

Feature-rich multiplex lexical networks reveal mental strategies of early language learning

Salvatore Citraro ^{1,2}, Michael S. Vitevitch ³, Massimo Stella ⁴⁺, Giulio Rossetti ²⁺

1 Department of Computer Science, University of Pisa Largo Bruno Pontecorvo, 3, Pisa

2 KDD-Lab, ISTI (CNR) G. Moruzzi, 1, Pisa

3 Department of Psychology, University of Kansas, USA

4 CogNosco Lab, Department of Computer Science, University of Exeter, UK

Corresponding author: m.stella AT exeter.ac.uk

+ These authors contributed equally.

Abstract

Knowledge in the human mind exhibits a dualistic vector/network nature. Modelling words as vectors is key to natural language processing, whereas networks of word associations can map the nature of semantic memory. We reconcile these paradigms - fragmented across linguistics, psychology and computer science - by introducing FEature-Rich MUltiplex LEXical (FERMULEX) networks. This novel framework merges structural similarities in networks and vector features of words, which can be combined or explored independently. Similarities model heterogenous word associations across semantic/syntactic/phonological aspects of knowledge. Words are enriched with multi-dimensional feature embeddings including frequency, age of acquisition, length and polysemy. These aspects enable unprecedented explorations of cognitive knowledge. Through CHILDES data, we use FERMULEX networks to model normative language acquisition by 1000 toddlers between 18 and 30 months. Similarities and embeddings capture word homophily via conformity, which measures assortative mixing via distance and features. Conformity unearths a language kernel of frequent/polysemous/short nouns and verbs key for basic sentence production, supporting recent evidence of children's syntactic constructs emerging at 30 months. This kernel is invisible to network core-detection and feature-only clustering: It emerges from the dual vector/network nature of words. Our quantitative analysis reveals two key strategies in early word learning. Modelling word acquisition as random walks on FERMULEX topology, we highlight non-uniform filling of communicative developmental inventories (CDIs). Conformity-based walkers lead to accurate (75%), precise (55%) and partially well-recalled (34%) predictions of early word learning in CDIs, providing quantitative support to previous empirical findings and developmental theories.

Introduction

The mental lexicon is the part of memory that stores information about a word's meanings, syntactic features, pronunciation and more [71, 77, 1]. Although often described as being like a mental dictionary [21, 27, 71], the mental lexicon is not static, and is instead a complex system, whose structure influences language processing and has been investigated across fields like psychology [71], linguistics [1, 13], computer science and artificial intelligence

[62, 7, 6]. Decades of multidisciplinary research have gathered evidence that words in the mental lexicon have a dual representation [27], analogous to the particle/wave duality of light in physics [4]. Psycholinguistics and distributional semantics posit that words in the lexicon possess both a networked organisation [17, 53, 72] and a vector-space nature [9, 25, 35, 36]. On the one hand, networks capture conceptual relationships (as links) between words (as nodes). On the other hand, vector-spaces identify alignment and distances between vectors, whose components represent word features. The network aspects of the mental lexicon started with seminal work by Quillian [53] and by Collins and Loftus [17]. These works showed how in a network of words linked through semantic associations, e.g. possessing a common attribute or overlapping in meaning, the length of the shortest path separating concepts was predictive of retrieval times from semantic memory and sentence understanding [17, 53]. The advent of network science has revived interest in this approach [13], with several recent works examining how the structure of semantic networks [66, 10, 59, 34, 33], phonological networks [72, 73], and their multiplex combination [63, 65, 37] influence language acquisition and processing.

In parallel, distributional semantics postulates that semantic memory possesses a vector space structure [25, 9, 50], where concepts are vectors whose components express either interpretable features [19] (e.g. possessing a semantic feature, being in a category or being acquired at a certain age) or latent aspects of language [32, 50, 39, 35] (e.g. overlap in meaning due to word co-occurrence in the same context). Although latent aspects of language limit the understanding of cognitive processing, models like Latent Semantic Analysis [35] and the Hyperspace Analogue to Language [39] were used extensively in cognitive inquiries of information processing, mainly due to their ability to extract semantic features without human intervention. More recently, transformer neural networks like BERT enabled vector representations for words depending on their context [9]. This enhancement revolutionised the field of natural language processing and predicted successfully semantic tasks like entity recognition or word meaning disambiguation [9, 32]. Although powerful predictors, these approaches provide relatively little access to the organisation of words in the human mind and can thus benefit from network models and interpretable distributional semantics [32]. Reconciling the non-latent, interpretable vector/network duality of words in the mental lexicon is the focus of this work.

We introduce FEature-Rich MUltiplex LEXical - *FERMULEX* - networks, which combines the vector-based and multiplex network aspects of words and their associations in the mental lexicon. Rather than merely building networks out of similarities between vectors of features [18], we view structure and feature similarities as two independent building blocks, whose contribution to represent words in the mind can be explored in parallel. Hence in *FERMULEX* networks, network structure remains and can be explored even when word similarities are switched off, and vice versa. This possibility does not exist in networks built from vector similarities (cf. [70]). We achieve this advancement by using the recent measure of conformity [56], an enhancement of assortative

mixing estimation through non-adjacent nodes.

We show that the dual network/vector representation of words is crucial for understanding key aspects of the mental lexicon that would go undetected by considering features - or networks - only. Using normative word learning norms [41] and phonological/semantic/syntactic [63] data in 1000 English toddlers, *FERMULEX* networks reveal a language kernel progressively built in the mental lexicon of toddlers and *undetectable* by either network core detection [30] or clustering in vector spaces [74]. This mental kernel contains general yet simple nouns and verbs that can build diverse sentences, with crucial relevance to children’s communication [26]. The identification of this kernel via *FERMULEX* provides quantitative evidence and modelling insights as to how can young children produce early sentences, as recently observed [26].

Modelling word acquisition as increasingly biased random walkers over the network/vectorial *FERMULEX* representation leads to more insights. We adopted this approach inspired by past work using random walkers over cognitive networks for investigating the mental lexicon [24]. We find that predicting word learning in the language kernel crucially depends on: (i) network/vectorial conformity [56] and (ii) the filling of communicative developmental inventories (CDIs) [22], i.e. lists of words sharing a semantic category and commonly used for measuring early cognitive development. We find that CDIs display a rich filling dynamic in word learning, which can be predicted by *FERMULEX* with accuracy, precision and recall up to 75%, 55% and 34%, respectively. These values are statistically significant with respect to a baseline random learner. Without conformity and CDI filling levels, in fact, predictions of word learning in the language kernel are equivalent to random guessing. Since the language kernel stores words crucial for producing early sentences, our results indicate that the documented ability for young toddlers to communicate via early sentences around month 30 [26] crucially depends on network, vector, and categorical aspects of the mental lexicon. Our approach with *FERMULEX* can encompass them all and thus represents a powerful tool for future cognitive research of various aspects of language.

1 Results

***FERMULEX* characterisation.** A combination of a multiplex network structure (Fig. 1, A) and a vector space of interpretable features (Fig. 1, B) results in a *FERMULEX* network (Fig. 1, C). Conformity [56] assesses structure-feature relationships on the aggregated topology. For each node and with respect to each feature, conformity quantifies the node assortative mixing, by extending this estimation to the non-adjacent but still reachable neighbors of a node. Studying conformity distributions, we can capture heterogeneous patterns between nodes.

Fig. 1, D sums up these patterns on the real data representing toddlers’ mental lexicon (see Methods for

FERMULEX Model

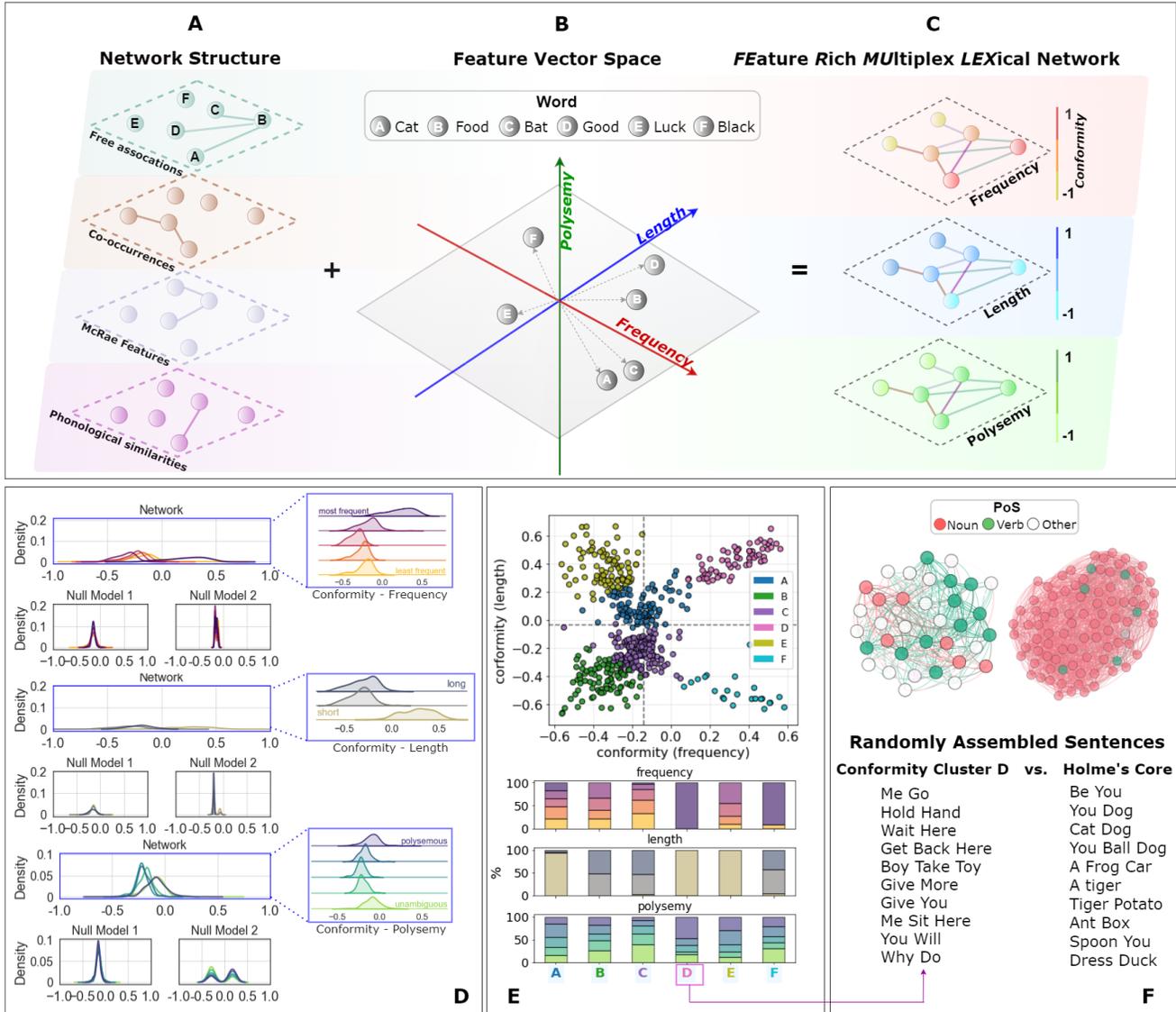


Figure 1: (A-C): Combining multiplex topology (A) and vector spaces (B) results in *FERMULEX* network (C); (D): kernel density estimates (KDEs) and ridgeline plots highlight conformity distribution for the frequency, length, and polysemy features in toddlers' mental lexicon and the randomised variants; (E): Above – two-dimensional scatter plot of conformity vector space, where each point is colored according to the cluster the point belongs to (K-means algorithm); Below – distribution of word features within each cluster, where a kernel language emerges, i.e. the cluster labeled as *D*; (F): content characterisation of the kernel compared to a competitor from a k-core decomposition.

