# Causal influence of linguistic learning on perceptual and conceptual processing: A brain-constrained deep neural network study of proper names and category terms

1  Causal influence of linguistic learning on perceptual and conceptual processing: A brain-constrained
2  deep neural network study of proper names and category terms

3

4  Phuc T. U. Nguyen[1]*, Malte R. Henningsen-Schomers[1,4], and Friedemann Pulvermüller[1,2,3,4]

5

6  [1]Dept. Of Philosophy and Humanities, Brain Language Laboratory, Freie Universität Berlin,
7  Habelschwerdter Allee 45, 14195 Berlin, Germany

8  [2]Berlin School of Mind and Brain, Berlin, Germany

9  [3]Einstein Center for Neurosciences, Berlin, Germany

10  [4]Cluster of Excellence "Matters of Activity. Image Space Material", Humboldt-Universität zu Berlin,
11  Berlin, Germany

12

13  * Corresponding authors: Phuc T. U. Nguyen (phuc.thu.uyen.nguyen@gmail.com) and Friedemann
14  Pulvermüller (friedemann.pulvermuller@fu-berlin.de )

15  Number of pages: 51

16  Number of Figures: 6

17  Number of tables: 3

18  Number of words

19  Abstract: 250/ 250

20  Significance statement: 120/120

21  Introduction: 646/650

22  Discussion: 1521/1500

36  Author contributions. M.H.-S. and F.P. conceived the study. M.H.-S., P.T.U.N and F. P. designed the
37  study, M.H.-S., and P.T.U.N implemented research. P.T.U.N analyzed data. P.T.U.N, M.H.-S., and
38  F.P. wrote the manuscript.

**Abstract**

39  
40       Language influences cognitive and conceptual processing, but the mechanisms through which
41  such causal effects are realized in the human brain remain unknown. Here, we use a brain-constrained
42  deep neural network model of category formation and symbol learning and analyze the emergent
43  model-internal mechanisms at the neural circuit level. In one set of simulations, the network was
44  presented with similar patterns of neural activity indexing instances of objects and actions belonging
45  to the same categories. Biologically realistic Hebbian learning led to the formation of instance-specific
46  neurons distributed across multiple areas of the network, and, in addition, to cell assembly circuits of
47  'shared' neurons responding to all category instances – the network correlates of conceptual
48  categories. In two separate sets of simulations, the network learned the same patterns together with
49  symbols for individual instances ('*proper names*') or symbols related to classes of instances sharing
50  common features ('*category terms*'). Learning category terms remarkably increased the number of
51  shared neurons in the network, thereby making category representations more robust while reducing
52  the number of neurons of instance-specific ones. In contrast, proper-name learning prevented
53  substantial reduction of instance-specific neurons and blocked the overgrowth of category-general
54  cells. Representational Similarity Analysis further confirmed that the neural activity patterns of
55  category instances became more similar to each other after category-term learning, relative to both
56  learning with proper names and without any symbols. These network-based mechanisms for concepts,
57  proper names and category terms explain why and how symbol learning changes object perception and
58  memory, as revealed by experimental studies.

### Significance Statement

59  
60       How do verbal symbols for specific individuals (*Micky Mouse*) and object categories (*house
61  mouse*) causally influence conceptual representation and processing? Category terms and proper
62  names have been shown to respectively promote category formation and instance learning, potentially
63  by respectively directing attention to category-critical and object-specific features. Yet the
64  mechanisms underlying these observations at the neural circuit level remained unknown. Using a
65  mathematically precise deep neural network model constrained by properties of the human brain, we
66  show category-term learning strengthens and solidifies conceptual representations, whereas proper
67  names support object-specific mechanisms. Based on network-internal mechanisms and unsupervised
68  correlation-based learning, this work offers neurobiological explanations for causal effects of symbol
69  learning on concept formation, category building and instance representation in the human brain.

70                                                           **Introduction**
71              Most signs and symbols are used to speak about objects and actions. This led philosophers and
72    logicians to propose that the referential link between symbol and world is most essential for meaning
73    and semantics (Wittgenstein, 1922; Frege, 1948). Yet there are quite different relationships between
74    symbols and their related real-world entities. One most essential difference exists between 'proper
75    names' used to speak about a single object or individual (e.g., "Mickey Mouse") and 'category terms,'
76    which can refer to members of an entire class or conceptual category (e.g., "house mouse"). Such
77    differences between referential symbols are well-described at the semantic level, but not understood in
78    terms of their underlying mechanisms in mind and brain.

79              The need for mechanistic neurobiological models of symbols and their meaning comes from
80    reports about causal influences of language on perception, attention, and memory. It had long been
81    speculated and recently been confirmed that, when human subjects learn words for objects, language
82    may help humans to attend to and distinguish between them (Majid et al., 2004; Whorf and Carroll,
83    2007; Miller et al., 2018; Vanek et al., 2021). Experimental research in infants showed that learning
84    'labels' for objects increases their attention to these objects (Baldwin and Markman, 1989), which
85    further establishes an attention-catching function of language. However, this general insight requires
86    further specification to capture the different effects of category terms and proper names. In particular,
87    learning a new symbol for a category of objects makes infants attend to the shared features of these
88    objects and facilitates their learning of the conceptual category (Gelman and Markman, 1986, 1987;
89    Plunkett et al., 2008); the latter even holds if the objects show little perceptual similarity (Graham et
90    al., 2013). On the other hand, the category building function of language is absent when object-
91    specific proper names are learned. In this case, the infant's attention is directed not towards the
92    common category features of objects, but to idiosyncratic and object-specific features instead (Scott
93    and Monesson, 2009; LaTourrette and Waxman, 2020). In summary, category-term learning directs
94    attention to shared features of objects (Waxman and Booth, 2001; Dewar and Xu, 2007; Althaus and
95    Mareschal, 2014; Althaus and Plunkett, 2016), whereas unique proper-name learning highlight
96    idiosyncratic and object-specific features (Best et al., 2010; Barnhart et al., 2018; Pickron et al., 2018;
97    LaTourrette and Waxman, 2020). These specific and replicable effects of proper names and category
98    terms on perception and attention have been explained in terms of different 'strategies' applied by the
99    learner. A neurobiological explanation of why these specific effects occur is still missing.

100             Why and how can proper names and category terms direct attention to specific versus shared
101   features of category members? To develop a mechanistic explanation, we used a brain-constrained
102   deep neural network designed according to the area structure and connectivity of major areas relevant
103   for language and conceptual processing (Garagnani et al., 2007; Tomasello et al., 2018; Pulvermüller
104   et al., 2021). Six "areas" of the model simulated processes in superior temporal and inferior frontal
105   perisylvian language areas and six extrasylvian model areas simulated inferior temporo-occipital
106   visual 'where' processing stream and dorsolateral prefrontal and motor cortices (Figure 1A). In the No
107   symbol (NoS) condition, the model learned activity patterns each representing 1 of 60 instances of
108   objects or actions belonging to 10 different categories. In learning-with-symbols conditions, the model
109   learned additional activity patterns representing word forms of proper names (PN) or category terms
110   (CT) (Figure 1B-C, 2A). After learning, the model was tested by activating previously trained instance
111   patterns of each category and, in addition, new patterns for novel instances belonging to the same
112   categories (Figure 2B). We documented the neural and cognitive effects of proper names and category
113   terms on instance and category learning in the model. In-depth analyses of the emerging activation
114   patterns and representations were provided by using Representational Similarity Analysis (RSA)
115   (Kriegeskorte et al., 2008) and by classifying neurons into instance-specific and category-general ones.

<center>**Materials and Methods**</center>

**Participants**

The current work does not contain experiments with human participants or animal subjects.

*Neurobiological constraints*

In contrast to many neural network models, the brain-constrained model applied aimed at biological plausibility by applying a range of structural and functional constraints (Tomasello et al., 2018; Henningsen-Schomers and Pulvermüller, 2022; for review, see Pulvermüller et al., 2021) realizing:

(1) neurophysiological dynamics of spiking pyramidal cells (Connors et al., 1982; Matthews, 2001),

(2) synaptic weights under the modification of unsupervised Hebbian-type learning (i.e., synaptic plasticity and learning were modified according to the biologically plausible unsupervised Hebbian principles that incorporated both long-term potentiation and long-term depression) (Artola and Singer, 1993),

(3) local and global activity regulation (Braitenberg, 1978; Yuille and Geiger, 1995) based on local and area-specific inhibition mechanisms (Knoblauch and Palm, 2002),

(4) excitatory and inhibitory within-area local connectivity (including sparse, random, and initially weak excitatory links whose probability falls off with distance) (Kaas, 1997; Braitenberg and Schüz, 1998),

(5) between-area global connectivity built on neuroanatomical evidence, and

(6) built-in uncorrelated white noise in neurons of (a) all areas during training and testing mimicked spontaneous baseline neuronal firing and (b) additional noise in neurons of areas not stimulated by patterns during training, which simulated uncorrelated sensory or motor activity unrelated to instances or symbols (Rolls and Deco, 2010).

Table 2 supplies the model specifications and parameters chosen in this current work.

**Model description**

We applied a brain-constrained deep neural network model including spiking model neurons and twelve model areas to model sensorimotor, conceptual and linguistic mechanisms in the left-hemispheric language-dominant fronto-temporo-occipital regions of the human brain, as described in previous studies by Tomasello et al., 2018; Henningsen-Schomers and Pulvermüller, 2022.

*Anatomical architecture and connectivity*

To distinguish between sub-parts of neural networks from their target cortical structures of the real human brain, all model areas are marked by an asterisk before (e.g., *A1, *V1). The architecture modelled three areas representing the ventral visual system (i.e., primary visual cortex (*V1), temporo-occipital area (*TO), anterior-temporal area (*AT)) and three areas representing the dorsolateral action system (i.e., dorsolateral fronto-central motor (*M1$_L$), premotor cortex (*PM$_L$), prefrontal cortex (*PF$_L$)). These formed the extrasylvian region for sensorimotor processing where semantic information was stored. Another 6 areas of the perisylvian region for word-form processing housed articulatory-phonological and acoustic-phonological information. These areas involved the three areas of the auditory system (i.e., primary auditory cortex (*A1), auditory belt (*AB), parabelt areas (*PB)) and three inferior frontal articulatory and prefrontal areas (i.e., inferior primary motor cortex (*M1$_i$), premotor cortex (*PM$_i$), prefrontal cortex (*PF$_i$)), respectively. Between-area connections were reciprocal and connected next neighbor areas, second-next neighbors (see Schomers, 2017) and long-distance corticocortical links supported by neuroanatomical evidence in the literature (Table 1).

In the current neural network model, the fundamental information processing units are artificial neuron-like elements, or cells. Each modelled area comprised two layers of 625 e-cells and 625 i-cells that mimicked an (excitatory) pyramidal spiking neuron and a cluster of (inhibitory)

164  interneurons hosted within the same cortical column in the cortical area. A more elaborate description
165  of the firing behavior of such neurons could be found in Garagnani et al. (2017), Tomasello et al.
166  (2018), Henningsen-Schomers and Pulvermüller (2022).

167  **Activity patterns applied to the networks**
168         60 'grounding patterns' were defined as sensorimotor activation patterns thought to represent
169  specific sensory-motor experiences of 60 different objects or "instances". Groups of 6 instances
170  overlapped in their neuronal grounding patterns and were taken as representations of different
171  instances of the same concept (e.g., different robots). Note that the images of robots and cat faces for
172  category members are to be taken purely for illustrative purposes here – the actual training patterns of
173  the models consisted of sets of activated neurons with no systematic relationship to images of robots
174  or cat faces. A category comprised three trained instances and three novel instances not presented
175  during training; all six instance patterns were used for network testing (Figure 2A-B). Each category
176  instance was neuronally coded as a set of perceptual and motor neuron activations in the primary
177  visual and hand-motor areas of the brain-constrained network. These instance-related grounding
178  patterns were activated either on their own or together with additional patterns of neuronal activation
179  in the network's articulatory and auditory cortices, which were thought to implement symbol forms,
180  that is, verbal labels or spoken word forms. These "word form patterns" were used either as proper
181  names and therefore specifically with only one grounding pattern, or as category terms and therefore
182  the same word form pattern co-occurred with all 3 trained grounding patterns of one category. To
183  control the effect of non-linguistic factors, a third class of trained grounding patterns was learnt
184  without concordant auditory-articulatory activation. Thus, we generated three classes of simulated
185  stimulation patterns: (i) instance-related grounding patterns applied to $*V1/*M1_L$ (Figure 1B-left), (ii)
186  category term patterns to $*A1/*M1_i$, (Figure 1B-middle) and (iii) proper name patterns to $*A1/*M1_i$
187  (Figure 1B-right). Sensorimotor experiences of instances were simulated with conceptual grounding
188  patterns, (i), and symbol-related auditory-articulatory activity were simulated using word form
189  patterns, (ii) and (iii).

