

Dear Author,

Here are the proofs of your article.

- You can submit your corrections **online**, via **e-mail** or by **fax**.
- For **online** submission please insert your corrections in the online correction form. Always indicate the line number to which the correction refers.
- You can also insert your corrections in the proof PDF and **email** the annotated PDF.
- For fax submission, please ensure that your corrections are clearly legible. Use a fine black pen and write the correction in the margin, not too close to the edge of the page.
- Remember to note the **journal title**, **article number**, and **your name** when sending your response via e-mail or fax.
- **Check** the metadata sheet to make sure that the header information, especially author names and the corresponding affiliations are correctly shown.
- **Check** the questions that may have arisen during copy editing and insert your answers/ corrections.
- **Check** that the text is complete and that all figures, tables and their legends are included. Also check the accuracy of special characters, equations, and electronic supplementary material if applicable. If necessary refer to the *Edited manuscript*.
- The publication of inaccurate data such as dosages and units can have serious consequences. Please take particular care that all such details are correct.
- Please **do not** make changes that involve only matters of style. We have generally introduced forms that follow the journal's style. Substantial changes in content, e.g., new results, corrected values, title and authorship are not allowed without the approval of the responsible editor. In such a case, please contact the Editorial Office and return his/her consent together with the proof.
- If we do not receive your corrections **within 48 hours**, we will send you a reminder.
- Your article will be published **Online First** approximately one week after receipt of your corrected proofs. This is the **official first publication** citable with the DOI. **Further changes are, therefore, not possible.**
- The **printed version** will follow in a forthcoming issue.

Please note

After online publication, subscribers (personal/institutional) to this journal will have access to the complete article via the DOI using the URL: [http://dx.doi.org/\[DOI\]](http://dx.doi.org/[DOI]).

If you would like to know when your article has been published online, take advantage of our free alert service. For registration and further information go to: <http://www.springerlink.com>.

Due to the electronic nature of the procedure, the manuscript and the original figures will only be returned to you on special request. When you return your corrections, please inform us if you would like to have these documents returned.

Metadata of the article that will be visualized in OnlineFirst

Please note: Images will appear in color online but will be printed in black and white.

ArticleTitle	A Multimodal Connectionist Architecture for Unsupervised Grounding of Spatial Language	
Article Sub-Title		
Article CopyRight	Springer Science+Business Media New York (This will be the copyright line in the final PDF)	
Journal Name	Cognitive Computation	
Corresponding Author	Family Name	Vavrečka
	Particle	
	Given Name	Michal
	Suffix	
	Division	Department of Cybernetics
	Organization	Czech Technical University
	Address	Karlovo náměstí 13, Prague, Czech Republic
	Email	vavrecka@fel.cvut.cz
Author	Family Name	Farkaš
	Particle	
	Given Name	Igor
	Suffix	
	Division	Department of Applied Informatics
	Organization	Comenius University
	Address	Mlynská dolina, 84248, Bratislava, Slovakia
	Email	
Schedule	Received	30 August 2012
	Revised	
	Accepted	12 March 2013
Abstract	<p>We propose a bio-inspired unsupervised connectionist architecture and apply it to grounding the spatial phrases. The two-layer architecture combines by concatenation the information from the visual and the phonological inputs. In the first layer, the visual pathway employs separate 'what' and 'where' subsystems that represent the identity and spatial relations of two objects in 2D space, respectively. The bitmap images are presented to an artificial retina and the phonologically encoded five-word sentences describing the image serve as the phonological input. The visual scene is hence represented by several self-organizing maps (SOMs) and the phonological description is processed by the Recursive SOM that learns to topographically represent the spatial phrases, represented as five-word sentences (e.g., 'blue ball above red cup'). Primary representations from the first-layer modules are unambiguously integrated in a multimodal second-layer module, implemented by the SOM or the 'neural gas' algorithms. The system learns to bind proper lexical and visual features without any prior knowledge. The simulations reveal that separate processing and representation of the spatial location and the object shape significantly improve the performance of the model. We provide quantitative experimental results comparing three models in terms of their accuracy.</p>	
Keywords (separated by '-')	Unsupervised learning - Self-organizing map - Symbol grounding - Spatial phrases	
Footnote Information		

A Multimodal Connectionist Architecture for Unsupervised Grounding of Spatial Language

Michal Vavrečka · Igor Farkaš

Received: 30 August 2012 / Accepted: 12 March 2013
© Springer Science+Business Media New York 2013

Abstract We propose a bio-inspired unsupervised connectionist architecture and apply it to grounding the spatial phrases. The two-layer architecture combines by concatenation the information from the visual and the phonological inputs. In the first layer, the visual pathway employs separate ‘what’ and ‘where’ subsystems that represent the identity and spatial relations of two objects in 2D space, respectively. The bitmap images are presented to an artificial retina and the phonologically encoded five-word sentences describing the image serve as the phonological input. The visual scene is hence represented by several self-organizing maps (SOMs) and the phonological description is processed by the Recursive SOM that learns to topographically represent the spatial phrases, represented as five-word sentences (e.g., ‘blue ball above red cup’). Primary representations from the first-layer modules are unambiguously integrated in a multimodal second-layer module, implemented by the SOM or the ‘neural gas’ algorithms. The system learns to bind proper lexical and visual features without any prior knowledge. The simulations reveal that separate processing and representation of the spatial location and the object shape significantly improve the performance of the model. We provide quantitative experimental results comparing three models in terms of their accuracy.

Keywords Unsupervised learning · Self-organizing map · Symbol grounding · Spatial phrases

M. Vavrečka (✉)
Department of Cybernetics, Czech Technical University,
Karlovo náměstí 13, Prague, Czech Republic
e-mail: vavrecka@fel.cvut.cz

I. Farkaš
Department of Applied Informatics, Comenius University,
Mlynská dolina, 84248 Bratislava, Slovakia

Introduction

The question of how to acquire, represent and use knowledge in the learning agent is fundamental in artificial intelligence and cognitive science research. Within the modern perspective, fueled by growing empirical evidence, we are looking for a system that, through interaction with the environment, learns the internal representations. These should store the constant attributes and regularities of the environment, giving rise to forming concepts, which become connected to the symbolic level (language). This approach to the representation of meaning differs from the classical symbolic (designer) approach based on formal principles [28, 33], which are subjected to the symbol grounding problem (Harnad 1990).

Harnad (1990) proposed a hybrid architecture based on discrimination and identification, where the former process is considered a subsymbolic (non-arbitrary) representation of perceptual inputs, while the latter assigns (non-arbitrary) concepts to (arbitrary) symbols. Harnad used neural networks for the subsymbolic representations and the classical architecture for symbol operations. In the overview of grounding architectures, Taddeo and Floridi [41] introduced the zero semantical commitment condition as a criterion for valid solution to the symbol grounding problem, completely avoiding the designer approach. This criterion, however, appears unsatisfiable not only in artificial, but in living systems as well [46].

In the past two decades, there was a number of different approaches and models of the symbol grounding (e.g., [1, 4, 7, 11, 23, 37, 40, 42]). These models typically ground linguistic symbols by linking them with agent’s sensorimotor behavior, or with objects and their features. Other approaches, instead, focus on the social symbol grounding where the symbols become grounded by simulating the cultural evolution in a population of agents (e.g., [38, 52]).

69 With respect to learning paradigms, we can distinguish
70 two types of connectionist models that link subsymbolic
71 (conceptual) knowledge with (linguistic) symbols. The
72 supervised approach is based on error correction learning
73 in which input patterns are linked with symbolic targets
74 (labels). The models listed above typically have this fea-
75 ture. Both inputs and outputs are assumed to be provided
76 by the environment, and the error information is used to
77 find the desired mapping between them.

78 On the other hand, the unsupervised approach treats both
79 perceptual stimuli and symbols equally as inputs, to be
80 associated (typically) by Hebbian-like learning. This implies
81 a different way of incorporating the symbolic (lexical) level.
82 The target signal (here the lexical level) only functions as an
83 additional input rather than being the source for error-based
84 learning. The unsupervised models are typically based on
85 self-organizing maps (SOM) that organize (high-dimen-
86 sional) input vectors according to their similarities [19]. For
87 instance, the DevLex model [20] also consists of two self-
88 organizing networks, one for lexical symbols (phonological
89 representations) and the other for conceptual (semantic)
90 representations, that are bidirectionally connected. They can
91 activate each other, but there is no additional layer for mul-
92 timodal representations, as opposed to the model proposed
93 here, and some other models of grounding (e.g., [4, 36]).

