

## Distributional Measures of Semantic Abstraction

Sabine Schulte im Walde<sup>1,\*</sup> and Diego Frassinelli<sup>2</sup>

<sup>1</sup>*Institute for Natural Language Processing, University of Stuttgart, Germany*

<sup>2</sup>*Department of Linguistics, University of Konstanz, Germany*

Correspondence\*:

Sabine Schulte im Walde

schulte@ims.uni-stuttgart.de

### 2 ABSTRACT

3 This article provides an in-depth study of distributional measures for distinguishing between  
4 degrees of *semantic abstraction*. Abstraction is considered a "central construct in cognitive  
5 science" (Barsalou, 2003) and a "process of information reduction that allows for efficient storage  
6 and retrieval of central knowledge" (Burgoon *et al.*, 2013). Relying on the distributional hypothesis  
7 (Harris, 1954; Firth, 1957), computational studies have successfully exploited measures of  
8 contextual co-occurrence and neighbourhood density to distinguish between conceptual semantic  
9 categorisations. So far, these studies have modeled semantic abstraction across lexical-semantic  
10 tasks such as ambiguity; diachronic meaning changes; abstractness vs. concreteness; and  
11 hypernymy (Sagi *et al.*, 2009; Hoffman *et al.*, 2013; Santus *et al.*, 2014; Schlechtweg *et al.*,  
12 2017; Naumann *et al.*, 2018; Frassinelli *et al.*, 2017). Yet, the distributional approaches target  
13 different conceptual types of semantic relatedness, and as to our knowledge not much attention  
14 has been paid to apply, compare or analyse the computational abstraction measures across  
15 conceptual tasks. The current article suggests a novel perspective that exploits variants of  
16 distributional measures to investigate semantic abstraction in English in terms of the abstract–  
17 concrete dichotomy (e.g., *glory–banana*) and in terms of the generality–specificity distinction  
18 (e.g., *animal–fish*), in order to compare the strengths and weaknesses of the measures regarding  
19 categorisations of abstraction, and to determine and investigate conceptual differences.

20 In a series of experiments we identify reliable distributional measures for both instantiations of  
21 lexical-semantic abstraction and reach a precision higher than 0.7, but the measures clearly differ  
22 for the abstract–concrete vs. abstract–specific distinctions and for nouns vs. verbs. Overall,  
23 we identify two groups of measures, (i) frequency and word entropy when distinguishing  
24 between more and less abstract words in terms of the generality–specificity distinction, and  
25 (ii) neighbourhood density variants (especially target–context diversity) when distinguishing  
26 between more and less abstract words in terms of the abstract–concrete dichotomy. We conclude  
27 that more general words are used more often and are less surprising than more specific words,  
28 and that abstract words establish themselves empirically in semantically more diverse contexts  
29 than concrete words. Finally, our experiments once more point out that distributional models of  
30 conceptual categorisations need to take word classes and ambiguity into account: results for  
31 nouns vs. verbs differ in many respects, and ambiguity hinders fine-tuning empirical observations.

32 **Keywords:** lexical-semantic abstraction, abstractness, concreteness, generality, specificity, hypernymy, vector spaces

## 1 INTRODUCTION

33 Over the years, interdisciplinary research on lexical semantics has seen multiple definitions of conceptual  
34 abstraction. For example, Barsalou (2003) considers abstraction as a “*central construct in cognitive*  
35 *science*” regarding categorical organisation in memory, and distinguishes between various types of  
36 abstraction. Burgoon *et al.* (2013) provide an extensive list and descriptions of past definitions of abstraction  
37 across research fields and research studies, and summarise the common core of abstraction types as  
38 “*a process of information reduction that allows for efficient storage and retrieval of central knowledge*  
39 *(e.g., categorization)*”. Among the various types of abstraction described by Barsalou (2003) and Burgoon  
40 *et al.* (2013), we find two types that have repeatedly been connected to each other across disciplines, i.e.,  
41 abstraction in terms of the abstract–concrete dichotomy (e.g., *glory* is more abstract than *banana*), and  
42 abstraction in terms of the generality–specificity distinction (e.g., *animal* is more abstract than *fish*). For  
43 example, one of the earliest datasets that collected abstractness ratings generated by humans was performed  
44 by Spreen and Schulz (1966), who in turn exploited two previously suggested tasks for abstractness  
45 ratings on a scale, to quantify abstractness (a) in contrast to concreteness in the sense of “*not perceived*  
46 *through senses*”, and (b) in contrast to specificity in the sense of “*general, generic*”. While the sense  
47 perception in task (a) was adopted as the standard task for collecting abstractness ratings in the following  
48 decades, these two categorisations demonstrate alternative instantiations of semantic abstraction, which  
49 were once more targeted in recent empirical studies. Theijssen *et al.* (2011) investigated annotations  
50 regarding (a) vs. (b) for noun senses in a corpus and for noun labels in dative alternations, and Bolognesi  
51 *et al.* (2020) correlated degrees of abstraction in collections of human-annotated concreteness vs. generality.  
52 Both studies were performed for English nouns and relied on existing norms of concreteness ratings  
53 (Coltheart, 1981; Brysbaert *et al.*, 2014, respectively) and the hierarchical organisation of hypernymy in  
54 WordNet (Miller and Fellbaum, 1991; Fellbaum, 1998b).

55 In a similar manner but with yet different distinctions, we also find various instantiations of abstraction  
56 across sub-fields of computational lexical-semantic research. Relying on the distributional hypothesis  
57 that words which are similar in meaning also occur in similar linguistic distributions (Harris, 1954;  
58 Firth, 1957), these studies successfully exploited distributional measures of contextual co-occurrence and  
59 neighbourhood density to distinguish between conceptual semantic categorisations. For example, Sagi  
60 *et al.* (2009) applied a measure of neighbourhood density to quantify diachronic lexical semantic change;  
61 Hoffman *et al.* (2013) proposed semantic diversity as a measure of lexical semantic ambiguity; Santus *et al.*  
62 (2014) utilised the information-theoretic measure entropy to distinguish hypernyms from their hyponyms;  
63 Frassinelli *et al.* (2017) and Naumann *et al.* (2018) applied variants of neighbourhood density and entropy  
64 to distinguish between abstract and concrete words. While these studies address different lexical-semantic  
65 tasks, all tasks have in common that they involve and model some notion of semantic abstraction, i.e.,  
66 diachronic innovative and reductive meaning change; lexical ambiguity; abstractness vs. concreteness in  
67 word meaning; and hypernymy. Yet, as to our knowledge, not much attention has been paid to the shared  
68 common meta-level task of quantifying abstraction across computational approaches, except for Rimell  
69 (2014) and Schlechtweg *et al.* (2017) using hypernymy measures for semantic entailment and diachronic  
70 change, respectively. Furthermore, a closer look into distributional neighbourhood variants reveals that the  
71 types of applied neighbourhoods are conceptually different, exploiting similarity between context words  
72 (Sagi *et al.*, 2009; Hoffman *et al.*, 2013; Naumann *et al.*, 2018) vs. exploiting similarity between nearest  
73 neighbours (Frassinelli *et al.*, 2017). In sum, most researchers involved in the respective sub-fields are  
74 not necessarily aware of each other, such that up to now we do not find a comprehensive application and  
75 comparison of distributional abstraction measures across semantic abstraction tasks.

76 The current article aims to fill this critical gap and provides a series of empirical studies that investigate  
77 conceptual categories of abstraction through variants of distributional measures. Focusing on the two types  
78 of abstraction originally suggested by Spreen and Schulz (1966), and brought back to attention by Theijssen  
79 *et al.* (2011) and Bolognesi *et al.* (2020), we distinguish abstraction in terms of the abstract–concrete  
80 dichotomy and in terms of the generality–specificity distinction. More specifically, we apply a selection of  
81 distributional measures to distinguish between English (i) abstract and concrete words and (ii) hypernyms  
82 and their hyponyms. As resources for our target words, we rely on the concreteness ratings in Brysbaert  
83 *et al.* (2014) and hypernymy relations in *WordNet* (Fellbaum, 1998b). Furthermore, we distinguish between  
84 noun and verb targets, given that lexical representations of word classes differ in their semantic abstraction  
85 regarding both concreteness and hypernymy (Miller and Fellbaum, 1991; Frassinelli and Schulte im Walde,  
86 2019; Schulte im Walde, 2020). The specific measures we apply are variants of neighbourhood densities  
87 (context-based and neighbour-based), the distributional inclusion measure *WeedsPrec* (Weeds *et al.*, 2014)  
88 and the information-theoretic measure *entropy* (Santus *et al.*, 2014; Shwartz *et al.*, 2017). The underlying  
89 distributional vector spaces are induced from the ENCOW web corpus (Schäfer and Bildhauer, 2012).

90 Overall, we thus suggest a novel perspective that brings together and effectively exploits empirical  
91 computational measures across two types of lexical-semantic abstraction. In this way, our studies enable  
92 us to compare the strengths and weaknesses of the distributional measures regarding categorisations of  
93 abstraction, and to determine and investigate conceptual differences as captured by the measures. In the  
94 remainder of this article, Section 2 introduces previous research perspectives and studies on the two types  
95 of semantic abstraction we focus on, both from a cognitive and from a computational perspective. Section 3  
96 then describes the data and methods we use in our study, before Section 4 provides the actual experiments  
97 and results which are then discussed in Section 5.

## 2 RELATED WORK

98 In the following, we introduce previous research perspectives and studies on the two types of semantic  
99 abstraction we focus on, i.e., abstraction in terms of the abstract–concrete dichotomy and in terms of  
100 the generality–specificity distinction. In this vein, Section 2.1 looks into abstraction from a cognitive  
101 perspective, while Section 2.2 provides an overview of computational models of abstraction. In Section 2.3  
102 we describe previous empirical investigations across the two types of abstraction. From a terminological  
103 perspective, we will use the word “concepts” when referring to mental representations, and “words” when  
104 referring to the corresponding linguistic surface forms humans are exposed to. Given the distributional  
105 nature of our studies, we will always refer to words as the targets of our analyses.

### 106 2.1 Cognitive Perspectives on Abstraction

107 Barsalou (2003) considers abstraction as a “*central construct in cognitive science*” regarding the  
108 organization of categories in the human memory. He attributes six different senses to abstraction:  
109 (i) abstracting a conceptual category from the settings it occurs in; (ii) generalising across category  
110 members; (iii) generalising through summary representations which are necessary for the behavioural  
111 generalisations in (ii); (iv) sparse schematic representations; (v) flexible interpretation; and (vi) abstractness  
112 in contrast to concreteness. Barsalou’s classification illustrates that the term “semantic abstraction” as  
113 well as its featural and inferential implications for memory representations are vague in that different  
114 instantiations go along with different representations; he himself focuses on summary representations (iii).

115 Burgoon *et al.* (2013) provide an extensive list and description of past definitions of abstraction across  
116 research fields and research studies, and state that, at the meta level, the term abstraction is referred to  
117 as “*a process of information reduction that allows for efficient storage and retrieval of central knowledge*  
118 (*e.g., categorization*)”. For their own study, they define abstraction as “*as a process of identifying a set*  
119 *of invariant central characteristics of a thing*”, and in what follows they compare existing definitions of  
120 abstraction regarding their roots, developments, antecedents, consequences, and methods for studying.

121 The distinction of the two abstraction types adopted in the current study comes from Spreen and Schulz  
122 (1966) indicating that the “*definition of abstractness or concreteness in previous studies shows that at*  
123 *least two distinctly different interpretations can be made*”, and pointing back to previous collections with  
124 judgements on generality by Gorman (1961) and judgements on concreteness as well as generality by  
125 Darley *et al.* (1959). Spreen and Schulz themselves collected ratings on both abstractness–concreteness and  
126 abstractness–specificity (among others) for 329 English nouns, and found a correlation of 0.626 between  
127 the ratings of the two abstraction variables. The two-fold distinction of abstraction outlined in the work by  
128 Spreen and Schulz (1966) is also included in the various instantiations of abstraction in Barsalou (2003)  
129 and Burgoon *et al.* (2013). In the following, we describe the lines of research involved in the representation  
130 and processing of abstract vs. concrete concepts and then those involved in general vs. specific concepts.

### 131 2.1.1 Abstract vs. Concrete Concepts

132 The most influential proposal about the processing, storing and comprehension of abstract concepts in  
133 contrast to concrete concepts can be traced back to Paivio (1971). He suggested the *dual-route theory* where  
134 a verbal system is primarily responsible for language aspects of linguistic units (such as words), while a  
135 non-verbal system, in particular imagery, is primarily responsible for sensory-motor aspects. Even though,  
136 in the meantime, a range of alternative as well as complementary theories have been suggested, Paivio’s  
137 theory offers an explanation why concrete concepts (which are supposedly accessed via both routes) are  
138 generally processed faster in lexical memory than abstract concepts (which are supposedly accessed only  
139 via the non-verbal system) across tasks and datasets, cf. Pecher *et al.* (2011) and Borghi *et al.* (2017) for  
140 comprehensive overviews.

141 Further than the dual-route theory, cognitive scientists have investigated other dimensions of abstractness.  
142 Most notably, Schwanenflugel and Shoben (1983) suggested the *context availability theory* where they  
143 compared the processing of abstract and concrete words in context and demonstrated that in appropriate  
144 contexts neither reading times nor lexical decision times differ, thus emphasising the role of context in  
145 conditions of abstractness. In addition, a number of properties have been pointed out where abstract and  
146 concrete concepts differ. (i) There is a strong consensus and experimental confirmation that concrete  
147 concepts are more *imaginable* than the abstract ones, and that it takes longer to generate images for abstract  
148 than for concrete concepts (Paivio *et al.*, 1968; Paivio, 1971; Paivio and Begg, 1971, i.a.). (ii) Abstract  
149 concepts are considered to be more *emotionally valenced* than concrete concepts (Kousta *et al.*, 2011;  
150 Vigliocco *et al.*, 2014; Pollock, 2018). (iii) *Free associations* to abstract concepts are assumed to differ from  
151 free associations to concrete concepts in terms of the number of types, but at the same time associations to  
152 concrete concepts have been found weaker and more symmetric than for abstract concepts (Crutch and  
153 Warrington, 2010; Hill *et al.*, 2014). (iv) Based on a *feature generation task*, features of abstract concepts  
154 are less property- and more situation-related than features of concrete words (Wiemer-Hastings and Xu,  
155 2005). (v) Accordingly, an appropriate embedding into *situations* has been identified as crucial for abstract  
156 vs. concrete meaning representations (Barsalou and Wiemer-Hastings, 2005; Hare *et al.*, 2009; Pecher  
157 *et al.*, 2011; Frassinelli and Lenci, 2012; Recchia and Jones, 2012).