190         For visualization and better conceptual understanding of the use of activity patterns, see
191  Figure 1B-C. Instances belonging to the same category were simulated by similar grounding patterns,
192  following Henningsen-Schomers and Pulvermüller (2022).: within-category instances had grounding
193  patterns that shared 50% of their feature neurons and differed from each other in the other half;
194  grounding patterns simulating instances from different categories had no neuronal overlap. For each
195  grounding pattern (i), a subset of twelve out of 625 potential cells per area were randomly chosen,
196  consisting of six unique neurons and six shared neurons. Shared neurons simulated features
197  characterizing all instances patterns of a category; they simulated shared conceptual features of all
198  category members (category-critical feature, e.g., members of the first category are robots in the same
199  height and are equipped with one camera, one speaker, two antennae, a power button, two metal legs,
200  and a pair of shoes; members of the second category are cats and have round-shaped head, eyes, nose,
201  mouth, ears, and whiskers (Figure 1B-left). Unique neurons simulated the 'idiosyncratic', fully
202  instance-specific visuomotor features; each of the corresponding feature neurons was only available in
203  one instance pattern (e.g., robots vary in the body shape and color, orientation of antennas, leg forms,
204  position of power button, shoes color). In sum, each category possessed 36 unique neurons from its six
205  exemplars and six shared neurons. For word form patterns, category term patterns (ii) of within-
206  category instances consisted of the same twelve neurons, which were co-activated with each of the 3
207  learnt grounding patterns of a category (e.g., to simulate the artificial words *fos* for all instances of the
208  robot category, and *coxt* for all instances of the cat category) (Figure 1B-middle); each proper name
209  pattern (iii) comprised twelve neurons, which were co-activated with one specific grounding pattern
210  (e.g., *xub*, *vit*, *hek* for the three instances of the robot category, respectively) (Figure 1B-right). The
211  choice of cells for pattern generation was pseudorandomized and constrained by the following criteria:
212  First, within-category neurons had to be non-adjacent to each other. This prevented coactivation
213  merely due to close distance. Second, no grounding patterns from two different categories shared any

214    neuron. Last, for each instance, the grounding patterns in *V1 and *M1$_L$ followed the same principles
215    but were not identical. The same rules applied for the grounding patterns in *A1 and *M1$_i$.

**Experiment design**

217    The current simulations involved three phases: model initialization, training phase, and testing phase,
218    which were carried out on the high-performance computing system of Freie Universität Berlin
219    (Bennett et al., 2020). During training, there were 3 different stimulation conditions, (1) where
220    grounding patterns were learnt without symbol (No symbol or control condition), (2) where all
221    grounding patterns of each category were presented together with the same word form pattern
222    (Category term condition), and (3) where each grounding pattern was co-presented with its own
223    specific word form pattern (Proper name condition). Thus, during learning, a stimulation pattern
224    included two activation patterns (to *V1 and *PF$_L$) when it was learnt outside symbol context (Figure
225    1C-top) or, a quadruplet including the two instance-related patterns plus two word form-related ones
226    (to A1 and PF$_i$) when learnt in symbol context (Figure 1C-bottom). Each test trial began with the
227    presentation of a grounding pattern of an instance (projected to the two sensorimotor model areas V1
228    and M1$_L$).

*Model initialization*

230        One crucial step prior to training was model initialization, which randomized all synaptic links
231    (and their corresponding weights) between within-area cells and between cells from connected areas.
232    Twelve sets of such synaptic links and weights (i.e., 12 different instantiations of the randomly
233    initialized neural network) were chosen, each set was then triplicated (cf. Schomers et al., 2017), and
234    each of these three copies entered one of the three training conditions – either No symbol, Category
235    term or Proper name. The use of distinct model instantiations can be seen as analogous to a within-
236    subject study design with 12 subjects. We chose to implement 3 separate sets of simulations for the 3
237    conditions to avoid any possible interference effects between concepts and symbols that may emerge
238    during training. Note, for example, that the relatively large representations that formed for category
239    terms might have interfered with further learning or may even have suppressed the activation of
240    conceptual representations without symbols. This configuration yielded a controlled 'within-subject'
241    design with training condition being a three-level repeated-measures factor (*No symbol*, *Category
242    term*, and *Proper name*). For the additional simulations performed to balance the number of word form
243    presentations, there were 4 levels.

*Training phase*

245        The neural network model was repeatedly presented with 30 instances from ten categories. To
246    mimic visuo-motor percepts associated with an instance, the extrasylvian primary sensorimotor areas,
247    *V1 and *M1$_L$, were each presented with its grounding pattern (i) for 16 time steps. Following the
248    experiment by LaTourrette and Waxman (2020) where instances were called either by a consistent
249    label or by distinct labels each, our within-category trained instances were either paired with the same
250    category term, by their distinct proper names, or they were not labeled at all. To mimic symbols in the
251    Category term and Proper name conditions, we presented to the primary perisylvian areas *A1 and
252    *M1$_i$ word form pattern (ii) and (iii), respectively, for 16 time steps (Figure 1C-bottom, 2A). Hence, in
253    different 'learning trials', the word form patterns of category terms were co-presented with one of
254    three different grounding patterns from one category, whereas those of proper names co-occurred with
255    only one specific grounding pattern. There were no word form patterns presented in the baseline No
256    symbol condition to control for the effect of either type of linguistic labels compared to learning
257    without ones (Figure 1C-top, 2A).

258        Because activity at the end of a trial might affect learning in the next trial, the network was
259    allowed to deactivate after each stimulated learning trial. To this end, we separated every two
260    consecutive pattern stimulations by a waiting interval during which only the uncorrelated white noise
261    mimicking spontaneous baseline neuronal firing was supplied to all areas (see principle 6 in Model
262    description – Neurobiological constraints). The goal was to reset the global network (i.e., all excitatory
263    and inhibitory cells displayed a membrane potential of zero) before a new grounding pattern was

264 inputted to the neural network model. This interstimulus interval (ISI) was terminated only after the
265 network activity had returned to its baseline value (thresh = 0.18, see table 2). As a result, the training
266 order was not influential in this experiment.

267 To balance learning conditions (NoS, CT, PN), each experiential grounding pattern
268 representing an instance was presented 2,000 times in one set of simulations. However, because each
269 category term pattern was co-presented with 3 different instance patterns, whereas proper name
270 patterns co-occurred with only one, this design leads to an imbalance of the number of learning trials
271 during which individual word form patterns were presented (3 times higher for category term than for
272 proper name presentations) (Figure 2C-top). Therefore, a second evaluation of learning trials was
273 performed and analyzed for which the number of word form pattern activations was balanced. In this
274 case, there were 1,000 learning trials in the Category term condition (CL_1x; each instance was
275 presented together with a category term in 1,000 training trials, resulting in a total of 3,000 training
276 trials per category terms) and 3,000 trials in the Proper name condition (PN_3x; each instance was
277 presented together with a proper name in 3,000 training trials, resulting in a total of 3,000 training
278 trials per proper name). For the control No symbol conditions, two comparison values were calculated,
279 after 1,000 (NoS_1x) and 3,000 (NoS_3x) trials (i.e., each instance was presented without symbol in
280 1,000 and 3,000 training trials, respectively) (Figure 2C-bottom). These different sub-designs are
281 summarized graphically in Figure 2C.

282 *Testing phase*
283 In the current experiment, we implemented a version of an old-new recognition task with the
284 use of new instances. For each of the ten categories we presented to the neural network six testing
285 instances: three trained instances and three novel instances (Figure 2B). In total, we used 30
286 previously learnt instances and 30 new instances. However, no actual old-new pairing took place
287 because we presented trained and novel instances to the neural network in separate test trials.

288 Memory performance of the network model was assessed in the absence of linguistic cues,
289 i.e., without stimulating the perisylvian primary areas *A1 or *M1i. To stimulate the experience of
290 individual instances, the extrasylvian primary areas *V1 and $*M1_L$ were activated for 2 time steps
291 with pure (i.e., free of any white noise) grounding patterns (i) and subsequentially deactivated towards
292 the baseline for 28 time steps. We recorded network responses 30 time steps from the onset of this
293 stimulation. Global resetting between two consecutive trials was conducted in the same manner as the
294 training phase. Hence, the test order was not of interest.

**Data analysis**
296 Grounding pattern production, data processing, and data analysis were performed using
297 Python 3.9.7, matplotlib 3.4.3 (Hunter, 2007), NumPy 1.20.3 (Harris et al., 2020), pandas 1.3.4
298 (Reback et al., 2022), SciPy 1.7.1 (Virtanen et al., 2020), seaborn 0.11.2 (Waskom, 2021). In the
299 current work, statistical significances were based on a conservative p-value threshold of .005
300 suggested by Di Leo and Sardanelli (2020). We used rstatix 0.7.0 (Kassambara, 2021) in the R
301 software environment (R Core Team, 2021) for statistical analyses.

302 When testing stimuli were presented to the primary sensorimotor areas, some of the 625
303 excitatory neurons per area fired in response to their conceptual grounding patterns. As described in
304 Procedure, we recorded all their responses during 30 time steps from stimulation. Let $\phi(e, t)$ denote
305 the output of an excitatory cell $e$ at time $t$, such that $\phi$ only takes up the value 0 or 1 and $t$ only allows
306 discrete values up to 30 (corresponding to thirty possible simulation time steps); let $\tau_{Favg} = 5$ be a
307 time constant, the estimated instantaneous firing rate $\omega_E(e, t)$ of cell $e$ at time $t$ can be calculated
308 based on the following equation:

309
$$\tau_{Favg} \cdot \frac{d\omega_E(e,t)}{dt} = -\omega_E(e,t) + \phi(e,t) \quad \text{Eq. (1)}$$

310     Solving Eq. (1) for $\omega_E(e, t)$ returns the cell's latest spiking activity (firing rate). We estimated
311 the mean firing rate based on $t = t_{30}$ and used this value for the subsequent RSAs. For details about
312 relevant calculation steps, see *Appendix* of Henningsen-Schomers and Pulvermüller, (2022).

313     Previous research found that several of the extra-sylvian areas targeted by the deep neural
314 model (including, for example, *V1 and *AT) are important for processing instance- and concept-
315 related information (see, for example, Binder et al., 2005; Martin, 2007; Ralph et al., 2017;
316 Henningsen-Schomers et al., 2023). Therefore, the current data analyses and statistical testing focused
317 on the extrasylvian region of the deep neural network. This decision was motivated by the main aim of
318 addressing possible causal influences of symbol learning on the perceptual processing of instances of
319 concepts and on conceptual processing itself.

### *Representational similarity analysis*

321     The estimated mean firing rate of 625 neurons in response to a testing instance reflected how
322 this instance was represented in neural network. To understand how differently the neural network
323 represented within- and between-category instances, we calculated the dissimilarity in firing patterns
324 for every pair of the 60 instances. Pairwise dissimilarities computed in terms of Euclidean distance
325 were organized in a $60 \times 60$ representational dissimilarity matrix (RDM) (Figure 3A): Each cell in the
326 matrix reflected the dissimilarity between the firing patterns of two instances. In total, there were 36
327 RDMs across three training conditions and twelve areas.

328     We defined two classes of pairwise dissimilarities, including between-category dissimilarity
329 ($Dissim_B$) and within-category dissimilarity ($Dissim_W$). A second way to define similarity types is
330 based on the type of instances under study, that is, dissimilarity between two trained instances
331 ($Dissim_{TT}$), between two novel instances ($Dissim_{NN}$), and between a trained and a novel instance
332 ($Dissim_{TN}$). For example, within-category dissimilarity could be classified as either dissimilarity
333 among trained instances $1 - 3$ ($Dissim_{W-TT}$), among novel instances $4 - 6$ ($Dissim_{W-NN}$) or between
334 trained and novel instances ($Dissim_{W-TN}$) (Figure 3A).

335     **Category learning.** Category learning was ~~defined as~~ evaluated through the ability to (1)
336 distinguish differences between categories and (2) group together category members. We assessed
337 how different types of symbols impacted upon category learning performance based on (1) the
338 dissimilarity between two between-category trained instances ($Dissim_{B-TT}$), and (2) the dissimilarity
339 between two within-category trained instances ($Dissim_{W-TT}$) (Figure 3A). Successful category
340 learning occurred when two instances from two distinct categories were considered as dissimilar (high
341 $Dissim_{B-TT}$) and/or when two within-category instances were considered as similar (low
342 $Dissim_{W-TT}$). If, as previously claimed, applying category terms invites one to encode the
343 commonalities among instances and thereby facilitates categorization, the deep neural network should
344 represent within-category instances similarly while highlighting the dissimilarities between instances
345 of different categories. In the Category term condition, we expected between-category dissimilarities
346 to be greater than within-category dissimilarities $Dissim_{B-TT_{CT}} > Dissim_{W-TT_{CT}}$. By contrast, we
347 proposed two scenarios for the Proper name condition. In the first scenario, if proper names focus the
348 neural network models on encoding only unique features and inhibit the encoding of category-critical
349 features, no traces of category learning will be observable, and the representations of individual
350 instances will be highly dissimilar regardless of their categorical membership ($Dissim_{B-TT_{PN}} \approx$
351 $Dissim_{W-TT_{PN}}$). However, because within-category instances shared 50% of their activated neurons
352 in the extrasylvian primary areas *V1 and *M1$_L$, the neural network could base on such similarities to
353 form category representation. In this second scenario, proper names are not sufficient to override
354 category learning; the neural network would house not only the unique representations of the instances
355 but also the commonalities of those belonging to the same category. Like the Category term condition,
356 the test data would as well show signs of category learning ($Dissim_{B-TT_{PN}} > Dissim_{W-TT_{PN}}$).
357 Taking into account such intrinsic perceptuomotor similarities among instances from the same
358 category, category learning was evaluated not only across symbol (i.e., category term or proper name)

359  learning conditions but also in control conditions (i.e., training without symbols). For example, a
360  superior causal influence of category terms on category learning performance would be expressed
361  through a significantly higher $Dissim_{B-TT_{CT}}$ and lower $Dissim_{W-TT_{CT}}$ relative to training with proper
362  names and also relative to training without symbols.