94 Dorffner et al. [1] have proposed unsupervised binding
95 between two primary (symbolic and conceptual) layers
96 mediated by the central linking layer. The linking layer
97 (which could be seen as the bimodal layer) interconnects
98 the two primary layers via its localist units that link both
99 representations (i.e., one unit connects one word-concept
100 pair of primary representations). First, one set of links
101 (weights to the linking layer) is trained using a competitive
102 mechanism exploiting the winner-take-all approach. Then,
103 the winner's weights toward the other layer are updated
104 according to the outstar rule [12]. Similarly, to DevLex,
105 these mappings were aimed at simulating word compre-
106 hension (the form to meaning) and word production
107 (meaning to the form).

108 Among the unsupervised approaches, there emerged an
109 alternative to link both the perceptual and symbolic
110 information (treated as an input) with multimodal repre-
111 sentations at the output. The example of this architecture is
112 the unsupervised feature-based model that was used to
113 account for early category formation in young infants [9].
114 Interestingly, this approach postulates the unsupervisory
115 role of linguistic labels that can effect categorization during
116 the acquisition process, which has also been supported by
117 experimental evidence.

118 The idea of unsupervised binding of two modalities
119 (as inputs) has also been applied in recent generative
120 probabilistic models such as the deep belief net (DBN).
121 The DBN was successfully trained to classify the isolated

hand-written digits, so the visual inputs were linked with
categorical labels [14]. The linking was established via the
training on image-label pairs, using the higher (bimodal)
layer that learned the joint distribution of input pairs.

Our model is similar to that of Gliozzi et al. [9] by treating
the information from two modalities as input. It differs from
it, however, by higher complexity and the task. Our model
was designed for grounding the spatial phrases rather than
object names (typical for early language learning). We test
our model in the area of spatial cognition, similarly to Regier
[34], who created a supervised neural network model con-
sisting of several modules to ground the spatial phrases.
Regier's model was able to ground both static spatial rela-
tions (e.g., left, right) and dynamic relations (e.g., around,
through). However, in Regier's model, the symbolic repre-
sentational level was considered prior and fixed. On the
contrary, we focus on unsupervised learning of spatial rela-
tions of two objects in 2D space, by linking perceptual
information and linguistic description, where neither level is
considered prior and fixed. The neural architecture we pro-
pose satisfies the requirement that the artificial system
(agent) should learn its own functions and representations
[53].

In this paper, we describe the 'experimental trajectory'
of our work whose aim was to design a bio-inspired model
at a reasonable level of abstraction. We converged to a
model that processes visual input separately using 'what'
and 'where' pathways, which is also a feature of biological
systems [45]. The motivation for our model was to
experimentally test whether it is possible (without errors)
to bind location, color and shape of two objects (Visual
Feature-Binding) without any prior knowledge and without
external information. The model also proposes a solution to
the (unsupervised) symbol grounding that can be consid-
ered as a temporal synchrony [6]. In this process, the
sequences of symbols (words), describing the spatial layout
of two objects and their identity, processed in the phono-
logical layer are grounded (bound) to proper features from
the visual subsystem (shape, color and location). Our
model exploits the simplification, being the fixed sentence
structure that facilitates the thematic role assignment in the
model(s).

The benefit of modularity in the model (including that for
separation of 'what' and 'where' information) was already
emphasized in earlier works. For instance, Jacobs et al. [15]
proposed a supervised approach to designing a modular
system, composed of competing expert networks and the
gating network, that could simultaneously learn two differ-
ent tasks. They applied their model to the learning of the
'where' and 'what' information (using simple bitmap ima-
ges) and pointed to the advantages of this modular feedfor-
ward architecture compared to the standard multi-layer
perceptron. In our models, the units also compete for inputs,

175 albeit using the principles of self-organization. Unlike
176 Jacobs et al., the modules are assumed to be given in our
177 models (the competition does not occur at the level of
178 modules, but rather the level of individual units).

179 The rest of the paper is organized as follows. In section
180 “The Models”, we introduce the architecture(s) of our
181 models in a greater detail. Section “Results” presents results
182 from four series of simulations. Section “Discussion” con-
183 tains the discussion about our final model and its relation to
184 other models. Section “Conclusion” concludes the paper.

185 The Models

186 In our model, the representation process takes advantage of
187 the unimodal layers of units. The visual layers represent
188 spatial location, shape and color of objects and the phono-
189 logical layer represents sentences. The multimodal level
190 integrates the outputs of these unimodal layers. In contrast to
191 the classical approaches that postulate the abstract symbolic
192 level as fixed and prior (defined by the designer), our model
193 offers possibility to learn and modify the phonological layer,
194 visual layer and, consequently, the multimodal level. The
195 schema of the system is depicted in Fig. 1.

196 In the simulations, we compare different versions of the
197 visual subsystem, analyzing the distinction between ‘what’

and ‘where’ pathways. The results help us to decide whether
198 this simplification is important for enhancing the overall
199 model performance. The visual system of our model is
200 therefore tested in three different configurations (see Fig. 1;
201 Table 1): a single SOM that learns to capture both ‘what’ and
202 in Model 3 information (Model 1), two separate SOMs for
203 ‘what’ and ‘where’ information (Model 2), and two separate
204 SOMs with reduced ‘where’ representations (Model 3).
205

206 In the last simulation, we compare two different types of
207 multimodal integration. Inspired by the biological evidence
208 about topographic organization of sensory and motor brain
209 areas, we assume that primary unimodal layers are topo-
210 graphically organized. Although examples of this organiza-
211 zing principle exist in higher areas as well [22], it
212 remains an empirical question whether topographically
213 organized responses are a general principle of the brain at
214 higher levels of organization. In the multimodal layer, we
215 hence compare the SOM and ‘neural gas’ (NG; [25])
216 algorithms as representatives of topographic and non-
217 topographic approaches, respectively. Both algorithms are
218 unsupervised, based on the competition among units, but
219 NG uses a flexible neighborhood function, as opposed to
220 the fixed neighborhood in SOM (that enforces topography).
221 The goal was to experimentally investigate the effect of the
222 neighborhood function in the multimodal layer. We used
223 the modified SOM Toolbox [50] for all simulations.

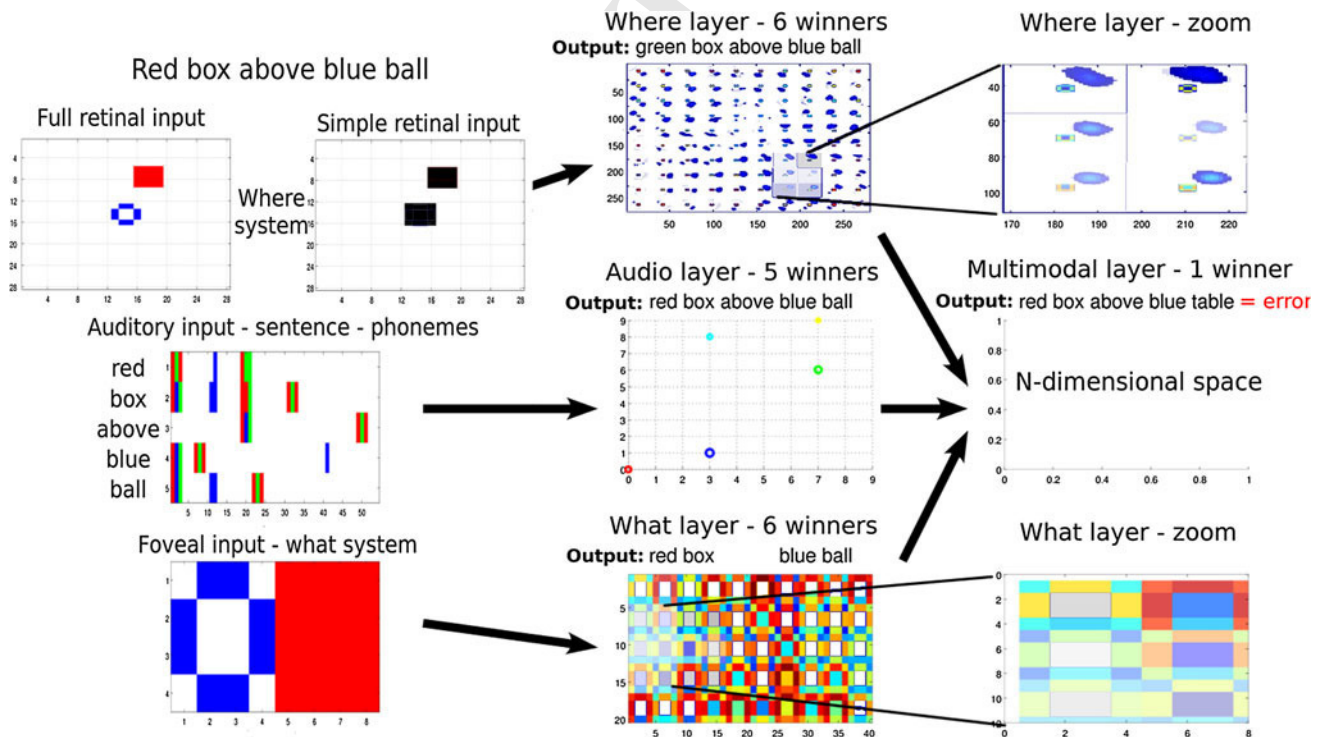


Fig. 1 Multimodal connectionist architecture for grounding spatial phrases. The phonological layer represents sentences and the visual layers represent spatial location, shape and color of objects. The multimodal level integrates the outputs of these unimodal layers