details on network layers and vectors of word features). Conformity with respect to frequency highlights an assortative mixing pattern but limited only to highly frequent words, i.e. only words occurring many times in child-directed speech tend to connect with each other in children's *FERMULEX* network. This effect is absent in lower-frequency words and it was not detected in single-layer semantic networks of adults [69]. Conformity of

word length highlights an assortative mixing pattern of very short words only. These two effects are expected to be related as shorter words tend to be more frequent in language [65].

Interestingly, conformity quantifies that polysemous words are likely to connect to each other to a smaller extent than most frequent and shortest words. This indicates an organisation of concepts where unambiguous/less polysemous words are linked to ambiguous/more polysemous words. This heterogeneous mixing by polysemy could be beneficial in providing context and differentiating among possible meanings of a polysemous word, as suggested by previous studies [12, 69]. If all ambiguous words were grouped together, sense disambiguation could not rely on links including less polysemous/unambiguous words and this homogeneity would ultimately violate the frequency-meaning law [23].

The above assortative mixing patterns are not a consequence of feature/distance distributions, because reshuffling node labels (*Null Model 1*) and rewiring links (*Null Model 2*) disrupt the heterogeneous mixing behaviour among classes (see Methods and SI). Hence, the above patterns indicate the presence of a core-periphery organisation in the dualistic multiplex/feature-rich structure of the mental lexicon: A set of highly frequent/shorter/polysemous words linked with each other creates a network core highlighted by conformity and invisible to previous inquiries [63, 62]. This preliminary evidence calls for further analysis of the core.

Fig. 1, E introduces an analysis of the core performed on: (i) dualistic network/vector and (ii) individual aspects of words in the mental lexicon of toddlers (see Methods and SI). We aim to find a language core that contains groups of words sharing similar structure-feature relationships. Among the six optimal clusters found (see Methods and SI), groups A and B (blue and gold) contain mostly short words. Cluster F (cyan) contains highly frequent words. Cluster D contains short, highly frequent and a relevant portion of polysemous words. Sets of clustered words with such features are known as language kernels in cognitive network science [10, 65, 66]. Language kernels facilitate communication through a small set of simple words suitable for expressing ideas in multiple contexts [10]. The conformity core (cluster D) satisfies this definition. In fact, 13% of the core is made of nouns, 33% of verbs and the other 54% include adjectives, adverbs and pronouns, which make it more likely to assemble syntactically well-formed sentences by random sampling compared to other word clusters (cf. Fig. 1 F). Identifying a network core via k-core decomposition [30] shows almost no meaning organisation and more expressions that are syntactically unrelated (see two random samples in Fig. 1, F). Analogously, K-Modes [31] attribute-only clusters are unable to form syntactically coherent bigrams. See SI for an analysis centered on computing the internal syntactic coherence of the cores. These comparisons provide unprecedented evidence showing a syntactically advantageous organisation of words in early children’s lexicon. This phenomenon goes undetected unless both the network and vector nature of words in the mind is considered.

Topology and cognitive relevance of the conformity core *FERMULEX*. We further compare the con-

formity core with the k-core decomposition [30] (where similarities are switched off) and with the most relevant K-Modes cluster (where network structure is switched off). Interestingly, the conformity core appears to be a synthesis of the other two potential language kernels. Fig. 2, C characterises the three cores with several qualitative functions assessing intra-cluster connectivity and inter-cluster distinctiveness (cf. Methods and the SI). The K-Modes core contains a rich set of short, highly frequent and polysemous words compared to the conformity core. The conformity core contains a more homogeneous set of words, which is crucial for syntactic sentences mixing specific and more general concepts [10, 51, 12]. The structural k-core has high transitivity, but the conformity core has a more *cliquish* configuration due to higher hub dominance score [75]. Cliquishness was recently shown to correlate with better recall from memory [68] due to the concentration of spreading activation in the clique [34]. These recent studies suggest that the higher cliquishness found here for the conformity core might be beneficial for language processing in toddlers. The conformity core also displays high values of conductance and cut ratio: this language kernel possess a dense internal structure but it is also strongly connected to the rest of the graph as well, considerably more than the other competitors. In other words, the conformity core is strongly internally connected (more than k-core) and homogeneous with respect to the features (more than k-mode). This higher connectivity might reflect an advantage in accessing and producing items from the language kernel in view of activation spreading models of the mental lexicon [34, 33, 13, 37].

Normative word learning as random walks on *FERMULEX*. To investigate how the conformity core and the whole *FERMULEX* structure emerge over time, we adopt a random walk framework. Random walks on cognitive network structures have successfully modelled phenomena like Zipf’s law [23] or semantic priming [24]. Here, we use structure-feature biased random walks to explore normative language learning, as reported in Fig. 2.

The simplest idea is to limit the walk to network structure only (*Graph Walk 1*). To explore the interplay between topology and features of words, we can weigh network links with the similarity between vectors representing adjacent words (*Graph Walk 2*). Let us consider an example. In Figure 2, at t_2 , *Graph Walk 1* should choose to learn either *cat* or *daddy* after the current word *mommy*. Because of network/vectorial similarities, *Graph Walk 2* will select *daddy* as the *next-to-be-learned* word. We can expand the set of next-to-be-learned candidate words: *Graph Walk 3* encodes this parallel word learning process by considering as potential candidates all neighbors of already learned words. With reference to Fig. 2 A, at t_3 , *Graph Walk 2* can only move to and learn *friend*, while *Graph Walk 3* can also activate and learn *cat* after *mommy*. Focus is given to considering how these models can predict the assembly over time of: (i) the conformity core, and of (ii) Communicative Development Inventories [22] (CDIs), which are commonly used by psycholinguists to measure a child’s communicative, receptive and expressive abilities. CDIs are clusters of words from the same semantic category - e.g. a list of words all relative to *time* - and thus represent a portion of the whole vocabulary available to children [42].

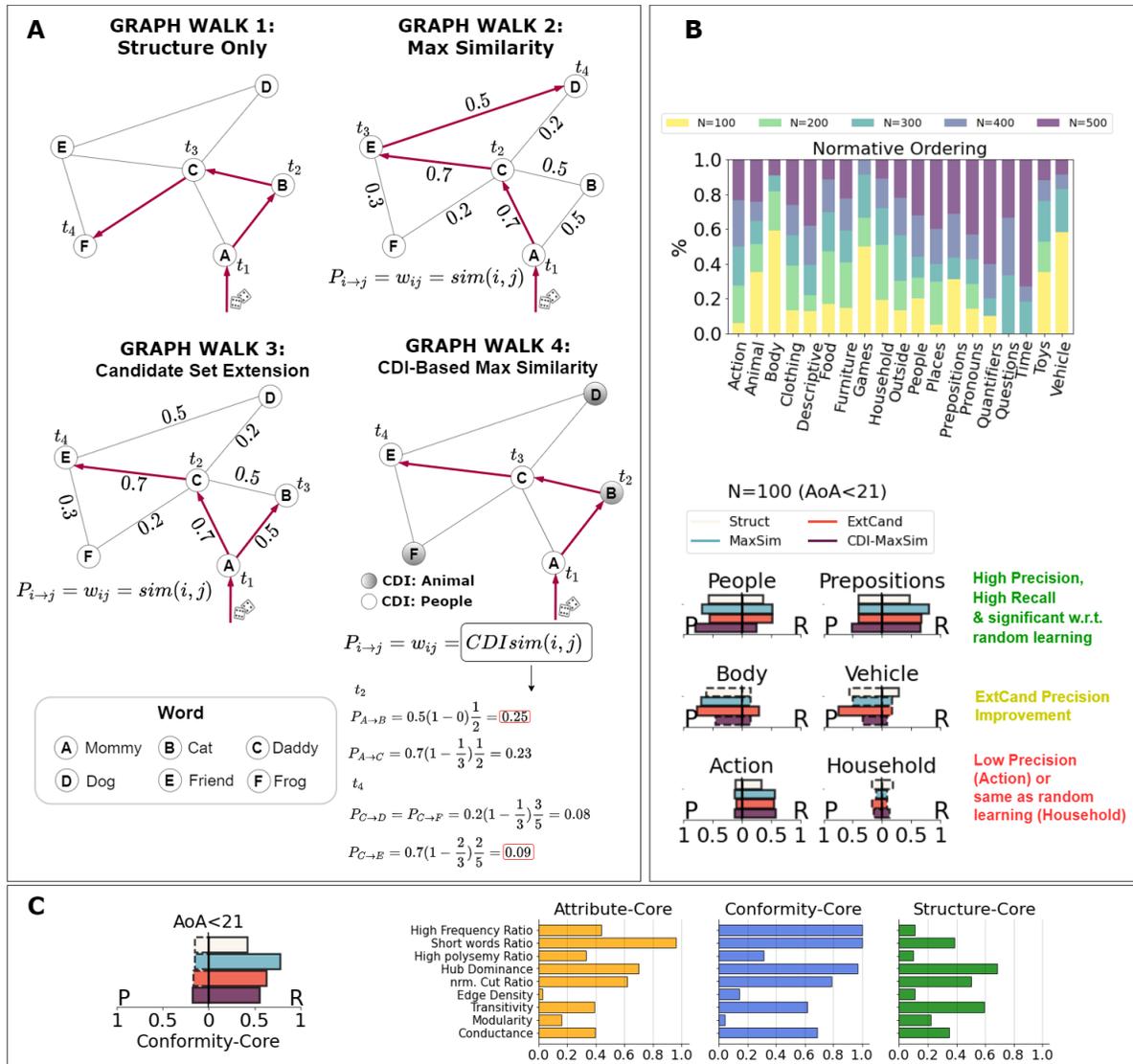


Figure 2: (A): Random walks combining progressively structure and vector information (Graph Walk 1-3) and CDIs integration (Graph Walk 4); (B): Above – CDIs filling in CHILDES normative learning; bars show that CDIs are not uniformly filled over time, e.g. more than half of *Body* and *Vehicle* categories are learned during early stage acquisition, whereas *Questions* and *Time* emerge later; Below – precision-recall evaluation over selected CDIs; solid bars identify statistically significant scores compared to a random learning baseline; (C): Left – precision-recall evaluation with respect to early acquisition of kernel words; Right – kernel characterisations using several quality measures.