363      **Generalization.** Assuming the neural network had encoded the commonalities between
364  within-category trained instances and formed category knowledge with the help of these shared
365  features, they might have as well represented novel instances as members of that category when
366  exposed to the category-critical features in these novel instances. Generalization performance would
367  then be reflected by how similarly within-category trained instances and within-category novel
368  instances stimulated the deep neural network. To evaluate the generalization performance of the neural
369  network on novel instances, pairwise dissimilarities between two trained instances ($Dissim_{W-TT}$) as
370  well as between a trained and a novel instance ($Dissim_{W-TN}$) were extracted. In the testing phase the
371  chance was low that the neural network readily applied category knowledge earned from thousands of
372  training trials onto a novel instance in the first and only exposure. In the case of poor generalization
373  performance, the activation pattern of within-category novel instances would be dissimilar from that of
374  the within-category trained instances (i.e., increasing $Dissim_{W-TN}$). Our criterion for a successful
375  generalization after learning with symbols was that $Dissim_{W-TN}$ should be as low as $Dissim_{W-TT}$
376  ($Dissim_{W-TN} \approx Dissim_{W-TT}$). In other words, their absolute dissimilarity difference $DissimDiff =$
377  $|Dissim_{W-TN} - Dissim_{W-TT}|$ must remain lower than when the deep neural network was trained
378  without symbols.

379  *Cell assembly analysis*
380      Motivated by the notion of cell assemblies (see, e.g., Hebb, 1949; Braitenberg, 1978; Fuster,
381  2005), that is, strongly interlinked sets of neurons forming as a consequence of correlated neuronal
382  activity and potentially carrying a main role in cognitive brain processing, we conducted cell assembly
383  analyses to discover possible neuronal correlates of grounding instances, concepts and symbols along
384  with instance-specific and category-critical neurons after repeated exposure to instances and their
385  category terms or proper names. We extracted cell assemblies (CAs) activated by each of the 60
386  grounding patterns used as testing instances based on the criterion described in previous work
387  (Garagnani and Pulvermüller, 2016; Henningsen-Schomers and Pulvermüller, 2022). Grounding
388  patterns in the testing phase tended to coactivate several excitatory neurons (e-cells) in an area, with at
389  least one being maximally responsive (non-response was under the threshold 0.01). To be part of a
390  CA, the firing rate of a given e-cell had to exceed 75% of the firing rate of the maximally responsive
391  cell of the same area. We then computed the number of unique, instance-specific and overlapping,
392  conceptual neurons among CAs for trained instances of the same category: neurons were classified
393  according to whether they were activated by just one grounding patterns or whether they responded to
394  two or three instances (thus being pair or triple-shared between the learnt instances of a concept).
395  Unique neurons were conceptualized as neurons which encoded specific, 'idiosyncratic' features of an
396  instance; shared neurons could be understood as those that encoded common features shared by at
397  least two instances and thus characteristic of their category. The specialized encoding of category-
398  critical features could be indicated by a higher proportion of shared neurons per area, while traces of
399  instance-specific features would be reflected by a larger proportion of unique neurons.

400      Representations are transformed through different levels of processing, i.e., from the primary
401  areas to secondary areas, and the central "connector hub" areas of the model. We quantified such
402  transformation as the change (i.e., gain/loss) in the number of unique and shared CA-cells in the
403  extrasylvian central areas (AT, PF$_L$) comparative to the extrasylvian primary areas (V1, M1$_L$). Gains in
404  a type of neuron, for example, shared neurons, are indicative of intensive encoding of concept related
405  commonalities on the course of processing, while loss of shared neurons in the central areas implies
406  that reduced encoding of idiosyncratic features and hence instance-related information. Percentage
407  gain was calculated as the difference between the number of neurons in the central and primary areas,
408  as a percentage with respect to the number of neurons in the primary areas:

409
$$Gain = \frac{n_{central} - n_{primary}}{n_{primary}} \times 100$$

410       **Representations of category-critical features.** A range of previous neurocomputational
411 studies show, that, when brain-like networks learn concepts and word meanings, they form cell
412 assemblies that are spread out across sensorimotor and more central areas of the network. The density
413 of shared semantic neurons in the most central connector hubs is greatest due to their high connectivity
414 degree and thus ample convergence of activity on these areas, resulting in especially strong activation,
415 in particular for shared semantic neurons (for discussion, see Garagnani et al., 2017; Tomasello et al.,
416 2018). Relative to instance-specific neurons, shared semantic neurons are activated more frequently
417 during semantic learning, which predicts that these will recruit the largest number of additional cell
418 assembly; these would therefore be semantic, too, and primarily located in the central hub regions. If a
419 labeling condition specifically invites the neural network to encode category-relevant features, we
420 expect (1) more shared neurons than unique neurons in the extrasylvian areas and (2) a greater gain in
421 shared neurons in the central semantic areas compared to the primary areas. Category learning might
422 still occur even in the presence of proper names because within-category similarities also characterize
423 sensorimotor experiences. If such information is sufficient, there should be traces of shared neurons in
424 the central, multimodal areas as well. Additionally, category terms should activate shared neurons
425 more than proper names.

426       **Representations of instance-specific features.** When a neural network represents instances
427 as unique entities, it shall reveal specific traces of each instance in the extrasylvian areas, especially in
428 the semantic hubs. In an extreme case where category learning is hindered and the neural network only
429 encodes the uniqueness of instances, there should be (1) more unique than shared neurons in the
430 extrasylvian areas and (2) a gain only in unique neurons in the central areas with respect to the primary
431 areas. Importantly, instances with proper names are expected to activate significantly more unique
432 neurons than categorically labelled instances.

433       We gather from all twelve model instantiations the CAs in response to all 30 trained instances
434 of 10 categories and classify CA-cells by their uniqueness to each instance (versus sharedness). To
435 facilitate readers' understanding about the results, we offer an interactive illustration of these CAs on
436 our web application at (https://phucthuun.shinyapps.io/CL_PN/). This web application enables one to
437 compare the differential effects of category terms versus proper names in representing category-
438 critical and instance-specific features of within-category and across-category instances.

439

440                                                                    **Results**

441     **Representational similarity analysis**

442          Figure 3B gives a first impression of the instance and category learning performance after
443     2,000 training trials. In the Category term condition, instances from the same category activated the
444     neural network similarly, whereas instances from different categories led to substantially more
445     dissimilar activation patterns across the different areas of the network (i.e., firing patterns were highly
446     dissimilar, as color-coded by dark blue and pink). Category knowledge was reflected in a relatively
447     reduced dissimilarity (light blues), which appears as homogenous within each category, contrasting
448     with those between categories, especially in the central areas (semantic hubs). Training the deep neural
449     network without the aid of symbols or with proper names reduced the networks' ability to distinguish
450     instances between categories: activity pattern dissimilarities between instances from different
451     categories were much more substantial in the Category term condition than in the Proper name
452     condition (color-coded with shades of intermediate blue). In contrast, within-category similarities and
453     generalization performance in the Category term condition were superior, as indicated by the more
454     homogeneous (light) blue shade across all 6 instances (trained and not-trained) from the same
455     category, relative to the other two conditions, where different shades of light blue are clearly visible.

456     *Category learning*

457          To evaluate category learning performance after 2,000 learning trials, within-category
458     dissimilarity ($Dissim_{W-TT}$) and between-category dissimilarity between activity patterns elicited by
459     grounding patterns of trained instances ($Dissim_{B-TT}$) were used. Figure 4A describes a global
460     tendency of the deep neural network, across its twelve areas and three training conditions, to identify
461     within-category instances as more similar and between-category instances as more dissimilar to each
462     other. This feature is explained by the grounding patterns presented, which were similar across
463     category instances, but not between. However, between-category dissimilarity is relatively enhanced
464     in central areas, a feature not explained by the stimulations. In the next step, dissimilarity values were
465     averaged for the six extrasylvian areas. The two-factorial repeated measures $(3 \times 2)$ – ANOVA with
466     training condition (No symbol/Category term/Proper name) and dissimilarity type
467     ($Dissim_{W-TT}/Dissim_{B-TT}$) confirmed the main effect of both factors ($F(2,22) = 2777.647, p <$
468     $.001, \eta^2 = 0.982$ and $F(1,11) = 11155.611, p < .001, \eta^2 = 0.996$, respectively) as well as their
469     interaction effect ($F(2,22) = 6113.987, p < .001, \eta^2 = 0.986$) on the dissimilarity between instances
470     within these extrasylvian areas. Figure 4B illustrates category-related activation performance of the
471     deep neural network in the extrasylvian areas of the three learning conditions: the neural network
472     successfully grouped together instances from the same category while distinguishing between
473     instances from the same vs. from two different categories. Pairwise comparisons with Bonferroni
474     correction were computed to observe the effect of training condition on each level of dissimilarity type
475     and vice versa. The results showed that $Dissim_{W-TT}$ was significantly lower than $Dissim_{B-TT}$ in all
476     three conditions ($ps < .001$); same category-membership was thus manifest as relatively enhanced
477     activation similarity in all conditions and across areas. The $Dissim_{W-TT}$ in the Category term
478     condition ($M$=0.229, $SD$=0.005) and the Proper name condition ($M$=0.264, $SD$=0.004) was
479     significantly smaller (i.e., greater similarity) than that in the control No symbol condition ($M$=0.29,
480     $SD$=0.006), and they were also significantly different from each other, with greatest similarities after
481     category term labelling ($ps < .001$). Relative to the control No symbol condition, the deep neural
482     network responded similarly to trained instances coming from the same category when it was trained
483     with symbols and such performance was above baseline. Importantly, the benefit of category terms
484     was superior to both training without symbols and with proper names. Likewise, the deep neural
485     network returned the highest $Dissim_{B-TT}$ ($M$=1.48, $SD$=0.018) for the Category term condition ($ps <$
486     $.001$), while $Dissim_{B-TT}$ in the Proper name condition ($M$=0.706, $SD$=0.01) was not significantly
487     different from that in the No symbol condition ($M$=0.749, $SD$=0.045) ($p = 0.01$), after application of
488     the Bonferroni-corrected significance threshold of $.005$. Compared to the No symbol condition,
489     training with proper names only gradually hindered the discrimination of between-category instances
490     but left the separation of within-category instances unaffected. By contrast, both aspects of category

491    learning were present with the aid of category terms, reduced within- and enhanced between-category
492    similarities.

493         The simulations performed to control for the number of word form presentations during
494    learning were evaluated using a two-factorial repeated measures ($4 \times 2$) – ANOVA with training
495    condition (now 4 levels, NoS_1x/NoS_3x/CT_1x/PN_3x) and dissimilarity type
496    ($Dissim_{W-TT}/Dissim_{B-TT}$). This confirmed the main effect of both factors ($F(1.67,18.35) =$
497    $1113.758, p < .001, \eta^2 = 0.964$ and $F(1,11) = 7485.295, p < .001, \eta^2 = 0.993$, respectively) as
498    well as their interaction effect ($F(1.65,18.10) = 1961.497, p < .001, \eta^2 = 0.973$) on the
499    dissimilarity between instances within extra-sylvian areas. Pairwise comparisons with Bonferroni
500    correction were computed to observe the effect of training condition on each level of dissimilarity type
501    and vice versa. In essence, $Dissim_{B-TT}$ in the Category term condition was significantly higher than
502    that in the Proper name and both No symbol control conditions ($ps < .001$) (Figure 4C); category-
503    term learning increased the dissimilarity across conceptual categories relative to no-symbol learning
504    and proper-name learning. The reverse effect, greater dissimilarity values for proper names than
505    category terms, was found within categories. These observations were therefore valid even when
506    proper names were 'shown' to the model three times more than category terms during learning.