Table 1 Summary of the visual features of the 3 models used in experiments. Each model uses the same phonological subsystem (RecSOM) and can be combined with the SOM or the NG module in the multimodal layer

Model	Visual input	Visual system
1	full	single SOM
2	'where'- full	'where' SOM
	'what' -2 objects' features	'what' SOM
3	'where' -blobs	'where' SOM
	'what' -2 objects' features	'what' SOM

224 Visual Input

225 The visual scenes consist of the trajector and the base objects
 226 in different spatial configurations. The base position is fixed
 227 in the center of the scene (the center of the retina) and the
 228 trajector position is located in one (or at the boundary
 229 between two) of the spatial quadrants relative to the base.
 230 The positions along the main semiaxes are linguistically
 231 referred to as up, down, left and right, but perceptually, the
 232 trajector position is fuzzy and random. The scene size (arti-
 233 ficial retina) contains 28×28 pixels and both objects consist
 234 of 4×4 pixels (Fig. 2a). The color of each pixel is
 235 encoded by the activity level, scaled to values between 0 and
 236 1 (0 = white, 0.33 = red, 0.66 = green, 1 = blue). Inputs to the
 237 visual SOM (Model 1) or just 'where' subsystem (Model 2
 238 and 3) are the 784-dimensional vectors. Each dimension
 239 represents the color information for Model 1 and 2, or
 240 monochrome activity (0 = white, 1 = black) for Model 3.
 241 Inputs to the 'what' subsystem (Model 2 and 3) are the
 242 32-dimensional vectors. The 'what' system incorporates a
 243 simple attentional mechanism and represents the foveal
 244 input of two consequently observed objects. Two visual
 245 fields (each with 4×4 receptors) simultaneously project
 246 visual information about the trajector and the base in a fixed
 247 position to the unimodal 'what' system. This subsystem
 248 represents color and shape of pairs of objects (trajector and
 249 base) irrespective of their spatial position.

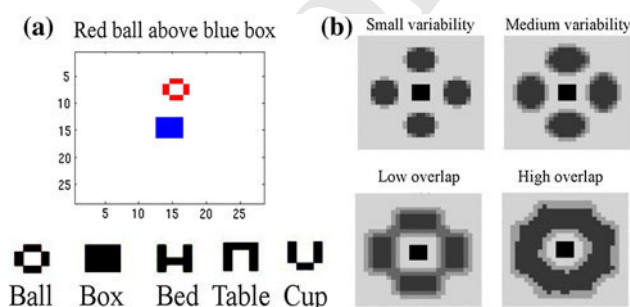


Fig. 2 a Example of a visual input scene and the monochrome visual 'vocabulary,' b Superimposed visual inputs with varying levels of spatial fuzziness

We trained all models with an increasing combinatorial 250
 complexity, starting with simple inputs with two colors, 251
 two object types and four spatial relations, up to more 252
 complex inputs consisting of three colors (red, green and 253
 blue), five object types (box, ball, table, cup and bed) and 254
 four spatial relations (above, below, left and right). The 255
 most complex scenario with two different objects in the 256
 scene amounts to 840 input configurations. The corre- 257
 sponding training set results in 42,000 examples (with 50 258
 instances per input configuration). We also present stimuli 259
 with increasing fuzziness in the spatial location to inves- 260
 tigate the relation between fuzziness and the error in the 261
 visual and multimodal layer. The fuzziness stands for 262
 variability of the trajector center with regard to the center 263
 of the spatial quadrant ranging from 2 to 8 pixels. The two 264
 conditions with the highest degree of fuzziness yield 265
 overlapping inputs (as seen in Fig. 2b). 266

Visual Subsystem 267

The sensory input of the visual subsystem is captured by an 268
 artificial retina that serves as an input to the primary visual 269
 layer. Visual layer consists of the SOM(s) that learn the 270
 nonlinear mapping of input vectors to output units in the 271
 topography-preserving manner (i.e., similar inputs are 272
 mapped to neighboring units in the map). The SOM per- 273
 forms standard computations in each iteration. After the 274
 presentation of a randomly chosen (rescaled) input vector 275
 \mathbf{x} , the output y_i of a unit i in the SOM is first computed as 276

where $\|\cdot\|$ denotes the Euclidean norm (it will also be used 278
 in forthcoming equations), and then, the k -WTA (winner- 279
 take-all) rule is applied. According to k -WTA, k most 280
 active units are proportionally kept active (with the activity 281
 of the best matching unit scaled to 1), and all other units are 282
 clamped to 0. In the models, we empirically found the 283
 optimal value $k = 6$. The motivation for this type of output 284
 representation rests in introducing some overlaps between 285
 similar patterns to facilitate generalization. 286

The output vectors of all unimodal modules are con- 287
 catenated (including the phonological input) and serve as 288
 the input vector to the multimodal layer. For all visual 289
 maps, standard computations are performed to update 290
 weights. Then, the best matching unit (winner) c is calcu- 291
 lated according to 292

$$c = \arg \min_i \{ \|\mathbf{x}(t) - \mathbf{w}_i(t)\| \},$$

the weights in the winner's neighborhood are updated as 294

$$\mathbf{w}_i(t+1) = \mathbf{w}_i(t) + \mu h_{ci}(t) [\mathbf{x}(t) - \mathbf{w}_i(t)],$$

where μ is the learning rate and $h_{ci}(t)$ is the neighborhood 296
 kernel around the winner c , with the neighborhood radius 297

298 linearly shrinking over time. Let us now take a more
299 detailed look at these layers and their inputs.

300 In Model 1, the single SOM was tested whether it could
301 learn to differentiate various positions of two objects, as
302 well as object types and their color. In Model 2, we used
303 separate SOMs for spatial locations (abstraction of the
304 ‘where’ system) and a separate SOM for color and shape of
305 objects (abstraction of the ‘what’ system). Model 3
306 employs the same ‘what’ and ‘where’ systems as Model 2,
307 but uses different inputs to the ‘where’ system consisting of
308 two monochromatic boxes (rather than concrete object
309 shapes in color) in the particular spatial position. The
310 dimension of all visual layers was fixed for all models,
311 namely $\dim(\mathbf{y}^{\text{what}}) = 25 \times 25$ neurons for the ‘what’ sys-
312 tem and $\dim(\mathbf{y}^{\text{where}}) = 30 \times 30$ neurons for the ‘where’
313 system. The similar size of matrices were estimated from
314 previous simulations [47], and they also stem from the
315 number of combinations in the most complex scenario (840
316 combinations in the ‘where’ system and 210 in the ‘what’
317 system). All SOM maps have a hexagonal neighborhood
318 function and the lattices have a toroid topology. The
319 overview of the characteristics of the three models is
320 summarized in Table 1.

321 Phonological Input

322 Phonological input (English sentences) was encoded as high-
323 dimensional patterns representing word forms using PatPho,
324 a generic phonological pattern generator that fits every word
325 (up to three syllables) onto a template according to its vowel-
326 consonant structure [21]. PatPho uses the concept of a syl-
327 labic template: a word representation is formed by combi-
328 nations of syllables in a metrical grid, and the slots in each
329 grid are made up by bundles of features that correspond to
330 consonants and vowels. Word representations can hence be
331 compared according to their phonological similarities. In our
332 case of 5-word sentences, each sentence consists of five
333 54-dimensional vectors with component values in the
334 interval (0,1).

335 Phonological Subsystem

336 The phonological input is fed (one vector at a time) to the
337 RecSOM [51], a recurrent SOM architecture, that uses a
338 detailed representation of the context information (the whole
339 output map activation) and has been demonstrated to be able
340 to learn to represent much richer dynamical behavior [44], in
341 comparison with other recurrent SOM models [13]. Rec-
342 SOM learns to represent the input (words) in the temporal
343 context (hence, capturing the sequential information). Rec-
344 SOM output, in terms of the map activation, feeds to the
345 multimodal layer, being integrated (by vector concatenation)
346 with the visual pathway. Like SOM, RecSOM is trained by a

347 competitive, Hebbian-like learning algorithm. As a property
348 of the RecSOM, its units become the sequence detectors after
349 training, topographically organized according to the suffix
350 (the most recent words).