158 Hand in hand with defining and investigating hypotheses about dimensions of abstract and concrete  
159 concepts, a number of data collections have been created. To name just a prominent subset of the  
160 large number of existing resources, Spreen and Schulz (1966) collected ratings of concreteness and  
161 specificity (among others) for 329 English nouns (see above); Paivio *et al.* (1968) collected ratings for  
162 925 English nouns on concreteness, imagery and meaningfulness; Coltheart (1981) put together the  
163 *MRC Psycholinguistic Database*, mostly comprising pre-existing information for almost 100,000 English  
164 words including concreteness, imageability, familiarity as well as frequency, semantic, syntactic and  
165 phonological information; Warriner *et al.* (2013) extended the *ANEW* norms from Bradley and Lang  
166 (1999) with 1,034 English words to almost 14,000, capturing emotion-relevant norms of valence, arousal  
167 and dominance; a similar collection for 20,000 English words regarding the same variables but using  
168 best–worst scaling instead of ratings has been done by Mohammad (2018); Brysbaert *et al.* (2014) created  
169 the so far largest human-generated collection containing concreteness ratings for 40,000 English words.  
170 The work by Connell and Lynott differs slightly on the variable depth, by focusing on the individual  
171 perception modalities and interoception (Lynott and Connell, 2009, 2013; Lynott *et al.*, 2020). While the  
172 vast amount of abstractness/concreteness datasets has been created for English, we also find collections for  
173 other languages, such as those for 2,654/1,000 nouns in German (Lahl *et al.*, 2009; Kanske and Kotz, 2010,  
174 respectively); 16,109 Spanish words (Algarabel *et al.*, 1988); 417 Italian words (Della Rosa *et al.*, 2010);  
175 and 1,659 French words (Bonin *et al.*, 2018). While traditional collections have been pen-and-paper-based,  
176 the collections from the last decade have moved towards crowd-sourcing platforms. As alternative to  
177 human-generated ratings, previous research suggested semi-automatic algorithms to create large-scale  
178 norms (Mandera *et al.*, 2015; Recchia and Louwerse, 2015; Köper and Schulte im Walde, 2016; Köper and  
179 Schulte im Walde, 2017; Aedmaa *et al.*, 2018; Rabinovich *et al.*, 2018).

## 180 2.1.2 General vs. Specific Concepts

181 Differently to the above distinction of semantic abstraction in terms of degrees of concreteness as opposed  
182 to abstractness, where concepts may be judged more or less abstract in comparison to otherwise semantically  
183 unrelated concepts (e.g., *banana–glory*), semantic abstraction in terms of generality is typically established  
184 in contrast to a semantically related concept (e.g., *animal–fish*). The lexical-semantic relation of interest  
185 here is hypernymy, where the more general concept represents the hypernym of the more specific hyponym.

186 An enormous body of work discusses hypernymy next to further semantic relations in the mental lexicon.  
187 For example, a seminal description of lexical relations can be found in Cruse (1986), who states that  
188 lexical relations “*reflect the way infinitely and continuously varied experienced reality is apprehended and*  
189 *controlled through being categorised, subcategorised and graded along specific dimensions of variation*”.  
190 Murphy (2003) focuses on the representation of semantic relations in the lexicon and discusses synonymy,  
191 antonymy, contrast, hyponymy and meronymy, across word classes. Most of her discussions concern  
192 linguistic vs. meta-linguistic representations of relations, reference of relations to words vs. concepts, and  
193 lexicon storage. The most extensive resource that systematically explores and compares types of lexical-  
194 semantic relations across word classes is established by the taxonomy of the Princeton *WordNet*, where  
195 hypernymy represents a key organisation principle of semantic memory (Fellbaum, 1990; Gross and Miller,  
196 1990; Miller *et al.*, 1990). Miller and Fellbaum (1991) provide a meta-level summary of relational structures  
197 and decisions. As basis for the *WordNet* organisation, they state that “*the mental lexicon is organised*  
198 *by semantic relations. Since a semantic relation is a relation between meanings, and since meanings*  
199 *can be represented by synsets, it is natural to think of semantic relations as pointers between synsets*”.  
200 The semantic relations in *WordNet* include the paradigmatic relations synonymy, hypernymy/hyponymy,

201 antonymy, and meronymy. For nouns, WordNet implements a hierarchical organisation of synsets (i.e., sets  
202 of synonymous word meanings) relying on hypernymy relations. Verbs are considered the most complex  
203 and polysemous word class; they are organised on a verb-specific variant of hypernymy, i.e., *troponymy*:  $v_1$   
204 *is to*  $v_2$  *in some manner*, that operates on semantic fields instantiated through synsets. Troponymy itself is  
205 conditioned on entailment and temporal inclusion.

## 206 2.2 Computational Models of Abstraction

207 Across both types of semantic abstraction, computational models have been suggested to automatically  
208 characterise or distinguish between more and less abstract words. They have been intertwined with cognitive  
209 perspectives to various degrees.

### 210 2.2.1 Abstract vs. Concrete Words

211 A common idea in this research direction is the exploitation of corpus-based co-occurrence information  
212 to infer textual distributional characteristics of cognitive semantic variables, including abstractness as  
213 well as further variables such as emotion, imageability, familiarity, etc. These models are large-scale  
214 data approaches to explore the role of linguistic information and textual attributes when distinguishing  
215 between abstract and concrete words. A subset of these distributional approaches is strongly driven by  
216 a cognitive perspective, thus aiming to explain the organisation of human semantic memory and lexical  
217 processing effects by the contribution of linguistic attributes. Common techniques for organising the  
218 textual information are semantic vector spaces such as Latent Semantic Analysis (LSA) (Salton *et al.*,  
219 1975), the Hyperspace Analogue to Language (HAL) (Burgess, 1998), and more recent variants of  
220 standard Distributional Semantic Models (DSMs) (Baroni and Lenci, 2010; Turney and Pantel, 2010), in  
221 combination with measures of distributional similarity and clustering approaches (Glenberg and Robertson,  
222 2000; Vigliocco *et al.*, 2009; Bestgen and Vincze, 2012; Troche *et al.*, 2014; Mander *et al.*, 2015; Recchia  
223 and Louwerse, 2015; Lenci *et al.*, 2018). Finally, our own studies provide preliminary insights into co-  
224 occurrence characteristics of abstract and concrete words with respect to linguistic parameters such as  
225 window size, parts-of-speech and subcategorisation conditions (Frassinelli *et al.*, 2017; Naumann *et al.*,  
226 2018; Frassinelli and Schulte im Walde, 2019). Overall, these studies agree on tendencies such that concrete  
227 words tend to have less diverse but more compact and more strongly associated distributional neighbours  
228 than abstract words.

### 229 2.2.2 General vs. Specific Words

230 From a computational perspective, hypernymy –which we take as instantiation to represent degrees  
231 of generality vs. specificity– is central to solving a number of NLP tasks such as automatic taxonomy  
232 creation (Hearst, 1998; Cimiano *et al.*, 2004; Snow *et al.*, 2006; Navigli and Ponzetto, 2012) and textual  
233 entailment (Dagan *et al.*, 2006; Clark *et al.*, 2007). An enormous body of computational work has applied  
234 variants of lexico-syntactic patterns in order to distinguish hypernymy among word pairs from other lexical  
235 semantic relations (Hearst, 1992; Pantel and Pennacchiotti, 2006; Yap and Baldwin, 2009; Schulte im  
236 Walde and Köper, 2013; Roth and Schulte im Walde, 2014; Nguyen *et al.*, 2017, i.a.). More closely  
237 related to the current study, Shwartz *et al.* (2017) provide an extensive overview and comparison of  
238 unsupervised distributional methods. They distinguish between families of distributional approaches, i.e.,  
239 *distributional similarity measures* (assuming asymmetric distributional similarities for hypernyms and  
240 their hyponyms regarding their contexts, e.g., Santus *et al.* (2016)), *distributional inclusion measures*

241 (comparing asymmetric directional overlap of context words, e.g., Weeds and Weir (2005); Kotlerman  
242 *et al.* (2010); Lenci and Benotto (2012)) and *distributional informativeness measures* (assuming different  
243 degrees of contextual informativeness, e.g., Rimell (2014); Santus *et al.* (2014)). Across modelling systems,  
244 most approaches model hypernymy between nouns; hypernymy between verbs has been addressed less  
245 extensively from an empirical perspective (Fellbaum, 1990; Fellbaum and Chaffin, 1990; Fellbaum, 1998a).

## 246 2.3 Empirical Models Across Types of Abstraction

247 In addition to interdisciplinary empirical research targeting concreteness or hypernymy that has been  
248 mentioned above, we find at least two empirical studies at the interface of cognitive and computational  
249 linguistics that brought together our two target types of abstraction beforehand, Theijssen *et al.* (2011)  
250 and Bolognesi *et al.* (2020). Similarly to the current work, Theijssen *et al.* used the observation in Spreen  
251 and Schulz (1966) defining abstraction in terms of concreteness and specificity as their starting point.  
252 They provide two empirical experimental setups to explore and distinguish between the abstraction types  
253 in actual system implementations, (1) based on existing annotations of noun senses in a corpus, and  
254 (2) based on human judgements on labelling nouns in English dative alternations. As resources they used  
255 the MRC database (Coltheart, 1981) and WordNet. Overall, they found cases where concreteness and  
256 specificity overlap and cases where the two types of abstraction diverge. Bolognesi *et al.* looked into the  
257 same two types of abstraction to correlate degrees of abstraction in the concreteness norms by Brysbaert  
258 *et al.* (2014) and in the WordNet hierarchy, and to investigate interactions between the four groups of  
259 more/less concrete  $\times$  more/less specific English nouns from the two resources. Their studies illustrate that  
260 concreteness and specificity represent two distinct types of abstraction.

261 Further computational approaches zoomed into statistical estimation of contextual diversity/neighbourhood  
262 density, in order to distinguish between degrees of semantic abstraction across types of abstraction. For  
263 example, McDonald and Shillcock (2001) applied the information-theoretic measure *relative entropy*  
264 to determine the degree of informativeness of words, where word-specific probability distributions over  
265 contexts were compared with distributions across corresponding sets of words. The contextual diversity  
266 measure by Adelman *et al.* (2006) is comparably more simple: they determined the number of documents  
267 in a corpus that contain a word. More recently, Danguécan and Buchanan (2016), Reilly and Desai (2017)  
268 and our own work in Naumann *et al.* (2018) explored variants of neighbourhood density measures for  
269 abstract and concrete words, i.e., the number of (different) context words and the distributional similarity  
270 between context words. Additional approaches to determine contextual diversity/neighbourhood density  
271 have arisen from other fields of research concerned with semantic abstraction, i.e., regarding ambiguity  
272 and diachronic meaning change (Sagi *et al.*, 2009; Hoffman *et al.*, 2013; Hoffman and Woollams, 2015).  
273 Overall, these studies demonstrated that contextual density/diversity differs for more vs. less abstract words  
274 and across types of abstraction, even though the applications of the measures were rather diverse.

## 3 MATERIALS AND METHODS

### 275 3.1 Abstraction Data: Concreteness and Hypernymy

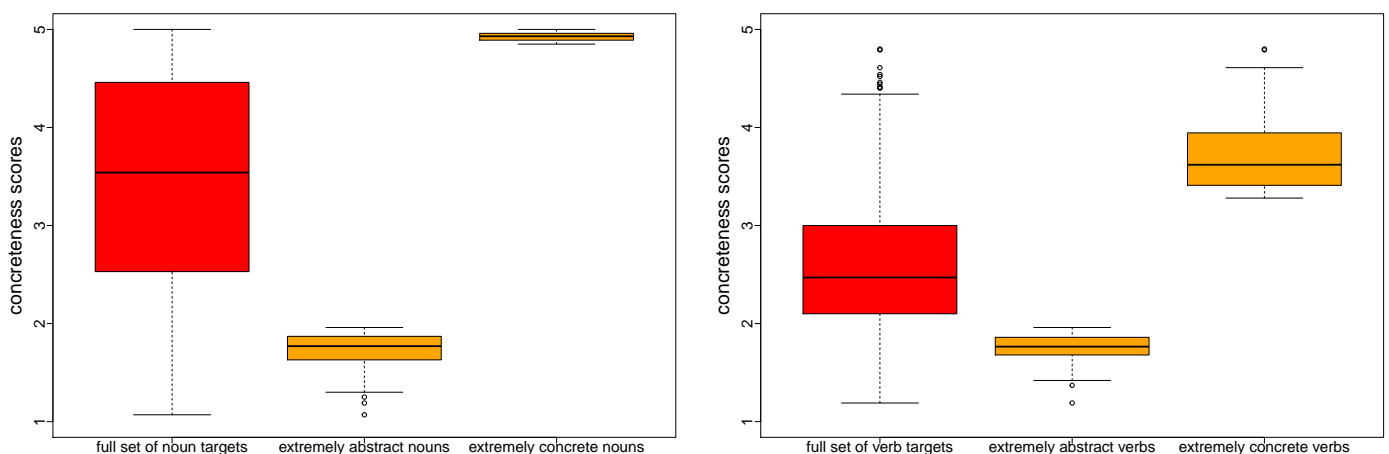
276 In the following, we introduce the resources we used for creating variants of abstraction data for our  
277 distributional experiments in Section 4. As motivated above, we distinguish semantic abstraction in terms  
278 of the abstract–concrete and the generality–specificity distinctions.

## 279 3.1.1 Concreteness Targets

280 Regarding abstraction in terms of the abstract–concrete dichotomy (henceforth referred to as **concreteness**  
 281 condition), we rely on the concreteness ratings for approximately 40,000 English words and two-word  
 282 expressions from Brysbaert *et al.* (2014). The ratings were collected via Amazon Mechanical Turk by  
 283 asking at least 25 participants to judge the concreteness vs. abstractness of the targets on a 5-point rating  
 284 scale from 1 (abstract) to 5 (concrete) regarding how strongly the participants thought the meanings of the  
 285 targets can(not) be experienced directly through their five senses. The overall targets’ scores of abstractness  
 286 vs. concreteness are represented by the mean values. For example, the concrete word *banana* received the  
 287 highest possible average rating of 5.0 because it is strongly perceived by human senses, while the abstract  
 288 word *glory* received a rather low average rating of 1.45.

289 The ratings had been collected for the targets out-of-context and without any further word-class  
 290 disambiguating information. In a post-processing step, Brysbaert *et al.* added part-of-speech (POS)  
 291 and frequency information from the SUBTLEX-US corpus (Brysbaert *et al.*, 2012). We repeated their  
 292 post-processing step, however relying on the ENCOW corpus data we also use in our studies (see below for  
 293 details), i.e., we automatically assigned each target its most frequently occurring POS tag in the ENCOW.