507    ***Generalization***
508         To evaluate the generalization performance of the deep neural network on novel instances,
509    pairwise dissimilarities between two trained instances ($Dissim_{W-TT}$) as well as between a trained and
510    a novel instance ($Dissim_{W-TN}$) were used. Figure 5A illustrates the tendency of the deep neural
511    network to represent two trained instances of the same category as more dissimilar, whereas the
512    representations of a novel and a trained instance from the same category were less dissimilar (lighter-
513    shaded columns were mostly higher than darker-shaded columns). In the six extrasylvian areas, a
514    $3 \times 2$ – ANOVA was computed with training condition (No symbol/Category term/Proper name) and
515    type of within-category dissimilarity ($Dissim_{W-TT}/Dissim_{W-TN}$) as repeated measures factors. Both
516    the main effects of training condition ($F(2,22) = 465.217, p < .001, \eta^2 = 0.956$) and dissimilarity
517    type ($F(1,11) = 7711.618, p < .001, \eta^2 = 0.939$) were significant. For these two factors, there was
518    also a significant interaction ($F(2,22) = 635.788, p < .001, \eta^2 = 0.707$) (Figure 5B). The
519    Greenhouse-Geisser sphericity correction to the violated sphericity assumption ($p = .024$) for training
520    conditions ($p[GG] = 2.38 \times 10^{-11}$) confirmed this result. Two-sided pairwise comparisons with
521    Bonferroni correction showed that $Dissim_{W-TN}$ in the Category term ($M = 0.214, SD = 0.004$) and
522    in the Proper name conditions ($M = 0.220, SD = 0.003$) were significantly lower than that in the
523    control No symbol condition ($M = 0.249, SD = 0.004$) ($ps < .001$), but they did not differ
524    significantly from each other ($p = .01$) (Figure 5B). $Dissim_{W-TN}$ was significantly lower than
525    $Dissim_{W-TT}$ in all three conditions ($ps < .001$) (Figure 5B), which means that within-category
526    trained instances were represented as less similar to each other than when each of them was compared
527    with a novel instance from the same category. In other words, trained instances resulted in neuronal
528    response patterns that were more similar to those caused by novel instances than those caused by
529    trained instances from the same category, a finding easily explained by the lack of learning of the
530    idiosyncratic features of novel instances. A further set of pairwise comparisons using Bonferroni
531    correction revealed that the absolute $DissimDiff$ in the No symbol condition ($M = 0.041, SD =$
532    $0.016$) was significantly higher than $DissimDiff$ in the Category term condition ($M = 0.016, SD =$
533    $0.012$) ($p < .001$) but not significantly different from that in the Proper name condition ($M =$
534    $0.044, SD = 0.02$) ($p = .009$). In other words, category-term learning resulted in the most similar
535    processing of learnt and not-learnt instances and thus to the greatest degree of generalization.

536         Results from the additional simulations controlling for the number of word form presentations
537    during learning (i.e., four training conditions NoS_1x, NoS_3x, CT_1x, PN_3x, see Methods) also
538    confirmed that generalization was maximal for novel members of categories for which category term
539    had been learned (Figure 5C). The mere exposure to instances or learning proper names showed little
540    generalization relative to category learning.

541     These results investigating brain-constrained neural network correlates of conceptual
542 generalization sit well with well-known observations that language-learning children often generalize
543 – or even overcategorize – category terms to novel items. In case of overgeneralization to an item,
544 subsequent learning may establish a novel category to which the item belongs. While our results offer
545 a mechanistic perspective on generalization, a detailed simulation of overgeneralization and
546 reclassification learning is left for future study.

547 **Cell assembly analysis**
548     Figure 6A illustrates the tendency of the deep neural network to encode fewer unique neurons
549 (U-shaped function across areas) and more shared neurons (inverted U-shaped function) in the extra-
550 sylvian central areas than in the extra-sylvian primary areas. In the first step, the number of unique
551 neurons and shared neurons activated by each instance were calculated and averaged across two
552 training conditions. The repeated measures $3 \times 2$ – ANOVA with training condition (No
553 symbol/Category term/Proper name) and neuron type (unique/shared) confirmed the significant main
554 effects ($F(2,22) = 902.098, p < .001, \eta^2 = 0.926$ and $F(1,11) = 13966.410, p < .001, \eta^2 =$
555 $0.998$, respectively), and a significant interaction involving both factors ($F(2,22) = 5027.907, p <$
556 $.001, \eta^2 = 0.985$). The supplementary $2 \times 2$ – ANOVA with training condition with symbols
557 (Category term/Proper name) and neuron type (unique/shared) returned comparable results with 2
558 significant main effects ($F(1,11) = 1009.255, p < .001, \eta^2 = 0.951$ and $F(1,11) = 23994.328, p <$
559 $.001, \eta^2 = 0.998$, respectively), and a significant interaction involving both factors ($F(1,11) =$
560 $4593.789, p < .001, \eta^2 = 0.986$). Pairwise comparisons with Bonferroni correction revealed that
561 category terms made the neural network reactivate more shared neurons ($M = 11.242, SD = 0.127$)
562 than unique neurons ($M = 2.861, SD = 0.051$) ($p < .001$). This also applied for training with proper
563 names (shared neurons: $M = 7.963, SD = 0.222$; unique neurons: $M = 3.89, SD = 0.064$) and
564 training without symbol (shared neurons: $M = 8.029, SD = 0.194$; unique neurons: $M = 4.493, SD =$
565 $0.08$) ($ps < .001$) (Figure 6B). Compared to this control condition, the number of unique instance-
566 specific neurons was moderately reduced by proper names, but radically so by category terms ($p <$
567 $.001$), whereas the number of shared, conceptual-category neurons remained unchanged after proper-
568 name learning ($p = .447$), but increased dramatically with category term acquisition ($p < .001$). The
569 latter is clear evidence for a facilitatory effect of language, more specifically, of category-term
570 learning, on conceptual category formation in brain-constrained deep neural networks.

571     With respect to the gain/loss of neurons in the extrasylvian central areas relative to the primary
572 ones, our repeated-measure $3 \times 2$ – ANOVA with two factors training condition (No symbol/Category
573 term/Proper name) and neuron type (unique/shared) confirmed both main effects on the percentage
574 change of neurons and their interaction to be significant ($F(2,22) = 55.17837, p < .001, \eta^2 =$
575 $0.5519424, F(1,11) = 6471.54090, p < .001, \eta^2 = 0.9954$, and $F(2,22) = 1484.43893, p <$
576 $.001, \eta^2 = 0.966$, respectively). According to the subsequent pairwise t-tests, the deep neural
577 networks gained shared neurons but lost unique neurons in the central areas, which held true for all
578 conditions ($ps < .001$) (Figure 6D, upward dotted lines represent positive gains in shared neurons and
579 downward solid lines mean negative gains in unique neurons). On the three levels of training
580 condition, the gain in shared neurons and the loss in unique neurons in the Category term condition
581 were significantly larger than that in the Proper name and No symbol conditions ($ps < .001$) (Figure
582 6D). Proper names did not significantly increase the gain in shared neurons ($p = .1$) but led only to a
583 moderate loss of unique neurons, as compared to the control training condition ($ps < .001$). These
584 results further confirm that training with category terms magnified both the gain in shared semantic
585 neurons in central areas and the loss of unique instance-specific neurons there. The simulations
586 performed for balancing the number of word form presentations during proper-name and category-
587 term learning also confirmed these observations (Figure 6C, E). Therefore, the overgrowth of shared
588 neurons in category-term learning does not depend on an abundant number of word form presentations
589 and cannot be explained by adding word form information to instance-related information.

590        Both RSA and CA-analyses were also conducted for the whole model architecture (6
591 extrasylvian and 6 perisylvian model areas). The data replicated the results indicating category
592 learning (Figure 4-1, Table 4-1), generalization (Figure 5-1, Table 5-1), and representations of
593 category-critical as well as instance-specific features (Figure 6-1, Table 6-1).

**Discussion**

594
595         When sensorimotor patterns simulating the processing of similar objects or actions from
596  different categories were presented, the brain-constrained network applied in the current study showed
597  successful conceptual category learning. Category learning outside symbol context was manifest in
598  greater similarities of activity patterns elicited by different instances of the same category as compared
599  with between-category pattern similarities. Importantly, compared with training of instances per se,
600  concurrent learning of category instances and symbols had a substantial effect on both categorial and
601  instance-specific processes. Category-term learning led to an additional increase in dissimilarities
602  between activity patterns across conceptual categories, while making category members substantially
603  more similar to each other. In contrast, proper-name learning did not change between-category
604  similarities and led to a relatively minor similarity increase between members of the same category.
605  The model gave evidence of generalization to novel members of learned categories and showed that
606  such generalization was maximal for novel members of categories for which category terms had been
607  learned. Meticulous analyses of neuronal activity patterns suggest that the enhancement of within-
608  category similarities and between-category dissimilarities in context of category symbols is due to an
609  increase in the number of cells responding to all category members. Likewise, relative persistence of
610  instance-specific neurons with proper-name learning underlies the maintained activation differences
611  between category instances observed in this case. All observed effects regarding pattern dissimilarities
612  and neuronal microstructure were greatly pronounced in the central 'connector hub' areas of the brain-
613  constrained model applied, as compared with primary areas. Table 3 summarizes major observations
614  in the current data and the corresponding learning aspects these observations reflect.

615  **Relationship to experimental and neurocomputational research**
616         Our results can be used to address observations delivered by neurocognitive and
617  neurobehavioral experiments. Neuropsychological evidence highlights the role of the prefrontal cortex
618  in categorical representation (for review see Kéri, 2003). Prefrontal areas (PF$_i$ and PF$_l$) are part of the
619  four central areas of our model, where conceptual neurons constituting category representations
620  emerged most numerously. This is explained by the high degree of convergence of neural activity on
621  these areas, which are not only located in the centre of the model architecture but also show the
622  highest connectivity degrees. Due to ample activity converging on these connector hub areas, their
623  frequently activated shared semantic neurons can most efficiently recruit other neurons, which
624  therefore take on similar response properties (Doursat and Bienenstock, 2006). This mechanism may
625  contribute to why these areas act as 'semantic hubs' and house neurons reflecting category
626  membership (e.g., PF and AT, see Miller et al., 2002; Seger and Miller, 2010; Garagnani and
627  Pulvermüller, 2016; Tomasello et al., 2017). On the other hand, the higher density of instance-specific
628  neurons in the primary visual/motor model area relative to the centre is evidence for exemplar learning
629  in the sensorimotor cortices (Bowman et al., 2020; Kéri, 2003) – a type of category learning that is
630  based on the representations of specific category instances (Nosofsky, 1988) and should be
631  independent of signs and symbols. Here, solid evidence for category formation was obtained even in
632  the control condition where only sensorimotor patterns were presented to the model without symbols.
633  In line with neural data (Freedman et al., 2001; Seger and Miller, 2010), experimental evidence shows
634  that perceptuomotor similarities among category members are sufficient to trigger category learning in
635  preverbal infants (Sloutsky and Fisher, 2004; de Heering and Rossion, 2015) and animals (Güntürkün
636  et al., 2018; Pusch et al., 2023).

637         When learning conceptual instances in context of category terms, infants show most
638  pronounced category building and an attention bias towards shared features of category members
639  (Waxman and Markow, 1995; Dewar and Xu, 2007; Althaus and Mareschal, 2014). In contrast,
640  encountering proper names for individual instances focuses their attention relatively more to object-
641  specific features (Barnhart et al., 2018; Pickron et al., 2018; La Tourette & W, 2020). In the current
642  network model, symbol association raises the number of neurons involving in the processing of a
643  given sensorimotor pattern. This can be interpreted as biased attention to the object or action for which
644  the pattern codes and thus explains why label learning generally increases attention to object features.

645 Furthermore, as category-term learning increases the number of category-critical shared semantic
646 neurons in the network, at the cost of reducing the number of instance-specific ones, the pre-observed
647 greater attention to shared features has a direct model correlate, along with the label-related tendency
648 to build stronger category representations. Infants' attentional focusing on instance-specific features of
649 objects is in line with the relative preservation of instance-specific neurons in the model of proper-
650 name learning. Thus, the opposing effects of proper name and category-term learning, which,
651 respectively, drive attention towards instance-specific and category general features of objects, are
652 captured by the current model.

653         A range of neurocomputational studies previously explored the putative brain basis of
654 cognitive processes (e.g., Deco and Rolls, 2005; Rolls and Deco, 2015; Palm, 2016), including
655 conceptual category learning and the influence of language on object perception (Rogers and
656 McClelland, 2014; Henningsen-Schomers and Pulvermüller, 2022). For example, Westermann and
657 Mareschal (2014) demonstrated, using a fully distributed parallel processing model, that learning a
658 category label made the neural patterns of category members more similar to each other, whereas
659 different categories moved away from each other in representational space. Our RSA in models
660 mimicking cortical area structure and connectivity, along with within-area excitatory and inhibitory
661 connectivity, achieved the same result. In addition, we determined the neuron-level mechanisms and
662 contributions of different model areas to this result and, in particular, revealed the model-central
663 connector hub areas as the loci where the differences between categorical and instance-specific
664 mechanisms as well as those between the shared- vs. specific-feature promoting roles of instance-
665 specific and category labels are most pronounced. As to our knowledge, the contrast between activity
666 patterns and neuronal correlates of proper names and category terms has not been addressed by
667 previous computational work.