351 Formally, each neuron $i \in \{1, 2, \dots, N\}$ in RecSOM has
352 two associated weight vectors: $\mathbf{w}_i \in \mathcal{R}^n$ – linked with an n -
353 dimensional input $\mathbf{s}(t)$ (in our case, the current word, with
354 dimension $n = 54$) feeding the network at time t , and the
355 weight vector $\mathbf{c}_i \in \mathcal{R}^N$ – linked with the context $\mathbf{y}(t-1) =$
356 $[y_1(t-1), y_2(t-1), \dots, y_N(t-1)]$ containing the unit acti-
357 vations $y_i(t-1)$ from the previous time step. The output of a
358 unit i at time t is $y_i(t) = \exp(-d_i(t))$, where

$$d_i(t) = \alpha \|\mathbf{s}(t) - \mathbf{w}_i\|^2 + \beta \|\mathbf{y}(t-1) - \mathbf{c}_i\|^2.$$

360 Here, $\alpha > 0$ and $\beta > 0$ are the model parameters that,
361 respectively, influence the effect of the input and the
362 context upon the neurons profile. Their suitable values are
363 usually found heuristically (in our model, we use $\alpha =$
364 $\beta = 0.1$). Both weight vectors are updated using the same
365 form of a SOM learning rule

$$\mathbf{w}_i(t+1) = \mathbf{w}_i(t) + \gamma h_{ci}(\mathbf{s}(t) - \mathbf{w}_i(t)),$$

$$\mathbf{c}_i(t+1) = \mathbf{c}_i(t) + \gamma h_{ci}(\mathbf{y}(t-1) - \mathbf{c}_i(t)),$$

367 where $c = \arg \min_i \{d_i(t)\}$, is the winner index at time t , and
368 $0 < \gamma < 1$ is the learning rate. (The winner can be equiva-
369 lently defined as the unit c with the highest activation $y_c(t)$:
370 $c = \arg \max_i \{y_i(t)\}$). The neighborhood function h_{ci} is a
371 Gaussian (of width σ) on the distance $d(i, c)$ of units i and c in
372 the map: $h_{ci} = \exp(-d(c, i)^2 / \sigma^2)$. The neighborhood width
373 σ linearly decreases in time to allow the formation of topo-
374 graphic representation of input sequences. After training, all
375 RecSOM units become sensitive to particular sentences,
376 ordered topographically according to sentence endings. The
377 output vector is composed of five consecutive winners rep-
378 resenting particular words in the sentence. The activations of
379 winning units are slowly decayed in time (decreased by value
380 0.1 at each step) toward the end of a sentence. This function
381 allows to represent the order of winners in the sequence,
382 hence differentiating between similar phonetic features in a
383 sentence (e.g., ‘red ball above red table’ or ‘blue ball above
384 red ball’). The size of RecSOM was set to $N = 20 \times 20$
385 neurons for all models based on results from previous
386 simulations.

Multimodal Layer 387

388 The multimodal layer is the core of the system, since it
389 learns to identify unique categories and represent them.
390 The main task for this layer is to process the output from
391 the unimodal layers and to find and learn the categories by
392 mapping different sources of information (visual and
393 phonological) that refer to the same objects in the external

394 world. Input vectors $\mathbf{m}(t)$ for the multimodal layer
 395 are taken as concatenated unimodal activation vectors
 396 (the ‘where’ and ‘what’ components are not separated in
 397 Model 1) using the above-mentioned k -WTA mechanism,
 398 explained in “[Visual Subsystem](#)”, i.e.,

$$\mathbf{m}(t) = [\mathbf{y}^{\text{where}}(t); \mathbf{y}^{\text{what}}(t); \mathbf{y}^{\text{phono}}(t)].$$

400 The multimodal module receives a 1,300-dimensional
 401 input in Model 1 and a 1925-dimensional input in Model 2
 402 and 3. Unlike sparse localized output codes ($k = 6$) used at
 403 the unimodal layer (to facilitate generalization), the output
 404 representation in the multimodal layer with the WTA
 405 mechanism is chosen to be localist ($k = 1$) for better
 406 interpretation of results and the error calculation.

407 We tested two unsupervised algorithms in the multi-
 408 modal layer, SOM and NG, that differ in the neighborhood
 409 function. The size of the multimodal layer was set to allow
 410 a distinct localist representation of all 840 object combi-
 411 nations in the most complex data set, so we used 841
 412 neurons (arranged in a 29×29 grid in case of SOM).

413 For clarity, we explain the NG algorithm briefly here. NG
 414 shares a number of features with the SOM. In each iteration t ,
 415 an input vector $\mathbf{m}(t)$ is randomly chosen from the training
 416 dataset. Subsequently, we compute $d_i(t) = \|\mathbf{m}(t) - \mathbf{z}_i\|$
 417 for all units, and then, we sort the units according to their
 418 increasing distances d_i , using indices $l = 0, 1, \dots, N - 1$
 419 (where $l(0)$ corresponds to the current winner’s index). We
 420 then update all weight vectors \mathbf{z}_i according to

$$\mathbf{z}_i(t + 1) = \mathbf{z}_i(t) + \eta \exp(-l(i)/\lambda)(\mathbf{m}(t) - \mathbf{z}_i(t))$$

422 with η being the learning rate and λ the so-called neigh-
 423 borhood range. We used $\eta = 0.5$ and $\lambda = n/2$ where n is
 424 the number of neurons. Both parameters are reduced with
 425 increasing t . It is known that after sufficiently many
 426 adaptation steps, the feature vectors cover the data space
 427 with minimum representation error [25]. Mathematically,
 428 the adaptation step of the NG can be interpreted as the
 429 gradient descent on a cost function.

430 Quantification of the Model Accuracy

431 To quantify the model accuracy, we designed the following
 432 procedure for computing the classification error. After the
 433 model has been trained, we again make a single sweep
 434 through the training set, in order to label all neurons,
 435 reflecting their responsiveness to each of the five input fea-
 436 tures (base color, base shape, spatial location, trajectory color,
 437 trajectory and shape). We attach five counter arrays $c_f^{(i)}(j)$
 438 to each neuron, initialized to zeros, each consisting of $n(f_j)$
 439 slots, corresponding to the number of different (possible)
 440 values of the feature f_j (depending on the task complexity),
 441 i.e., $j = 1, 2, \dots, n(f_j)$. For each training input pattern, we
 442 find the winner (as in the SOM algorithm) whose five counter

values are increased by one (i.e., for each current feature 443
 value). After sweeping through the training set, we assign 444
 unique feature labels to all neurons by applying the ‘maxi- 445
 mum response principle,’ according to which each neuron 446
 becomes a representative of only the most frequent value of 447
 the given feature (for which that neuron became the winner 448
 most often), i.e., $f_j^{(i)} = \arg \max_j \{c_f^{(i)}(j)\}$. 449

Then, we can measure the model accuracy, as the per- 450
 centage of correctly classified test inputs. The feature of the 451
 testing pattern is considered to be correctly classified, if it 452
 matches the winner’s representative feature. The calcula- 453
 tion of the classification error rate is first made for each 454
 feature separately and then also for the whole scene-sen- 455
 tence input (overall error) that requires that all features in 456
 the testing input be correctly classified. 457

In the case of the sequential RecSOM, in addition to the 458
 classification error, we also compute the confusion error. It 459
 occurs if the same neuron wins more than once during a 460
 sentence, most typically in case of multiple occurrences of 461
 the same word in a sentence, e.g., in ‘red box above red 462
 ball,’ or ‘red ball below blue ball.’ So whenever the same 463
 winner occurs twice, we increase the error counter by one. 464
 The confusion error stands for the percentage of examples 465
 with the same winner and it helps to detect erroneous cases 466
 not revealed by the classification error. 467

468 Results

We present results corresponding to the three models as 469
 described in “[The Models](#)”, tracking our ‘experimental 470
 trajectory,’ along which we eventually converged to the 471
 architecture with SOM maps in visual subsystems and NG 472
 in the multimodal layer. We trained each model for 100 473
 epochs and tested it with a novel set of inputs. For each run, 474
 the data set was randomly split to the training and testing 475
 subsets using the 70:30 ratio. 476

477 Model 1

In Model 1, the single SOM in the visual system is tested 478
 whether it can learn to represent all visual features simul- 479
 taneously. We observe a high error in this system for the 480
 trajectory features, because trajectory positions overlap in the 481
 specific area. Errors for trajectory color (37 %) and trajectory 482
 shape (65 %) are rather high. Although the spatial location 483
 of the trajectory is fuzzy, the error for this feature in the test 484
 set is the lowest (14 %). Low errors also result for base 485
 color (18 %) and base shape (28 %). We also test whether 486
 the level of fuzziness (shown in Fig. 2b) affects the error in 487
 the SOM map. All features except for the spatial location 488
 are not sensitive to the fuzziness level, as the errors vary 489

490 within a 3 % range. On the other hand, the error for spatial
491 location correlates with the fuzziness starting from 3 % for
492 fixed position of the trajector to 14 % for highly overlap-
493 ping spatial locations.

