t$  to a language cluster (Winner-takes-all, see Eq. (15),  $\eta_t^n = \text{argmax}_i(P(w_t^n | L_i))$ )
    for  $k \in \{\text{size, color, orientation, texture, shape}\}$  do
         $v_t^k \leftarrow$  assign each data point  $\mathbf{x}_t^k$  to a visual cluster (Winner-takes-all, see Eq. (15),  $v_t^k = \text{argmax}_j(P(\mathbf{x}_t^k | K_j^k))$ )
        for  $j = 1 : J^k$  do
            for  $i = 1 : I$  do
                 $T_{ij}^k \leftarrow$  how many times did visual class  $i$  co-occur with language class  $j$  ( $T_{ij}^k = \sum_{t: v_t^k == j} \sum_n \delta(\eta_t^n, i)$ , where  $\delta$  is Kronecker delta)
            end for
        end for
        end for
         $[km, im, jm] \leftarrow \text{argmax}_k \text{argmax}_i \text{argmax}_j T_{ij}^k$  (the visual cluster  $K_{jm}^{km}$  is mapped to the language cluster  $L_{im}$ )
         $\mathbf{X}_{del}^{km} \leftarrow$  data points assigned to  $K_{jm}^{km}$  and  $L_{im}$ 
         $\theta_{new,jk}^k \leftarrow N(\mathbf{x}_{del}^k)$  learn Gaussian on the to be deleted data
         $\mathbf{X}^{km} \leftarrow \mathbf{X}^{km} \setminus \mathbf{X}_{del}^{km}$  delete data assigned to  $K_{jm}^{km}$  and  $L_{im}$ 
         $J^{km} \leftarrow J^{km} - 1$ 
        relearn  $K_j^{km} \sim N(\mathbf{m}^k, \mathbf{S}^k)$  on new data  $\mathbf{x}^{km}$  with  $J^{km}$  clusters
    end while
cluster visual data using new  $\theta_{new}^k$  parameters (cluster centers  $\mathbf{m}^k$ , covariance matrices  $\mathbf{S}^k$ ) and perform one-step mapping (Model 1)

```

---

773 directly associate words from the sentence with individual  
774 visual features, and therefore, associations to all visual features  
775 must be taken into account. The whole algorithm is described  
776 in Algorithm 2 in the form of a pseudocode.

777 Algorithm 2 can be further extended when we incorpo-  
778 rate the language model and corresponding probabilities of  
779 sequences of individual visual features (see [58]).

780

## REFERENCES

- 781 [1] L. W. Barsalou, "Grounded cognition," *Annu. Rev. Psychol.*, vol. 59,  
782 pp. 617–645, Jan. 2008.
- 783 [2] F. B. Klein. (2016). *GWR and GNG Classifier—File Exchange—MATLAB Central*. [Online]. Available:  
784 <http://uk.mathworks.com/matlabcentral/fileexchange/57798-gwr-and-gng-classifier>
- 785 [3] M. R. Brent and J. M. Siskind, "The role of exposure to isolated words in  
786 early vocabulary development," *Cognition*, vol. 81, no. 2, pp. B33–B44,  
787 2001.
- 788 [4] C. Büchel, C. Price, and K. Friston, "A multimodal language region in  
789 the ventral visual pathway," *Nature*, vol. 394, no. 6690, pp. 274–277,  
790 1998.
- 791 [5] A. Angelosi, A. Greco, and S. Harnad, "From robotic toil to symbolic  
792 theft: Grounding transfer from entry-level to higher-level categories,"  
793 *Connection Sci.*, vol. 12, no. 2, pp. 143–162, 2000.
- 794 [6] S. Coradeschi, A. Loutfi, and B. Wrede, "A short review of symbol  
795 grounding in robotic and intelligent systems," *Künstliche Intelligenz*,  
796 vol. 27, no. 2, pp. 129–136, 2013, doi: 10.1007/s13218-013-0247-2.
- 797 [7] C. Crick and B. Scassellati, "Controlling a robot with intention derived  
798 from motion," *Topics Cogn. Sci.*, vol. 2, no. 1, pp. 114–126, 2010.
- 799 [8] J. C. Culham and K. F. Valyear, "Human parietal cortex in action,"  
800 *Current Opin. Neurobiol.*, vol. 16, no. 2, pp. 205–212, 2006.
- 801 [9] M. Daoutis and N. Mavridis, "Towards a model for grounding semantic  
802 composition," in *Proc. Artif. Intell. Simulat. Behav. (AISB)*, 2014.
- 803 [10] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood  
804 from incomplete data via the EM algorithm," *J. Roy. Stat. Soc. B (Methodol.)*, vol. 39, no. 1, pp. 1–38, 1977.
- 805 [11] J. Donahue *et al.*, "Long-term recurrent convolutional networks for  
806 visual recognition and description," in *Proc. IEEE Conf. Comput. Vis.  