294 If this POS did not represent an overall proportion of at least 95% of all POS tags of that target or  
 295 if our most-frequent POS was not identical to the POS tag assigned by Brysbaert *et al.*, we discarded  
 296 the target in order to minimise POS ambiguity among targets. We also discarded target words with an  
 297 ENCOW frequency below 10,000. Our final concreteness set of targets contains 5,448 nouns and 1,280  
 298 verbs. Henceforth, we will refer to this selection of datapoints as the **full concreteness** collection. We also  
 299 created target subsets of the 500 most concrete and the 500 most abstract nouns, and ditto for the 200  
 300 most concrete/abstract verbs. We will refer to these subsets as the **concreteness extremes** subsets. Figure 1  
 301 illustrates the distributions of concreteness scores across the full and extreme target sets; the underlying  
 302 files are provided in the supplement.



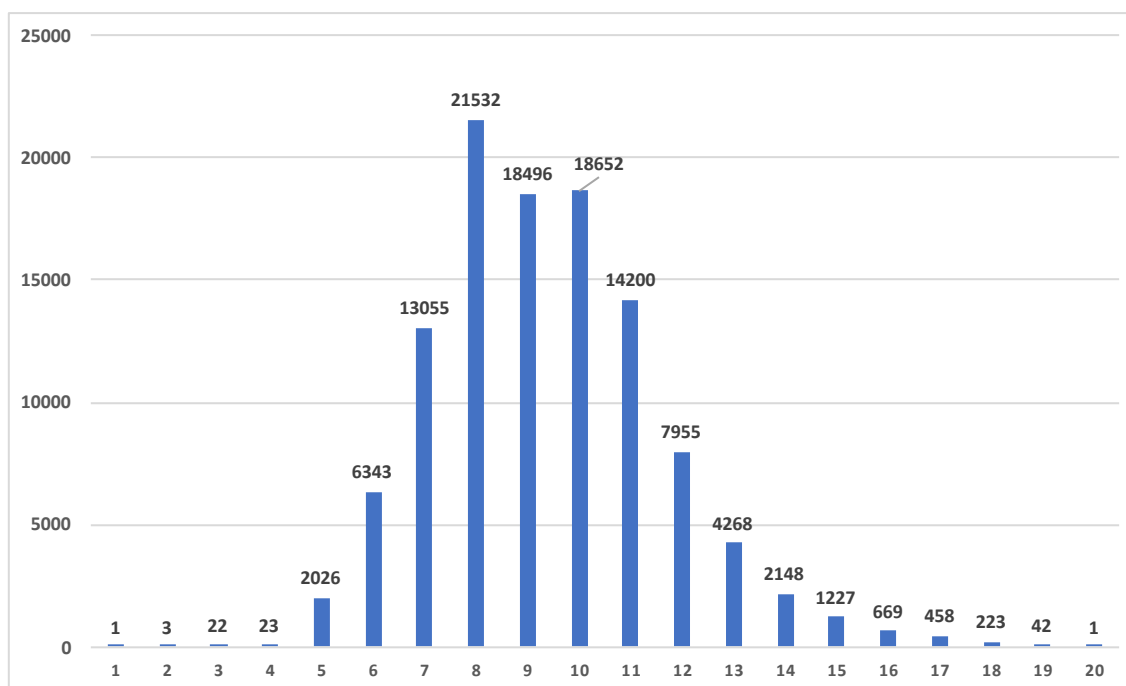
**Figure 1.** Distributions of concreteness scores on a 5-point rating scale from 1 (abstract) to 5 (concrete) for our full concreteness sets of 5,448/1,280 nouns/verbs and for the 500/200 most extreme abstract and concrete nouns/verbs.



## 303 3.1.2 Hypernymy Targets

304 Regarding abstraction in terms of generality (henceforth referred to as **hypernymy** condition), we rely  
 305 on WordNet, a standard lexical semantic taxonomy for English developed at Princeton University (Miller  
 306 and Fellbaum, 1991; Fellbaum, 1998b) that was also used by previous work on the generality–specificity  
 307 abstraction distinction (Theijssen *et al.*, 2011; Bolognesi *et al.*, 2020). The lexical database was inspired  
 308 by psycholinguistic research on human lexical memory and organises English nouns, verbs, adjectives  
 309 and adverbs into classes of synonyms (*synsets*), which are connected by lexical and conceptual semantic  
 310 relations. Words with several senses are assigned to multiple synsets. As mentioned above, WordNet  
 311 implements a hierarchical organisation of noun synsets relying on hypernymy relations, and verbs are  
 312 organised by a verb-specific variant of hypernymy, i.e., *troponymy*:  $v_1$  is to  $v_2$  in some manner, which itself  
 313 is conditioned on entailment and temporal inclusion.

314 We extracted all noun and verb synset pairs from WordNet version 3.0 that are in a hyponym–hypernymy  
 315 relation and paired all nouns/verbs from the respective subsets (such as *trout–fish* and *swim–move*, where  
 316 the first word in the pairs is the semantically more specific hyponym and the second word in the pairs is the  
 317 semantically more general hypernym), resulting in a total of 295,963/67,586 word pairs for nouns/verbs.  
 318 We then discarded any pairs containing multiword targets (such as *edible fruit*) as well as targets starting  
 319 with a capital letter (mostly proper names such as *Xhosa*) or starting with a number, leaving a total of  
 320  $\approx 110,000/47,500$  noun/verb pairs containing  $\approx 38,000/8,500$  different nouns/verbs. Figure 2 shows the  
 321 number of synsets per level in the noun hierarchy, with level 1 representing the top-most and therefore  
 322 most general synset  $\{entity\}$ . For verbs this analysis is not straightforward, as many synsets do not have  
 323 a hypernym, and the top levels are not consistently connected downwards (also see Richens (2008) on  
 324 “anomalies in the WordNet verb hierarchy”); this is the reason why some hypernymy-level-related analyses  
 325 in Section 4 will not be performed for verbs.



**Figure 2.** Number of synsets per hypernymy level in the WordNet noun hierarchy, with level 1 representing the top-most and therefore most general synset  $\{entity\}$ .

326 **3.2 Vector Space Variants**

327 The basis for our experiments is represented by the POS-tagged version of the sentence-shuffled English  
 328 COW corpus ENCOW16AX<sup>1</sup>, containing  $\approx 10$  billion words (Schäfer and Bildhauer, 2012; Schäfer, 2015).  
 329 From the corpus, we extracted co-occurrences (i.e., context words) for all nouns and verbs in the corpus  
 330 by applying a standard range of co-occurrence options: We relied on 2-word and 20-word symmetric  
 331 windows (left+right) across the lemmatised version of the corpus and distinguished between (a) taking  
 332 only co-occurring noun context words into account (henceforth: N space) and (b) taking all co-occurring  
 333 nouns, verbs and adjectives into account (henceforth: N-V-A space), when creating our noun–context and  
 334 verb–context matrices. The windows were applied within-sentence because the corpus is sentence-shuffled  
 335 for copyright reasons, such that going beyond sentence border is not meaningful. Furthermore, to reduce  
 336 noise in the co-occurrence data, we restricted the corpus lemmas to words starting with at least two letters;  
 337 by using a co-occurrence frequency cut-off of 50; and by discarding the most frequent content words:  
 338 *people, time, year* (nouns); *be, do, have* (verbs); and *other, more, many, such, same, few, most* (adjectives),  
 339 given that high-frequency words are notorious hubs and popular nearest neighbours in the vector spaces  
 340 (Radovanović *et al.*, 2010; Dinu *et al.*, 2015; Köper *et al.*, 2016, i.a.). The raw co-occurrence frequency  
 341 counts were weighted by the association measure *local mutual information (lmi)*, cf. Evert (2005). LMI  
 342 is an information-theoretic association measure that compares observed frequencies  $O$  with expected  
 343 frequencies  $E$ , taking marginal frequencies into account:  $LMI = O \times \log \frac{O}{E}$ , with  $E$  representing the  
 344 product of the marginal frequencies over the sample size.<sup>2</sup>

345 Our co-occurrence matrices are general-purpose and not prone to our specific resource-induced targets,  
 346 which is required by some abstraction measures (see following Section 3.3). Table 1 shows the sizes of  
 347 our vector space matrix variants in numbers of targets and dimensions, i.e., context words. Table 2 shows  
 348 co-occurrence frequencies and lmi scores for a sample noun, i.e., *fish*, and a selection of its context words  
 349 within a window of  $\pm 20$  words.

target POS	window size	dimension POS	# targets	# dimensions
N	2	N	22,017	22,017
		N-V-A	24,279	40,571
	20	N	29,721	29,721
		N-V-A	30,748	51,249
V	2	N	6,259	16,373
		N-V-A	6,544	28,736
	20	N	7,338	25,254
		N-V-A	7,530	43,329

**Table 1.** Sizes of vector space variants in terms of numbers of target types and dimension types in the co-occurrence (context) matrices.

350 **3.3 Abstraction Measures**

351 The following subsections introduce our selection of distributional methods to measure abstraction both  
 352 in terms of the abstractness–concreteness dichotomy and in terms of the generality–specificity distinction.

<sup>1</sup> <https://corporafromtheweb.org/encow16/> provides details on corpus version and toolchains.

<sup>2</sup> See <http://www.collocations.de/AM/> for a detailed description of association measures.

context word & POS		frequency	lmi
water	NN	56,049	133,387.53
tank	NN	39,118	150,223.00
catch	V	37,003	117,624.73
eat	V	31,558	87,119.87
small	ADJ	30,864	45,470.63
big	ADJ	24,835	37,067.61
chip	NN	19,407	72,473.17
oil	NN	18,404	41,075.69
salmon	NN	8,983	38,461.76
tropical	ADJ	6,629	23,600.64
serve	V	6,571	4,433.21
eye	NN	4,052	1,701.02

**Table 2.** Example context words for the target noun *fish* within a window of  $\pm 20$  words, accompanied by co-occurrence frequencies and local mutual information (lmi) scores.

### 353 3.3.1 Neighbourhood Densities

354 Our main focus regarding vector space measures of abstraction lies on variants of neighbourhood densities.  
 355 As described in Section 2, previous work has applied such measures to a number of tasks involving semantic  
 356 abstraction (not necessarily using the identical term “neighbourhood density”), such as lexical semantic  
 357 ambiguity (Hoffman *et al.*, 2013), lexical semantic change (Sagi *et al.*, 2009), hypernymy (Santus *et al.*,  
 358 2014) and lexical concreteness (Frassinelli *et al.*, 2017; Naumann *et al.*, 2018). The underlying assumption  
 359 of the empirical models across tasks is that the neighbourhood density of more abstract words is lower  
 360 than the neighbourhood density of less abstract (i.e., more specific/concrete) words, because conceptual  
 361 connections between abstract words and their semantically associated words are more diverse/variable  
 362 and less meaning-specific than conceptual connections between more specific/concrete words and their  
 363 semantically associated words.

364 In this vein, neighbourhood density measures score the variability of contexts in which words occur in  
 365 different ways. They either (i) measure neighbourhood density by relying on **context words**, assuming that  
 366 more abstract words co-occur with a larger variety of context words, or they (ii) measure neighbourhood  
 367 density by relying on **neighbour words**, assuming that more abstract words have a larger variety of  
 368 distributionally similar words. As mentioned above, these types of neighbourhood densities are conceptually  
 369 rather different, exploiting similarity between context words vs. exploiting similarity between nearest  
 370 neighbours. In addition, neighbourhood density measures differ with respect to involving (or not involving)  
 371 the respective target words in the calculation. Finally, all variants of measures need to define the number  $k$   
 372 of context/neighbour words that are taken into account, i.e., how many words are involved as “strongest”  
 373 context/neighbour words. The four variants are defined and computed as follows.

374 CC The neighbourhood density of a target word  $t$  is defined as the average vector-space distance **between**  
 375 **the  $k$  strongest context words** of  $t$ .

376 TC The neighbourhood density of a target word  $t$  is defined as the average vector-space distance **between**  
 377  **$t$  and its  $k$  strongest context words**.

378 NN The neighbourhood density of a target word  $t$  is defined as the average vector-space distance **between**  
 379 **the  $k$  nearest neighbours** of  $t$ .

380 TN The neighbourhood density of a target word  $t$  is defined as the average vector-space distance **between**  
 381  **$t$  and its  $k$  nearest neighbours**.

382 The strongest context words are determined on the basis of the local mutual information strength of  
 383 co-occurrence (see previous Section 3.2). Vector-space distance between words in order to determine  
 384 nearest neighbours is computed by calculating the *cosine* of the angle between the respective word vectors.  
 385 See Table 7 in the Appendix 1 for examples of strongest context and neighbour words regarding a selection  
 386 of target nouns and verbs.

### 387 3.3.2 Contextual Entropy

388 For measuring the contextual entropy of a target word we rely on standard word entropy, which has  
 389 been suggested as an asymmetric method for hypernymy prediction by Shwartz *et al.* (2017), inspired by  
 390 a previous second-order co-occurrence variant (Santus *et al.*, 2014). The underlying assumption is that  
 391 more abstract words are more uncertain (and therefore receive a higher entropy value) than less abstract  
 392 (i.e., more specific/concrete) words. For each target word  $w$  in our vector spaces we calculated the word  
 393 entropy  $H(w)$ , taking all of  $w$ 's context words  $c$  from our vector spaces into account, see Equation (1). The  
 394 computation requires per-target probabilities over context words, which we calculated based on the raw  
 395 target–context co-occurrence frequencies.

$$H(w) = - \sum_c p(c|w) \cdot \log_2(p(c|w)) \quad (1)$$

### 396 3.3.3 Weeds Precision

397 Weeds Precision (WeedsPrec) represents an asymmetric method suggested by Weeds *et al.* (2004) that  
 398 quantifies the weighted inclusion of the features of word  $w_1$  in the features of word  $w_2$ . In our case the  
 399 features refer to the words' context words  $c$ . The underlying assumption is that more context words  $c$  of the  
 400 more specific hyponym are among its hypernym's context words than there are context words of the more  
 401 general hypernym among its hyponym's context words. If  $WeedsPrec(w_1, w_2) > WeedsPrec(w_2, w_1)$ ,  
 402 then  $w_1$  is predicted as the hyponym and  $w_2$  as the hypernym, and vice versa, see Equation (2). For example,  
 403 one would expect more context words of the hyponym *cat* also as context words of its hypernym *animal*  
 404 (such as *eyes, fur, tail*) than vice versa, because the hypernym also co-occurs with words relevant for other  
 405 animals (such as *flapper* for *fish*) that are however not relevant for *cats*.