CDIs are not filled uniformly under normative learning. In the CHILDES data [42], toddlers are found not to learn CDIs uniformly over time (cf. Fig. 2, B). This means that some semantic domains of toddlers’ lexicon are filled earlier during normative learning. However, the above random walkers do not include information about the semantic category a word belongs to. *Graph Walk 4* proposes a CDI-based similarity integrating information about CDIs’ availability and attractiveness. In the figure, at t_2 , *Graph Walk 4* moves from *mommy* to *cat*,

because *Animal*-CDI is relatively emptier than *People*-CDI, i.e. *People* already contains *mommy*. However, at t_4 , the model learns *friend* from *daddy*, because the *feature similarity equation* term is stronger than the CDI-based ones (see Methods and SI).

Toddler’s language kernel rises from CDI density and network/vector dualities. Figure (2, C - left) reports precision and recall in reconstructing the conformity core early on during cognitive development. Performance metrics statistically higher than random learning (significance of 0.05, see SI B) are highlighted with full bars. Non-significant results are visualised as dashed bars. The normative growth of children’s language kernel was captured with a precision higher than random learning only by our most advanced model, combining CDI density, multiplex network structure and feature similarities. This provides strong evidence that semantic spheres and their filling over time provide insights additional to network/vector duality for capturing how early production of syntactically coherent sentences is achieved [26]. Compared to other CDIs (see next section), our walkers achieved a relatively lower precision in predicting the assembly of the conformity core. This indicates that the language kernel *does not* emerge all at once during early cognitive development, unlike other kernels highlighted in older children [65]. The emergence of the conformity core is thus a gradual phenomenon, that is not strongly biased by similarities and cannot thus be captured by biased random walks only.

Random walks highlight different strategies at work in different CDIs. Random walks produce word ordering lists that we evaluate with respect to CHILDES normative ordering, i.e. the order in which most children produced words over time (Fig. 2, B - above). Random learning is used as a baseline to test whether walks considering word topology and feature predict more words as correctly learned over time. See Methods/SI for details of our statistical approach.

Table 1 presents a coarse-grained evaluation of the walkers (cf. Discussion). Fig. 2, B sums up results with respect to CDIs focusing on the very early stage of acquisition, which corresponds to $N = 100$ words learned before 21 months [63]. The selected CDIs are captured differently by the models. CDIs like *People* and *Prepositions* are predicted with higher-than-random precision and recall for all Graph Walk models. *CDI-MaxSim* precision is slightly better than in the other models. Interestingly, the two most filled CDIs in this stage of acquisition, i.e. *Body* and *Vehicle*, are predicted with high precision but low recall (cf. Methods/SI). This means that the few words predicted are the expected ones, but the models cannot fill the CDIs. *ExtCand* precision is higher. Not all CDIs can be predicted in this way, e.g. *Action* and *Household*. Furthermore, model performances for *Household* are not distinguishable from a random learning, i.e. all bars are dotted. The high recall but the low precision of *Action* is poorly relevant: less of 0.1% of the CDI is covered in this stage of acquisition (however, cf. the SI, where *Action* category is well captured in other stages).

		Accuracy Relevant CDIs	Precision Relevant CDIs	Recall Relevant CDIs
AoA < 21	Random Learning	0.67 —	0.17 —	0.19 —
	Struct	0.70 0.26	0.40 0.64	0.30 0.58
	MaxSim	0.76 0.26	0.37 0.70	0.34 0.52
	ExtCand	0.65 0.21	0.55 0.76	0.30 0.58
	CDI-MaxSim	0.75 0.42	0.25 0.58	0.34 0.47
21 < AoA < 23	Random Learning	0.71 —	0.17 —	0.19 —
	Struct	— 0.00	0.24 0.64	0.24 0.71
	MaxSim	0.82 0.36	0.28 0.57	0.25 0.64
	ExtCand	0.83 0.26	0.24 0.42	0.24 0.50
	CDI-MaxSim	0.66 0.10	0.26 0.71	0.26 0.71
23 < AoA < 24	Random Learning	0.69 —	0.17 —	0.19 —
	Struct	0.73 0.21	0.19 0.42	0.21 0.52
	MaxSim	0.75 0.36	0.17 0.42	0.23 0.42
	ExtCand	0.73 0.31	0.20 0.21	0.21 0.52
	CDI-MaxSim	0.69 0.21	0.19 0.52	0.23 0.52
24 < AoA < 26	Random Learning	0.70 —	0.17 —	0.19 —
	Struct	0.73 0.31	0.20 0.44	0.22 0.61
	MaxSim	0.72 0.31	0.21 0.38	0.26 0.44
	ExtCand	0.71 0.42	0.18 0.33	0.22 0.50
	CDI-MaxSim	0.72 0.31	0.23 0.38	0.22 0.44
AoA > 26	Random Learning	0.61 —	0.24 —	0.24 —
	Struct	0.68 0.78	0.32 0.72	0.36 0.61
	MaxSim	0.70 0.84	0.33 0.77	0.35 0.66
	ExtCand	0.64 0.52	0.28 0.44	0.29 0.66
	CDI-MaxSim	0.79 0.31	0.33 0.77	0.41 0.66

Table 1: Model performances over each bin of acquisition. *Relevant CDI fraction* is the ratio of statistically significant precision/recall values against a random learning model.

2 Discussion

This work introduces a cutting-edge combination of network [17, 66, 71] and vector [25, 35] aspects of knowledge in the human mind, which historically run in parallel when modelling language and its cognitive processes [13].

Using data from 1000 toddlers between 18 and 30 months from the CHILDES project [42], our *FERMULEX* network revealed a core of words facilitating word production [51] and invisible to methods based on network structure [63, 65, 30] or vector similarities only. This core was detected via conformity [15], a metric extending assortative mixing estimation in a multi-scale, node-centric fashion. Our numerical experiments identified this core as a set of highly frequent, short, polysemous and well-connected nouns and verbs, i.e. a language kernel containing concepts versatile enough to communicate via basic sentences (cf. [10]) and whose access via spreading activation is facilitated by network connectivity [34, 13, 68]. Revealing the presence of such a core through our analyses provides for the first time quantitative support of recent empirical findings showing that typical learners can start combining words in basic sentences after 30 months of age [51]. The kernel persisted even when co-occurrences from child-directed speech were ignored (see SI): the conformity core emerged from an interplay between semantic/phonological associations and psycholinguistic norms in the mental lexicon of linguistic knowledge.

To investigate the assembly over time of such a crucial core of linguistic knowledge, we implemented artificial models of word learning as biased random walks over *FERMULEX*, inspired by past approaches using walkers to investigate the mental lexicon [24, 23]. We found that the conformity core does not emerge suddenly over time, differently from other language kernels modelled as viable component in other studies [65]. Instead, the conformity core is progressively built in ways that are captured only by combining the network and vector aspects of words together with CDI filling rates. This finding quantitatively stresses that the conformity core - containing building blocks for producing syntactically coherent words - emerges from strategies dependent on semantic categories, which are partly captured by CDIs [42].

We also used the same random walkers for capturing how different CDIs filled over time through normative learning, giving unprecedented focus [63] to learning strategies for individual aspects of children’s knowledge. In our analyses, different CDIs are found to fill at different times over developmental stages, further emphasizing that language learning is not a uniformly random process. Inventories relative to food and action themes are found to be predicted well by our model, confirming recent independent studies [14, 52] that these salient familiar themes are crucial for predicting early language acquisition.

Notice also that words in some CDIs might be learned according to context-specific strategies [60, 16], so that a single, general word-learning strategy might not fit all cases. For instance, according to the *Pervasiveness*

Hypothesis by Clerkin and colleagues [16], toddlers would tend to learn earlier words more frequently occurring across several daily contexts. This visual prevalence/occurrence would be crucially missing from CDIs like *Household* or *Action*, which were in fact poorly reproduced by our model. These negative findings indicate the presence of local strategies for learning words in physical settings that are at work in toddlers but missing from the current instance of *FERMULEX*.

For inventories like *Body* or *Vehicle*, a combination of network structure and feature similarities corresponded to a significant boost in precision over predictions from random learning. This is quantitative evidence for combining network and vector aspects of the mental lexicon. A further boost in precision was found when the random walker was allowed to backtrack. This indicates that some components of the mental lexicon are not built sequentially, without appending words to the most recent lexical item, as assumed in attachment kernel models [64], but rather filling gaps in the whole vocabulary available to children, as shown also by other approaches with persistent homology and gap filling [61].

Interestingly, recency in word acquisition is found to be more a powerful strategy for reconstructing the filling of CDIs like *People* or *Prepositions*, where our most elaborate random walker based on recency beats the back-tracking one. Our quantitative results open the way for further discussion and interpretation in light of psychological studies behind early language learning.

This first conception of *FERMULEX* has some key limitations, which can be addressed in future research. For example, our approach considers only normative learning, i.e. how most children learn words over time [63]. This learning dynamic might be different from how individual children with different language learning skills might learn words over time [7]. Future research should thus test the presence of the language kernel and its time-evolution dynamics in a longitudinal cohort of children. Since the occurrence of the language kernel characterises normative learning in a large population of 1000 and more toddlers [42] and it supports the production of early sentences observed in normative talkers [26], we expect for the kernel to be present in normative learners but also to be disrupted or incomplete in late talkers [5]. If supported by data, then the language kernel revealed here could become a crucial early predictor of delayed language development in young children. Another limitation is that our predictions do not treat learning as the outcome of a statistical process, where words are learned with certain probabilities. Rather we model word learning as a binary learned/not learned process. We chose to follow this approach for model parsimony and indicate the addition of statistical learning [55] within the *FERMULEX* framework as an exciting future research direction. Future enhancements of random-walk models should account also for distinctiveness in addition to similarity. The recent work by Siew [58] indicates that global feature distinctiveness, i.e. how many different semantic features are possessed by a word, correlates with earlier acquisition. Hence, random walkers accounting for switches between distinctiveness

and similarity might enhance prediction results and represent an exciting future research direction. Another important approach for future research might be casting language acquisition as a percolation problem, which has been explored in feature-rich networks only recently [2]. An important limitation of our study is that it adopts CDIs for modelling language learning, however these inventories are not grounded in theories from cognitive psychology [22] but were rather created *ad-hoc* by psycholinguists. Future instances of *FERMULEX* networks should rely on word learning data that is more representative across semantic and syntactic categories.