807 Pattern Recognit.*, Boston, MA, USA, 2015, pp. 2625–2634.
- 808 [12] E. Elhamifar and R. Vidal, "Robust classification using structured sparse  
809 representation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.  
810 (CVPR)*, Colorado Springs, CO, USA, 2011, pp. 1873–1879.
- 811 [13] E. Elhamifar and R. Vidal, "Sparse subspace clustering: Algorithm, the-  
812 ory, and applications," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35,  
813 no. 11, pp. 2765–2781, Nov. 2013.
- 814 [14] A. Fazly, A. Alishahi, and S. Stevenson, "A probabilistic computational  
815 model of cross-situational word learning," *Cogn. Sci.*, vol. 34, no. 6,  
816 pp. 1017–1063, 2010.
- 817 [15] M. Fleischman and D. Roy, "Grounded language modeling for automatic  
818 speech recognition of sports video," in *Proc. ACL*, Columbus, OH, USA,  
819 2008, pp. 121–129.
- 820 [16] M. C. Frank, N. D. Goodman, and J. B. Tenenbaum, "Using speakers'  
821 referential intentions to model early cross-situational word learning,"  
822 *Psychol. Sci.*, vol. 20, no. 5, pp. 578–585, 2009.
- 823 [17] V. Gligozzi, J. Mayor, J. F. Hu, and K. Plunkett, "Labels as features  
824 (not names) for infant categorization: A neurocomputational approach,"  
825 *Cogn. Sci.*, vol. 33, no. 4, pp. 709–738, 2009.
- 826 [18] E. Goldberg, N. Driedger, and R. I. Kittredge, "Using natural-language  
827 processing to produce weather forecasts," *IEEE Expert*, vol. 9, no. 2,  
828 pp. 45–53, Apr. 1994.
- 829 [19] P. Gorniak and D. Roy, "Speaking with your sidekick: Understanding  
830 situated speech in computer role playing games," in *Proc. AIIDE*,  
831 Marina Del Rey, CA, USA, 2005, pp. 57–62.
- 832 [20] S. Harnad, "The symbol grounding problem," *Phys. D Nonlin.  
833 Phenomena*, vol. 42, nos. 1–3, pp. 335–346, 1990.
- 834 [21] A. Karpathy and L. Fei-Fei, "Deep visual-semantic alignments for gen-  
835 erating image descriptions," in *Proc. IEEE Conf. Comput. Vis. Pattern  
836 Recognit.*, Boston, MA, USA, 2015, pp. 3128–3137.
- 837 [22] J. Kennedy, P. Baxter, and T. Belpaeme, "Comparing robot embodiments  
838 in a guided discovery learning interaction with children," *Int. J. Soc.  
839 Robot.*, vol. 7, no. 2, pp. 293–308, 2015.
- 840 [23] P. Lamere *et al.*, "The CMU SPHINX-4 speech recognition system," in  
841 *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, vol. 1,  
842 Hong Kong, 2003, pp. 2–5.
- 843 [24] P. Li, I. Farkas, and B. MacWhinney, "Early lexical development in  
844 a self-organizing neural network," *Neural Netw.*, vol. 17, nos. 8–9,  
845 pp. 1345–1362, 2004.
- 846 [25] E. M. Markman, "Constraints children place on word meanings," *Cogn.  
847 Sci.*, vol. 14, no. 1, pp. 57–77, 1990.
- 848 [26] S. Marsland, J. Shapiro, and U. Nehmzow, "A self-organising  
849 network that grows when required," *Neural Netw.*, vol. 15, nos. 8–9,  
850 pp. 1041–1058, 2002.
- 851 [27] N. Mavridis, "Grounded situation models for situated conversational  
852 assistants," Ph.D. dissertation, Massachusetts Inst. Technol., Cambridge,  
853 MA, USA, 2007.
- 854 [28] N. Mavridis, "A review of verbal and non-verbal human–robot  
855 interactive communication," *J. Robot. Auton. Syst.*, vol. 63, pp. 22–35,  
856 Jan. 2015.
- 857 [29] B. McMurray, J. S. Horst, and L. K. Samuelson, "Word learning emerges  
858 from the interaction of online referent selection and slow associative  
859 learning," *Psychol. Rev.*, vol. 119, no. 4, pp. 831–877, 2012.
- 860 [30] G. Metta *et al.*, "The iCub humanoid robot: An open platform for  
861 research in embodied cognition," in *Proc. ACM 8th Workshop Perform.  
862 Metrics Intell. Syst.*, Gaithersburg, MD, USA, 2008, pp. 50–56.
- 863 [31] G. Metta, P. Fitzpatrick, and L. Natale, "YARP: Yet another robot  
864 platform," *Int. J. Adv. Robot. Syst.*, vol. 3, no. 1, pp. 43–48, 2006.