406 The computation requires raw target–context co-occurrence frequencies  $|w_{ic}|$ . Next to the original  
 407 weighted, token-based version of WeedsPrec in Equation (2) we also apply a non-weighted, type-based  
 408 version (WeedsPrec') where we compute *whether* a context word is included in a specific vector, rather  
 409 than *how often* it is included, see Equation (3).

$$WeedsPrec(w_1, w_2) = weeds-token = \frac{\sum_{c \in (\vec{w}_1 \cap \vec{w}_2)} |w_{1c}|}{\sum_{c \in \vec{w}_1} |w_{1c}|} \quad (2)$$

$$WeedsPrec'(w_1, w_2) = weeds-type = \frac{\sum_{c \in (\vec{w}_1 \cap \vec{w}_2)} 1}{\sum_{c \in \vec{w}_1} 1} \quad (3)$$

## 4 DISTRIBUTIONAL ABSTRACTION EXPERIMENTS

410 In this section we report our empirical experiments on distributional models of abstraction. Subsection 4.1  
411 describes the setup of the experiments, and subsection 4.2 presents the results of distinguishing between  
412 degrees of abstraction in terms of concreteness and hypernymy.

### 413 4.1 Abstraction Experiments: Setup

414 **Main experiments:** The nature of our target datasets differs with respect to the underlying type of  
415 abstraction. For this reason, we defined a common strategy to make the results comparable across datasets:  
416 As a major point of comparison we rely on **pairs** of target words, which combine abstract with concrete  
417 words, and hypernyms with their hyponyms. For the hypernymy pairs, the two words are directly provided  
418 by the resource: we paired each word in a synset with each word in the superordinated synset(s), see  
419 Section 3.1; for the concrete–abstract pairs, we followed our previous work (Naumann *et al.*, 2018;  
420 Frassinelli and Schulte im Walde, 2019) and took our collection of extremes with 500+500 nouns and  
421 200+200 verbs to create 250,000/40,000 concrete–abstract noun/verb word pairs. Note that Figure 1 already  
422 included the distributions of concreteness scores for these extreme target subsets.

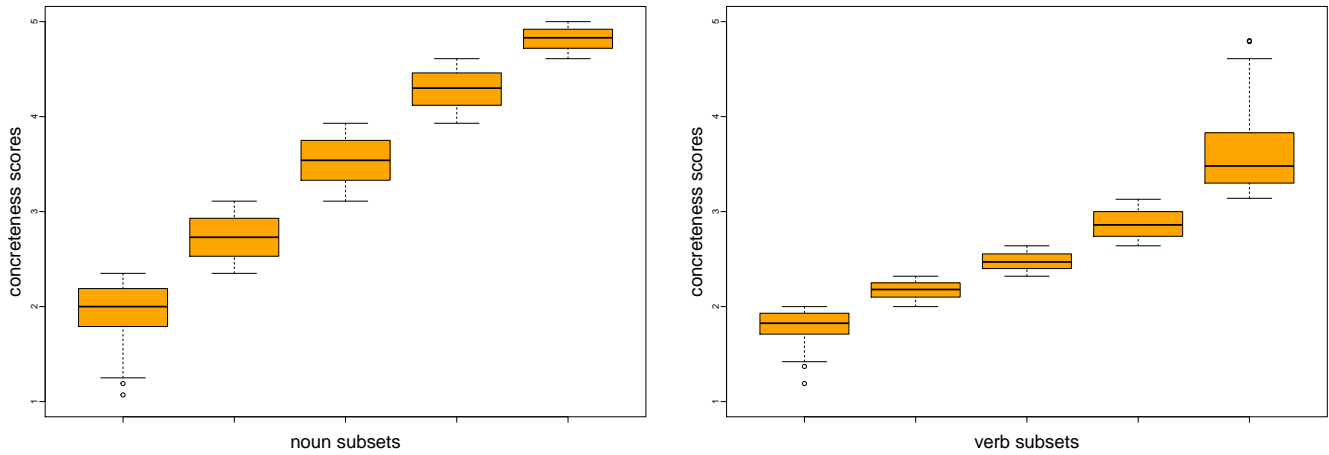
423 The task for our measures regarding target pairs was to identify the more abstract word in each pair.  
424 The results are computed by determining precision (which in this setup is identical to accuracy), i.e., the  
425 proportion of empirically identified abstract words that were indeed the more abstract words in the pairs.  
426 We focus on precision here because the differences of our vector spaces regarding the proportions of target  
427 words they cover (i.e., their recall) is only marginal. We nevertheless include the numbers of retrieved  
428 distinctions per measure and target space in the full results in Appendix 2.

429 In addition to this first set of experiments where we compared all of our abstraction measures on noun  
430 and verb concreteness and hypernymy pairs across vector spaces, we then focused on specific aspects in  
431 the experimental paradigm, as follows.

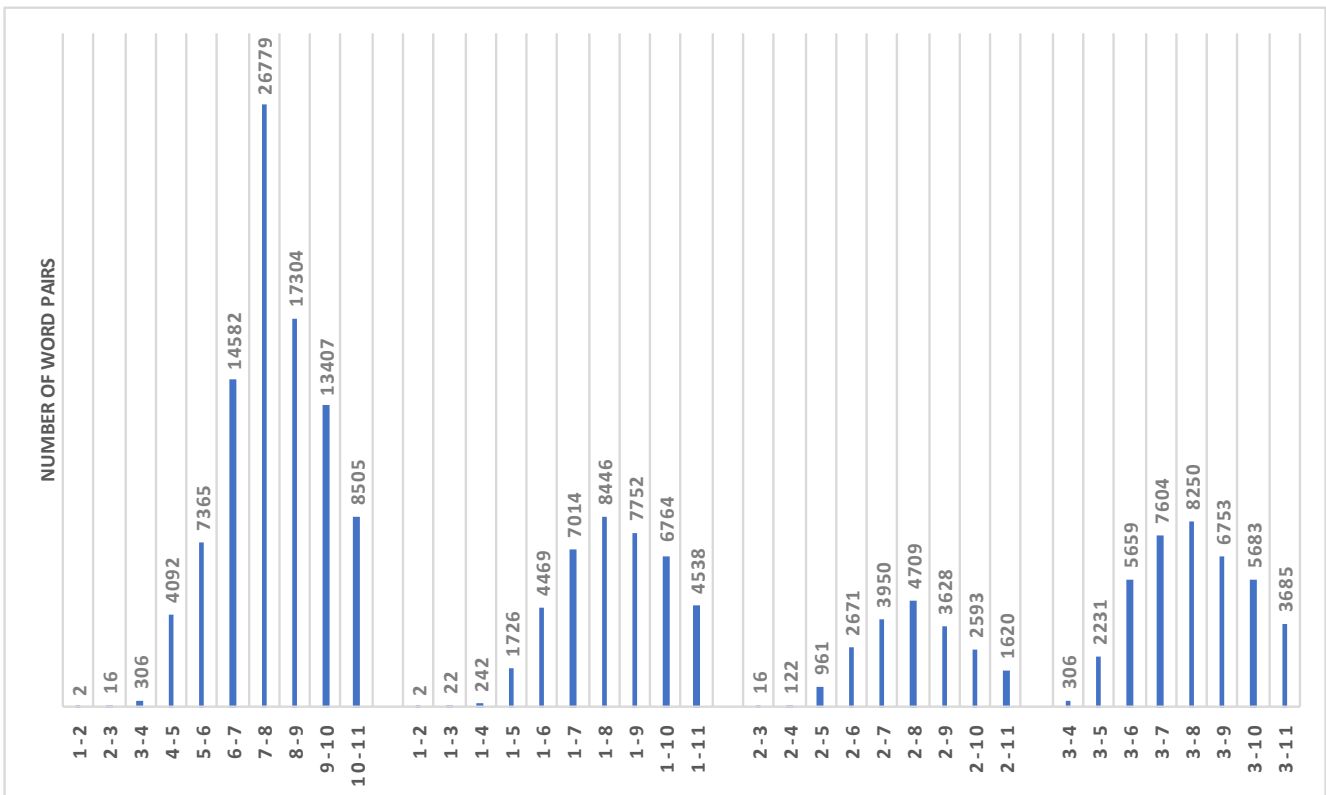
432 **Strength of abstraction:** We hypothesised that the measures are more or less successful with respect to  
433 how “different” the concrete and abstract words are in their degrees of concreteness (again, for noun and  
434 verb targets), and how “different” the hypernyms and hyponyms are in their degrees of specificity (for nouns  
435 only, cf. Section 3.1). Similarly to the previous experiments, this setup also relies on concrete–abstract and  
436 hyponym–hypernym pairs but the target sets were created in a different way.

437 For concreteness, we took our full concreteness dataset (see Section 3.1) and divided the 5,448/1,280  
438 nouns/verbs (separately for each word class) into five equally-sized subsets, after having sorted them by  
439 their concreteness scores. Figure 3 shows the distributions of concreteness scores across the five 20%  
440 dataset proportions. Then we created pairs using the targets in subset 1 and the targets in subset 2 (i.e.,  
441 pairing the 20% most abstract words with each of the targets in the second 20% most abstract words), for  
442 each of the targets in subset 1 with each of the targets in subset 3, etc., resulting in a total of 1,187,010  
443 pairs per range combination for nouns, and 65,536 pairs per range combination for verbs. In this way, we  
444 compare distinctions for pairs that are more or less similar in their degrees of concreteness, rather than  
445 the most extreme subsets. Note, in this respect, that the sizes of the boxes in Figure 3 indicate that we are  
446 facing a large number of very concrete nouns, while for verbs the majority is located in the range [2; 3].  
447 For hypernymy, we took into account the hierarchical levels of nouns when creating pairs, by pairing  
448 the top-level noun in the hierarchy (*entity*) with all second-level nouns, then with all third-level nouns,  
449 etc., and by pairing all second-level nouns with all third-level nouns, then with all forth-level nouns, etc.

450 Figure 4 shows the numbers of pairs after combining words from synsets of specific hierarchical hypernymy  
 451 levels. Note that we go down to level 11 in the WordNet hierarchy for this specific analysis. In the actual  
 452 experiments we will however disregard the level combinations with <100 pairs (i.e., 1-2, 1-3, 2-3).



**Figure 3.** Concreteness ranges of noun and verb subsets (each containing 20% of respective total data).



**Figure 4.** Numbers of word pairs in synset combinations across hierarchical levels.

453 **Correlations and interactions between measures:** We zoomed into correlations and interactions of  
454 abstraction distinctions across measures, in order to see whether the actual decisions of the measures are  
455 more or less strongly correlated with corpus frequency and with each other, and how they interact and  
456 complement each other. For this set of experiments we only used the concreteness targets (both nouns  
457 and verbs), which provide scores on a scale, differently to the pair-wise organised hierarchical hypernymy  
458 targets (which we could organise into hypernymy-based chains of levels but this would add a level of  
459 interpretation to the actual human categorisations that we do not judge appropriate). In addition, we used  
460 the 329 noun targets from Spreen and Schulz (1966) which are rated on a scale for both concreteness and  
461 specificity. For this set of experiments we exploit Spearman's rank-order correlation coefficient  $\rho$  (Siegel  
462 and Castellan, 1988) and regression models.

463 We now describe how we apply the abstraction measures to the pair-wise distinction between degrees  
464 of abstraction in concrete–abstract pairs and hyponym–hypernym pairs. For measuring contextual word  
465 entropy and WeedsPrec, we follow a straightforward one-step procedure: Relying on one of our vector-  
466 space matrices, we compute the extent of feature inclusion (WeedsPrec) regarding both words' dimensions,  
467 and we compute the word entropy for both words; the comparison of the respective two values then decides  
468 which word in a word pair is predicted as the more/less abstract one, see Section 3.3. For measuring  
469 neighbourhood density, two-step procedures are required: Regarding the CC and TC variants, we first need  
470 to identify the  $k$  strongest context words (i.e., co-occurrence dimensions) for each target word, and then  
471 compute the respective average cosine distances between the strongest context words (CC) or between the  
472 target and the strongest context words (TC). Regarding the NN and TN variants, we first need to identify the  
473  $k$  nearest neighbour words for each target word, and then compute the respective average cosine distances  
474 between the strongest neighbour words (NN) or between the target and the strongest neighbour words (TN).  
475 For all four neighbourhood density variants we rely on one of our vector-space matrices in the first step  
476 (i.e., N vs. N-V-A dimensions), and in step two we again face the same choice between the vector-space  
477 matrix variants. See Appendix 1 for a selection of noun and verb targets and their strongest context and  
478 neighbour words.

## 479 4.2 Abstraction Experiments: Results

### 480 4.2.1 Main Experiments

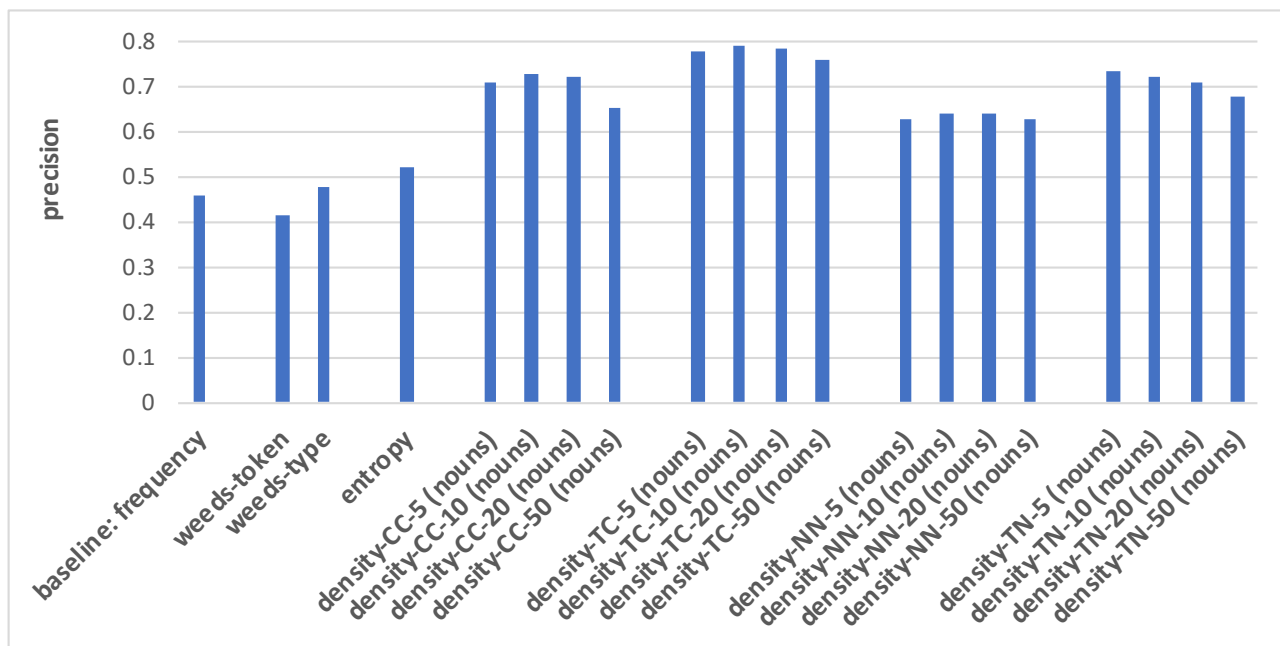
481 Figures 5–8 present the results when distinguishing between degrees of abstraction across measures in  
482 terms of precision, i.e., the proportion of abstract words suggested by the measures that were indeed the  
483 more abstract words in the pairs. As baseline we use frequency, assuming that a word in a word pair is  
484 more abstract if it is more frequent. The weighted vs. non-weighted variants of WeedsPrec are referred to  
485 as “weeds-token” vs. “weeds-type”, respectively. For neighbourhood density we report results for 5, 10,  
486 20 and 50 contexts/neighbours across our four variants CC, TC, NN and TN, and we distinguish between  
487 taking into account only nouns or only verbs (depending on the target POS)<sup>3</sup> as contexts/neighbours vs.  
488 *all* nouns, verbs and adjectives (N-V-A). We only show results using the N-V-A vector spaces induced  
489 from a co-occurrence window of 20 words, and the density variants that take only single-POS words as  
490 contexts/neighbours into account, because these generally provided the best results; the full result tables  
491 are available in Appendix 2.