**Model explanation**

669         The present simulations offer explanations of the observed phenomena based on neuroscience
670 principles. Of special relevance here are the biological learning mechanisms applied, which include
671 unsupervised Hebbian synaptic strengthening of connections between co-activated neurons and
672 weakening of links between cells firing independently of each other. This principle explains why
673 category labels primarily interlink with the shared neurons of instance representations belonging to the
674 same category. The reason lies in the highest correlation values, as instance-specific neurons are silent
675 when the category term is used together with other category instances. This implies some weakening
676 of connections between the category terms' and the instance-specific neurons, based on the 'anti-
677 Hebbian' "neurons out-of-sync delink" rule. The opposite difference applies to proper names, whose
678 neural correlates strongly connect to instance-specific neurons but weaken their links with the
679 category-critical shared neurons whenever a different category member co-occurs with its own and
680 thus different name. Effects are most clearly present in the central areas of the network where the
681 neural correlates of words and entities are equally manifest so that their correlation structure can easily
682 be mapped.

**Limitations and future direction**

684         The current simulations use idealized instance and category learning conditions. The
685 activation patterns representing conceptual instances and word forms were chosen to be non-
686 overlapping, except for the neurons coding for shared features. These are idealizations considering
687 both the features of word forms and those of objects and actions could be shared across categories (cf.
688 phonological e.g., "cat"-"hat" or perceptual color/shape similarities). Such similarities are irrelevant to
689 category membership and hence were omitted to keep the simulation well-controlled. Secondly, only a
690 small number of conceptual features were realized and a small set of shared features determined
691 concept membership. This situation may hold for some concrete terms but not for others and certainly
692 not for abstract concepts (Henningsen-Schomers et al., 2022). Furthermore, proper names and
693 category terms were acquired by different networks to allow straightforward separation and evaluation
694 of the mechanistic side of different label types – although label types are normally co-present in the

695   same mind and brain. In future, it is desirable to complement this work by simulations of more
696   realistic conceptual categories and to build one model in which interaction/interference effects
697   between different learning conditions are possible.

698   **Conclusion**
699           The current study strived to meet the need for a mechanistic model of symbols and their
700   meaning within a neurobiological computational framework by addressing specific features of proper
701   names (Mickey Mouse) and category symbols (house mouse). Developmentalists and linguists have
702   long been proposing that category terms and proper names distinctively impact infants' locus of
703   attention towards category-shared and instance-specific object and action features, respectively. By
704   simulating concept and instance learning in a deep neural network with neurobiologically realistic
705   architecture and brain-like connectivity, we demonstrate that learning these two different symbol types
706   had opposing effects on the emergent neuronal cell assemblies representing and processing instances
707   of a category and the shared conceptual features of that category, which can explain pre-observed
708   differences in perceptual, attentive and memory processes related to the specific and shared features of
709   category instances. These explanations were based on unsupervised Hebbian associative learning
710   mechanism binding neurons involved in correlated processing of instance-specific category-general
711   information. The current work could thus not only replicate but also offer underlying neuronal
712   mechanisms and causal neurobiological explanations for well-established observations in cognitive
713   science.

714 **References**
715 Althaus N, Mareschal D (2014) Labels Direct Infants' Attention to Commonalities during Novel
716      Category Learning. PLOS ONE 9:e99670.

717 Althaus N, Plunkett K (2016) Categorization in infancy: labeling induces a persisting focus on
718      commonalities. Developmental Science 19:770–780.

719 Arikuni T, Watanabe K, Kubota K (1988) Connections of area 8 with area 6 in the brain of the
720      macaque monkey. Journal of Comparative Neurology 277:21–40.

721 Artola A, Singer W (1993) Long-term depression of excitatory synaptic transmission and its
722      relationship to long-term potentiation. Trends in Neurosciences 16:480–487.

723 Baldwin DA, Markman EM (1989) Establishing World-Object Relations: A First Step. Child
724      Development 60:381–398.

725 Barnhart WR, Rivera S, Robinson CW (2018) Effects of Linguistic Labels on Visual Attention in Children
726      and Young Adults. Front Psychol 9:358.

727 Bauer RH, Fuster JM (1978) The effect of ambient illumination on delayed-matching and delayed-
728      response deficits from cooling dorsolateral prefrontal cortex. Behavioral Biology 22:60–66.

729 Bauer RH, Jones CN (1976) Feedback training of 36 – 44 Hz EEG activity in the visual cortex and
730      hippocampus of cats: Evidence for sensory and motor involvement. Physiology & Behavior
731      17:885–890.

732 Bennett L, Melchers B, Proppe B (2020) Curta: A General-purpose High-Performance Computer at
733      ZEDAT, Freie Universität Berlin. :5 S.

734 Best C, Robinson C, Sloutsky V (2010) The Effect of Labels on Visual Attention: An Eye Tracking Study.
735      Proceedings of the Annual Meeting of the Cognitive Science Society 32 Available at:
736      https://escholarship.org/uc/item/0wn1j6px [Accessed October 9, 2022].

737 Binder JR, Westbury CF, McKiernan KA, Possing ET, Medler DA (2005) Distinct Brain Systems for
738      Processing Concrete and Abstract Concepts. Journal of Cognitive Neuroscience 17:905–917.

739 Bowman CR, Iwashita T, Zeithamova D (2020) Tracking prototype and exemplar representations in
740      the brain across learning Behrens TE, Barense M, Barense M, Tompary A, eds. eLife 9:e59360.

741 Braitenberg V (1978) Cell Assemblies in the Cerebral Cortex. In: Theoretical Approaches to Complex
742      Systems (Heim R, Palm G, eds), pp 171–188 Lecture Notes in Biomathematics. Berlin,
743      Heidelberg: Springer.

744 Braitenberg V, Schüz A (1998) Cortex: Statistics and Geometry of Neuronal Connectivity. Berlin,
745      Heidelberg: Springer Berlin Heidelberg. Available at: http://link.springer.com/10.1007/978-3-
746      662-03733-1.

747 Bressler SL, Coppola R, Nakamura R (1993) Episodic multiregional cortical coherence at multiple
748      frequencies during visual task performance. Nature 366:153–156.

749 Catani M, Jones DK, Donato R, ffytche DH (2003) Occipito-temporal connections in the human brain.
750      Brain 126:2093–2107.

751    Catani M, Jones DK, Ffytche DH (2005) Perisylvian language networks of the human brain. Annals of
752        Neurology 57:8–16.

753    Chafee MV, Goldman-Rakic PS (2000) Inactivation of Parietal and Prefrontal Cortex Reveals
754        Interdependence of Neural Activity During Memory-Guided Saccades. Journal of
755        Neurophysiology 83:1550–1566.

756    Connors BW, Gutnick MJ, Prince DA (1982) Electrophysiological properties of neocortical neurons in
757        vitro. Journal of Neurophysiology 48:1302–1320.

758    de Heering A, Rossion B (2015) Rapid categorization of natural face images in the infant right
759        hemisphere Culham JC, ed. eLife 4:e06564.

760    Deacon TW (1992) Cortical connections of the inferior arcuate sulcus cortex in the macaque brain.
761        Brain Research 573:8–26.

762    Deco G, Rolls ET (2005) Neurodynamics of Biased Competition and Cooperation for Attention: A
763        Model With Spiking Neurons. Journal of Neurophysiology 94:295–313.

764    Dewar K, Xu F (2007) Do 9-month-old infants expect distinct words to refer to kinds? Developmental
765        Psychology 43:1227–1238.

766    Di Leo G, Sardanelli F (2020) Statistical significance: p value, 0.05 threshold, and applications to
767        radiomics—reasons for a conservative approach. Eur Radiol Exp 4:18.

768    Distler C, Boussaoud D, Desimone R, Ungerleider LG (1993) Cortical connections of inferior temporal
769        area TEO in macaque monkeys. J Comp Neurol 334:125–150.

770    Doursat R, Bienenstock E (2006) Neocortical Self-Structuration as a Basis for Learning. 5th
771        International Conference on Development and Learning (ICDL 2006).

772    Dum RP, Strick PL (2002) Motor areas in the frontal lobe of the primate. Physiol Behav 77:677–682.

773    Dum RP, Strick PL (2005) Frontal Lobe Inputs to the Digit Representations of the Motor Areas on the
774        Lateral Surface of the Hemisphere. J Neurosci 25:1375–1386.

775    Eacott MJ, Gaffan D (1992) Inferotemporal-frontal Disconnection: The Uncinate Fascicle and Visual
776        Associative Learning in Monkeys. European Journal of Neuroscience 4:1320–1332.

777    Freedman DJ, Riesenhuber M, Poggio T, Miller EK (2001) Categorical Representation of Visual Stimuli
778        in the Primate Prefrontal Cortex. Science 291:312–316.

779    Frege G (1948) Sense and Reference. The Philosophical Review 57:209.

780    Fuster JM (2005) Cortex and Mind. Oxford University Press. Available at:
781        https://academic.oup.com/book/1939 [Accessed April 10, 2023].

782    Fuster JM, Bauer RH, Jervey JP (1985) Functional interactions between inferotemporal and prefrontal
783        cortex in a cognitive task. Brain Research 330:299–307.

784    Fuster JM, Jervey JP (1981) Inferotemporal Neurons Distinguish and Retain Behaviorally Relevant
785        Features of Visual Stimuli. Science 212:952–955.

786    Garagnani M, Lucchese G, Tomasello R, Wennekers T, Pulvermüller F (2017) A Spiking
787        Neurocomputational Model of High-Frequency Oscillatory Brain Responses to Words and

788 Pseudowords. Frontiers in Computational Neuroscience 10 Available at:
789 https://www.frontiersin.org/article/10.3389/fncom.2016.00145 [Accessed March 5, 2022].

790 Garagnani M, Pulvermüller F (2016) Conceptual grounding of language in action and perception: a
791 neurocomputational model of the emergence of category specificity and semantic hubs
792 Barbas H, ed. Eur J Neurosci 43:721–737.

793 Garagnani M, Wennekers T, Pulvermüller F (2007) A neuronal model of the language cortex.
794 Neurocomputing 70:1914–1919.

795 Gelman SA, Markman EM (1986) Categories and induction in young children. Cognition 23:183–209.

796 Gelman SA, Markman EM (1987) Young Children's Inductions from Natural Kinds: The Role of
797 Categories and Appearances. Child Development 58:1532–1541.

798 Gierhan SME (2013) Connections for auditory language in the human brain. Brain and Language
799 127:205–221.

800 Graham S, Keates J, Vukatana E, Khu M (2013) Distinct Labels Attenuate 15-Month-Olds' Attention to
801 Shape in an Inductive Inference Task. Frontiers in Psychology 3 Available at:
802 https://www.frontiersin.org/articles/10.3389/fpsyg.2012.00586 [Accessed July 24, 2022].

803 Güntürkün O, Koenen C, Iovine F, Garland A, Pusch R (2018) The neuroscience of perceptual
804 categorization in pigeons: A mechanistic hypothesis. Learn Behav 46:229–241.

805 Guye M, Parker GJM, Symms M, Boulby P, Wheeler-Kingshott CAM, Salek-Haddadi A, Barker GJ,
806 Duncan JS (2003) Combined functional MRI and tractography to demonstrate the
807 connectivity of the human primary motor cortex in vivo. NeuroImage 19:1349–1360.

808 Harris CR et al. (2020) Array programming with NumPy. Nature 585:357–362.

809 Hebb DO (1949) The Organization of Behavior: A Neuropsychological Theory. Oxford, England: Wiley.

810 Henningsen-Schomers MR, Garagnani M, Pulvermüller F (2022) Influence of language on perception
811 and concept formation in a brain-constrained deep neural network model. Philosophical
812 Transactions of the Royal Society B: Biological Sciences 378:20210373.

813 Henningsen-Schomers MR, Garagnani M, Pulvermüller F (2023) Influence of language on perception
814 and concept formation in a brain-constrained deep neural network model. Phil Trans R Soc B
815 378:20210373.

816 Henningsen-Schomers MR, Pulvermüller F (2022) Modelling concrete and abstract concepts using
817 brain-constrained deep neural networks. Psychological Research 86:2533–2559.

818 Hunter JD (2007) Matplotlib: A 2D Graphics Environment. Comput Sci Eng 9:90–95.

819 Kaas JH (1997) Topographic Maps are Fundamental to Sensory Processing. Brain Research Bulletin
820 44:107–112.

821 Kaas JH, Hackett TA (2000) Subdivisions of auditory cortex and processing streams in primates.
822 Proceedings of the National Academy of Sciences 97:11793–11799.

823 Kassambara A (2021) rstatix: Pipe-Friendly Framework for Basic Statistical Tests. Available at:
824 https://CRAN.R-project.org/package=rstatix [Accessed July 25, 2022].