494 The phonological RecSOM layer performs better compar-
495 ed to the visual layer because the phonetic features,
496 being sequentially fed to the system, are not fuzzy. There
497 are 0 % errors for base color, base shape, trajector color
498 and spatial term. Error for the trajector shape is 1 %. On
499 the other hand, there is a 22 % confusion error, which
500 results in the confusion of the neuron response (see the
501 illustrative Fig. 3 with only 6×6 neurons) and increases
502 the error in the multimodal layer. There are only 4 winning
503 neurons in case of confusion (the same neuron wins twice
504 within one sentence). We observe that this problem can
505 partially be eliminated using the decayed winner activation
506 of winners (as described in “Phonological Subsystem”). It
507 would be possible to solve this problem more reliably by
508 excluding the winner from competition until the rest of the
509 sentence (as, e.g., done in [16]).

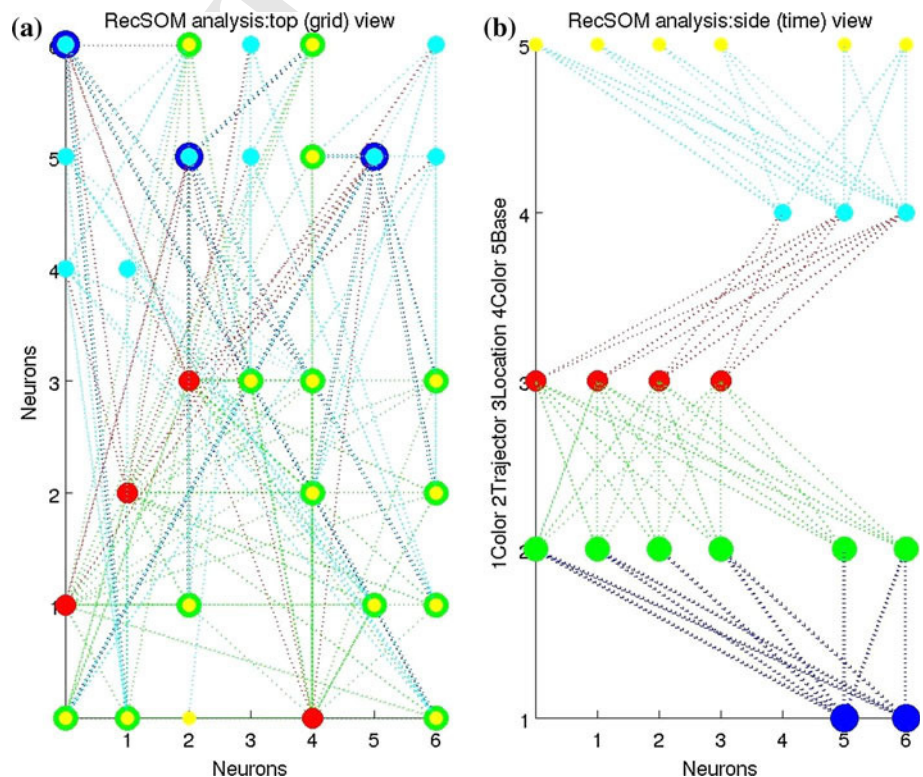
510 The performance of the multimodal layer heavily
511 depends on the effectiveness of unimodal layers. The errors
512 for the representation of trajector color (8 %), base color
513 (1 %) and base shape (2 %) are low. On the other hand,
514 there are high errors for both the trajector shape (46 %) and

spatial term (25 %). This is due to poor performance of the
visual layer. The overall error of the system reaches 68 %.

Model 2

Model 2 processes ‘what’ and ‘where’ information using
separate SOMs, and we identify a difference in accuracy
between the two systems. The ‘what’ system outperforms
the ‘where’ system, as documented by low errors for base
color (1 %), base shape (8 %), trajector color (0 %) and
trajector shape (5 %). We did not test the performance of
the ‘what’ system for the spatial term simply because that
information was not made available to this system. The
errors for the ‘where’ and phonological systems are identi-
cal to Model 1, because these layers receive the same
input as in Model 1. Notably, the additional ‘what’ layer
changed the performance of the multimodal layer. Errors
for base color (2 %) and base shape (4 %) in the multi-
modal layer remain the same as in Model 1, but lower
errors are observed for trajector color (1 %) and trajector
shape (5 %). On the other hand, the system exhibits a much
higher error for the spatial term (71 %) compared to Model
1 (25 %). The multimodal SOM layer is probably not able
to merge the information from three unimodal layers. The
overall error is 75 %, caused by the problem with the

Fig. 3 Unit responses in the phonological layer of Model 1. If the same neuron responds to the same feature (e.g., *shape*) of the trajector and the base (shown by *overlapping dots*), it will increase the error for the whole scene/sentence as well. **a** Visualization of the RecSOM grid (time is represented *bottom-up* by the size of the *dot*; **b** The time course of sentence processing (y-axis) in the *bottom-up* direction



538 representation of the spatial term. A more detailed analysis
539 is explained in the Discussion section.

540 Model 3

541 The simplification of inputs to the ‘where’ system is
542 achieved by using monochromatic bounding boxes instead
543 of object shapes and colors. This expectedly led to a lower
544 error (8.3 % in the most complex and fuzzy scenario)
545 compared to full retinal images (see Fig. 4). We do not
546 compare the results for object features (shape and color),
547 because there is no information about them provided to the
548 ‘where’ system in Model 3. The analysis of the SOM
549 structure revealed a better organization of specific clusters
550 in favor of bounding box inputs for the spatial term rep-
551 resentation. These results lead us to the conclusion that it is
552 possible to simplify the information projected to the
553 ‘where’ system to optimize the speed and effectiveness of
554 the model. However, the simplification of the ‘where’
555 inputs does not affect the performance of the multimodal
556 layer. There are similar results for the object features,
557 spatial term (70 %) and also the overall error (74 %).
558 Therefore, we tested the NG algorithm in the multimodal
559 layer in further simulations trying to improve the accuracy.

560 Comparison of SOM and NG in Multimodal Layer

561 We compare the effectiveness of the SOM and NG algo-
562 rithms in the multimodal layer for all three models. We
563 observe a different type of clustering in the unimodal layers
564 that are transferred to the multimodal layer, where the SOM
565 is not able to adapt to the concatenated outputs from uni-
566 modal layers, apparently due to neighborhood constraints
567 (Model 1SOM and 2SOM). The results of the NG algorithm

(Model 1NG and 2NG) for the same input data confirm this
568 hypothesis. The multimodal layer based on NG is able to
569 correctly map all the object features except spatial term
570 without any problem. There is a 0 % error for both simpli-
571 fied inputs (Model 3NG) and also for full retinal projections
572 to the ‘where’ system (Model 2NG). The errors for the
573 multimodal NG module and the single SOM in the visual
574 layer (Model 1NG) are as follows: 1 % for base color, 2 %
575 for base shape, 6 % for trajector color and 26 % for trajector
576 shape. These results are significantly better than those for
577 the multimodal SOM. Surprisingly, we observe the lowest
578 error for the representation of the spatial term in the multi-
579 modal layer for NG algorithm and a single SOM visual layer
580 (Model 1NG). There is a 12 % error compared to 24 % for
581 Model 2NG (see Fig. 5) and 32 % for Model 3NG (see
582 Table 2). The SOM algorithm leads to higher errors of the
583 spatial term for both models, namely 25 % (Model 1SOM),
584 70 % (Model 2SOM) and 73 % (Model 3SOM). These
585 results are contradictory, because Model 2SOM and 3SOM
586 with separate ‘what’ and ‘where’ systems perform better for
587 all features except the spatial term (see Discussion). Pre-
588 liminary results of this comparison were also presented in
589 Vavrečka, Farkaš and Lhotská [49].

591 The comparison of the overall accuracy (overall error) is
592 shown in Fig. 6 and Table 2. The best results are obtained
593 for ‘what’ and ‘where’ subsystems and the NG algorithm in
594 the multimodal layer (Model 2NG). There is a 25 % error
595 compared to 70 % overall error for the multimodal SOM in
596 the most complex scenario. Hence, the better, albeit not
597 perfect, results are achieved with NG by sacrificing the
598 topographic organization of responses in the multimodal
599 layer.