<sup>3</sup> When taking into account a single POS for context/neighbour words, as context words we use nouns for both noun and verb targets, and as nearest neighbours we use same-POS neighbour words (i.e., noun nearest neighbours for noun targets and verb nearest neighbours for verb targets).

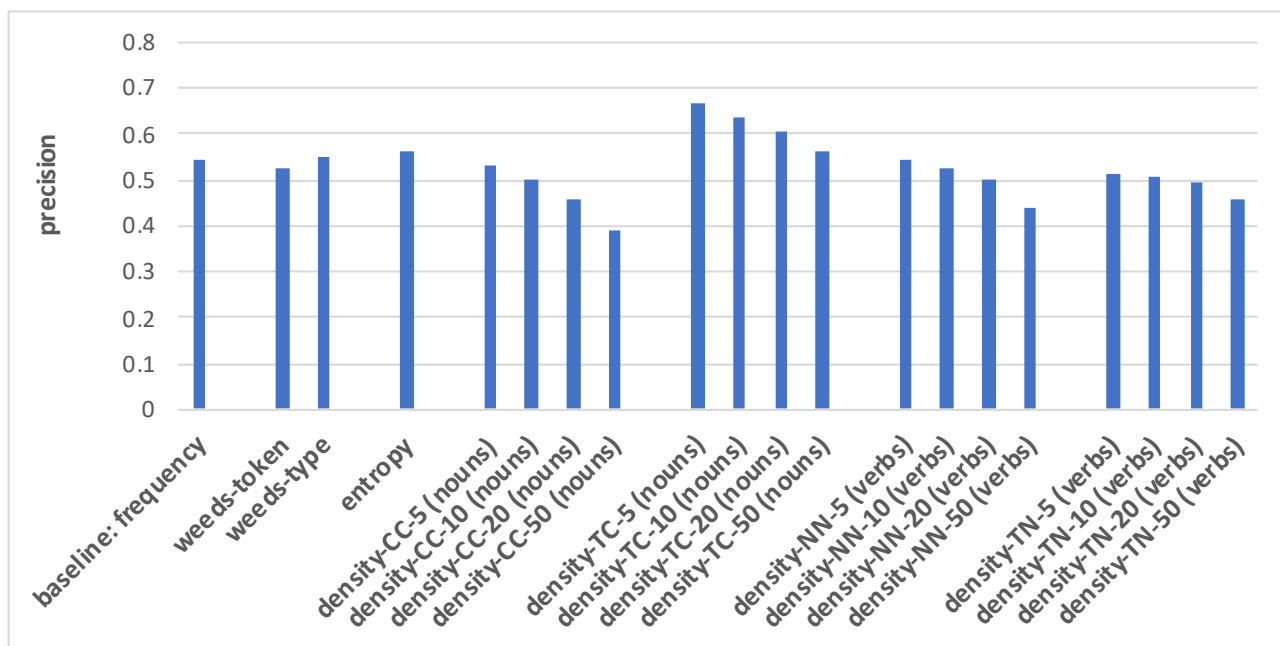
492 For both noun and verb targets, distinguishing between degrees of concreteness in Figures 5 and 6  
493 is performed best when applying the neighbourhood density measure TC: the strength of distributional  
494 similarity between a target word and its strongest context words distinguishes between the most abstract  
495 and the most concrete words with a precision of up to 0.79 for nouns and 0.67 for verbs, respectively.  
496 This means that the distributionally most similar context words in relation to a target are most indicative  
497 of the target's concreteness, and the higher this average vector-space similarity is, the more concrete  
498 are the target words. The next-best variants differ across the two POS types of our targets: for noun  
499 targets, the density measures are generally better than the baseline, weeds-token/-type and entropy, with  
500 density-NN representing the worst of the four density variants; for verb targets, the other density variants  
501 are at most en par with the baseline, weeds-token/-type and entropy, and overall the density variants are  
502 worse than for nouns, while the other measures perform better distinctions than for nouns. I.e., the baseline,  
503 weeds-token and entropy achieve 0.46/0.42/0.53 for nouns and 0.54/0.54/0.57 for verbs; for nouns the  
504 frequency baseline is even below the random baseline of 0.5. An additional insight from the figures is  
505 that in the vast majority of cases the strongest five or ten contexts/neighbours are the most indicative of  
506 their degrees of concreteness: in most cases the results worsen when more contexts/neighbours are taken  
507 into account. Including as contexts/neighbours only nouns/same-POS words (as in Figures 5 and 6, cf.  
508 footnote 3) vs. nouns, verbs and adjectives (see "all" in the full result tables in the Appendix) does not  
509 seem to strongly influence the qualities of the distinctions.

510 The prediction of hypernymy in Figures 7–8 provides a totally different pattern of results. For both noun  
511 and verb targets the best results are achieved by the frequency baseline (0.73/0.71), entropy (0.72/0.71),  
512 and the WeedsPrec variants: 0.72/0.73 for weeds-token and 0.73/0.71 for weeds-type, in comparison to the  
513 best density variants (for noun targets and density-NN-5: 0.52; for verb targets and density-NN-10: 0.56).  
514 Overall, most of the density-based results hardly beat the random baseline (0.5). Furthermore, the tendency  
515 that the density-based distinction results decrease when taking more context/neighbour words into account  
516 is visible only in some variants, and also not as clearly as in the results for distinguishing between degrees  
517 of concreteness.

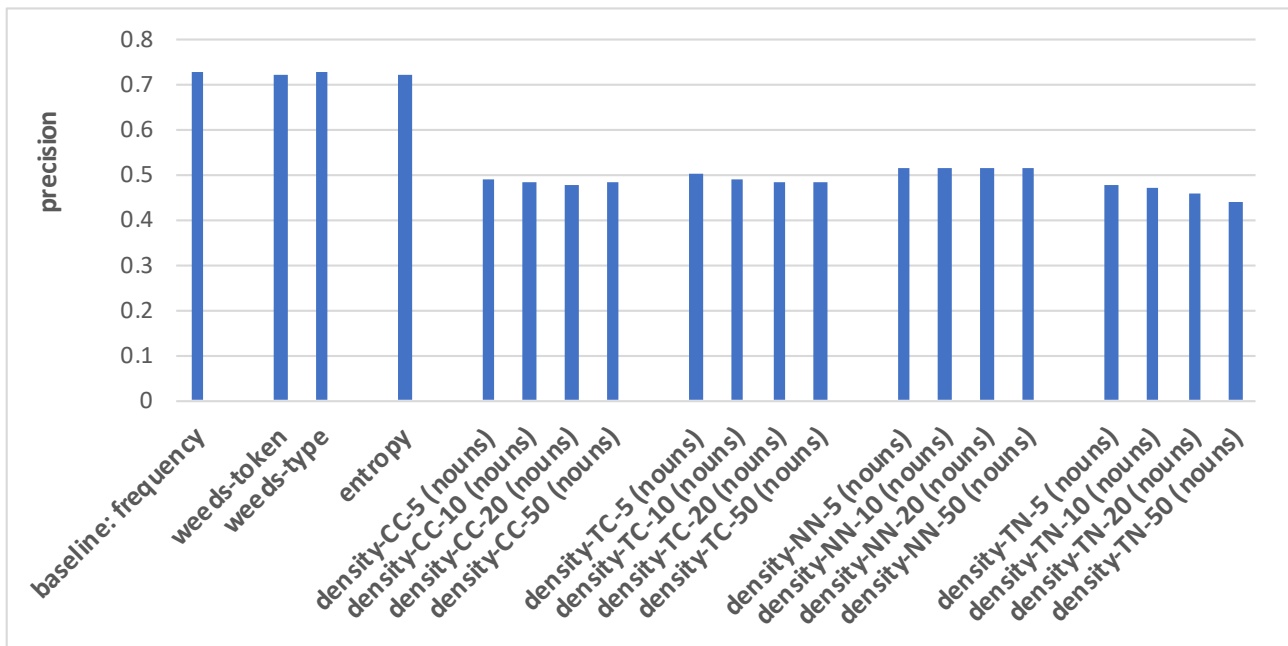




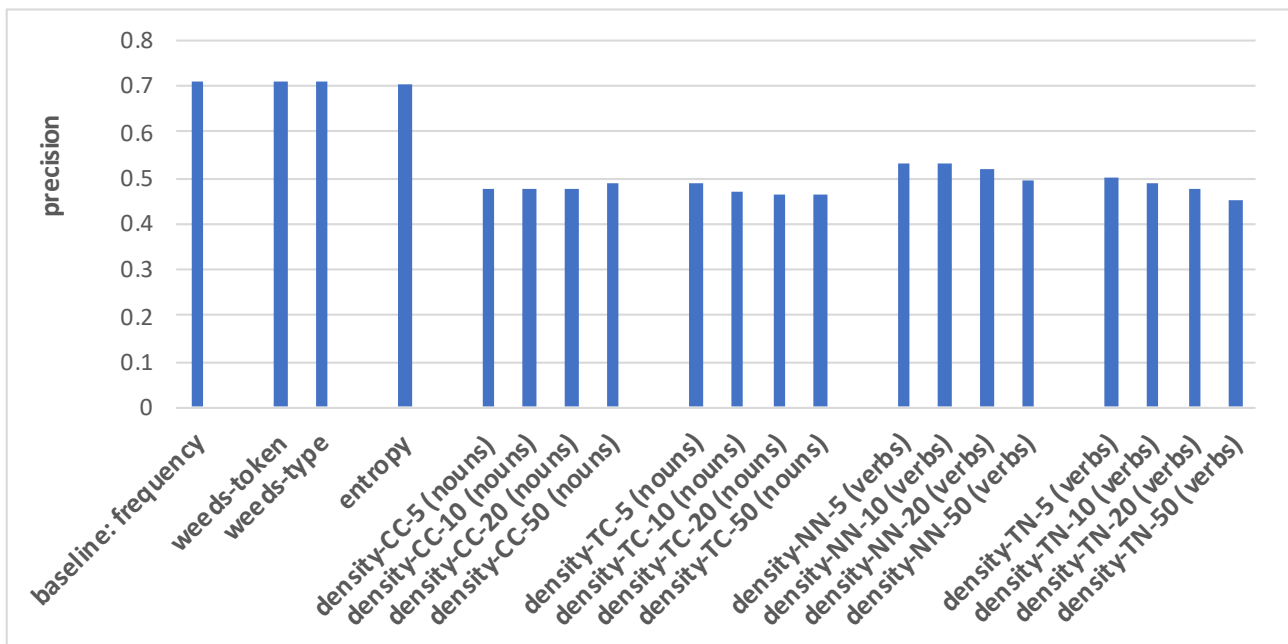
**Figure 5.** Pair-wise precision results for concreteness of nouns relying on an N-V-A vector space. Densities take only nouns as context/neighbour words into account.



**Figure 6.** Pair-wise precision results for concreteness of verbs relying on an N-V-A vector space. Densities take only nouns as context/neighbour words into account.



**Figure 7.** Pair-wise precision results for hypernymy of nouns relying on an N-V-A vector space. Densities take only nouns as context/neighbour words into account.



**Figure 8.** Pair-wise precision results for hypernymy of verbs relying on an N-V-A vector space. Densities take only nouns as context/neighbour words into account.

## 518 4.2.2 Strength of Abstraction

519 Following the main set of experiments we now zoom into the role of differences in results according to  
520 the strengths of concreteness and the levels of hypernymy. We hypothesise that the measures are more  
521 or less successful with respect to how “different” the concrete and abstract words are in their degrees of  
522 concreteness, and how “different” the hypernyms and hyponyms are in their degrees of specificity. We once  
523 more compare the baseline, weeds-token/-type, and entropy; for the neighbourhood variants we present the  
524 results relying on the 10 strongest context/neighbour words, because these proved rather successful and  
525 stable in the main experiments, and here we are not interested in the best results but rather in tendencies  
526 across subsets.

527 Figure 9 shows the results<sup>4</sup> across four sets of combinations of concreteness degrees for nouns. Note that  
528 we use the interval [0.4; 0.8] for precision values on the y-axis, for better visibility of trends and differences  
529 in results. The left-most set of results compares the distinctions between the most abstract and the second  
530 most abstract 20% of the targets, then the second and the third most abstract 20% of the targets, etc. So in  
531 this first set, the distances between concreteness degrees are identical (i.e., we use adjacent levels), but the  
532 concreteness ranges of the involved subsets differ. We can see that for the best three measures (densities TC,  
533 CC and TN) there is a slight upward trend which only drops for a mid-range comparison (subsets 3–4), even  
534 though we always look at adjacent levels. The four measures frequency, entropy and weeds-token/-type are  
535 better for mid-range nouns than for extremely abstract/concrete nouns but overall obtain lower precision  
536 values than the above three density variants. Density-NN shows the most idiosyncratic pattern of results,  
537 with mid-range precision values.

538 When comparing the results for nouns with increasing differences in concreteness degrees (see second,  
539 third and forth sets of results, using reference labels 1, 2, and 3), we can clearly see that for the four  
540 density variants the task becomes easier (and, accordingly, the results of the best measures improve) with  
541 stronger differences in concreteness scores. The overall best result (0.77) is obtained when distinguishing  
542 between nouns in levels 1 vs. 5, which represents the strongest difference in concreteness scores and  
543 is therefore similar to the previous extreme-range distinctions in the main experiments. The measures  
544 frequency, entropy and weeds-token/-type also show a slight increase in precision values but then drop for  
545 every comparison involving the most extreme concrete nouns (i.e., set 5).