825    Kéri S (2003) The cognitive neuroscience of category learning. Brain Research Reviews 43:85–109.

826    Knoblauch A, Palm G (2002) Scene segmentation by spike synchronization in reciprocally connected
827          visual areas. I. Local effects of cortical feedback. Biol Cybern 87:151–167.

828    Kriegeskorte N, Mur M, Bandettini P (2008) Representational similarity analysis - connecting the
829          branches of systems neuroscience. Frontiers in Systems Neuroscience 2 Available at:
830          https://www.frontiersin.org/articles/10.3389/neuro.06.004.2008 [Accessed April 8, 2023].

831    LaTourrette AS, Waxman SR (2020) Naming guides how 12-month-old infants encode and remember
832          objects. PNAS 117:21230–21234.

833    Lu M-T, Preston JB, Strick PL (1994) Interconnections between the prefrontal cortex and the
834          premotor areas in the frontal lobe. Journal of Comparative Neurology 341:375–392.

835    Majid A, Bowerman M, Kita S, Haun DBM, Levinson SC (2004) Can language restructure cognition?
836          The case for space. Trends in Cognitive Sciences 8:108–114.

837    Makris N, Pandya DN (2009) The extreme capsule in humans and rethinking of the language circuitry.
838          Brain Struct Funct 213:343–358.

839    Malcolm P. Young, Jack W. Scanneil, Gully A. P. C. Burns, Colin Blakemore (1994) Analysis of
840          Connectivity: Neural Systems in the Cerebral Cortex. Reviews in the Neurosciences 5:227–
841          250.

842    Martin A (2007) The Representation of Object Concepts in the Brain. Annu Rev Psychol 58:25–45.

843    Matthews GG (2001) Neurobiology: molecules, cells, and systems, 2nd ed. Malden, MA: Blackwell
844          Science.

845    Meyer JW, Makris N, Bates JF, Caviness VS, Kennedy DN (1999) MRI-Based topographic parcellation
846          of human cerebral white matter. Neuroimage 9:1–17.

847    Miller EK, Freedman DJ, Wallis JD (2002) The prefrontal cortex: categories, concepts and cognition
848          Parker A, Derrington A, Blakemore C, eds. Phil Trans R Soc Lond B 357:1123–1136.

849    Miller TM, Schmidt TT, Blankenburg F, Pulvermüller F (2018) Verbal labels facilitate tactile
850          perception. Cognition 171:172–179.

851    Nosofsky RM (1988) Exemplar-based accounts of relations between classification, recognition, and
852          typicality. Journal of Experimental Psychology: Learning, Memory, and Cognition 14:700–708.

853    Palm G (2016) Neural Information Processing in Cognition: We Start to Understand the Orchestra,
854          but Where is the Conductor? Frontiers in Computational Neuroscience 10 Available at:
855          https://www.frontiersin.org/articles/10.3389/fncom.2016.00003 [Accessed May 14, 2023].

856    Pandya DN (1995) Anatomy of the auditory cortex. Rev Neurol (Paris) 151:486–494.

857    Pandya DN, Barnes CL (1987) Architecture and connections of the frontal lobe. In: The frontal lobes
858          revisited, pp 41–72. New York, NY, US: The IRBN Press.

859    Pandya DN, Yeterian EH (1985) Architecture and Connections of Cortical Association Areas. In:
860          Association and Auditory Cortices (Peters A, Jones EG, eds), pp 3–61 Cerebral Cortex. Boston,
861          MA: Springer US. Available at: https://doi.org/10.1007/978-1-4757-9619-3_1 [Accessed
862          October 25, 2022].

863 Parker A, Gaffan D (1998) Interaction of frontal and perirhinal cortices in visual object recognition
864          memory in monkeys. European Journal of Neuroscience 10:3044–3057.

865 Parker GJM, Luzzi S, Alexander DC, Wheeler-Kingshott CAM, Ciccarelli O, Lambon Ralph MA (2005)
866          Lateralization of ventral and dorsal auditory-language pathways in the human brain.
867          NeuroImage 24:656–666.

868 Paus T, Castro-Alamancos MA, Petrides M (2001) Cortico-cortical connectivity of the human mid-
869          dorsolateral frontal cortex and its modulation by repetitive transcranial magnetic
870          stimulation. European Journal of Neuroscience 14:1405–1411.

871 Petrides M, Pandya DN (2009) Distinct Parietal and Temporal Pathways to the Homologues of Broca's
872          Area in the Monkey Ungerleider L, ed. PLoS Biol 7:e1000170.

873 Pickron CB, Iyer A, Fava E, Scott LS (2018) Learning to Individuate: The Specificity of Labels
874          Differentially Impacts Infant Visual Attention. Child Dev 89:698–710.

875 Plunkett K, Hu J-F, Cohen LB (2008) Labels can override perceptual categories in early infancy.
876          Cognition 106:665–681.

877 Pulvermüller F, Tomasello R, Henningsen-Schomers MR, Wennekers T (2021) Biological constraints
878          on neural network models of cognitive function. Nat Rev Neurosci 22:488–502.

879 Pusch R, Clark W, Rose J, Güntürkün O (2023) Visual categories and concepts in the avian brain. Anim
880          Cogn 26:153–173.

881 R Core Team (2021) R: A Language and Environment for Statistical Computing. Available at:
882          https://www.R-project.org/ [Accessed July 25, 2022].

883 Ralph MAL, Jefferies E, Patterson K, Rogers TT (2017) The neural and computational bases of
884          semantic cognition. Nat Rev Neurosci 18:42–55.

885 Rauschecker JP, Scott SK (2009) Maps and streams in the auditory cortex: nonhuman primates
886          illuminate human speech processing. Nat Neurosci 12:718–724.

887 Rauschecker JP, Tian B (2000) Mechanisms and streams for processing of "what" and "where" in
888          auditory cortex. Proceedings of the National Academy of Sciences 97:11800–11806.

889 Reback J et al. (2022) pandas-dev/pandas: Pandas 1.4.3. Available at:
890          https://zenodo.org/record/3509134 [Accessed July 25, 2022].

891 Rilling JK (2014) Comparative primate neuroimaging: insights into human brain evolution. Trends in
892          Cognitive Sciences 18:46–55.

893 Rilling JK, Glasser MF, Jbabdi S, Andersson J, Preuss TM (2012) Continuity, Divergence, and the
894          Evolution of Brain Language Pathways. Front Evol Neurosci 3:11.

895 Rilling JK, Glasser MF, Preuss TM, Ma X, Zhao T, Hu X, Behrens TEJ (2008) The evolution of the
896          arcuate fasciculus revealed with comparative DTI. Nat Neurosci 11:426–428.

897 Rilling JK, van den Heuvel MP (2018) Comparative Primate Connectomics. BBE 91:170–179.

898 Rizzolatti G, Luppino G (2001) The Cortical Motor System. Neuron 31:889–901.

899  Rogers TT, McClelland JL (2014) Parallel Distributed Processing at 25: Further Explorations in the
900      Microstructure of Cognition. Cogn Sci 38:1024–1077.

901  Rolls ET, Deco G (2010) The Noisy BrainStochastic Dynamics as a Principle of Brain Function. Oxford
902      University Press. Available at: https://academic.oup.com/book/7413 [Accessed October 9,
903      2022].

904  Rolls ET, Deco G (2015) Networks for memory, perception, and decision-making, and beyond to how
905      the syntax for language might be implemented in the brain. Brain Research 1621:316–334.

906  Romanski L m., Bates J f., Goldman-Rakic P s. (1999a) Auditory belt and parabelt projections to the
907      prefrontal cortex in the Rhesus monkey. Journal of Comparative Neurology 403:141–157.

908  Romanski LM (2007) Representation and Integration of Auditory and Visual Stimuli in the Primate
909      Ventral Lateral Prefrontal Cortex. Cerebral Cortex 17:i61–i69.

910  Romanski LM, Tian B, Fritz J, Mishkin M, Goldman-Rakic PS, Rauschecker JP (1999b) Dual streams of
911      auditory afferents target multiple domains in the primate prefrontal cortex. Nat Neurosci
912      2:1131–1136.

913  Saur D, Kreher BW, Schnell S, Kümmerer D, Kellmeyer P, Vry M-S, Umarova R, Musso M, Glauche V,
914      Abel S, Huber W, Rijntjes M, Hennig J, Weiller C (2008) Ventral and dorsal pathways for
915      language. Proceedings of the National Academy of Sciences 105:18035–18040.

916  Schomers MR (2017) Establishing action-perception circuits as a neural basis for meaning-carrying
917      linguistic symbols – the role of frontal speech motor areas and fronto-temporal connectivity.
918      Available at: http://doi.org/.

919  Schomers MR, Garagnani M, Pulvermüller F (2017) Neurocomputational Consequences of
920      Evolutionary Connectivity Changes in Perisylvian Language Cortex. J Neurosci 37:3045–3055.

921  Scott LS, Monesson A (2009) The Origin of Biases in Face Perception. Psychol Sci 20:676–680.

922  Seger CA, Miller EK (2010) Category Learning in the Brain. Annu Rev Neurosci 33:203–219.

923  Seltzer B, Pandya DN (1989) Intrinsic connections and architectonics of the superior temporal sulcus
924      in the rhesus monkey. Journal of Comparative Neurology 290:451–471.

925  Sloutsky VM, Fisher AV (2004) Induction and Categorization in Young Children: A Similarity-Based
926      Model. Journal of Experimental Psychology: General 133:166–188.

927  Thiebaut de Schotten M, Dell'Acqua F, Valabregue R, Catani M (2012) Monkey to human comparative
928      anatomy of the frontal lobe association tracts. Cortex 48:82–96.

929  Tomasello R, Garagnani M, Wennekers T, Pulvermüller F (2017) Brain connections of words,
930      perceptions and actions: A neurobiological model of spatio-temporal semantic activation in
931      the human cortex. Neuropsychologia 98:111–129.

932  Tomasello R, Garagnani M, Wennekers T, Pulvermüller F (2018) A Neurobiologically Constrained
933      Cortex Model of Semantic Grounding with Spiking Neurons and Brain-Like Connectivity.
934      Frontiers in Computational Neuroscience 12 Available at:
935      https://www.frontiersin.org/article/10.3389/fncom.2018.00088 [Accessed February 12,
936      2022].

937     Ungerleider LG, Gaffan D, Pelak VS (1989) Projections from inferior temporal cortex to prefrontal
938         cortex via the uncinate fascicle in rhesus monkeys. Exp Brain Res 76:473–484.

939     Vanek N, Sóskuthy M, Majid A (2021) Consistent verbal labels promote odor category learning.
940         Cognition 206:104485.

941     Virtanen P et al. (2020) SciPy 1.0: fundamental algorithms for scientific computing in Python. Nat
942         Methods 17:261–272.

943     Wakana S, Jiang H, Nagae-Poetscher LM, van Zijl PCM, Mori S (2004) Fiber Tract–based Atlas of
944         Human White Matter Anatomy. Radiology 230:77–87.

945     Waskom M (2021) seaborn: statistical data visualization. JOSS 6:3021.

946     Waxman SR, Booth AE (2001) Seeing Pink Elephants: Fourteen-Month-Olds' Interpretations of Novel
947         Nouns and Adjectives. Cognitive Psychology 43:217–242.

948     Waxman SR, Markow DB (1995) Words as Invitations to Form Categories: Evidence from 12- to 13-
949         Month-Old Infants. Cognitive Psychology 29:257–302.

950     Webster MJ, Bachevalier J, Ungerleider LG (1994) Connections of Inferior Temporal Areas TEO and TE
951         with Parietal and Frontal Cortex in Macaque Monkeys. Cerebral Cortex 4:470–483.

952     Westermann G, Mareschal D (2014) From perceptual to language-mediated categorization.
953         Philosophical Transactions of the Royal Society B: Biological Sciences 369:20120391.

954     Whorf BL, Carroll JB (2007) Language, thought, and reality: selected writings, 28. print. Cambridge,
955         Mass: The MIT Press.

956     Wittgenstein L (1922) Tractatus logico-philosophicus. London: Routledge & Kegan Paul.

957     Yeterian EH, Pandya DN, Tomaiuolo F, Petrides M (2012) The cortical connectivity of the prefrontal
958         cortex in the monkey brain. Cortex 48:58–81.

959     Young M, Scannell JW, Burns G (1995a) The analysis of cortical connectivity. New York : Austin:
960         Springer ; R.G. Landes.

961     Young MP, Scannell JW, Burns G (1995b) The Analysis of Cortical Connectivity, 1st ed. Springer Berlin,
962         Heidelberg. Available at: https://link.springer.com/book/9783540604013 [Accessed October
963         25, 2022].

964     Yuille AL, Geiger D (1995) Winner-Take-All Mechanisms. In: Handbook of Brain Theory and Neural
965         Networks (Arbib MA, ed), pp 1–1056. MIT Press.