3 Methods

Multiplex Layers. We modelled word learning as a cognitive process acting on a mental representation of linguistic knowledge. Structure in this representation is given by a multiplex lexical network, where nodes represent words that are replicated and connected across different semantic and phonological levels of the network [63].

Only layers of relevance for word learning acquisition were considered [63], namely: (i) free associations, indicating memory recall patterns between words from semantic memory [20], (ii) co-occurrences in child-directed speech [63, 42], (iii) feature-sharing norms, indicating which concepts shared at least one semantic feature from the McRae dataset [44] and (iv) phonological similarities [72], representing which words differed by the addition/substitution/deletion of one phoneme only. Hills and colleagues showed that the words with larger degrees in free association networks were also more likely to be acquired at earlier ages, a phenomenon known also as *lure of the associates* (cf. also [28]). A subsequent study by Carlson and colleagues [11] found a similar effect also in phonological networks built from child-directed speech [72]. Investigations of co-occurrence and feature sharing networks by Beckage and Colunga reported that highly connected words were distinct trademarks of early word production in typical talkers [6]. Importantly, these four aspects of knowledge in the human mind produced network representations that were irreducible [63]. Layers represented different connectivity patterns among words and could thus not be aggregated or erased without decreasing structural information about the system in terms of Von Neumann graph entropy.

Normative age of acquisition. Network models of language acquisition often use normative datasets that follow the development of language production in toddlers [28]. The most prominent data source is CHILDES (Child Language Data Exchange System), a multi-language corpus of the TalkBank system established by MacWhinney and Snow, storing data about language acquisition in toddlers between age 16 and 36 months [42].

We used CHILDES data to rank words in the order they are learned by most English toddlers. By considering the fraction of children producing a certain word in a given month, within each month, words were assigned a

production probability. Month after month, a rank in descending order of production probability was constructed as a proxy for normative learning of most toddlers, as done in previous studies [29, 6, 63].

Features. This study selected word features shown in previous research to influence early language acquisition, namely frequency in child-directed speech [42, 63], word length [29, 7] and polysemy [12]. Polysemy scores indicated the numbers of meanings relative to a given word in WordNet [45], a proxy to word polysemy successfully used in quantitative studies of early word learning [63]. Due to the highly-skewed distribution of variables (e.g., Zipf’s law for word frequency [76]), we regularised data by recasting it from numerical to categorical, as to avoid biases in computing conformity [56]. We grouped each variable into discrete bins, fine tuning bin boundaries so as to obtain non-empty bins featuring the same order of magnitude of entries. This fine-tuning led to splitting words in quintiles for both word frequency and polysemy and in tertiles for length.

Conformity. We characterise the interplay between structure and features through conformity [56], which estimates the mixing patterns of nodes in a feature-rich network, i.e. a categorical node-attributed network. This measure can find heterogeneous behaviour among all nodes of a network. Conformity enables a multi-scale strategy by leveraging node distances for computing label-similarities between a target node and other nodes. A distance damping parameter α is needed for decreasing the impact of label-similarities over longer network distances between the target node and its connected neighbors. Based on previous investigations [56], we adopt a value of $\alpha = 2$ giving more emphasis to closer neighbours in a given network topology. See the SI or [56] for a formal description of the measure and the motivation behind its choice in this work.

When analysing conformity, we need to test whether the measured values are a trivial consequence of structural (or attributive) patterns or rather come from a non-trivial interplay between the two. To characterise this, we resort to two null models: (i) random re-shuffling the node attribute labels while maintaining network topology (*Null Model 1*, Fig. 1, D, [65]), and (ii) randomly rewiring of links while preserving the node degree and attribute labels (*Null Model 2*, Fig. 1, 1 D). In other words, let us consider this question: Are two labels at the endpoint of an edge significant for the distribution of conformity or can we observe similar patterns by randomly rewiring the attributive or structural model components? While rewiring labels or connectivity patterns, respectively, we keep the other component fixed. For building *Null Model 1*, a random label permutation is enough to disrupt correlations between structure and features. For building *Null Model 2*, we used a configuration model [46] to obtain a degree preserving graph randomisation, that is, given N nodes and any arbitrary degree sequence $\{k_i\} = (k_1, k_2, k_N)$, we place k_i stubs on each node i in the graph; then we match each stub with another one until all stubs are matched. The conformity distributions of the null models in Fig. 1, D refer to the average node scores from 100 randomised instances of *FERMULEX* network.

All conformity distributions are analysed through kernel density estimates (KDEs) and ridgelines (1, D); in

particular, these last ones get a better picture of mixing heterogeneity between the class labels on the original toddlers' lexicon.

Core: Definition and Evaluation. For finding a potential language core, we model each word as a vector of conformity scores. This results in a vector space where classic clustering algorithms as K-Means [40] can be run. We reveal a relevant set of words among the six optimal clusters identified by K-Means through the elbow method. The SI provides methodological details about this configuration.

A set of several quality functions are proposed to characterise the language core. We focus on modularity, conductance, cut ratio, internal edge density, hub dominance and transitivity [47]. Modularity, conductance and cut ratio focus on the links within and outside a community: They measure how well-separated a cluster is from the rest of the network. Edge density, transitivity and hub dominance characterise the internal structure of the core. In particular, transitivity and hub dominance characterise it in terms of triadic closure and *cliquishness* level, i.e. the creation of subgraphs where each node is fully connected to others. See the SI for their formal description. All in all, these network metrics are used to characterise the structure of the different cores found via conformity (in *FERMULEX*), via core-detection on the network structure only [30] and via K-Modes on feature embeddings only [31]. Notice that these measures, combined, provide info about the distinctiveness and connectedness of a given component/cluster in a network.

Graph Walks. We aim to model early word acquisition by progressively combining the network and vector components of *FERMULEX*. To achieve this goal, the core idea is to generate a word rank that is progressively filled according to the different graph walk strategies, each one incorporating specific assumptions. In this work we compare four alternative random walk models each one having a unique rationale on how to weigh links and/or to determine the set of candidates for the next to-be-learned word. In particular:

- *Struct* (Graph Walk 1): Words are connected by unweighted links, hence the next word is chosen according to the underlying structure only. Similarly, the set of candidates is chosen from the adjacent neighborhood of the current word;
- *MaxSim* (Graph Walk 2): Edges are weighted according to the pairwise similarity between nodes' features. Jaccard similarity is used (cf. SI), and frequency, length and polysemy are all considered. The same strategy of *Struct* is used for the set of candidates;
- *ExtCand* (Graph Walk 3): The same strategy of *MaxSim* is used for weighing links; the set of candidates is chosen from the adjacent neighborhood of all the words already learned;
- *CDI-MaxSim* (Graph Walk 4): Links are weighted according to a CDI-based pairwise similarity between

the attributes of nodes as well as the availability and attractiveness (cf. SI), and it needs to be updated at each iteration. The same strategy of *Struct* and *MaxSim* is used for the set of candidates.

Struct and *MaxSim* are biased random walks considering, respectively, topology or similarity between words (i.e., the network structure or the vector space) while *ExtCand* and *CDI-MaxSim* aim for a more holistic approach.

ExtCand visit strategy is designed to mime non-sequential word learning in children (cf. [7]), where the word acquired at step $t + 1$ could be similar to any word already learned before, thus enabling an interplay between exploration and exploitation of CDIs. When the last word determines the topology of similar candidates for the next acquisition step, resembling a Markovian process [29], the walker possesses a bias to remain within the same CDI. By considering as to-be-learned candidates all previously learned words, the walker has a chance of backtracking and acquiring more words within the CDI sharing tightly similar concepts.

CDI-MaxSim, the CDI-based model relies on pairwise similarity between two words modulated by additional information on the filling of CDIs they belong to. For additional details and a formal description of the pairwise similarity function adopted refer to the SI.

Graph Walk Evaluation. Accuracy, precision and recall are used to evaluate the goodness of ranks' prediction, as commonly done in statistics and machine learning. Accuracy is defined as the number of correct predictions, i.e. true positives or TP, divided by the total number of predictions. In this domain, TPs are words belonging to a CDI that are learned by a random walker in a specific bin of age of acquisition. Precision is the fraction of relevant elements among all the retrieved ones including non-relevant elements, i.e. false positives or FP. In this domain, FPs are words that fill a CDI as expected in a particular age of acquisition bin, but they are not the exact same words considered in normative learning. For instance, *dog* might contribute to increase FPs because it belongs to the *Animals* CDI but the normative learning contemplated *cat* instead of *dog*. Finally, recall is the fraction of relevant elements that are retrieved. Missing relevant elements (false negatives or FNs) are CDI's words that are not retrieved by a random walker in a particular bin of age of acquisition. The above definitions imply that there can be predictions with high recall and low precision, because there are many words that satisfy the semantic category roughly represented by the CDI (e.g. guessing as learned names of animals) but different from the specific words learned during normative acquisition (e.g. other names of animals). This interplay spans from the specific characterisation of random-walk predictions and it is accounted for in the Results and Discussion sections. See the SI for a complete formalization of the measures, and toy examples.

Appendix: Supporting Information

A Conformity

A.1 Definition

Conformity provides a multi-scale strategy to estimate local homophily in complex networks, overcoming classic measures as Newman’s assortativity [48] that only produce a global, averaged score. Conformity is not the only multi-scale strategy in the literature. Some valuable variants of local Newman’s assortativity also exist [49]. The reason behind the choice of conformity rather than other measures is because conformity provides node similarities grounded in the real distances between nodes. Other approaches, for instance, only leverage random walks as a proxy of information about paths of all possible lengths [49, 3].

To the best of our knowledge, no other works have used node-centric homophily estimation in a mental lexicon. These measurements were typically applied in social network analysis [43] or mobility data [8].