- 865 [32] M. Mishkin, L. G. Ungerleider, and K. A. Macko, "Object vision and  
866 spatial vision: Two cortical pathways," *Trends Neurosci.*, vol. 6, no. 10,  
867 pp. 414–417, 1983.
- 868 [33] P. Monaghan, K. Mattock, R. A. I. Davies, and A. C. Smith,  
869 "Gavagai is as Gavagai does: Learning nouns and verbs from cross-  
870 situational statistics," *Cogn. Sci.*, vol. 39, no. 5, pp. 1099–1112, 2015,  
871 doi: 10.1111/cogs.12186.
- 872 [34] N. Otsu, "A threshold selection method from gray-level histograms,"  
873 *IEEE Trans. Syst., Man, Cybern.*, vol. SMC-9, no. 1, pp. 62–66,  
874 Jan. 1979.
- 875 [35] W. V. Quine, "On the reasons for indeterminacy of translation," *J.  
876 Philos.*, vol. 67, no. 6, pp. 178–183, 1970.
- 877 [36] T. Regier, "The emergence of words: Attentional learning in form and  
878 meaning," *Cogn. Sci.*, vol. 29, no. 6, pp. 819–865, 2005.
- 879 [37] D. Roy and N. Mukherjee, "Towards situated speech understanding:  
880 Visual context priming of language models," *Comput. Speech Lang.*,  
881 vol. 19, no. 2, pp. 227–248, 2005.
- 882 [38] D. K. Roy, "Learning visually grounded words and syntax for a scene  
883 description task," *Comput. Speech Lang.*, vol. 16, no. 3, pp. 353–385,  
884 2002.
- 885 [AQ3]

AQ5

- 888 [39] R. G. Schwartz and B. Y. Terrell, "The role of input frequency in lexical  
889 acquisition," *J. Child Lang.*, vol. 10, no. 1, pp. 57–64, 1983.  
 890 [40] J. M. Siskind, "A computational study of cross-situational techniques  
891 for learning word-to-meaning mappings," *Cognition*, vol. 61, nos. 1–2,  
892 pp. 39–91, 1996.  
 893 [41] K. Smith, A. D. M. Smith, R. A. Blythe, and P. Vogt, *Cross-Situational  
894 Learning: A Mathematical Approach* (LNCS 4211). Heidelberg,  
895 Germany: Springer, 2006, pp. 31–44.  
 896 [42] C. E. Snow, "Mothers' speech to children learning language," *Child  
897 Dev.*, vol. 43, no. 2, pp. 549–565, 1972.  
 898 [43] G. Spitsyna, J. E. Warren, S. K. Scott, F. E. Turkheimer, and R. J. Wise,  
899 "Converging language streams in the human temporal lobe," *J. Neurosci.*,  
900 vol. 26, no. 28, pp. 7328–7336, 2006.  
 901 [44] F. Stramandinoli, D. Marocco, and A. Cangelosi, "The grounding of  
902 higher order concepts in action and language: A cognitive robotics  
903 model," *Neural Netw.*, vol. 32, pp. 165–173, Aug. 2012.  
 904 [45] Y. Sugita and J. Tani, "Learning semantic combinatoriality from the  
905 interaction between linguistic and behavioral processes," *Adapt. Behav.*,  
906 vol. 13, no. 1, pp. 33–52, 2005.  
 907 [46] M. Taddeo and L. Floridi, "Solving the symbol grounding problem: A  
908 critical review of fifteen years of research," *J. Exp. Theor. Artif. Intell.*,  
909 vol. 17, no. 4, pp. 419–445, 2005.  
 910 [47] A. Taniguchi, T. Taniguchi, and A. Cangelosi, "Multiple categorization  
911 by iCub: Learning relationships between multiple modalities and  
912 words," in *Proc. IROS Workshop Mach. Learn. Methods High Level  
913 Cogn. Capabilities Robot.*, 2016.  
 914 [48] V. Tikhanoff *et al.*, "An open-source simulator for cognitive robotics  
915 research: The prototype of the iCub humanoid robot simulator," in *Proc.  
916 ACM 8th Workshop Perform. Metrics Intell. Syst.*, Gaithersburg, MD,  
917 USA, 2008, pp. 57–61.  
 918 [49] V. Tikhanoff, A. Cangelosi, and G. Metta, "Integration of speech and  
919 action in humanoid robots: iCub simulation experiments," *IEEE Trans.  
920 Auton. Mental Develop.*, vol. 3, no. 1, pp. 17–29, Mar. 2011.  
 921 [50] M. Tomasello and N. Akhtar, "Two-year-olds use pragmatic cues to  
922 differentiate reference to objects and actions," *Cogn. Dev.*, vol. 10, no. 2,  
923 pp. 201–224, 1995.  