600 The last analysis is dedicated to the comparison of SOM
601 (Model 3SOM) and NG (Model 3NG) algorithms in the

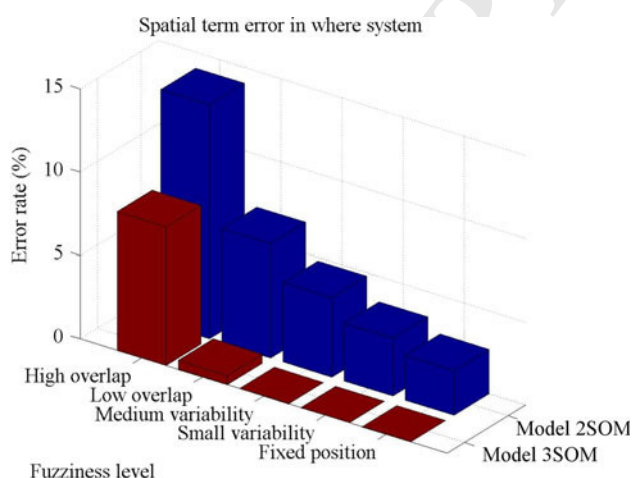


Fig. 4 Visualization of spatial term errors in the ‘where’ layer for full retinal inputs (blue) and for bounding box inputs (red) as a function of the fuzziness level of trajector spatial location

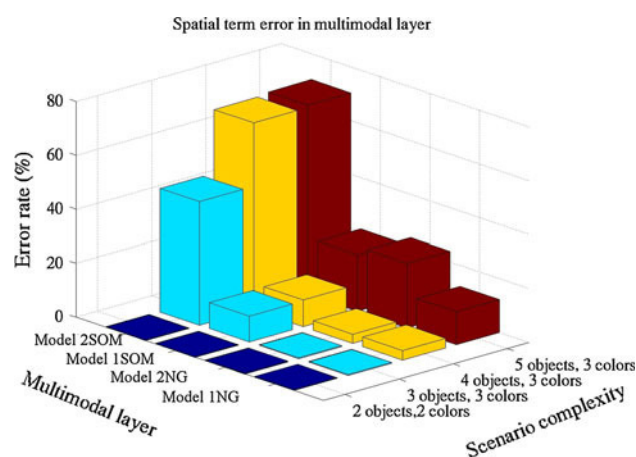
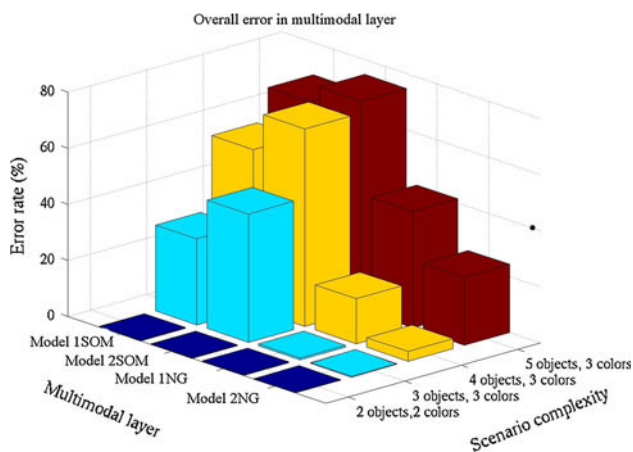


Fig. 5 Comparison of the errors in the multimodal layer for the representation of the spatial term. Model 1NG (NG in the multimodal layer and a single SOM in the visual system) performs best

Table 2 Summary of error rates for specific layers and models

Subsystem	Model	TrajCol	TrajShape	SpatTerm	BaseCol	BaseShape	Overall
Where	1,2SOM; 1,2NG	39.3	68.2	14.2	19.2	30.5	91.7
	3SOM; 3NG	-	-	8.3	-	-	-
What	2,3SOM; 2,3NG	0.4	5.3	-	0.9	0.8	-
Phono	1,2,3SOM; 1,2,3NG	0.0	1.2	0.2	0.0	0.0	12.3
Multimodal	1SOM	8.3	46.0	24.6	0.5	2.0	68.3
	2SOM	0.9	5.4	70.3	1.9	3.8	74.7
	3SOM	0.9	4.1	72.7	1.2	1.6	75.3
	1NG	5.6	26.4	12.3	0.6	1.7	41.5
	2NG	0.0	0.3	24.0	0.0	0.0	24.3
	3NG	0.0	0.0	31.5	0.0	0.0	31.5

**Fig. 6** Errors in the multimodal layer for whole scene (overall) representation. Model 2NG based on ‘what’ and ‘where’ visual system and NG in multimodal layer performs best

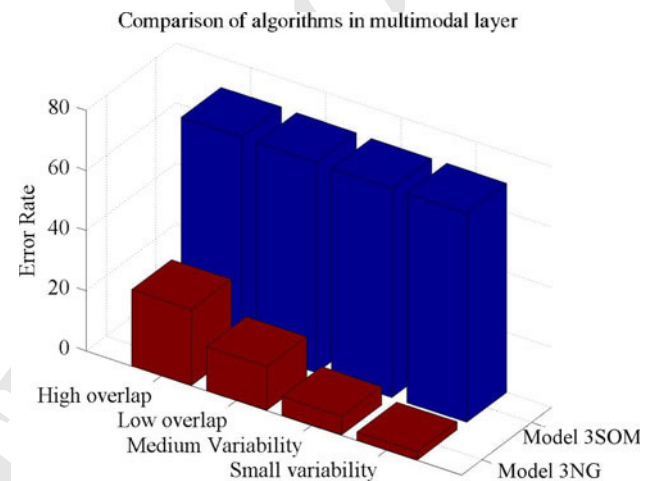
602 multimodal layer that have to process different levels of
 603 spatial fuzziness. Fig. 7 reveals a lower error for NG at all
 604 levels of fuzziness and the high errors for SOM regardless
 605 of the fuzziness level (70 %). Hence, the multimodal SOM
 606 is unable to unambiguously represent neither fuzzy nor
 607 distinct inputs.

608 Discussion

609 We analyze the presented models in the context of theo-
 610 retic assumptions, especially the perceptual theory of
 611 cognition and conceptual approaches to knowledge repre-
 612 sentation. We also discuss various aspects of our model, its
 613 relation to other models and the features of Visual Feature-
 614 Binding and temporal synchrony [6].

615 Architecture

616 We should also compare our architecture with the system
 617 for the representation of spatial relations developed by

**Fig. 7** Errors in the multimodal layer for SOM (Model 3SOM) and NG (Model 3NG) algorithms as a function of the fuzziness level of the trajectors’ spatial location (see Fig. 2b)

Regier [34]. The main difference lies in the unsupervised 618
 manner of our architecture compared to the supervised 619
 approach adopted by Regier. His system is composed of 620
 specific modules for the calculation of angle between tra- 621
 jector and base, an object’s intersection and dynamic 622
 properties in motion inputs. It resembles the designer’s 623
 approach described in Ziemke [53] as there is modular 624
 architecture engineered for the specific task. Our system 625
 is more generic and biologically inspired as the subsystems 626
 copy the information processing in human brain (unsu- 627
 pervised learning, ‘what’ and ‘where’ pathways, multi- 628
 modal integration). The advantage of Regier’s system is the 629
 ability to represent dynamic spatial relations (around, 630
 through, etc.). On the other hand, our unsupervised archi- 631
 tecture based on RecSOM [51] in a visual subsystem and 632
 the growing-when-required networks [24] in the phono- 633
 logical and multimodal layer was able to process visual 634
 sequences (around, through, outside, over and under) and it 635
 reached 88 % accuracy [48]. 636

637 In our model, the representations take advantage of the
 638 two or three unimodal layers of units. The phonological layer
 639 represents unique labels (linguistic terms), whereas the
 640 visual ‘where’ subsystem represents fuzzy information
 641 about the spatial locations of objects in the external world
 642 and the ‘what’ subsystem captures shapes and colors of
 643 objects in a fixed foveal position. The multimodal level
 644 integrates the outputs of these unimodal layers. The ground-
 645 ed meaning is simultaneously represented by all layers
 646 (phonological, visual and multimodal), making this
 647 approach resemble the theory of Peirce [30] who defined
 648 basic components of a sign—representamen and interpretant.
 649 Our model represents the sign hierarchically guaranteeing
 650 better processing and storing of representations, because the
 651 sign (the multimodal level) is modifiable from both modalities
 652 (the sequential ‘representamen’ via the phonological level and
 653 the parallel ‘interpretant’ via the visual level). This feature
 654 makes the units in the higher layer bimodal (i.e., they can be
 655 stimulated by any of the primary layers) and their activation
 656 can be forwarded for further processing. Bimodal (and multimodal)
 657 neurons are known to be ubiquitous in the association areas of
 658 the brain [39]. The multimodal layer is formed by exploiting
 659 the concept of self-organized conjunctive representations that
 660 have been hypothesized to exist in the brain with the purpose
 661 of binding the features such as various perceptual properties of
 662 objects [26]. We adhere to the view that conjunctive neurons,
 663 as an alternative to mechanisms of temporal synchrony, are the
 664 plausible connectionist approach for addressing the binding
 665 problem [29]. Here, we extend the concept of binding by linking
 666 the subsymbolic and symbolic information. Hence, each output
 667 unit learns to represent a unique combination of perceptual and
 668 symbolic information.