In the following, we report a concise description of conformity [56]. Given a node-attributed graph $G = (V, E, A)$, where V is the set of nodes, E the set of edges, and A the set of categorical node attributes, we calculate for a node $u \in V$ the conformity score $\psi(u, \alpha)$ with respect to an attribute $l \in A$. The damping parameter α allows to decrease the impact of label-similarities over longer network distances between the target node and all the reachable neighbors. To define conformity formally, we need a couple of support functions, namely the indicator $I_{u,v}$

$$I_{u,v} = \begin{cases} 1 & \text{if } l_u = l_v \\ -1 & \text{otherwise,} \end{cases} \quad (1)$$

that compares the attribute values of two nodes, and the similarity function f_{u,l_u}

$$f_{u,l_u} = \frac{|\{v | v \in \Gamma_u \wedge l_u = l_v\}|}{|\Gamma_u|}, \quad (2)$$

that computes the ratio of u ’s first-order neighbors that share the same attribute value l_u . Thus, given a real number α in $[0, +\infty)$, conformity of node $u \in V$ is defined as in the following:

$$\psi(u, \alpha) = \frac{\sum_{d \in D} \frac{\sum_{v \in N_{u,d}} I_{u,v} f_{v,l_v}}{|N_{u,d}| d^\alpha}}{\sum_{d \in D} d^{-\alpha}}, \quad (3)$$

where $D = \max(\{dist(i, j) | i, j \in V\})$, i.e. the maximum distance among all node pairs. The computed score is normalized to ensure that conformity lies in the range $[-1, 1]$.

A.2 Comments on null models

Some class labels exhibit more assortative mixing than others, e.g. shortest words in conformity by length or most frequent words in conformity by frequency (cf. Results). While reshuffling node labels or rewiring links, we intend to observe whether similar distributions can emerge trivially from random label permutations or random link configurations. The heterogeneous distributions do not emerge while measuring conformity on the ensemble of networks obtained from the two randomisation processes (cf. Results and Methods). In conformity by frequency and by length, null models distributions are mainly disassortative among all classes; hence, the *anomalous* behaviour of most frequent and shortest words is flattened by both randomisation processes. Conformity distribution with respect to polysemy slightly differs from the other two attributes. In particular, the null model that rewires links shows a bi-modal distribution. The explanation behind this behaviour could be similar to the explanation used to describe the quasi uniformly mixed pattern of polysemous words (cf. Results): ambiguous and unambiguous words must link in non-trivial patterns that are harder to break while randomising node connectivity.

A.3 Conformity vector space

Multi-dimensional conformity information is used in the language core analysis for finding a relevant set of words in language acquisition. We model each node as its vector of conformity scores, where conformity by frequency, by length and by polysemy are the vector components. This allows to build a vector space where classic clustering analysis can be performed. The difference between a clustering method on the features only (cf. K-Modes [31], Results) is that using vector of conformity scores we integrate structure-feature relationships, thus we aim to group words having similar mixing patterns across the features.

An optimal K-Means [40] instance is used to cluster words. For selecting the optimal number of clusters k , we leverage the elbow method, namely plotting the sum of squared errors in function of k and determining the point of inflection of the curve. $k = 6$ is identified as the optimal point and chosen as the number of centroids to initialize the algorithm.

B Core Evaluation

We report here a description of the quality functions used for characterising the language core(s). All the measures are implemented in the CDLib library [57], and other detailed information can be found in the library documentation. Let $G = (V, E)$ be a graph with $v \in V$ and $e \in E$, and C a partition of G with $c \in C$, with c composed by a subset of V and a subset of E . We aim to characterise a cluster/community/core c with the following measures:

- *Conductance*: the fraction of total edge volume that points outside the community:

$$conductance(c) = \frac{|c|}{2|e_c| + |c|},$$

where $|c|$ is the cardinality of the community, namely the number of community nodes, and e_c the number of community edges;

- *Edge density*: the internal density of the community set:

$$density(c) = \frac{|e_c|}{\frac{|c|(|c|-1)}{2}};$$

- *Hub dominance*: indicates the ratio of the degree of the most connected node in a community with respect to the theoretically maximal degree within the community, namely

$$Hub_dom(c) = \begin{cases} 1 & \text{iff } |c| = 1 \\ \frac{\max_{v \in c} k_v}{|c|-1} & \text{otherwise,} \end{cases},$$

where k_v is the degree of node v ;

- *Modularity*: measures the strength of the division of a network into sets of well-separated clusters or modules, and it is calculated as the sum of the differences between the fraction of edges that actually fall within a given community and the expected fraction if edges were randomly distributed:

$$Q = \frac{1}{(2e_c)} \sum_{vw} \left[A_{vw} - \frac{k_v k_w}{(2e)} \right] \delta(c_v, c_w),$$

where $A_{v,w}$ is the entry of the adjacency matrix for $v, w \in V$, k_v, k_w the degree of v, w and $\delta(c_v, c_w)$ is an indicator function taking value 1 iff v, w belong to the same community, 0 otherwise;

- *Normalized cut ratio*: is the fraction of existing edges (out of all possible edges) leaving the community:

$$cut_ratio(c) = \frac{|c|}{2|e_c|+|c|} + \frac{|c|}{2(|e|-|e_c|)+|c|};$$

- *Transitivity*: is the average clustering coefficient of community nodes with respect to their connection within the community itself:

$$CC(c) = \frac{1}{|c|} \sum_{v \in c} \frac{2\Delta}{k_v(k_v-1)}$$

where Δ is the number of triangles including node v in the community c .

C Core Analysis

C.1 Persistence without layers

A key result of this work is the identification of a language kernel with interesting structural properties and non-trivial content organisation. Randomly picked pairs/triads of words from this kernel can build simple, syntactically well-formed sentences (cf. Results). However, it can be observed that this core can appear just because there is the co-occurrence layer in the multiplex network, i.e. links between concepts co-occurring in child-directed speech. A strength of the *FERMULEX* model is the possibility to *switch off* a layer from the structural component. Removing completely a layer also allows us to observe the emergence of an interesting conformity core. Fig. 3 establishes that a language core emerges even when the 1000 co-occurrences links from the child-directed speech layer are removed. The language kernel (here, the cluster labeled as D) persists to show homogeneity across all the features. This gives more strength to the hypothesis that the kernel found with conformity stems from a more broad interplay between semantic and phonological layers. Finally, for a complete overview of the whole clustering result, the few differences we can notice without the co-occurrence layer is a more homogeneous distribution of word length within the clusters (cf. Fig 1, E), probably due to the removal of long words from the network.

C.2 Degree assortativity and hierarchical organisation

We may want to observe whether the removal of the conformity-core can disrupt global characteristics of the network. Several studies [38, 67, 54] identify complex global properties in the degree-degree assortativity and in the hierarchical organisation of lexical networks, measured through $Knn(k)$ and $C(k)$ curves, respectively. $Knn(k)$ curves show the average degree of neighbors of nodes with degree k [38]. Similarly, $C(k)$ curves show the average clustering coefficient of nodes with degree k [54]. If $Knn(k)$ increases with k , the network behaviour is

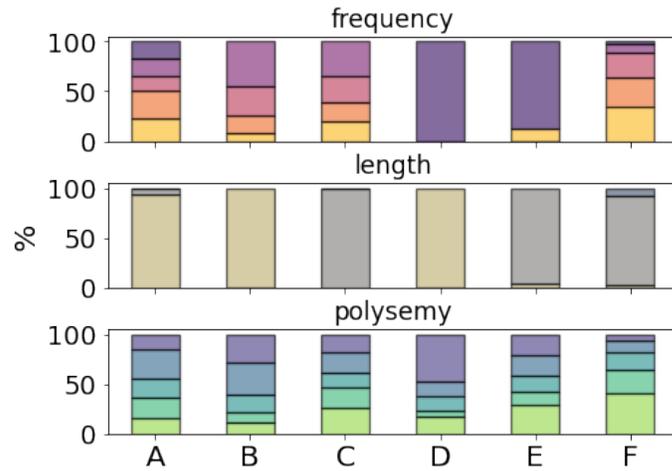


Figure 3: Conformity vector space characterisation without the co-occurrence layer.

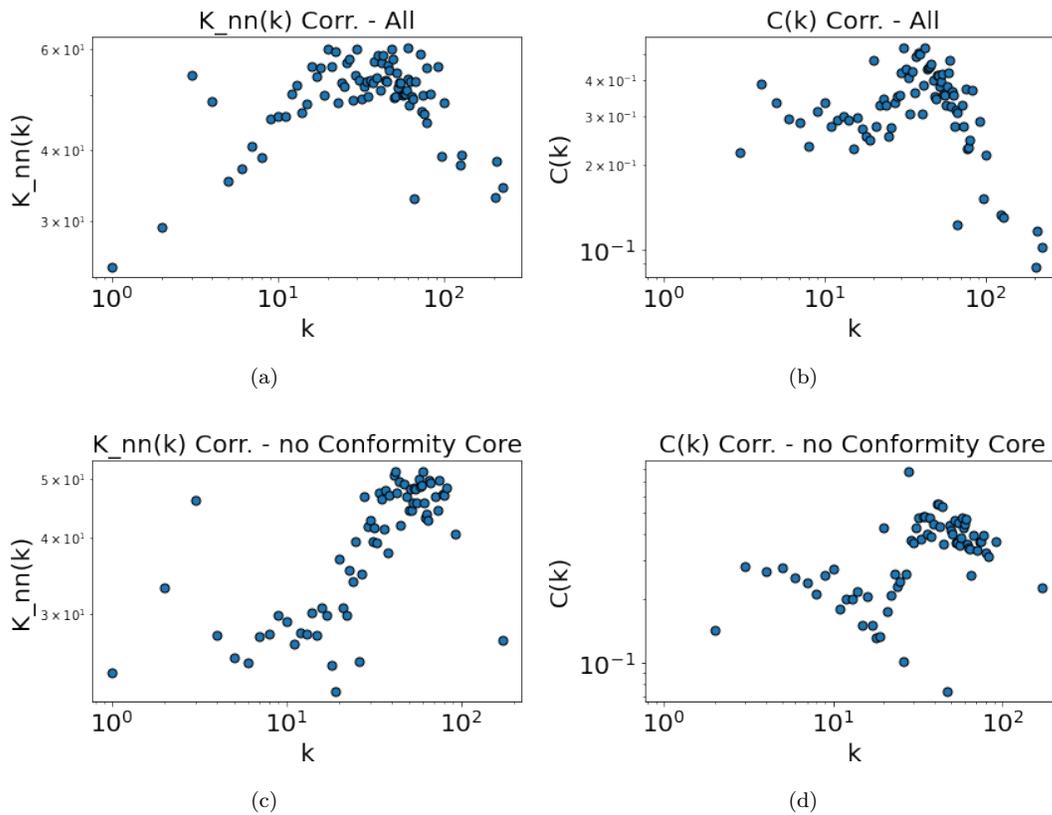


Figure 4: Degree-degree assortativity and hierarchical organisation in CHILDES mental lexicon studying $K_{nn}(k)$ and $C(k)$ curves, respectively, (a-b) with and (c-d) without the conformity core.

assortative by degree; if $K_{nn}(k)$ decreases with k , the network behaviour is disassortative. If $C(k)$ decreases with k , the network exhibits a hierarchical organisation, otherwise the network does not present this characteristic.