546 Regarding abstraction measures, our insights from the main experiments are confirmed: for distinguishing  
547 between degrees of noun concreteness, the neighbourhood density measure TC is the best and most  
548 consistent in all cases, density-TN and density-CC are the next-best measures, and density-NN as well as  
549 frequency, entropy and weeds-token/-type represent the least successful measures.

550 Figure 10 shows the results across four sets of combinations of concreteness degrees for verbs. Note  
551 that we now use the interval [0.4; 0.65] for precision values on the y-axis, for better visibility of trends and  
552 differences in results. The left-most set of results across concreteness ranges for adjacent subsets shows a  
553 less clear pattern than for nouns. Across measures, the best results are achieved for the most abstract and  
554 for the most concrete subset combinations (1–2 and 4–5) and drop for the middle range combinations (2–3  
555 and 3–4).

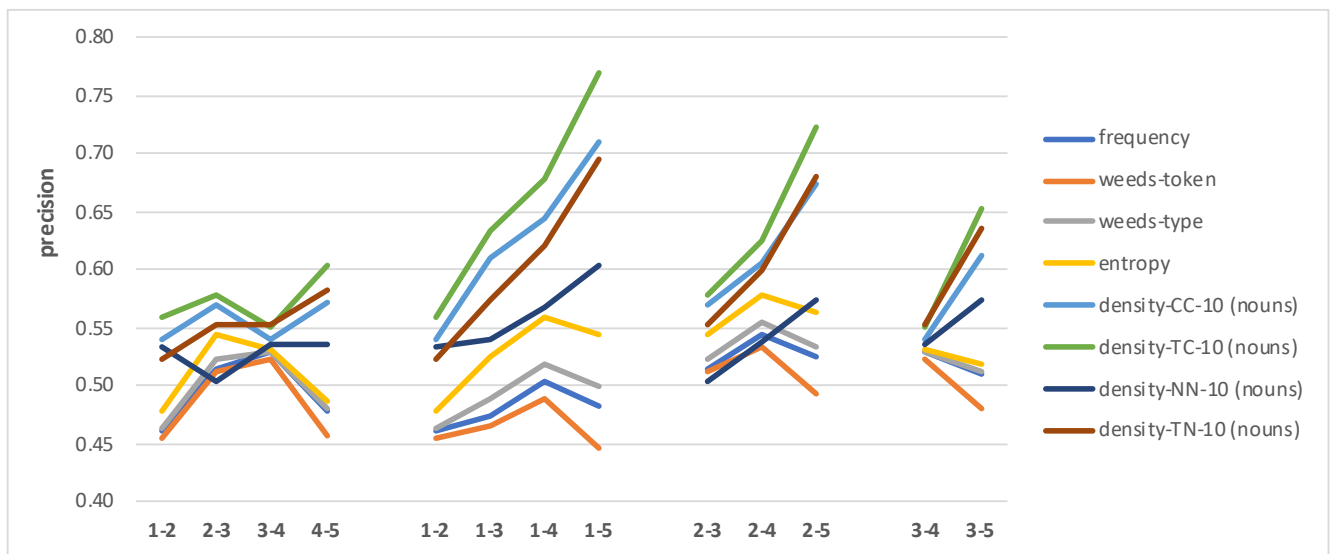
556 When comparing the results for verbs with increasing differences in concreteness degrees (see second,  
557 third and forth sets of results, again using reference labels 1, 2, and 3), we can see that the task is once more  
558 the easiest for the strongest differences in concreteness scores. But as for the adjacent-level comparisons

---

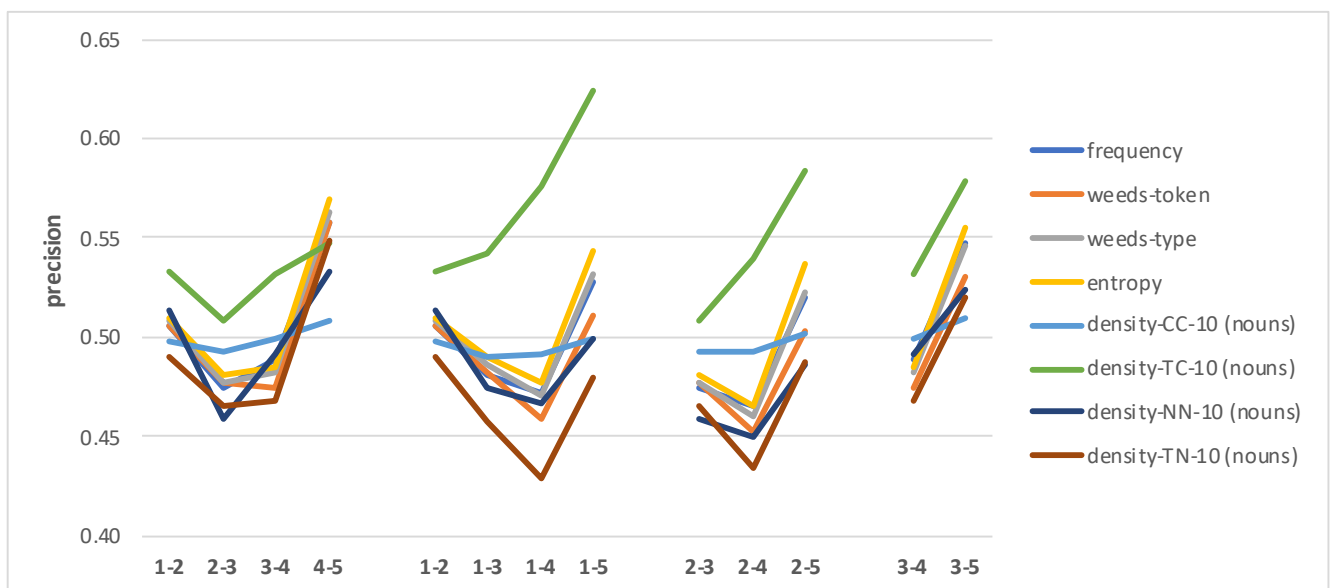
<sup>4</sup> Note that even though the precision scores are discrete, we use lines to illustrate the results, for better visibility and comparison.

559 for verb subsets, decisions involving the middle ranges are worse. Overall, the results are clearly below  
 560 those for nouns, with a best result of 0.62 obtained by density-TC when distinguishing between verbs in  
 561 levels 1 vs. 5.

562 Regarding abstraction measures, our insights from the main experiments are confirmed to some extent:  
 563 for distinguishing between degrees of verb concreteness, the neighbourhood measure density-TC is the  
 564 best in most cases, and frequency, entropy and weeds-token/-type are extremely similar to each other and  
 565 represent the next-best set of measures, however clearly below density-TC precision results and not much  
 566 above the other density variants. Density-CC seems to be least influenced by the degree of concreteness,  
 567 showing similar results across comparisons.

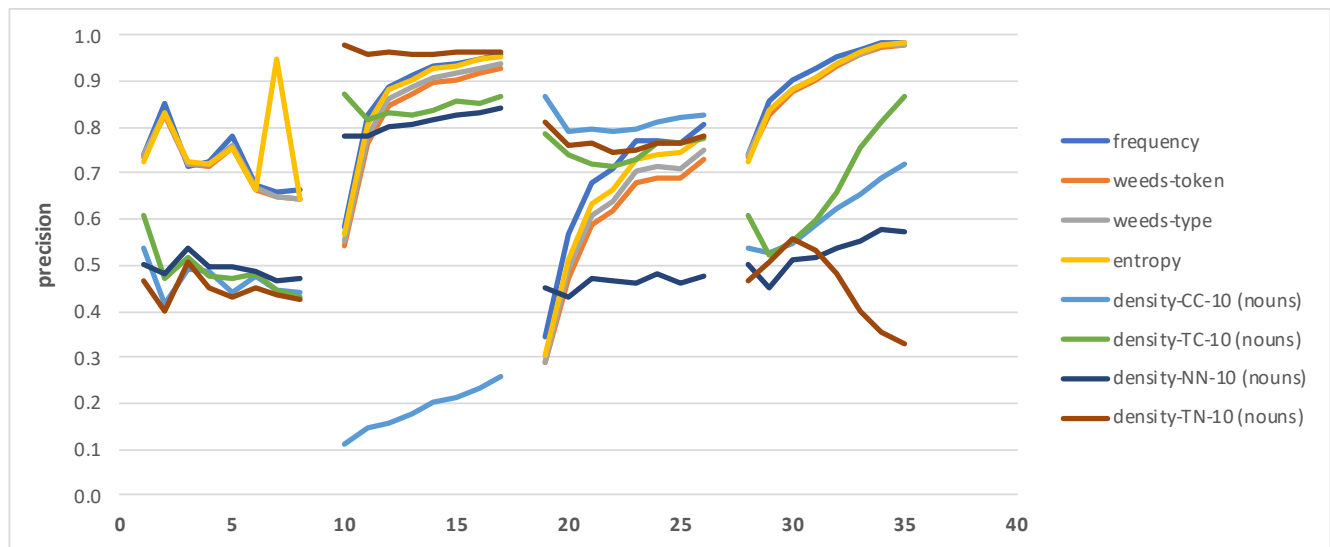


**Figure 9.** Results across combinations of concreteness ranges for nouns.



**Figure 10.** Results across combinations of concreteness ranges for verbs.

568 Figure 11 shows the results across four sets of combinations of hypernymy levels for nouns. Note that  
 569 in this case we use the full interval  $[0; 1]$  for precision values on the y-axis. The left-most set of results  
 570 compares the distinctions between pairs of related nouns from adjacent levels of hypernymy. Please  
 571 remember that we omit the combinations 1-2, 1-3 and 2-3 because these sets of pairs contain only 2, 16, and  
 572 22 pairs, respectively. Differently to the noun concreteness distinctions, there seems to be a slight downward  
 573 trend in precision. At the same time, there is more up and down across the level combinations, so the trends  
 574 are also less clear overall. What is clearly visible, on the contrary, is that frequency, entropy and weeds-  
 575 token/type are by far the best measures in this left-most set of distinctions for directly hypernymy-related  
 576 nouns across levels in the hierarchy (down to level 11).



**Figure 11.** Results across combinations of hypernymy levels for nouns.

577 Similarly, when comparing the results for related nouns with increasing differences in hypernymy levels  
 578 (see second, third and fourth sets of results, again using reference levels 1, 2, and 3), we can clearly see  
 579 that also here the task becomes easier (and, accordingly, the results improve) with stronger differences in  
 580 hypernymy levels. While this is clearly true for frequency, entropy and weeds-token/type, the patterns differ  
 581 more strongly for the density variants which mostly show less variability in results. Similarly to the main  
 582 results for hypernymy prediction, we once more observe that frequency, entropy and weeds-token/type  
 583 generally represent the best measures, while the density variants are worse.

#### 584 4.2.3 Correlations and Interactions between Measures

585 Overall, when looking at the distributions of frequency, entropy, weeds-token/-type and the neighbourhood  
 586 densities across types of abstraction and POS we see how subgroups of the measures are often extremely  
 587 similar to each other (and possibly interchangeable) in terms of predictive power. We now zoom into  
 588 correlations and interactions of abstractness distinctions across abstraction measures, in order to see whether  
 589 the actual scores provided by the measures are more or less strongly correlated with corpus frequency  
 590 and with each other, and how they interact and complement each other. For this set of experiments we  
 591 thus compare scores for words rather than binary decisions for word pairs, and as mentioned above we  
 592 use our concreteness targets (both nouns and verbs), which provide scores on a scale, and we use the 329  
 593 noun targets from Spreen and Schulz (1966) because those were rated on a scale for both concreteness  
 594 and specificity. We disregard the weeds-token/-type precision measures, as they would require setting  
 595 additional parameters in order to generate one score out of the two scores per pair.

596 **Correlations:** Figure 12 shows the correlations between noun concreteness scores, corpus frequency,  
 597 entropy and our four neighbourhood density variants (once more relying on  $k=10$ ). As before, the measures  
 598 use N-V-A spaces with a window of 20 words. First of all, we can see that the concreteness scores using  
 599 entropy are strongly correlated with corpus frequency ( $\rho=0.964$ ), while the density measures show no or  
 600 very low correlations with corpus frequency and entropy, so the density measures produce rather different  
 601 scores for abstraction in comparison to frequency and entropy. Among themselves, the density measures  
 602 show stronger agreement on their scores: regarding context densities, CC-10 and TC-10 correlate strongly  
 603 ( $\rho=0.814$ ); regarding nearest neighbour densities, NN-10 and TN-10, we find  $\rho=0.719$ . In contrast, we  
 604 see low correlations for NN-10 with CC-10/TC-10 ( $\rho < 0.3$ ), while for TN-10 we find medium-level  
 605 correlations of  $\rho \approx 0.5$  with the two context variants.

606 Figure 13 shows the correlations between verb concreteness scores, corpus frequency, entropy and our  
 607 four neighbourhood density variants ( $k=10$ ). As for the nouns, we find extremely high correlations between  
 608 corpus frequency and entropy; no correlations between these two measures and concreteness scores; strong  
 609 correlations for CC-10/TC-10 and NN-10/TN-10; moderate correlations between TN-10 and the context  
 610 variants; and low correlations between NN-10 and the context variants. Differently to the noun distinctions,  
 611 we do not find any correlation between any of the abstraction measures and concreteness.

612 Figures 14 and 15 look into correlations between abstraction ratings and abstraction measures for a subset  
 613 of 226 noun targets from Spreen and Schulz (1966). These 226 targets represent the intersection of the  
 614 nouns in Spreen and Schulz (1966) and our full concreteness subset Brysbaert *et al.* (2014). First of all,  
 615 Figure 14 shows the correlations between the concreteness and specificity ratings for these 226 noun targets  
 616 in the two norms. The two sets of concreteness ratings, which represent the main point of comparison,  
 617 strongly correlate ( $\rho=0.939$ ). Between the two sets of concreteness ratings and the specificity ratings we  
 618 find a lower but still meaningful correlation of  $\rho \approx 0.7$  for both resources. (Note that Spreen and Schulz  
 619 report a correlation of 0.626 between the concreteness and specificity ratings for their full set of 329 nouns.)

620 As in Figure 12, Figure 15 shows the correlations between noun concreteness scores, corpus frequency,  
 621 entropy and our four neighbourhood density variants (once more relying on  $k=10$ ) for the set of 226 nouns,  
 622 once more using N-V-A spaces with a window of 20 words. The overall picture is very much the same as  
 623 for our full set of 5,448 target nouns in Figure 12, for the concreteness ratings in Brysbaert *et al.* (2014) and  
 624 the concreteness and specificity ratings in Spreen and Schulz (1966), with one exception: frequency and  
 625 entropy show a moderate negative correlation with all abstraction rating sets:  $-0.47 < \rho < -0.41$  for both sets  
 626 of concreteness ratings, and  $-0.65 < \rho < -0.51$  for specificity ratings. The outcome of this last analysis is in  
 627 line with what we would have expected (but did not happen) to see in all three figures: generally, abstract  
 628 nouns are more frequent/entropic than concrete nouns, as we will also see below in the regression analysis,  
 629 so we expected a negative correlation between both frequency and entropy and the concreteness ratings.