966

967                                        **Tables**

968  **Table 1.** Connectivity structure of the modelled cortical areas with neuroanatomical evidence. Table
969  taken from Tomasello et al. (2018).

| Modelled areas | References |
| --- | --- |
| **Between-area connectivity (black arrows)** | |
| Perisylvian system | |
| A1, AB, PB | Pandya and Yeterian, 1985; Pandya, 1995; Rauschecker and Tian, 2000 |
| $PF_i$, $PM_i$, $M1_i$ | Pandya and Yeterian, 1985; Young et al., 1995 |
| Extra-sylvian system | |
| V1, TO, AT | Bressler et al., 1993; Distler et al., 1993 |
| $PF_L$, $PM_L$, $M1_L$ | Pandya and Yeterian, 1985; Arikuni et al., 1988; Lu et al., 1994; Rizzolatti and Luppino, 2001; Dum and Strick, 2002, 2005 |
| Between system | |
| AT, PB | Gierhan, 2013 |
| $PF_i$, $PF_L$ | Yeterian et al., 2012 |
| **Long distance cortico-cortical connections (purple arrows)** | |
| Perisylvian system | |
| $PF_i$, PB | Meyer et al., 1999 p.19; Romanski et al., 1999a, 1999b; Paus et al., 2001; Catani et al., 2005 p.200; Parker et al., 2005; Rilling et al., 2008; Makris and Pandya, 2009 |
| PB, $PM_i$ | Rilling et al., 2008; Saur et al., 2008 |
| AB, $PF_i$ | Romanski et al., 1999a, 1999b; Kaas and Hackett, 2000; Petrides and Pandya, 2009; Rauschecker and Scott, 2009 |
| Extra-sylvian system | |
| AT, $PF_L$ | Bauer and Jones, 1976 p.197; Fuster et al., 1985 p.198; Ungerleider et al., 1989; Eacott and Gaffan, 1992; Webster et al., 1994; Parker and Gaffan, 1998; Chafee and Goldman-Rakic, 2000 |
| AT, $PM_L$ | Bauer and Fuster, 1978; Fuster et al., 1985; Pandya and Barnes, 1987; Seltzer and Pandya, 1989; Chafee and Goldman-Rakic, 2000 |
| TO, $PF_L$ | Bauer and Jones, 1976; Fuster and Jervey, 1981; Fuster et al., 1985; Seltzer and Pandya, 1989; Makris and Pandya, 2009 |
| Between systems | |
| PB, $PF_L$ | Pandya and Barnes, 1987; Romanski et al., 1999a, 1999b |
| AT, $PF_i$ | Pandya and Barnes, 1987; Ungerleider et al., 1989; Webster et al., 1994; Romanski, 2007; Petrides and Pandya, 2009; Rilling, 2014 |
| **Second-next neighbor "jumping" links (blue arrows)** | |
| Perisylvian system (Rilling et al., 2008, 2012; Thiebaut de Schotten et al., 2012; Rilling and van den Heuvel, 2018) | |
| A1, PB | Pandya and Yeterian, 1985; Malcolm P. Young et al., 1994 |
| $PF_i$, $M1_i$ | Deacon, 1992; Young et al., 1995b; Guye et al., 2003 |
| Extra-sylvian system (Thiebaut de Schotten et al., 2012) | |
| V1, AT | Catani et al., 2003; Wakana et al., 2004 |
| $PF_L$, $M1_L$ | Deacon, 1992; Young et al., 1995a; Guye et al., 2003 |

970

971  **Table 2.** Parameter values used in the simulations. For details and more elaborate discussion on the
972  corresponding equations as well as their mathematical implementations, please see Henningsen-
973  Schomers et al. (2022).

| Eq. (1) | Time constant (excitatory cells) | $\tau = 2.5$ (time steps) |
|---------|----------------------------------|---------------------------|
|  | Time constant (inhibitory cells) | $\tau = 5$ (time steps) |
|  | Total input rescaling factor | $k_1 = 0.01$ |
|  | Noise amplitude | $k_2 = 7\sqrt{(24/\Delta t)}$ ($\Delta t = 0.5$ ms) |
|  | Global inhibition strength | $k_G = 0.80$ (time steps) |
| Eq. (3) | Spiking threshold | thresh $= 0.18$ |
|  | Adaptation strength | $\alpha = 8.0$ |
| Eq. (4) | Adaption time constant | $\tau_{ADAPT} = 10$ (time steps) |
| Eq. (5) | Rate-estimate time constant | $\tau_{Favg} = 30$ (time steps, training) |
|  |  | $\tau_{Favg} = 5$ (time steps, testing) |
| Eq. (6) | Global inhibition time constant | $\tau_{FGLOB} = 12$ (time steps) |
| Eq. (7) | Postsynaptic potential thresholds | $\vartheta_+ = 0.15$ (LTP) |
|  |  | $\vartheta_- = 0.14$ (LTD) |
|  | Presynaptic output activity required for any synaptic change | $\vartheta_{pre} = 0.05$ (LTP) |
|  | Learning rate | $\Delta w = 0.0008$ |

974

975

976    **Table 3.** Critical and significant observations and the corresponding aspects of learning.

| Analysis | Learning aspect | Observation |
|---|---|---|
| RSA | Category learning | Successful category learning in all learning conditions<br>$$Dissim_{B-TT} > Dissim_{W-TT}$$<br>Interaction effect of Symbol type and within/between categories<br>$$Dissim_{B-TT_{CT}} > Dissim_{B-TT_{PN}}; Dissim_{B-TT_{CT}} > Dissim_{B-TT_{PN}}$$<br>$$Dissim_{W-TT_{CT}} < Dissim_{W-TT_{PN}}; Dissim_{W-TT_{CT}} < Dissim_{W-TT_{NoL}}$$ |
| | Generalization | Symbol effect on dissimilarity differences within category<br>$$DissimDiff_{CT} < DissimDiff_{NoS}$$<br>$$DissimDiff_{CT} < DissimDiff_{PN}$$ |
| CA Analysis | Representations of category-critical features | Tendency to encode shared features in all learning conditions<br>$$n_S > n_U$$<br>Symbol effect on the number of shared neurons<br>$$n_{S_{CT}} > n_{S_{PN}}; n_{S_{CT}} > n_{S_{NoL}}$$<br>Gain in shared neurons in central area in all learning conditions<br>$$n_{S-central} > n_{S-primary}$$<br>Symbol effect on across-area gain of shared neurons<br>$$Gain_{S_{CT}} > Gain_{S_{PN}}; Gain_{S_{CT}} > Gain_{S_{NoL}}$$ |
| | Representations of instance-specific features | Symbol effect on the number of unique neurons<br>$$n_{U_{PN}} > n_{U_{CT}}; n_{U_{NoL}} > n_{U_{CT}}$$<br>Loss in unique neurons in central areas in all learning conditions<br>$$n_{S-central} > n_{S-primary}$$<br>Symbol effect on across-area loss of unique neurons<br>$$Loss_{S_{PN}} < Loss_{S_{CT}}$$ |

977    Abbreviation: $Dissim_{W-TT}/Dissim_{W-TN}$ = Dissimilarity between a trained instance and another
978    trained instance/novel instance of the same category; $Dissim_{B-TT}$ = Dissimilarity between two
979    trained instances from different categories; $DissimDiff = |Dissim_{W-TN} - Dissim_{W-TT}|$; $n_S$ =
980    number of shared neuron; $n_U$ = number of unique neuron; CT = Category term; PN = Proper name;
981    NoS = No symbol.

982

987 **Figure legends**

988 **Figure 1. A) Area structure and between-area connectivity of the neural network model**. **Left:**
989 The network model's 12 cortical areas in the left fronto-temporo-occipital lobes: inferior-frontal
990 articulatory (red) and superior temporal auditory systems (blue) of the perisylvian areas, and the lateral
991 frontal hand-motor system (yellow/orange/brown) and visual "what" stream (green) in the extrasylvian
992 cortex. **Right:** Connections among the 12 modelled brain areas: direct connections between adjacent
993 areas (black arrows), second nearest-neighbor areas (blue arrows), and long-distant links (purple
994 arrows). Figure modified from Tomasello et al. (2018). **B) Schematic illustrations of activity**
995 **patterns for instances of two categories.** The categories are illustrated with images of robots and cat
996 faces, but note that this is for illustrative purposes. The actual input to the model was not images, but
997 grounding patterns consisting of sets of activated neurons (see main text for details). Active neurons of
998 given activity patterns were either shared among instances of the same category (black) or unique to
999 each instance (color). Each model area included $25 \times 25$ excitatory neurons, i.e., 625 cells. **Left:** In
1000 grounding patterns (i) presented to *V1/*M1$_L$, 6 shared active neurons (black) code for the common
1001 perceptual-semantic features of the category "a" and 6 unique neurons (color) represent instance-
1002 specific perceptuomotor features from each of the category members. Member instances of one
1003 category activated the same six shared neurons while the instance from another category activated a
1004 different set of six shared neurons; each instance also activated six unique neurons. **Middle:** 12
1005 neurons (black) make up word form pattern for the category term; in the Category term condition,
1006 member instances co-activated with the same word form pattern (ii) in *A1/*M1$_i$. **Right:** 12 unique
1007 neurons (color) represent each proper name of an individual instance, which are activated 1-to-1 with
1008 these instances in the Proper name condition. Instances were co-activated with distinct different word
1009 form patterns (iii) in *A1/*M1$_i$ regardless of category. **C) Simulating no-symbol learning** (top),
1010 **category-term learning** (bottom-left), and **proper-name learning** (bottom-right) where no word
1011 form pattern, word form patterns (ii), and word form pattern (iii) were presented to *A1/*M1$_i$,
1012 respectively.

1013 **Figure 2.** Experiment design used for instance learning and conceptual grounding. **A) Training phase**
1014 with 30 object instances from ten categories. The categories are illustrated with images of robots and
1015 cat faces, but note that this is for illustrative purposes. The actual input to the model was not images,
1016 but grounding patterns consisting of sets of activated neurons (see main text for details). For each
1017 trained instance, the grounding pattern (i) was either presented to the network on its own (No symbol)
1018 or combined with a 'word form pattern' of type (ii) (Category term) or type (iii) (Proper name). **B)**
1019 **Testing phase** with a collection of the initially trained 30 instances and 30 novel instances from the 10
1020 original categories, resulting in 60 testing instances (i.e., 6 per category). **C) Training conditions in**
1021 **the main simulations** (top) **and control simulations** (bottom) differ in the number of training trials
1022 (tt) to match the number of instance representations and the number of word form representations,
1023 respectively.

1024 **Figure 3. A) Schematic extraction of a $60 \times 60$ Representational Dissimilarity Matrix** (RDM)
1025 which represents 12 instances from two different categories and the similarities between any instance
1026 pair. For illustration, we once again use the categories of robots and cat faces. The schematic
1027 dissimilarity matrix illustrates how between-category (cells outside the red boundaries) within-
1028 category dissimilarities (cells within the red boundaries) were calculated. Of interest are the (1) within-
1029 category dissimilarity among trained instances ($Dissim_{W-TT}$, lightest blue shade), (2) within-category
1030 dissimilarity between a trained and a novel instance ($Dissim_{W-TN}$, intermediate blue shade), and (3)
1031 between-category dissimilarity of two trained instances ($Dissim_{B-TT}$, darkest blue shade). The RDM
1032 is symmetric about its diagonal (grey) of zeros (representing the non-dissimilarity of each of the
1033 instances to itself). Only the upper half of the RDM is used for analysis and the lower half could be
1034 abandoned (black). **B) RDMs for each of the twelve model areas in three main simulations**: No
1035 symbol (top row), Category term (middle row), and Proper name (bottom row). Squares indicate the
1036 degrees to which network activity in the 12 network areas elicited by (12 out of 60) grounding patterns
1037 in the three learning conditions differed between each other within and between categories and are
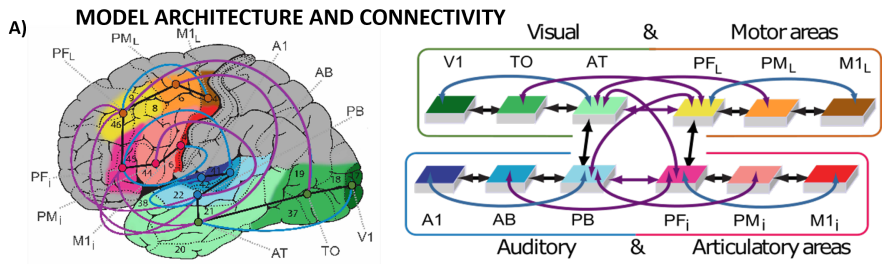
1038    color-coded from turquoise (no dissimilarity, $Dissim = 0$), blue, pink, and to dark red (high
1039    dissimilarity, $Dissim > 3$).