Fig. 4 shows that $Knn(k)$ and $C(k)$ curves drastically change when removing the language kernel found with conformity. The network is highly disassortative and hierarchically organised when all nodes are present, but switches to a highly degree-assortative behaviour and splits into two $C(k)$ branches when the core is removed ($r_{Knn} = -0.27$ with $N = V$, $r_{Knn} = 0.55$ with $N = V - V_{core}$; $r_C = -0.33$ with $N = V$, $r_C = 0.32$ with $N = V - V_{core}$). Both disassortativity and hierarchies may spawn from super-general concepts being embedded in the core. These nodes may attract many more links than other nodes, forming a hub-node structure that creates degree disassortativity and implies that the neighbors of hubs are not linked to each other. We suggest that these findings are coherent with the intra-cluster and inter-clusters analysis to characterise the core (cf. Results). The core is internally structured as a tight *clique*, i.e. it presents high transitivity and hub dominance values. Nevertheless, high conductance/cut-ratio values indicate that the core is highly connected to the rest of the graph. Thus, removing this set of words disrupts the global disassortative and hierarchical organisation of the system, i.e. the underlying structure that guarantees the system connectivity. Further analyses are needed if we aim to interpret these results from a cognitive perspective. Despite some relevant exceptions [69], the lack of a compact set of studies on the degree-disassortative behaviour and the hierarchical organisation of lexical networks makes it difficult to shed light on the underlying cognitive phenomena structuring these non-trivial topological patterns.

C.3 Comparison with other cores

A language kernel is defined as relatively small set of words enabling the creation of simple yet general and frequently used sentences, thus facilitating early communication [10]. This kernel was identified in the cluster found via conformity. Words in this kernel form simple and syntactically structured sentences, and are heterogeneous in part of speech composition. Conversely, the kernel found via k-core network decomposition [30] is unable to form syntactically coherent bigrams/trigrams (cf. Results, Fig. 1, F). For a more robust investigation, we report additional analysis here, including in the discussion the kernel found via K-Modes clustering as well [31].

We focus on the whole internal content organisation of the three cores, regardless of their underlying structure, that we characterized up to now. We aim to compute the ratios of the internal syntactic coherence for the cores. From each core we extract all the possible bigrams. This approach considers each link from the complete subgraph made of core-words, and allows us to count the frequency of each part of a speech pair. We find that the most frequent bigram in the structural-based and attribute-based cores is the *noun-noun* pair, 0.42% and 0.20%, respectively. Conversely, the complete subgraph from the conformity-core continues to present a more heterogeneous part of speech composition, where a prominent pair is not observed. In fact, similar frequencies are found for the *verb-noun* (0.1%), the *noun-adjective* (0.06%) and the *verb-adjective* (0.05%) pairs, which are the

Algorithm 1 Graph Walk

Require: Undirected Graph $G = (V, E)$, starting node n

- 1: Initialize weights on E
- 2: Initialize word ordering list T
- 3: Let n be the current visited word and add it to T
- 4: Initialize set of word candidates C
- 5: **while** $len(T) < |V|$ **do**
- 6: **if** C is not empty **then**
- 7: **for** each candidate c in C **do**
- 8: Compute the similarity between n and c
- 9: Let $max(c)$ be the current word
- 10: **if** $max(c)$ not in T **then**
- 11: add c to T
- 12: remove c from C
- 13: **else**
- 14: add randomly a not already learned $v \in V$ to C
- 15: **return** T

three most frequent bigrams in the conformity-core.

D Graph Walks

D.1 Pseudo-code

Algorithm 1 introduces a general schema of a graph walk to describe the four proposed variants. We impose an undirected graph as input, initializing the edge weights as desired (line 1), e.g. all weights are equal to 1 if we want to ignore feature similarity. Then, the word acquisition ordering is initialized, and a randomly selected node n is added to the rank at $t = 1$ (lines 2-3). The set of word candidates is initialized (line 4) to be filled according to the different strategies of each graph walk. Once having a starting word and the first set of candidates, the walk starts until the whole dataset is covered, e.g. each node has a position assigned in the rank (line 5). We iteratively select the new current word by computing the similarity between the current word and the word candidates, choosing the one which maximises similarity (lines 7-9); we re-compute similarity at each iteration because we might need to update this quantity, e.g. when this last one is based on CDI’s availability. The new current word is added to the rank only if it was not already learned, otherwise it is removed from the set of candidate (lines 10-12). If the set is empty (line 6), a randomly chosen and still not learned word is added to it for continuing the walk (lines 13-14); this is equivalent to adding an error ϵ , only when it is strictly necessary.

D.1.1 CDI-based similarity

The CDI-based model relies on pairwise similarity between two words i and j modulated by additional information on the CDIs they belong to, namely $i \in I$ and $j \in J$. Let A be the set of features relative to i and B the features of j . The probability for the random walker to acquire j after i is the factorisation of three terms:

$$P_{i \rightarrow j} = \begin{cases} J(i, j) & \text{iff } g(\text{cdi}(j)) \cdot h(i, \text{cdi}(j)) = 0 \\ J(i, j) \cdot g(j) \cdot h(i, \text{cdi}(j)) & \text{otherwise,} \end{cases}$$

The first term in $P_{i \rightarrow j}$ accounts for the similarity between sets A and B quantified via the Jaccard index, namely $\text{sim}(A, B) = |A \cap B| / |A \cup B|$ or the ratio of elements common to both sets of features A and B . The second term is the target CDI availability, namely the amount of words still available for acquisition in the CDI containing target word j :

$$g(j) = 1 - \frac{|\{w \in J | \text{isactive}\}|}{|J|}.$$

The more words are available for acquisition in the target CDI J , the higher $g(j)$ and thus the probability for the random walker to move to $j \in J$. The third term is the CDI attractiveness:

$$h(i, J) = \frac{|\{j \in \Gamma_i | j \in J\}|}{|\Gamma_i|},$$

where Γ_i is the set of adjacent neighbors of the word available for acquisition i . The more i 's neighbors are within the target CDI J , the higher $h(i, J)$ and thus the probability for the next word to be *attracted* in the CDI where many neighbors are already present.

D.2 Evaluation

Accuracy, precision and recall evaluate the performances of the random walks (cf. Results and Methods). The measures are built upon the confusion matrix of predictions, i.e. a matrix containing the number of correct predictions, namely true positives (TPs) and true negatives (TNs), and the number of incorrect predictions, namely false positives (FPs) and false negatives (FNs). Contextualizing these concepts in this domain, TPs are CDI's words correctly learned by a random walker in a selected AoA bin, while TNs are all other words that a graph walk correctly predict as not belonging to a CDI in that AoA bin; FPs are words that fill a CDI as expected, but they are not the exact same words considered in normative learning, while FNs are CDI's words that are not retrieved in that AoA bin.

Formally, accuracy is the number of TPs divided by the total number of predictions:

$$\text{Accuracy}(CDI, AoA) = \frac{TP}{TP + TN + FP + FN}.$$

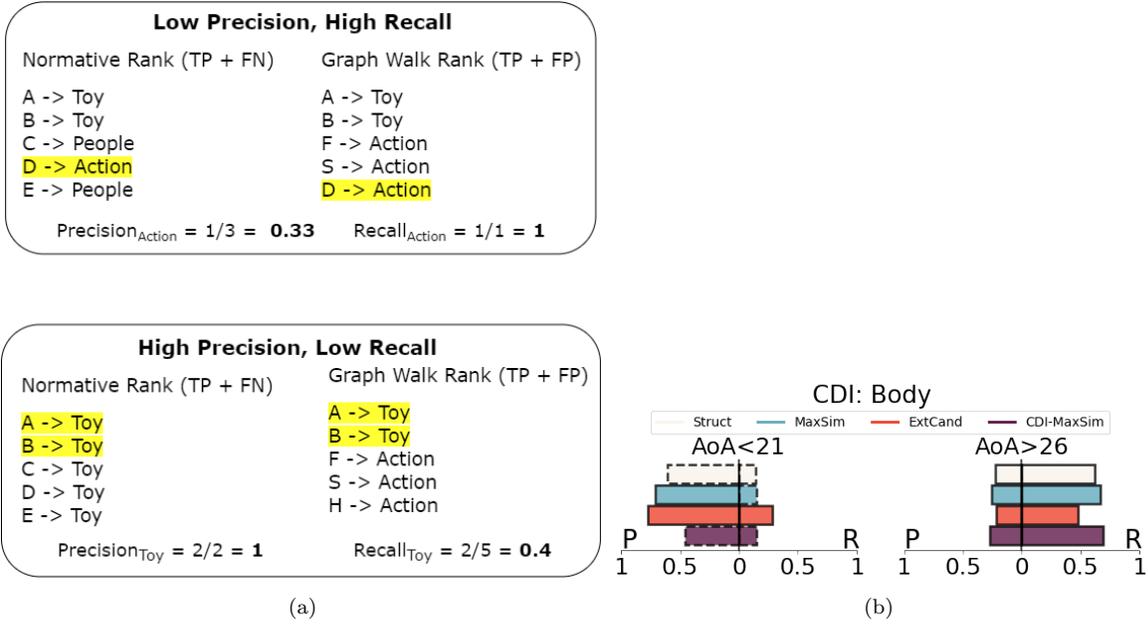


Figure 5: (a): Toy example focusing on the meaning of precision and recall of a CDI; (b): Precision and recall of *Body* CDI at two different stages of acquisition, i.e., $AoA_{\leq 21}$ months, namely the first 100 learned words ($N=100$), and $AoA_{>26}$, namely the last 129 learned words according to CHILDES dataset.

Accuracy can answer poorly to some questions as *how many (expected) CDI's words a graph walk can retrieve in a specific AoA bin?* Precision and recall address this better. Hence, precision is the fraction of relevant elements among the retrieved ones,

$$Precision(CDI, AoA) = \frac{TP}{TP+FP},$$

while recall is the fraction of relevant elements that are retrieved,

$$Recall(CDI, AoA) = \frac{TP}{TP+FN}.$$

Example n. 1. We aim to evaluate the performances of a random walk focusing on how the model is filling the *Animal*-CDI at a very early stage of acquisition, i.e. considering the words learned before 21 months. $Recall(Animal, < 21m)$ increases whatever animal-related word the model retrieves, e.g. *dog* and *frog*, but $Precision(Animal, < 21m)$ does not increase if *frog* is not learned before 21 months. FPs as *frog* are non relevant words; moreover, the model can miss FNs as *cat*, i.e. relevant words learned before 21 months.