630 Overall, the correlations for nouns and verbs (and for our targets and the subset of the targets from Spreen  
 631 and Schulz) show similar patterns regarding strong frequency–entropy correlations and tendencies in the  
 632 intra- and extra-density correlations. We however did not observe any meaningful correlation between the  
 633 abstraction measures and the concreteness scores of our verb targets, while we found correlations of  $\rho \approx 0.3$   
 634 between the abstraction measures and our noun ratings. This fits to our insights from the main experiments,  
 635 where the pair-wise distinctions for concreteness of verbs were worse than for nouns, and often similar to a  
 636 random baseline; nevertheless we reached precision scores of up to 0.79/0.67 for nouns/verbs, respectively.  
 637 For the much smaller set of 226 nouns from Spreen and Schulz (1966) the picture is similar to that for  
 638 our noun targets, but in addition frequency and entropy show a moderate negative correlation with both  
 639 concreteness and specificity ratings.

	frequency	entropy	density-CC-10	density-TC-10	density-NN-10	density-TN-10
concreteness	0.000	-0.076	0.263	0.335	0.126	0.248
frequency		0.964	0.089	0.136	0.033	0.189
entropy			-0.003	0.065	-0.019	0.095
density-CC-10				0.814	0.234	0.490
density-TC-10					0.255	0.552
density-NN-10						0.719

**Figure 12.** Spearman's  $\rho$  correlations between noun concreteness measures (N-V-A space).

	frequency	entropy	density-CC-10	density-TC-10	density-NN-10	density-TN-10
concreteness	-0.009	-0.032	-0.004	0.031	0.021	-0.046
frequency		0.970	0.002	0.141	-0.016	-0.048
entropy			0.085	0.180	-0.029	0.067
density-CC-10				0.694	0.217	0.350
density-TC-10					0.198	0.314
density-NN-10						0.749

**Figure 13.** Spearman's  $\rho$  correlations between verb concreteness measures (N-V-A space).

	S&S	
	concreteness	specificity
concreteness (B et al.)	0.939	0.687
concreteness (S&S)		0.704

**Figure 14.** Spearman's  $\rho$  correlations between the Spreen and Schulz and the Brysbaert *et al.* ratings for the subset of 226 nouns in the intersection.

	frequency	entropy	density-CC-10	density-TC-10	density-NN-10	density-TN-10
concreteness (B et al.)	-0.414	-0.454	0.255	0.336	0.027	0.224
concreteness (S&S)	-0.416	-0.468	0.257	0.349	0.023	0.231
specificity (S&S)	-0.506	-0.647	0.353	0.375	0.005	0.205
frequency		0.873	0.029	0.009	0.289	0.318
entropy			-0.239	-0.220	0.150	0.070
density-CC-10				0.819	0.203	0.475
density-TC-10					0.248	0.511
density-NN-10						0.764

**Figure 15.** Spearman's  $\rho$  correlations between ratings and measures for the subset of 226 nouns in the intersection of Spreen and Schulz and Brysbaert *et al.*

640 **Interactions:** The correlation analysis reported in Figure 12 shows a strong positive relationship for  
 641 nouns in the N-V-A space between frequency and entropy as well as between the density variants TC,  
 642 CC, TN and NN. For this reason, we must consider collinearity issues between the various predictors  
 643 (features) when modeling concreteness using linear regression models. In the following analyses, we will  
 644 model concreteness (as a continuous value ranging from 1 to 5) given different feature combinations. After  
 645 centering around the mean all the predictors, to test which triplet of variables best captures variability  
 646 in concreteness scores, we run eight independent models and select the one with the highest adjusted  
 647 R-squared value, as a measure of explained variance in the data. For an overview of the performance of  
 648 the eight models, see Table 3. The model including entropy, density-TC, and density-TN (highlighted  
 649 by bold font) is the one explaining the highest amount of variance in the concreteness scores (adjusted  
 650 R-squared: 13.4%) and does not show any collinearity problem ( $VIF < 1.64$ ). For this reason, we will  
 651 focus the following analysis on this model. The results discussed below are also fully in line with the results  
 652 in the other seven models from Table 3. As shown in Table 4, all three predictors (entropy, density-TC,  
 653 density-TN) are highly significant ( $p\text{-value} < 0.0001$ , after alpha correction because of multi-comparisons)  
 654 when modeling the concreteness of a noun. Words that are more concrete show: significantly lower entropy  
 655 scores, higher density-TC and higher density-TN; moreover, the interaction between the two density  
 656 measures indicates a positive overall effect. In the same table, we also report the "relative importance"  
 657 of each predictor (normalised to 100%) using the method developed by Lindeman *et al.* (1980). This  
 658 measure indicates the contribution of each predictor to the total amount of variance explained by the model.  
 659 Density-TC by itself explains 68.7% of the variance captured by the model, density-TN 20.7% and entropy  
 660 only 7.3%. The contribution of the various features is very stable across models and in line with what has  
 661 been discussed in the previous sections. When looking at all eight models, density measures involving  
 662 contextual information like density-TC and density-CC always contribute the most, as opposed to nearest  
 663 neighbour measures like density-NN and density-TN.

664 In Table 5, we see similar patterns to those emerged for nouns also for verbs. Once again, the model  
 665 including entropy, density-TC and density-TN is the one obtaining the highest R-squared value. However,  
 666 compared to nouns, the explained variance is extremely low (only 2%). When zooming in on the effect of  
 667 the single predictors on concreteness, Table 6 indicates some differences. The model shows only a strong  
 668 significant positive effect of density-TC ( $p < 0.0001$ ; after alpha correction) indicating that the contextual  
 669 density of concrete words is higher than the abstract one. For verbs, entropy ( $p = 0.008$ ), density-TN ( $p =$   
 670  $0.031$ ) and the interaction between the two density measures ( $p = 0.910$ ) do not reach significance. Once  
 671 more, density-TC is the feature with the strongest effect on concreteness scores, both for nouns and verbs.

Formula	Adj. R-squared
freq (ENCOW) + (density-TC $\times$ density-TN)	12.5%
freq (ENCOW) + (density-TC $\times$ density-NN)	11.9%
freq (ENCOW) + (density-CC $\times$ density-TN)	9.3%
freq (ENCOW) + (density-CC $\times$ density-NN)	8.1%
<b>entropy + (density-TC <math>\times</math> density-TN)</b>	<b>13.4%</b>
entropy + (density-TC $\times$ density-NN)	12.8%
entropy + (density-CC $\times$ density-TN)	9.9%
entropy + (density-CC $\times$ density-NN)	8.5%

**Table 3.** Comparison of model variants processing noun targets in the N-V-A space, and their explained variance (represented in terms of adjusted R-squared). The dependent variable is concreteness (1–5).



	Estimate	Std. Error	t-value	p-value	RI
<i>(Intercept)</i>	3.44	0.01	234.91	***	-
entropy	-0.11	0.01	-8.53	***	7.3%
density-TC	2.80	0.17	16.76	***	68.8%
density-TN	0.83	0.12	7.07	***	20.7%
density-TC × density-TN	4.45	0.86	5.20	***	2.3%

Significant codes: 0 '\*\*\*\*' 0.0001 '\*\*\*' 0.001 '\*\*' 0.006 ' ' 1

**Table 4.** Linear regression output for the best predictor combination for nouns in the N-V-A condition: entropy + (density-TC × density-TN). RI indicates the relative importance (normalised to 100%). The significance codes are all adjusted because of the 8 multi-comparisons.

Formula	Adj. R-squared
freq (ENCOW) + (density-TC × density-TN)	1.5%
freq (ENCOW) + (density-TC × density-NN)	1.2%
freq (ENCOW) + (density-CC × density-TN)	-0.2%
freq (ENCOW) + (density-CC × density-NN)	-0.2%
<b>entropy + (density-TC × density-TN)</b>	<b>2.0%</b>
entropy + (density-TC × density-NN)	1.6%
entropy + (density-CC × density-TN)	0.0%
entropy + (density-CC × density-NN)	0.0%

**Table 5.** Comparison of model variants processing verb targets in the N-V-A space, and their explained variance (represented in terms of adjusted R-squared). The dependent variable is concreteness (1-5).

	Estimate	Std. Error	t-value	p-value	RI
<i>(Intercept)</i>	2.58	0.02	140.42	***	-
entropy	-0.04	0.02	-2.67		18.5%
density-TC	1.21	0.25	4.84	***	72.4%
density-TN	-0.33	0.15	-2.16		9.0%
density-TC × density-TN	-0.16	1.40	-0.11		0.0%

Significant codes: 0 '\*\*\*\*' 0.0001 '\*\*\*' 0.001 '\*\*' 0.006 ' ' 1

**Table 6.** Linear regression output for the best predictor combination for verbs in the N-V-A condition: entropy + (density-TC × density-TN). RI indicates the relative importance (normalised to 100%). The significance codes are all adjusted because of the 8 multi-comparisons.

## 5 DISCUSSION

672 The previous section provided a series of vector-space experiments to investigate two conceptual  
673 categorisations of lexical-semantic abstraction (abstractness–concreteness and generality–specificity)  
674 through variants of distributional computational measures. The current section summarises, interprets  
675 and discusses the insights from the empirical experiments with respect to differences in the conceptual  
676 organisation of English nouns and verbs, and the roles of corpus frequency, distributional co-occurrence,  
677 distributional similarity and distributional neighbourhoods for mental distinctions between degrees of  
678 semantic abstraction.

679 Our experiments brought together a variety of distributional vector-space measures that had previously  
680 been applied to different tasks of lexical-semantic abstraction. We focused on the two types of semantic  
681 abstraction originally suggested by Spreen and Schulz (1966) and brought back to attention by Theijssen  
682 *et al.* (2011) and Bolognesi *et al.* (2020). They distinguished abstraction in terms of the abstract–concrete  
683 dichotomy (e.g., *glory* is more abstract than *banana*), and abstraction in terms of the generality–specificity  
684 distinction (e.g., *animal* is more abstract than *fish*). Assuming that a large-scale web corpus provides an  
685 adequate basis for general-language distributional information, we empirically explored corpus frequency  
686 and corpus co-occurrence as proxies to lexical-semantic meaning and lexical meaning relatedness. We  
687 thereby relied on the distributional hypothesis (Harris, 1954; Firth, 1957) indicating that words which are  
688 similar in meaning also occur in similar linguistic distributions.

689 In this vein, we induced variants of neighbourhood densities (context-based and neighbour-based),  
690 token- and type variants of the distributional, vector-based inclusion measure *WeedsPrec*, as well as word  
691 frequency and word entropy, in order to empirically capture noun and verb target words differing in their  
692 degrees of semantic abstraction. We applied these distributional measures to distinguish between degrees  
693 of abstraction regarding the abstract–concrete dichotomy as well as regarding the generality–specificity  
694 distinction. Overall, we identified reliable vector-space measures for both instantiations of lexical-semantic  
695 abstraction (reaching a precision higher than 0.7), but the measures clearly differed for concreteness vs.  
696 hypernymy and for nouns vs. verbs. In order to distinguish between more and less abstract words in  
697 terms of hypernymy, we found that word frequency computed on corpus data, word entropy, and the  
698 distributional inclusion measure (originally suggested for hypernymy) were the most salient predictors,  
699 while neighbourhood density measures could hardly beat the random baseline. In order to distinguish  
700 between more and less abstract words in terms of concreteness, the neighbourhood density measures were  
701 generally more successful than frequency, word entropy and distributional inclusion, especially when  
702 integrating only the strongest contexts/neighbours. Among the density measures the variant that considers  
703 the distributional similarity between a target word and its strongest context words (density-TC) seems  
704 the most appropriate and is also the one with the highest impact in the regression studies. This overall  
705 picture was similar for concreteness ratings for nouns and verbs, but (i) the precision scores for verbs were  
706 generally lower than for nouns and could hardly beat the random baseline, and (ii) frequency, entropy and  
707 weeds-token were not much different from (or even better than) the density variants CC, NN and TN.

708 As a side line of research we explored differences in distinctions between degrees of abstraction regarding  
709 variants of vector spaces in the experimental paradigm. While our main set of experiments did not go into  
710 depth regarding this variable, our full results in the Appendix demonstrate surprisingly clear differences  
711 regarding window size and parts-of-speech of vector dimensions: Results exploiting vector spaces induced  
712 from a co-occurrence window of  $\pm 20$  words (in comparison to only  $\pm 2$  words) and density variants  
713 taking only single-POS words as contexts/neighbours into account generally provided the best results.

714 Whether it was more profitable to rely on noun-only vs. N-V-A (nouns, verbs, adjectives) dimensions in  
715 the co-occurrence vectors depended on the target POS and type of abstraction: For noun concreteness  
716 the N-V-A spaces seemed more indicative, while for verb concreteness and noun and verb specificity the  
717 noun-only spaces were more salient.

718 When zooming into the role of measure-based distinctions according to the strengths of concreteness and  
719 the levels of hypernymy, i.e., hypothesising that the measures are more or less successful with respect to  
720 how “different” the concrete and abstract words are in their degrees of concreteness, and how “different”  
721 the hypernyms and hyponyms are in their degrees of specificity, our insights from the main experiments  
722 were largely confirmed and partially even strengthened: The stronger the differences in concreteness, the  
723 better the quality of distinctions in terms of precision. While this is true for both noun and verb targets, the  
724 picture was again clearer for nouns than for verbs; in the latter case, distinctions for target verbs involving  
725 the mid-range scale of concreteness were worse than those involving any of the extreme ranges. Taking  
726 into account that the concreteness ranges for verbs in the mid-range subsets are rather small ([2.0; 2.3] for  
727 subset 2; [2.3; 2.6] for subset 3; and [2.6; 3.1] for subset 4), this tendency is reasonable because concreteness  
728 scores from different subsets were still rather similar to each other. Also, mid-range concreteness scores  
729 are generally more difficult in their generation by humans and consequently noisier in their distributional  
730 representation (Pollock, 2018). Finally, verbs are generally more ambiguous than nouns, especially when  
731 their semantic properties have been evaluated out of context, and furthermore perception-based concreteness  
732 ratings might not be as appropriate for verbs as they are for nouns. Regarding abstraction measures, our  
733 zooming-in experiments confirmed that the target–context measure density-TC is the best one for predicting  
734 abstraction in terms of concreteness, while frequency, entropy and weeds-token/-type are the best ones for  
735 predicting abstraction in terms of hypernymy.