1040    **Figure 4.** Bar charts depicting dissimilarities between network activity elicited by trained grounding
1041    patterns after learning for each of the three training conditions. **A)** Main simulation: Within-category
1042    (W-TT) and between-category (B-TT) dissimilarity values across all 30 trained activity patterns were
1043    averaged for each of the twelve model areas. **B-C)** Within-category (W-TT) and between-category (B-
1044    TT) dissimilarities across the 30 trained items were averaged for extrasylvian model areas. The three
1045    training conditions of the main simulations (B) were No symbol (NoS, grey), Category term (CT,
1046    blue) and Proper name (PN, pink). The four training conditions of the control simulation (C) were No
1047    symbol with each instance presented over 1,000 (NoS_1x, blue-striped grey) or 3,000 trials (NoS_3x,
1048    pink-striped grey), Category term where each instance presented over 1,000 trials (CT_1x, blue) and
1049    Proper name where each instance presented over 3,000 trials (PN_3x, pink). Error bars represent 95%
1050    confidence intervals of the mean. Circles above the bars represent post hoc pairwise comparisons
1051    between a reference (circles with filled colored) and a corresponding mean (unfilled circles) after
1052    Bonferroni correction (critical $p$ value = 0.005). 10 comparisons relevant to the main effects of
1053    training condition and dissimilarity type and for their interaction are illustrated. Asterisks represent
1054    two-tailed p values: ** $p < .005$, *** $p < .001$, ns: not significant. The results were replicated in the
1055    whole model architecture (6 extrasylvian and 6 perisylvian model areas); see Figure 4-1 and Table 4-1.

1056    **Figure 5.** Bar charts depicting dissimilarities between network activity elicited by trained novel
1057    grounding patterns after learning for each of the three training conditions. **A)** Main simulation:
1058    Within-category dissimilarity values between any two trained instances (W-TT) and between trained
1059    and novel instances were averaged for each of the twelve model areas. **B&C)** Within-category
1060    dissimilarities between any two trained instances (W-TT) and between trained and novel instances (W-
1061    TN) were averaged for extrasylvian model areas. The three training conditions of the main simulations
1062    (B) were No symbol (NoS, grey), Category term (CT, blue) and Proper name (PN, pink). The four
1063    training conditions of the control simulation (C) were NoS_1x (blue-striped grey) or NoS_3x (pink-
1064    striped grey), CT_1x (blue) and PN_3x (pink). For further explanation, see Figure 4. The results were
1065    replicated in the whole model architecture (6 extrasylvian and 6 perisylvian model areas); see Figure
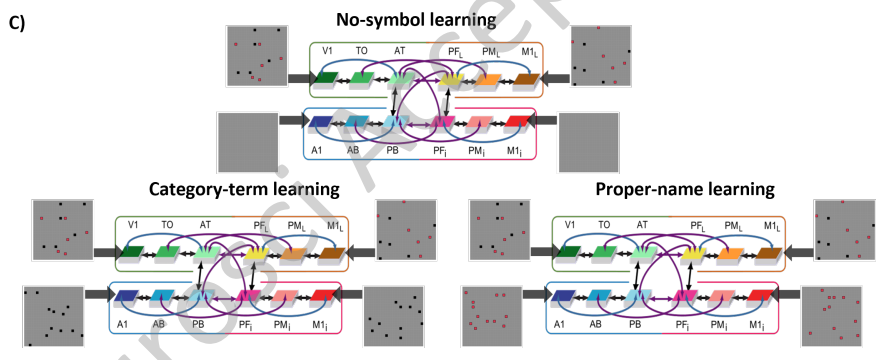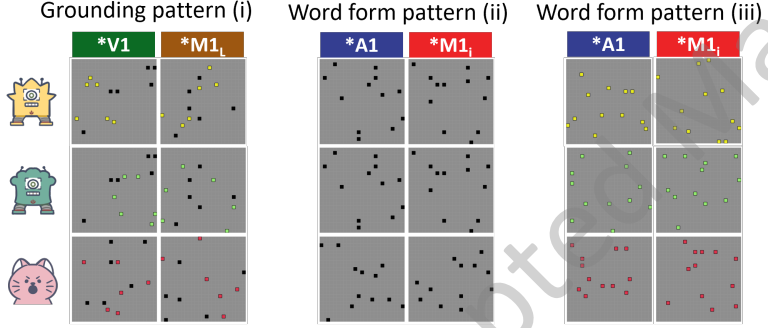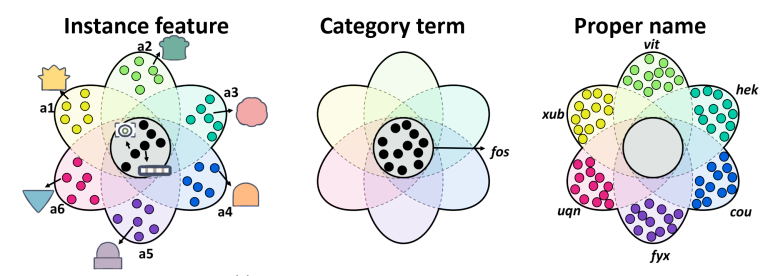1066    5-1 and Table 5-1.

1067    **Figure 6.** Bar charts depicting average numbers of instance specific ('unique') and category-general
1068    ("shared") neurons activated by grounding patterns of instances learnt in the three training conditions,
1069    No symbol (grey), Category term (blue) and Proper name (pink). **A)** Main simulation: The number of
1070    activated unique (U) and shared (S) neurons in response to each of the 30 trained instances was
1071    averaged across all twelve model areas. **B &C)** The number of activated neurons in response to the 30
1072    trained grounding patterns was averaged for each of the six extrasylvian areas. **D&E)** Changes in
1073    neuronal activation seen between extrasylvian primary areas, where stimulation was given, and the
1074    'higher' more central connector hub areas central to the architecture. Changes in the number of
1075    activated neurons in response to trained grounding patterns are shown for the three training conditions.
1076    Unique neurons are shown by solid lines with crossed ends, shared ones by broken lines with
1077    triangular ends. The three training conditions of the main simulations (B, D) were No symbol (NoS,
1078    grey), Category term (CT, blue) and Proper name (PN, pink). The four training conditions of the
1079    control simulation (C) were NoS_1x (blue-striped grey) or NoS_3x (pink-striped grey), CT_1x (blue)
1080    and PN_3x (pink). For further explanations see Figure 4. The results were replicated in the whole
1081    model architecture (6 extrasylvian and 6 perisylvian model areas); see Figure 6-1 and Table 6-1.

1082

**A)** MODEL ARCHITECTURE AND CONNECTIVITY

Visual & Motor areas

V1  TO  AT  PF_L  PM_L  M1_L

A1  AB  PB  PF_i  PM_i  M1_i

Auditory & Articulatory areas

**B)** ACTIVITY PATTERNS

Member instances of a category called *fos*:

An instance of a category called *cax*:

Instance feature

Category term

Proper name

Grounding pattern (i)

*V1  *M1_L

Word form pattern (ii)

*A1  *M1_i

Word form pattern (iii)

*A1  *M1_i

**C)**

No-symbol learning

V1  TO  AT  PF_L  PM_i  M1_L

A1  AB  PB  PF_i  PM_i  M1_i

Category-term learning

V1  TO  AT  PF_L  PM_L  M1_L

A1  AB  PB  PF_i  PM_i  M1_i

Proper-name learning

V1  TO  AT  PF_L  PM_L  M1_L

A1  AB  PB  PF_i  PM_i  M1_i

# EXPERIMENT DESIGN

**A)**

| Category (10) | | a | | | b | | | ... |
|---|---|---|---|---|---|---|---|---|
| **Training instance (30)** | | | | | | | | ... |
| **Activity pattern** | Grounding pattern (30) | (i) | (i) | (i) | (i) | (i) | (i) | (i) |
| | No symbol | | | | | | | |
| | Category term (10) | | (ii) *fos* | | | (ii) *cax* | | (ii) *...* |
| | Proper name (30) | (iii) *xub* | (iii) *vit* | (iii) *hek* | (iii) *dre* | (iii) *tla* | (iii) *tsu* | (iii) |

**B)**

| Category (10) | | a | | | | | | b | | | | | | ... |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Testing instance (60)** | | | | | | | | | | | | | | ... |
| **Activity pattern** | Grounding pattern (30) | (i) | (i) | (i) | (i) | (i) | (i) | (i) | (i) | (i) | (i) | (i) | (i) | (i) |
| | No symbol | | | | | | | | | | | | | |

training instances    novel instances    training instances    novel instances

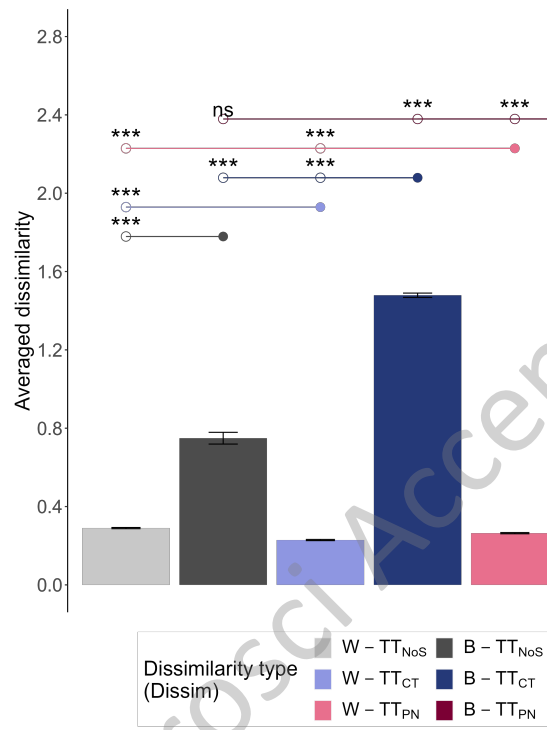**C) The number of training trials in the main and control simulations**

| **Main simulation** (matched for instance presentations) | | | |
|---|---|---|---|
| Training trials (tt) per instance | No symbol | Category term | Proper name |
| 2000 | **NoS** | **CT** (6000 tt/symbol) | **PN** (2000 tt/symbol) |

| **Control simulation** (matched for word form presentations) | | | |
|---|---|---|---|
| Training trials (tt) per instance | No Symbol | Category term | Proper name |
| 1000 | **NoS_1x** | **CT_1x** (3000 tt/symbol) | |
| 3000 | **NoS_3x** | | **PN_3x** (3000 tt/symbol) |

**A)**

Category a    Category b

Trained  (1) (2) (3)

Novel

Trained

Novel

**B)**

**No Symbol**

V1    TO    AT    PF$_L$    PM$_L$    M1$_L$

A1    AB    PB    PF$_i$    PM$_i$    M1$_i$

**Category Term**

V1    TO    AT    PF$_L$    PM$_L$    M1$_L$

A1    AB    PB    PF$_i$    PM$_i$    M1$_i$

**Proper Name**

V1    TO    AT    PF$_L$    PM$_L$    M1$_L$

A1    AB    PB    PF$_i$    PM$_i$    M1$_i$

**A**

| V1 | TO | AT | PF_L | PM_L | M1_L |

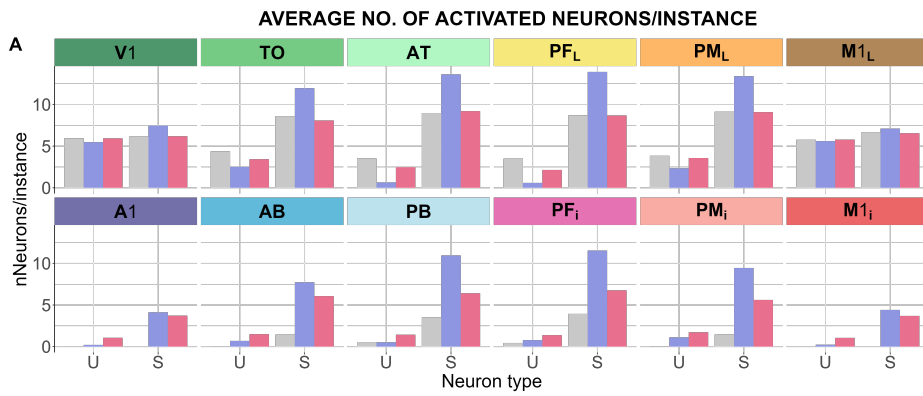| A1 | AB | PB | PF_i | PM_i | M1_i |

**B. Main simulation - Extra-sylvian areas**

**C. Control simulation - Extra-sylvian areas**

Dissimilarity type (Dissim)

W – TT_NoS    B – TT_NoS
W – TT_CT     B – TT_CT
W – TT_PN     B – TT_PN

Dissimilarity type (Dissim)

W – TT_NoS_1x    B – TT_NoS_1x
W – TT_NoS_3x    B – TT_NoS_3x
W – TT_CT_1x     B – TT_CT_1x
W – TT_PN_3x     B – TT_PN_3x

A

B. Main simulation - Extra-sylvian areas

C. Control simulation - Extra-sylvian areas

**AVERAGE NO. OF ACTIVATED NEURONS/INSTANCE**

A

**B. Main simulation**  **C. Control simulation**

**GAINS IN NO. OF ACTIVATED NEURONS**

**D. Main simulation**  **E. Control simulation**