Example n. 2. Fig. 5 (a) focuses on two possible extremes, i.e. when precision is low but recall is high (above), or precision is high but recall is low (below). In the first case, the normative learning contains only one *Action* word in the sliced AoA bin of five words, but the graph walk retrieves three *Action* words. Recall is maximised, i.e. the expected word D is retrieved; however, FPs as F and S decrease precision. In the second case, *Toy* words

only fill the sliced normative bin, but the graph walk correctly predicts two expected words out of five. Precision is maximised, i.e. A and B are expected words; however, FNs as C, D and E decrease recall.

Fig. 5 (b) sums up a real example on the CHILDES dataset. Fig. 6 reports the precision-recall bars of the conformity-core for each bin of age of acquisition (cf. Results, focus on the core words learned before 21 months only). Fig. 7 reports the complete precision-recall bars of each CDI and age of acquisition (cf. Results, focus on *Action*, *Body*, *Household*, *People*, *Prepositions* and *Vehicle* CDIs only, acquired before 21 months).

D.2.1 Random learning

We compare the graph walk performances against a learning model that assigns to all words a position in the rank randomly, regardless of any type of structural, vectorial or CDI-based information (cf. Results). In each precision-recall plot, solid bars specify whether the values are statistically significant with respect to this random word assignment. We use the following z-score for the test:

$$z = \frac{M_1 - M_2}{\sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{n}}}$$

where M_1 is the mean precision/recall value from n runs of the walk model, with σ_1 standard deviations, and M_2 is the mean precision/recall value from n runs of the random learning model, with σ_2 standard deviations. Thus, dotted bars indicate values that are not statistically different from the random learning distribution ($z > 0.05$ or precision/recall higher in the random learning than in the random walk model).

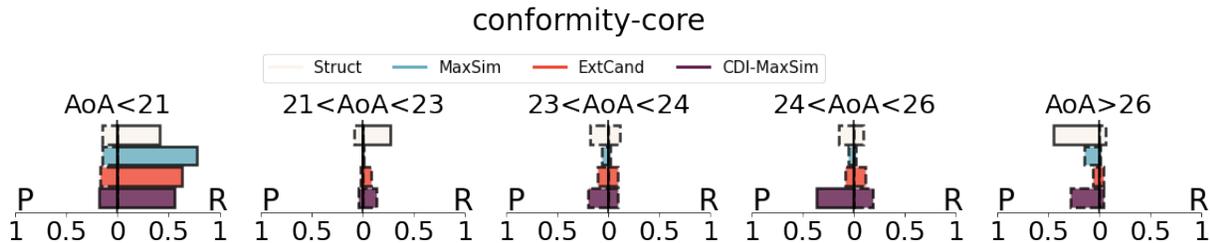


Figure 6: Precision-recall evaluation of the core over all bins of age of acquisition.

Acknowledgments

This work is supported by the European Union – Horizon 2020 Program under the scheme “INFRAIA-01-2018-2019 – Integrating Activities for Advanced Communities”, Grant Agreement n.871042, “SoBigData++: European Integrated Infrastructure for Social Mining and Big Data Analytics” (<http://www.sobigdata.eu>).

References

- [1] Jean Aitchison. *Words in the mind: An introduction to the mental lexicon*. John Wiley & Sons, 2012.
- [2] Oriol Artime and Manlio De Domenico. “Percolation on feature-enriched interconnected systems”. In: *Nature communications* 12.1 (2021), pp. 1–12.
- [3] Aleix Bassolas and Vincenzo Nicosia. “First-passage times to quantify and compare structural correlations and heterogeneity in complex systems”. In: *Communications Physics* 4.1 (2021), pp. 1–14.
- [4] Friedrich Beck. “Mind, brain, and dualism in modern physics”. In: *Psycho-Physical Dualism Today: An Interdisciplinary Approach*, New York: Rowman & Littlefield (2008), pp. 69–97.
- [5] Nicole Beckage, Linda Smith, and Thomas Hills. “Small worlds and semantic network growth in typical and late talkers”. In: *PloS one* 6.5 (2011), e19348.
- [6] Nicole M Beckage and Eliana Colunga. “Language networks as models of cognition: Understanding cognition through language”. In: *Towards a Theoretical Framework for Analyzing Complex Linguistic Networks*. Springer, 2016, pp. 3–28.
- [7] Nicole M Beckage and Eliana Colunga. “Network growth modeling to capture individual lexical learning”. In: *Complexity* 2019 ().
- [8] Eszter Bokányi et al. “Universal patterns of long-distance commuting and social assortativity in cities”. In: *Scientific reports* 11.1 (2021), pp. 1–10.
- [9] Gemma Boleda. “Distributional semantics and linguistic theory”. In: *Annual Review of Linguistics* 6 (2020), pp. 213–234.
- [10] Ramon Ferrer I Cancho and Richard V Solé. “The small world of human language”. In: *Proceedings of the Royal Society of London. Series B: Biological Sciences* 268.1482 (2001), pp. 2261–2265.
- [11] Matthew T Carlson, Morgan Sonderegger, and Max Bane. “How children explore the phonological network in child-directed speech: A survival analysis of children’s first word productions”. In: *Journal of memory and language* 75 (2014), pp. 159–180.
- [12] Bernardino Casas et al. “The polysemy of the words that children learn over time”. In: *Interaction Studies* 19.3 (2018), pp. 389–426.
- [13] Nichol Castro and Cynthia SQ Siew. “Contributions of modern network science to the cognitive sciences: revisiting research spirals of representation and process”. In: *Proceedings of the Royal Society A* 476.2238 (2020), p. 20190825.

-
- [14] Lucas M Chang and Gedeon O Deák. “Adjacent and Non-Adjacent Word Contexts Both Predict Age of Acquisition of English Words: A Distributional Corpus Analysis of Child-Directed Speech”. In: *Cognitive Science* 44.11 (2020), e12899.
- [15] Salvatore Citraro and Giulio Rossetti. “Identifying and exploiting homogeneous communities in labeled networks”. In: *Applied Network Science* 5.1 (2020), pp. 1–20.
- [16] Elizabeth M Clerkin et al. “Real-world visual statistics and infants’ first-learned object names”. In: *Philosophical Transactions of the Royal Society B: Biological Sciences* 372.1711 (2017), p. 20160055.
- [17] Allan M Collins and Elizabeth F Loftus. “A spreading-activation theory of semantic processing.” In: *Psychological review* 82.6 (1975), p. 407.
- [18] Cesar H Comin et al. “Complex systems: Features, similarity and connectivity”. In: *Physics Reports* 861 (2020), pp. 1–41.
- [19] Simon De Deyne and Gert Storms. “Word associations: Network and semantic properties”. In: *Behavior research methods* 40.1 (2008), pp. 213–231.
- [20] Simon De Deyne et al. “The “Small World of Words” English word association norms for over 12,000 cue words”. In: *Behavior research methods* 51.3 (2019), pp. 987–1006.
- [21] Jeffrey L Elman. “An alternative view of the mental lexicon”. In: *Trends in cognitive sciences* 8.7 (2004), pp. 301–306.
- [22] Larry Fenson et al. *MacArthur-Bates communicative development inventories*. Paul H. Brookes Publishing Company Baltimore, MD, 2007.
- [23] Ramon Ferrer-i-Cancho and Michael S Vitevitch. “The origins of Zipf’s meaning-frequency law”. In: *Journal of the Association for Information Science and Technology* 69.11 (2018), pp. 1369–1379.
- [24] Thomas L Griffiths, Mark Steyvers, and Alana Firl. “Google and the mind: Predicting fluency with PageRank”. In: *Psychological science* 18.12 (2007), pp. 1069–1076.
- [25] Fritz Günther, Luca Rinaldi, and Marco Marelli. “Vector-space models of semantic representation from a cognitive perspective: A discussion of common misconceptions”. In: *Perspectives on Psychological Science* 14.6 (2019), pp. 1006–1033.
- [26] Pamela A Hadley, Megan M McKenna, and Matthew Rispoli. “Sentence diversity in early language development: Recommendations for target selection and progress monitoring”. In: *American journal of speech-language pathology* 27.2 (2018), pp. 553–565.

-
- [27] Thomas T Hills and Yoed N Kenett. “Networks of the Mind: How Can Network Science Elucidate Our Understanding of Cognition?” In: *Topics in Cognitive Science* ().
- [28] Thomas T Hills and Cynthia SQ Siew. “Filling gaps in early word learning”. In: *Nature Human Behaviour* 2.9 (2018), p. 622.
- [29] Thomas T Hills et al. “Longitudinal analysis of early semantic networks: Preferential attachment or preferential acquisition?” In: *Psychological science* 20.6 (2009), pp. 729–739.
- [30] Petter Holme. “Core-periphery organization of complex networks”. In: *Physical Review E* 72.4 (2005), p. 046111.
- [31] Zhexue Huang. “Clustering large data sets with mixed numeric and categorical values”. In: *Proceedings of the 1st pacific-asia conference on knowledge discovery and data mining,(PAKDD)*. Citeseer. 1997, pp. 21–34.
- [32] Joshua Jackson et al. “From text to thought: How analyzing language can advance psychological science”. In: *Perspectives on Psychological Science* (2021).
- [33] Yoed N Kenett. “What can quantitative measures of semantic distance tell us about creativity?” In: *Current Opinion in Behavioral Sciences* 27 (2019), pp. 11–16.
- [34] Abhilasha A Kumar, Mark Steyvers, and David A Balota. “A Critical Review of Network-Based and Distributional Approaches to Semantic Memory Structure and Processes”. In: *Topics in Cognitive Science* (2021).
- [35] Thomas K Landauer and Susan T Dumais. “A solution to Plato’s problem: The latent semantic analysis theory of acquisition, induction, and representation of knowledge.” In: *Psychological review* 104.2 (1997), p. 211.
- [36] Alessandro Lenci. “Distributional models of word meaning”. In: *Annual review of Linguistics* 4 (2018), pp. 151–171.
- [37] Orr Levy et al. “Unveiling the nature of interaction between semantics and phonology in lexical access based on multilayer networks”. In: *Scientific reports* 11.1 (2021), pp. 1–14.
- [38] HaiTao Liu. “Statistical properties of Chinese semantic networks”. In: *Chinese Science Bulletin* 54.16 (2009), pp. 2781–2785.
- [39] Kevin Lund and Curt Burgess. “Producing high-dimensional semantic spaces from lexical co-occurrence”. In: *Behavior research methods, instruments, & computers* 28.2 (1996), pp. 203–208.