 924 [51] M. Vavrečka and I. Farkaš, "A multimodal connectionist architecture  
925 for unsupervised grounding of spatial language," *Cogn. Comput.*, vol. 6,  
926 no. 1, pp. 101–112, 2014.  
 927 [52] P. Vogt, "Language evolution and robotics: Issues on symbol grounding,"  
928 in *Artificial Cognition Systems*. Hershey, PA, USA: IGI Glob., 2006,  
929 p. 176.  
 930 [53] J. F. Werker and P. J. McLeod, "Infant preference for both male and  
931 female infant-directed talk: A developmental study of attentional and  
932 affective responsiveness," *Can. J. Psychol. Revue Can. De Psychol.*,  
933 vol. 43, no. 2, pp. 230–246, 1989.  
 934 [54] F. Xu and J. B. Tenenbaum, "Word learning as Bayesian inference,"  
935 *Psychol. Rev.*, vol. 114, no. 2, pp. 245–272, 2007.  
 936 [55] C. Yu and L. B. Smith, "Rapid word learning under uncertainty via  
937 cross-situational statistics," *Psychol. Sci.*, vol. 18, no. 5, pp. 414–420,  
938 2007.  
 939 [56] C. Yu and L. B. Smith, "Modeling cross-situational word-referent  
940 learning: Prior questions," *Psychol. Rev.*, vol. 119, no. 1, pp. 21–39,  
941 2012.  
 942 [57] D. Yurovsky, C. Yu, and L. B. Smith, "Competitive processes in cross-  
943 situational word learning," *Cogn. Sci.*, vol. 37, no. 5, pp. 891–921, 2013.  
 944 [58] K. Štepánová, "Hierarchical probabilistic model of language acquisi-  
945 tion," Ph.D. dissertation, Dept. Cybern., Czech Tech. Univ. Prague,  
946 Prague, Czech Republic, 2016.  
 947 [59] K. Štepánová *et al.*, "Where is my forearm? Clustering body parts from  
948 simultaneous tactile and linguistic input," in *Proc. Cogn. Artif. Life  
949 (KUZ XVII)*, 2017.



**Karla Štepánová** received the master's degree in condensed matter physics from Charles University, Prague, Czech Republic, and the Ph.D. degree in artificial intelligence and biocybernetics from Czech Technical University in Prague, Prague, in 2017.

She is a Researcher with the Czech Institute of Informatics, Robotics, and Cybernetics, Prague. She was a Visiting Researcher with Plymouth University, Plymouth, U.K., in 2016. Her current research interests include probabilistic models of cognition, unsupervised learning, language acquisition, and multimodal integration.



**Frederico Belmonte Klein** was born in Porto Alegre, Brazil, in 1982. He received the Electrical Engineering degree and the medical degree from the Federal University of Rio Grande do Sul, Porto Alegre, in 2006 and 2014, respectively. He is currently pursuing the Ph.D. degree in cognitive robotics with the University of Plymouth, Plymouth, U.K.



**Angelo Cangelosi** received the degree in psychology and cognitive science from the University of Rome La Sapienza, Rome, Italy, and the University of Genoa, Genoa, Italy.

He is a Professor of artificial intelligence and cognition and the Director of the Centre for Robotics and Neural Systems with Plymouth University, Plymouth, U.K. and was a Visiting Scholar with the University of California at San Diego, San Diego, CA, USA, and the University of Southampton, Southampton, U.K. His current research interests include language grounding and embodiment in humanoid robots, developmental robotics, human–robot interaction, and on the application of neuromorphic systems for robot learning.



**Michal Vavrečka** received the Ph.D. degree in general psychology from the Faculty of Social Studies, in 2008.

He is with the Czech Institution for Informatics Robotics and Cybernetics, Czech Technical University in Prague, Prague, Czech Republic. He develops cognitive architecture for symbol grounding, language acquisition, and knowledge representation.

## AUTHOR QUERIES

### AUTHOR PLEASE ANSWER ALL QUERIES

**PLEASE NOTE:** We cannot accept new source files as corrections for your paper. If possible, please annotate the PDF proof we have sent you with your corrections and upload it via the Author Gateway. Alternatively, you may send us your corrections in list format. You may also upload revised graphics via the Author Gateway.

AQ1: Please confirm/give details of funding source.

AQ2: Please provide the postal code for “Plymouth University Plymouth, U.K.”

AQ3: Please provide the page range for References [9], [47], and [59].

AQ4: Please provide the department name for Reference [27].

AQ5: Please confirm if the location and publisher information for Reference [41] is correct as set.

AQ6: Please specify the type of the degree and year attained by the author “A. Cangelosi.”

AQ7: Please provide the organization name and location for the Ph.D. degree obtained by “M. Vavrečka” biography.