670 Visual Feature-Binding

671 Our Models 2 and 3 propose the unsupervised solution to the
 672 Visual Feature-Binding, based on the integration of the ‘what’
 673 and ‘where’ pathways. With respect to the Visual Feature-
 674 Binding [6], the model is based on convergent hierarchical
 675 coding, also called combination coding [35]. The neurons react
 676 only to combinations of features, that is, to an object of a
 677 particular shape and color at a particular retinal position
 678 (localist representation). Hierarchical processing implies that
 679 increasingly complex features are represented by higher levels
 680 in the hierarchy. Complex objects and situations are constructed
 681 by combining simpler elements. On the other hand, the
 682 convergent hierarchical coding requires as many binding units
 683 as there are distinguishable objects. It should result in a
 684 combinatorial explosion for large-scale simulations. Our model
 685 is able to represent 840 combinations, but it can also suffer
 686 from the combinatorial explosion because we represent pairs of

688 objects instead of separate entities in the primary visual
 689 layers. In case of 10 objects, 5 colors in 4 spatial locations,
 690 we would need to represent 2450 object pairs in a primary
 691 ‘what’ system, instead of 50 separate objects. It is also
 692 possible to add a separate layer for the color processing, in
 693 which case there will only be 10 objects presented in the
 694 ‘what’ system. Alternatively, we could represent the features
 695 in the activity of a population of neurons distributed within
 696 and across levels of the cortical hierarchy as the distributed
 697 representation [8], although some authors have raised the
 698 question whether the combinatorial explosion is really a
 699 problem [10]. It is estimated that the number of objects,
 700 scenarios, colors and other features in the brain is approximately
 701 10 million items. It is obviously beyond the limits of recent
 702 cognitive systems, but it is below the number of neurons in
 703 the mammalian visual cortex, so the combination coding could
 704 be a sufficient method. It could also be possible to adopt
 705 Neural Modeling Fields [31], the unsupervised learning
 706 method based on Gaussian mixture models that arguably does
 707 not suffer from combinatorial complexity. The application of
 708 this theory to the area of symbol grounding resulted in 95 %
 709 accuracy of the system that learned the repertoire of 112
 710 actions [5].

Temporal Synchrony

711 Our model is able to map the words in the sentence with the
 712 fixed grammar to the objects in the environment without any
 713 prior knowledge (temporal synchrony). Previous models of
 714 symbol grounding [2–5] deal with the lexical level, but our
 715 model goes beyond words because it can represent sentences
 716 in RecSOM. The ability of temporal synchrony can be
 717 considered as an extension of the symbol grounding. Cangelosi
 718 et al. [2] recommend to ground-specific words at the first
 719 stage (sensorimotor toil) and then compositionally chain them
 720 at the grounded language level (symbolic theft). There are
 721 separate objects presented to their system within a training
 722 phase, grounding basic object features. Our approach can be
 723 considered an alternative to this theory. We also ground words
 724 in the first stage, but unlike the mentioned approach, we
 725 present sentences as linguistic inputs to be bound with proper
 726 features from the visual subsystem (shape, color and location).
 727 Compared to the classic sensorimotor toil experiments based
 728 on the grounding of two features, our system is able to ground
 729 5 features simultaneously, which speeds up the process of
 730 symbol grounding (faster acquisition of the grounded lexicon).
 731 Tikhanoff [43] proposed an architecture (and implemented it
 732 in iCub robot) that was able to understand basic sentences,
 733 but it was based on supervised learning. Our model is a
 734 proof of concept that unsupervised architectures can also
 735 find proper mapping between multiple visual and lexical
 736 features. We are able to build representations solely from
 737 sensory inputs, arguing

739 that the co-occurrence of inputs from the environment is a
740 sufficient source of information to create an intrinsic rep-
741 resentational system.

742 Performance

743 The analysis of the model behavior revealed that the tra-
744 jector shape and the spatial term representations are the most
745 difficult subtasks for visual unimodal systems. The difficulty
746 is caused by the variability and fuzziness of these inputs.
747 The correct representation of the trajectory shape requires a
748 separate unimodal ‘what’ system. The errors for (both SOM
749 and NG) Model 1, 2 and 3 confirm the necessity of the
750 ‘what’ system in the complex environment because we
751 observe a 60 % increase of errors in the model without a
752 separate ‘what’ system. On the other hand, the error for the
753 spatial term in Model 2 and 3 reflects some problems with an
754 increasing number of inputs from different subsystems to
755 the multimodal layer, because there is a lower error for
756 Model 1 compared to Model 2 and 3 (both SOM and NG).
757 The problem could reside in the number of dimensions. The
758 multimodal module receives a 1300-dimensional input in
759 Model 1 and a 1925-dimensional input in Model 2 and 3.
760 The increase of dimensionality together with a localist
761 unimodal output function may decrease the effectiveness for
762 the spatial term representation, although other features are
763 represented better in a high-dimensional space. This con-
764 tradiction has to be investigated in greater detail.

765 The results for specific algorithms in the multimodal
766 layer confirm our hypothesis that the SOM algorithm,
767 based on the fixed neighborhood function, is not able to
768 adapt to the joint distribution of the outputs from unimodal
769 layers. The SOM-based models aim at the topology-pre-
770 serving property for the input data, but they are weak with
771 regard to properly representing clusters with different non-
772 uniform data distributions [18]. On the other hand, the NG
773 algorithm is not subject to topographic constraints and,
774 thus, leads to better clusters. Our results are also in line
775 with Pezzulo and Calvi [32], who conclude that perceptual
776 symbols may not be topographically organized, although
777 some parts of the perceptual and motor areas show topo-
778 graphic hierarchical organization. Grounding models based
779 on topographically organized connectionist networks (e.g.,
780 [17]) to simulate the perceptual symbol system also exist,
781 but our results do not confirm this assumption for more
782 complex inputs.

783 The mapping in our models is actually a clustering pro-
784 cess that makes the system also vulnerable to errors in the
785 input space. Successful clustering presumes that at least one
786 modality provides distinct activation vectors for different
787 classes to drive the clustering process (i.e., the classes are
788 well separable in the corresponding input subspace). On the
789 other hand, the occurrence of both phonological and visual

fuzzy inputs is rare in the real world, so our system could be
considered a step toward solving the symbol grounding
problem (at least at this small scale).

Conclusion

We have created an unsupervised connectionist system that
is able to extract constant attributes and regularities from
the environment and link them with abstract symbols. The
meaning is non-arbitrarily represented at the conceptual
level that guarantees the correspondence of the internal
representational system with the external environment. We
can also conclude that it is advantageous to follow the
biologically inspired hypothesis about the processing of
visual information in separate subsystems. The question for
future research is to find a proper way of output coding
from the unimodal layers to increase system accuracy and
to scale up the model. The main advantage of our model is
the hierarchical representation of the sign components.

Acknowledgments This work has been supported by the research
program MSM 6840770012 of the CTU in Prague, SAIA scholarship
and GAČR Grant P407/11/P696 (M.V.) and by VEGA Grant 1/0439/
11 (I.F.).

References

1. Dorffner G, Hentze M., Thurner G. A connectionist model of categorization and grounded word learning. In: Koster C, Wijnen F, editors. Proceedings of the groningen assembly on language acquisition (GALA'95), 1996.
2. Cangelosi A, Greco A, Harnad S. From robotic toil to symbolic theft: grounding transfer from entry-level to higher-level categories. *Conn Sci.* 2000;12(2):143–62.
3. Cangelosi A, Parisi D. The processing of verbs and nouns in neural networks: insights from synthetic brain imaging. *Brain Lang.* 2004;89(2):401–08.
4. Cangelosi A, Riga T. An embodied model for sensorimotor grounding and grounding transfer: experiments with epigenetic robots. *Cogn Sci.* 2006;30(4):673–89.
5. Cangelosi A, Tikhanoff V, Fontanari JF, Hourdakis E. Integrating language and cognition: a cognitive robotics approach. *IEEE Comput Intell Mag.* 2007;2(3):65–70.
6. Feldman J. The neural binding problem(s). *Cogn Neurodyn.* 2012; doi:10.1007/s11571-012-9219-8.
7. Fontanari JF, Tikhanoff V, Cangelosi A, Ilin R, Perlovsky LI. Cross-situational learning of object-word mapping using neural modeling fields. *Neural Netw.* 2009;22:579–85.
8. Goldstein EB. *Wahrnehmungspsychologie.* Heidelberg: Spektrum Akademischer Verlag, 2002.
9. Gliozzi V, Mayor J, Hu J-F, Plunkett K. Labels as features (not names) for infant categorization: a neurocomputational approach. *Cogn Sci.* 2009;33(4):709–38.
10. Ghose GM, Maunsell J. Specialized representations in visual cortex: a role for binding? *Neuron* 1999;24:79–85.
11. Greco A., Caneva C. Compositional symbol grounding for motor patterns. *Front Neurobot.* 2010;4(111), doi:10.3389/fnbot.2010.00111.