-
- [40] James MacQueen et al. “Some methods for classification and analysis of multivariate observations”. In: *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*. Vol. 1. 14. Oakland, CA, USA. 1967, pp. 281–297.
- [41] Brian MacWhinney. *The CHILDES project: The database*. Vol. 2. Psychology Press, 2000.
- [42] Brian MacWhinney. *The CHILDES project: Tools for analyzing talk, Volume II: The database*. Psychology Press, 2014.
- [43] Miller McPherson, Lynn Smith-Lovin, and James M Cook. “Birds of a feather: Homophily in social networks”. In: *Annual review of sociology* (2001).
- [44] Ken McRae et al. “Semantic feature production norms for a large set of living and nonliving things”. In: *Behavior research methods* 37.4 (2005), pp. 547–559.
- [45] George A Miller. *WordNet: An electronic lexical database*. MIT press, 1998.
- [46] Michael Molloy et al. “A critical point for random graphs with a given degree sequence”. In: *The Structure and Dynamics of Networks*. Princeton University Press, 2011, pp. 240–258.
- [47] Mark Newman. *Networks*. Oxford university press, 2018.
- [48] Mark EJ Newman. “Mixing patterns in networks”. In: *Physical review E* 67.2 (2003), p. 026126.
- [49] Leto Peel, Jean-Charles Delvenne, and Renaud Lambiotte. “Multiscale mixing patterns in networks”. In: *Proceedings of the National Academy of Sciences* 115.16 (2018), pp. 4057–4062.
- [50] Jeffrey Pennington, Richard Socher, and Christopher D Manning. “Glove: Global vectors for word representation”. In: *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 2014, pp. 1532–1543.
- [51] Jan Pepper and Elaine Weitzman. *It takes two to talk: A practical guide for parents of children with language delays*. The Hanen Centre, 2004.
- [52] Ron Pomper and Jenny R Saffran. “Familiar object salience affects novel word learning”. In: *Child development* 90.2 (2019), e246–e262.
- [53] M Ross Quillian. “Word concepts: A theory and simulation of some basic semantic capabilities”. In: *Behavioral science* 12.5 (1967), pp. 410–430.
- [54] Erzsébet Ravasz and Albert-László Barabási. “Hierarchical organization in complex networks”. In: *Physical review E* 67.2 (2003), p. 026112.

-
- [55] Alexa R Romberg and Jenny R Saffran. “Statistical learning and language acquisition”. In: *Wiley Interdisciplinary Reviews: Cognitive Science* 1.6 (2010), pp. 906–914.
- [56] Giulio Rossetti, Salvatore Citraro, and Letizia Milli. “Conformity: A path-aware homophily measure for node-attributed networks”. In: *IEEE Intelligent Systems* 36.1 (2021), pp. 25–34.
- [57] Giulio Rossetti, Letizia Milli, and Rémy Cazabet. “CDLIB: a python library to extract, compare and evaluate communities from complex networks”. In: *Applied Network Science* 4.1 (2019), pp. 1–26.
- [58] Cynthia SQ Siew. “Global and Local Feature Distinctiveness Effects in Language Acquisition”. In: *Cognitive Science* 45.7 (2021), e13008.
- [59] Cynthia SQ Siew et al. “Cognitive network science: A review of research on cognition through the lens of network representations, processes, and dynamics”. In: *Complexity* 2019 ().
- [60] Serene Siow and Kim Plunkett. “Exploring the variable effects of frequency and semantic diversity as predictors for a word’s ease of acquisition in different word classes”. In: *Proceedings of the Annual Meeting of the Cognitive Science Society*. Vol. 43. 43. 2021.
- [61] Ann E Sizemore et al. “Knowledge gaps in the early growth of semantic feature networks”. In: *Nature human behaviour* 2.9 (2018), pp. 682–692.
- [62] Massimo Stella. “Modelling early word acquisition through multiplex lexical networks and machine learning”. In: *Big Data and Cognitive Computing* 3.1 (2019), p. 10.
- [63] Massimo Stella, Nicole M Beckage, and Markus Brede. “Multiplex lexical networks reveal patterns in early word acquisition in children”. In: *Scientific Reports* 7 (2017), p. 46730.
- [64] Massimo Stella and Markus Brede. “Patterns in the English language: phonological networks, percolation and assembly models”. In: *Journal of Statistical Mechanics: Theory and Experiment* 2015.5 (2015), P05006.
- [65] Massimo Stella et al. “Multiplex model of mental lexicon reveals explosive learning in humans”. In: *Scientific reports* 8.1 (2018), pp. 1–11.
- [66] Mark Steyvers and Joshua B Tenenbaum. “The large-scale structure of semantic networks: Statistical analyses and a model of semantic growth”. In: *Cognitive science* 29.1 (2005), pp. 41–78.
- [67] Akira Utsumi. “A complex network approach to distributional semantic models”. In: *PloS one* 10.8 (2015), e0136277.
- [68] Olga Valba and Alexander Gorsky. “K-clique percolation in free association networks. The mechanism behind the 7 ± 2 law?” In: *arXiv preprint arXiv:2110.09317* (2021).

-
- [69] Bram Van Rensbergen, Gert Storms, and Simon De Deyne. “Examining assortativity in the mental lexicon: Evidence from word associations”. In: *Psychonomic bulletin & review* 22.6 (2015), pp. 1717–1724.
- [70] Alexander Veremyev et al. “Graph-based exploration and clustering analysis of semantic spaces”. In: *Applied Network Science* 4.1 (2019), pp. 1–26.
- [71] Michael S Vitevitch. “Can network science connect mind, brain, and behavior?” In: *Network science in cognitive psychology*. Routledge, 2019, pp. 184–197.
- [72] Michael S Vitevitch. “What can graph theory tell us about word learning and lexical retrieval?” In: (2008).
- [73] Michael S Vitevitch et al. “Using complex networks to understand the mental lexicon”. In: *Yearbook of the Poznan Linguistic Meeting*. Vol. 1. 1. Sciendo. 2014, pp. 119–138.
- [74] Christopher Whelan, Greg Harrell, and Jin Wang. “Understanding the k-medians problem”. In: *Proceedings of the International Conference on Scientific Computing (CSC)*. The Steering Committee of The World Congress in Computer Science, Computer . . . 2015, p. 219.
- [75] Hiroto Yamaguchi et al. “Controlling Internal Structure of Communities on Graph Generator”. In: *2020 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*. IEEE. 2020, pp. 937–940.
- [76] George Kingsley Zipf. *Human behavior and the principle of least effort: An introduction to human ecology*. Ravenio Books, 2016.
- [77] Michael Zock. “Words in Books, Computers and the Human Mind”. In: *Journal of Cognitive Science* 16.4 (2015), pp. 355–378.

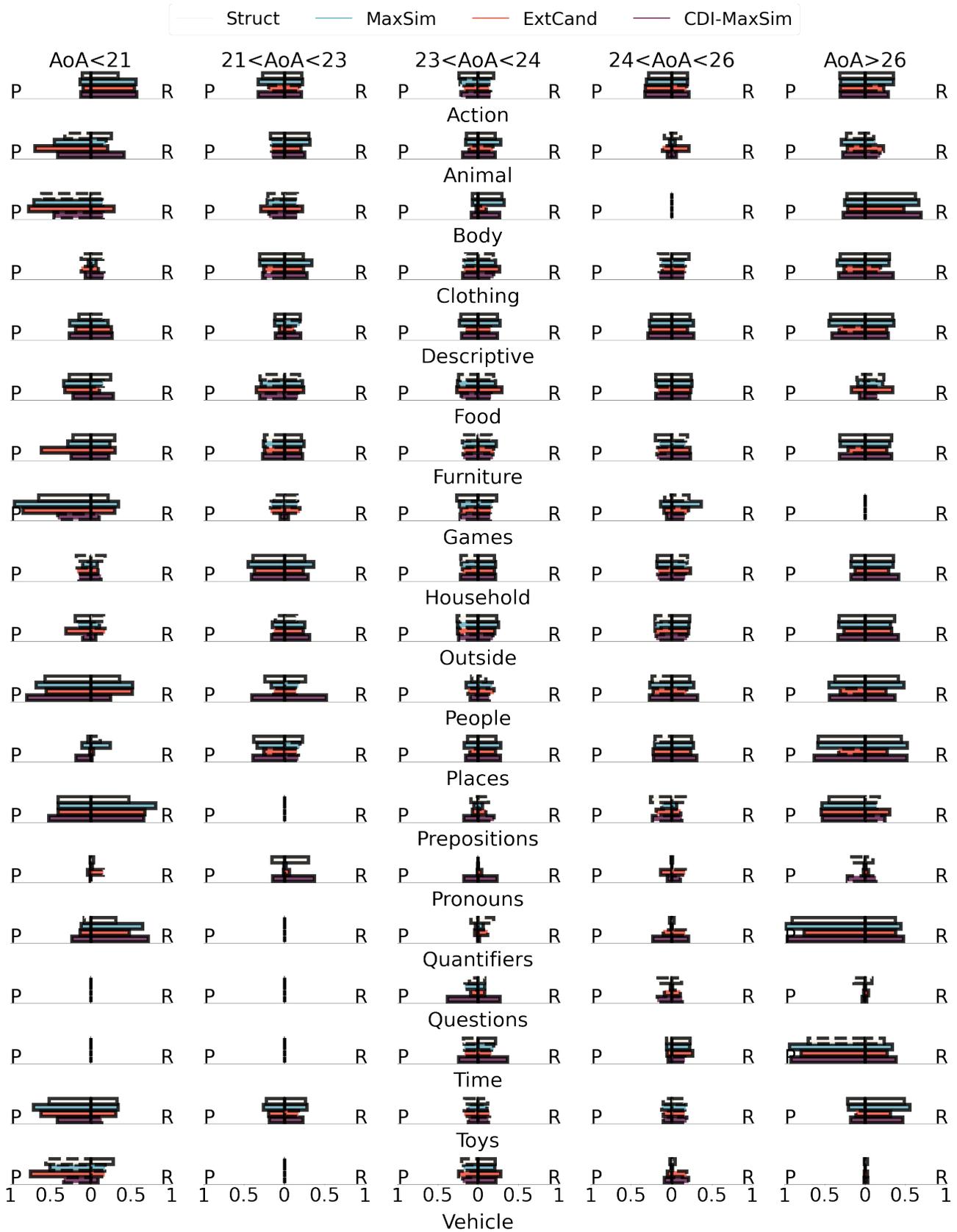


Figure 7: Precision-recall evaluation of all CDIs over all bins of age of acquisition.