- 843 12. Grossberg S. Competitive learning: from interactive activation to
844 adaptive resonance. *Cogn Sci* 1987;11(1):23–63.
- 845 13. Hammer B, Micheli A, Sperduti A, Strickert M. Recursive self-
846 organizing network models. *Neural Netw.* 2004;17(8–9):
847 1061–85.
- 848 14. Hinton G, Osindero S, Teh Y. A fast learning algorithm for deep
849 belief nets. *Neural Comput.* 2006;18:1527–54.
- 850 15. Jacobs RA, Jordan MI, Barto AG. Task decomposition through
851 competition in a modular connectionist architecture: the what and
852 vision tasks. *Cogn Sci.* 1991;15(2):219–50.
- 853 16. James DJ, Miikkulainen R. SardNet: A self-organizing feature
854 map for sequences. *Adv Neural Inf Process Syst* 1995;7:577–84.
- 855 17. Joyce D, Richards L, Cangelosi A, Coventry KR (2003) On the
856 foundations of perceptual symbol systems: specifying embodied
857 representations via connectionism. In: Detje F, Drner D, Schaub
858 H, editors. *The logic of cognitive systems. Proceedings of the*
859 *fifth international conference on cognitive modeling, Universi-*
860 *taetsverlag Bamberg*, pp. 147–52.
- 861 18. Kim B, Sang-Woo B, Minho L. Growing fuzzy topology adaptive
862 resonance theory models with a pushpull learning algorithm.
863 *Neurocomputing.* 2011;74(4):646–55.
- 864 19. Kohonen T. *Self-Organizing Maps*, 3rd edn. Berlin: Springer;
865 2001.
- 866 20. Li P, Farkaš I, MacWhinney B. Early lexical development in a
867 self-organizing neural network. *Neural Netw.* 2004;17(8–9):
868 1345–62.
- 869 21. Li P, MacWhinney B. PatPho: a phonological pattern generator
870 for neural networks. *Behav Res Methods Instrum Comput.* 2002;
871 34:408–15.
- 872 22. Malach R, Levy I, Hasson U. The topography of high-order
873 human object areas. *Trends Cogn Sci.* 2002;6(4):176–84.
- 874 23. Marocco D, Cangelosi A, Fischer K, Belpaeme T. Grounding
875 action words in the sensorimotor interaction with the world:
876 experiments with a simulated iCub humanoid robot. *Front Neuro-*
877 *robot.* 2010;4(7), doi:[10.3389/fnbot.2010.00007](https://doi.org/10.3389/fnbot.2010.00007).
- 878 24. Marsland S, Shapiro J, Nehmzow U. A self-organising network
879 that grows when required. *Neural Netw.* 2002;15(8–9):1041–58.
- 880 25. Martinetz T, Berkovich S, Schulten K. “Neural-gas” network for
881 vector quantization and its application to time-series prediction.
882 *IEEE Trans Neural Netw.* 1993;4(4):558–69.
- 883 26. Mel B, Fiser J. Minimizing binding errors using learned con-
884 junctive features. *Neural Comput.* 2000;12:247–78.
- 885 27. Miikkulainen R. Dyslexic and category-specific aphasic impair-
886 ments in a self-organizing feature map model of the lexicon. *Brain*
887 *Lang.* 1997;59:334–66.
- 888 28. Newell A, Simon HA. *Human problem solving*. Englewood
889 *Cliffs: Prentice-Hall*; 1972.
- 890 29. O’Reilly RC, Busby RS, Soto R. Three forms of binding and their
891 neural substrates: alternatives to temporal synchrony. In: Cleere-
892 mans A, editor. *The unity of consciousness: binding, integration,*
893 *and dissociation*. Oxford: Oxford University Press, 2003; 168–92.
- 894 30. Peirce, C.S. *Collected papers of Charles Sanders Peirce*. In
895 *Hartshorne C, editor. Harvard University Press*, 1931.
- 896 31. Perlovsky LI. *Neural networks and intellect: using model-based*
897 *concepts*. Oxford University Press, New York, 2001.
- 898 32. Pezzulo G, Calvi G. Computational explorations of perceptual
899 symbol systems theory. *New Ideas Psychol.* 2011;29:275–297.
- 900 33. Pylyshyn Z. *Computation and cognition: towards a foundation for*
901 *cognitive science*. Cambridge: MIT Press; 1984.
- 902 34. Regier T. *The human semantic potential: spatial language and*
903 *constrained connectionism*. Cambridge: MIT Press; 1996.
- 904 35. Riesenhuber M, Poggio T. Neural mechanisms of object recog-
905 nition. *Curr Opin Neurobiol.* 2002;12:162–168.
- 906 36. Roy D. Grounding words in perception and action: computational
907 insights. *Trends Cogn Sci.* 2005;9:389–396.
- 908 37. Roy D, Pentland A. Learning words from sights and sounds: a
909 computational model. *Cogn Sci* 2002; 26:113–146.
- 910 38. Steels L, Kaplan F. Situated grounded word semantics. In: *Pro-*
911 *ceedings of the 16th international joint conference on artificial*
912 *intelligence, vol 2. 1999; p. 862–67.*
- 913 39. Stein B, Meredith M. *Merging of the senses*. Cambridge: MIT
914 *Press*; 1993.
- 915 40. Sugita Y, Tani J. Learning semantic combinatoriality from the
916 interaction between linguistic and behavioral processes. *Adapt*
917 *Behav* 2005; 13(1):33–52.
- 918 41. Taddeo M, Floridi L. The symbol grounding problem: a critical
919 review of fifteen years of research. *J Exp Theor Artif Intell.*
920 2005;17(4):419–45.
- 921 42. Tikhonoff V, Cangelosi A, Fitzpatrick P, Metta G, Natale L, Nori F.
922 An open-source simulator for cognitive robotics research: the
923 prototype of the iCub humanoid robot simulator. In: *Performance*
924 *metrics for intelligent systems (PerMIS) workshop, 2008; p. 57–61.*
- 925 43. Tikhonoff, V. *Development of cognitive capabilities in humanoid*
926 *robots*. PhD thesis. School of Computing, Communications &
927 *Electronics, University of Plymouth, UK, 2009.*
- 928 44. Tiño P, Farkaš I, van Mourik J. Dynamics and topographic
929 organization in recursive self-organizing map. *Neural Comput.*
930 2006;18:2529–67.
- 931 45. Ungerleider LG, Mishkin M. Two cortical visual systems. In:
932 *Ingle DJ et al. editors. Analysis of visual behavior*. MIT Press,
933 *Cambridge*; 1982.
- 934 46. Vavrečka M. Symbol grounding in context of zero semantic
935 commitment (in Czech). In: Kelemen J, Kvasnička V, editors.
936 *Kognice a umělý život VII. (1st ed.) Opava : Slezská univerzita*
937 *2006; 365–377.*
- 938 47. Vavrečka, M. Grounding of spatial terms (in Czech). In: J. Kel-
939 emen J, Kvasnička V, editors. *Kognice a umelý život VII, Opava:*
940 *Slezsk univerzita, 2007; p. 365–77.*
- 941 48. Vavrečka M. Application of cognitive semantics in the model of
942 the spatial terms representation (in Czech). PhD thesis, Masaryk
943 *University in Brno, Czech Republic (2008).*
- 944 49. Vavrečka M, Farkaš I, Lhotská L. Bio-inspired model of spatial
945 cognition. In *Lecture notes in computer science 7062 LNCS (Part 1),*
946 *2011;443–450.*
- 947 50. Vesanto J, Himberg J, Alhoniemi E, Parhankangas J. Self-
948 *Organizing Map in Matlab: the SOM Toolbox*. In: *Proceedings of the*
949 *matlab DSP conference, 2000; p. 35–40.*
- 950 51. Voegtlin T. Recursive self-organizing maps. *Neural Netw* 2002;
951 15(8–9):979–91.
- 952 52. Vogt P, Divina F. Social symbol grounding and language evo-
953 lution. *Interact Stud.* 2007;8:31–52.
- 954 53. Ziemke T. Rethinking grounding. In: Riegler A, Peschl M, von
955 Stein A, editors. *Understanding representation in the cognitive*
956 *sciences*. New York: Plenum Press; 1999. p. 177–90.
- 957