736 A final study looked into correlations between concreteness and specificity ratings, the abstraction  
737 measure, and their interactions. These correlations confirmed that corpus frequency and word entropy  
738 measure abstraction in a similar way, and ditto for the context-based density measures CC and TC and  
739 the neighbour-based density measures NN and TN (while density-NN seems to differ most from the other  
740 density variants). Moreover, based on a series of regression studies, we confirmed that density-TC is the  
741 strongest option to quantify concreteness both for nouns and for verbs.

742 Bringing together our results across experiments, we can identify two groups of measures, (i) frequency  
743 and word entropy, whose distinctions are correlated and which are stronger than neighbourhood density  
744 measures when distinguishing between more and less abstract words in terms of the generality–specificity  
745 distinction, and (ii) the neighbourhood density variants, which are stronger than group (i) when  
746 distinguishing between more and less abstract words in terms of the abstractness–concreteness dichotomy.  
747 The distributional inclusion variants of WeedsPrec cluster together with frequency and entropy, and are  
748 clearly more useful for hypernymy than for concreteness. Regarding group (i), the relationship between  
749 frequency, word entropy and the lexical-semantic relation hypernymy has been demonstrated before  
750 (Shwartz *et al.*, 2017; Bott *et al.*, 2021), and our experiments confirmed this strong interaction across a  
751 variety of experimental conditions regarding strength of hypernymy. Regarding group (ii), we effectively  
752 and successfully exploited the usefulness of neighbourhood density measures that had previously been  
753 suggested and applied to different instantiations of lexical-semantic abstraction. At the same time we  
754 demonstrated that there are indeed conceptual differences between the measures that result in different  
755 distinction qualities for our two target types of abstraction.

756 Now let us look at these empirical results and insights from a conceptual perspective. First of all, we  
757 can induce from our results that lexical-semantic abstraction in terms of generality in the human lexicon

758 is mirrored by how often we use words, which itself is highly correlated with the words' entropy values.  
759 While this is neither surprising nor novel, one might not have expected such a clear picture over diverse  
760 settings regarding degrees of generality. I.e., more general words are used more often and are therefore  
761 also less surprising. The density measures do not seem appropriate to model the generality–specificity  
762 distinction, thus indicating that they do not capture degrees of semantic relatedness (which is taken into  
763 account by the vector similarity variants of WeedsPrec, for example). Secondly, we can induce from our  
764 results that contextual diversity/neighbourhood density is a strong indicator of lexical-semantic abstraction  
765 in terms of concreteness. Given that density-TC seems to represent the overall most salient measure, we  
766 may induce that abstract words establish themselves empirically in semantically more diverse contexts than  
767 concrete words, thus abstract concepts are lexically connected to more different concepts, while concrete  
768 concepts are lexically connected to less diverse but on the other hand semantically more strongly associated  
769 concepts, and these semantically most indicative associated words are predominantly represented by nouns.  
770 In this vein, lexical entries of abstract and concrete words may be refined with respect to their tendencies to  
771 co-occur with more or less highly distributionally similar, and consequently –according to the distributional  
772 hypothesis– also more or less semantically related words (nouns). The differences in the success of the  
773 abstraction measures regarding our two target types of semantic abstraction seems directly related to a core  
774 distinction: while words differing in their degree of concreteness are not necessarily semantically related  
775 (e.g., *glory–banana*), words differing in their degree of specificity (e.g., *animal–fish*) are, at least with  
776 regard to hypernymy in WordNet. Overall, our insights should generally be useful for computational models  
777 exploiting degrees of semantic abstraction, such as standard classification approaches and topic models,  
778 and similarly for more complex computational systems where the degree of contextual abstraction plays a  
779 role, such as figurative language detection, text simplification, summarisation, and machine translation.

780 Our experiments also point out once more that distributional measures, distributional similarity and  
781 distributional semantic relatedness differ across word classes. On the one hand, concreteness and hypernymy  
782 represent two lexical-semantic types of abstraction, and therefore their organisation is also defined in  
783 different ways in the respective resources. I.e., concreteness scores had been collected on a word-type  
784 basis, where participants were not provided a part-of-speech categorisation and part-of-speech tags were  
785 assigned post-hoc. Even though we applied a rather restrictive procedure to POS label identification and  
786 discarded ambiguous words, this basis is sub-optimal for any word-class-dependent analyses: we calculated  
787 Spearman's  $\rho$  correlation for the POS assignment based on SUBTLEX (Brysbaert *et al.*, 2012) and our  
788 ENCOW-based procedure, obtaining  $\rho=0.624$  for our noun targets and  $\rho=0.750$  for our verb targets, which  
789 we consider as rather low and pointing to an undesired disagreement in POS assignment. On the other hand,  
790 all our studies have been on a type-basis: vector spaces and concreteness ratings are type-based, and while  
791 WordNet does distinguish between word senses, we only indirectly used this option, because we utilised all  
792 senses in word pairs, but we did not distinguish between senses. This is more crucial for verbs than for  
793 nouns, which are notoriously more ambiguous. Overall, future work should therefore target contextualised,  
794 token-based distributional representations and sense-based abstraction ratings.

## 6 CONCLUSION

795 In this article, we provided a series of empirical studies that investigated conceptual categories of semantic  
796 abstraction through distributional variants of abstraction measures. We distinguished abstraction in terms of  
797 the abstract–concrete dichotomy and in terms of the generality–specificity distinction, and brought together  
798 a variety of distributional measures that had previously been applied to different tasks of lexical-semantic  
799 abstraction. We thus suggested a novel perspective that exploited empirical measures across two types of

800 semantic abstraction, in order to compare the strengths and weaknesses of the measures for categorisations  
801 of abstraction, and to determine and investigate conceptual differences as captured by the measures.

802 In a series of experiments we identified reliable vector-space measures for both instantiations of lexical-  
803 semantic abstraction (reaching a precision of  $>0.7$ ), and we demonstrated that the measures clearly  
804 differed for concreteness vs. hypernymy and for nouns vs. verbs. We could identify two groups of  
805 measures, (i) frequency, word entropy and weeds-token/-type when distinguishing between more and  
806 less abstract words in terms of the generality–specificity distinction, and (ii) the neighbourhood density  
807 variants (especially target–context diversity, with nouns providing the most salient context words) when  
808 distinguishing between more and less abstract words in terms of the abstractness–concreteness dichotomy.  
809 We concluded that more general words are used more often and are therefore also less surprising than  
810 more specific words, and that abstract words establish themselves empirically in semantically more diverse  
811 contexts than concrete words, i.e., abstract concepts are lexically connected to more different concepts,  
812 while concrete concepts are lexically connected to less diverse but at the same time semantically more  
813 strongly associated concepts.

814 Finally, we demonstrated the need to take word classes and ambiguity into account. On the one hand,  
815 results for nouns vs. verbs clearly differ, and both ratings and vector spaces should take semantic differences  
816 between word classes into account; on the other hand, ambiguity (which is more severe for verbs than for  
817 nouns) prevents from fine-tuning empirical observations and conclusions.

## 7 ACKNOWLEDGEMENTS

818 We would like to acknowledge our students Simone Beckmann Escandon, Maximilian Bräuninger, Christos  
819 Lontos and Daniela Naumann for their great help in performing some of the initial studies for this research,  
820 for being always up for interesting discussions and for their strenuous ability to survive the frequent intense  
821 discussions between the authors of this paper.

## REFERENCES

- 822 Adelman, J. S., Brown, G. D., and Quesada, J. F. (2006). Contextual Diversity, Not Word Frequency,  
823 Determines Word-Naming and Lexical Decision Times. *Psychological Science*, **17**(9), 814–823.
- 824 Aedmaa, E., Köper, M., and Schulte im Walde, S. (2018). Combining Abstractness and Language-specific  
825 Theoretical Indicators for Detecting Non-Literal Usage of Estonian Particle Verbs. In *Proceedings of the*  
826 *NAACL 2018 Student Research Workshop*, pages 9–16, New Orleans, LA, USA.
- 827 Algarabel, S., Ruiz, J. C., and Sanmartin, J. (1988). The University of Valencia's Computerized Word Pool.  
828 *Behavior Research Methods, Instruments, and Computers*, **20**(4), 398–403.
- 829 Baroni, M. and Lenci, A. (2010). Distributional Memory: A General Framework for Corpus-based  
830 Semantics. *Computational Linguistics*, **36**(4), 673–721.
- 831 Barsalou, L. W. (2003). Abstraction in Perceptual Symbol Systems. *Philosophical Transactions of the*  
832 *Royal Society London B*, **358**, 1177–1187.
- 833 Barsalou, L. W. and Wiemer-Hastings, K. (2005). Situating Abstract Concepts. In D. Pecher and R. Zwaan,  
834 editors, *Grounding Cognition: The Role of Perception and Action in Memory, Language, and Thinking*,  
835 chapter 7, pages 129–163. Cambridge University Press, New York.
- 836 Bestgen, Y. and Vincze, N. (2012). Checking and Bootstrapping Lexical Norms by Means of Word  
837 Similarity Indexes. *Behavior Research Methods*, **44**, 998–1006.

- 838 Bolognesi, M., Burgers, C., and Caselli, T. (2020). On Abstraction: Decoupling Conceptual Concreteness  
839 and Categorical Specificity. *Cognitive Processing*, **21**, 365–381.
- 840 Bonin, P., Meot, A., and Bugaiska, A. (2018). Concreteness Norms for 1,659 French Words: Relationships  
841 with other Psycholinguistic Variables and Word Recognition Times. *Behavior Research Methods*, **50**,  
842 2366–2387.
- 843 Borghi, A. M., Binkofski, F., Castelfranchi, C., Cimatti, F., Scorolli, C., and Tummolini, L. (2017). The  
844 Challenge of Abstract Concepts. *Psychonomic Bulletin*, **143**, 263–292.
- 845 Bott, T., Schlechtweg, D., and Schulte im Walde, S. (2021). More than just Frequency? Demasking  
846 Unsupervised Hypernymy Prediction Methods. In *Findings of the Association for Computational*  
847 *Linguistics: ACL-IJCNLP*, pages 186–192, Bangkok, Thailand (online).
- 848 Bradley, M. M. and Lang, P. J. (1999). Affective Norms for English Words (ANEW): Instruction Manual  
849 and Affective Ratings. Technical Report C-1, The Center for Research in Psychophysiology, University  
850 of Florida, Philadelphia, PA.
- 851 Brysbaert, M., New, B., and Keuleers, E. (2012). Adding Part-of-Speech Information to the SUBTLEX-US  
852 Word Frequencies. *Behavior Research Methods*, **44**, 991–997.
- 853 Brysbaert, M., Warriner, A. B., and Kuperman, V. (2014). Concreteness Ratings for 40 Thousand generally  
854 known English Word Lemmas. *Behavior Research Methods*, **64**, 904–911.
- 855 Burgess, C. (1998). From Simple Associations to the Building Blocks of Language: Modeling Meaning in  
856 Memory with the HAL Model. *Behavior Research Methods, Instruments and Computers*, **30**, 188–198.
- 857 Burgoon, E. M., Henderson, M. D., and Markman, A. B. (2013). There Are Many Ways to See the Forest  
858 for the Trees: A Tour Guide for Abstraction. *Perspectives on Psychological Science*, **8**(5), 501–520.
- 859 Cimiano, P., Schmidt-Thieme, L., Pivk, A., and Staab, S. (2004). Learning Taxonomic Relations from  
860 Heterogeneous Evidence. In *Proceedings of the ECAI Workshop on Ontology Learning and Population*,  
861 Valencia, Spain.
- 862 Clark, P., Murray, W. R., Thompson, J., Harrison, P., Hobbs, J., and Fellbaum, C. (2007). On the Role of  
863 Lexical and World Knowledge in RTE3. In *Proceedings of the Workshop on Textual Entailment and*  
864 *Paraphrasing*, pages 54–59, Prague, Czech Republic.
- 865 Coltheart, M. (1981). The MRC Psycholinguistic Database. *Quarterly Journal of Experimental Psychology*,  
866 **33A**, 497–505.
- 867 Cruse, D. A. (1986). *Lexical Semantics*. Cambridge Textbooks in Linguistics. Cambridge University Press,  
868 Cambridge, UK.
- 869 Crutch, S. J. and Warrington, E. K. (2010). The Differential Dependence of Abstract and Concrete  
870 Words upon Associative and Similarity-based Information: Complementary Semantic Interference and  
871 Facilitation Effects. *Cognitive Neuropsychology*, **27**, 46–71.
- 872 Dagan, I., Glickman, O., and Magnini, B. (2006). The PASCAL Recognising Textual Entailment Challenge.  
873 *Lecture Notes in Computer Science*, **3944**, 177–190.
- 874 Danguécan, A. N. and Buchanan, L. (2016). Semantic Neighborhood Effects for Abstract versus Concrete  
875 Words. *Frontiers in Psychology*, **7**(1034).
- 876 Darley, F. L., Sherman, D., and Siegel, G. M. (1959). Scaling of Abstraction Level of Single Words.  
877 *Journal of Speech and Hearing Research*, **2**(2), 161–167.
- 878 Della Rosa, P. A., Catricala, E., Vigliocco, G., and Cappa, S. F. (2010). Beyond the Abstract–Concrete  
879 Dichotomy: Mode of Acquisition, Concreteness, Imageability, Familiarity, Age of Acquisition, Context  
880 Availability, and Abstractness Norms for a Set of 417 Italian Words. *Behavior Research Methods*, **42**(4),  
881 1042–1048.

- 882 Dinu, G., Lazaridou, A., and Baroni, M. (2015). Improving Zero-Shot Learning by Mitigating the Hubness  
883 Problem. In *Proceedings of the International Conference on Learning Representations, Workshop Track*,  
884 San Diego, CA, USA.
- 885 Evert, S. (2005). *The Statistics of Word Co-Occurrences: Word Pairs and Collocations*. Ph.D. thesis,  
886 Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart.
- 887 Fellbaum, C. (1990). English Verbs as a Semantic Net. *Journal of Lexicography*, **3**(4), 278–301.
- 888 Fellbaum, C. (1998a). A Semantic Network of English Verbs. In Fellbaum (1998b), pages 69–104.
- 889 Fellbaum, C., editor (1998b). *WordNet – An Electronic Lexical Database*. Language, Speech, and  
890 Communication. MIT Press, Cambridge, MA, USA.
- 891 Fellbaum, C. and Chaffin, R. (1990). Some Principles of the Organization of Verbs in the Mental Lexicon.  
892 In *Proceedings of the 12th Annual Conference of the Cognitive Science Society of America*, pages  
893 420–427, Cambridge, MA, USA.
- 894 Firth, J. R. (1957). *Papers in Linguistics 1934-51*. Longmans, London, UK.
- 895 Frassinelli, D. and Lenci, A. (2012). Concepts in Context: Evidence from a Feature-Norming Study. In  
896 *Proceedings of the 34th Annual Meeting of the Cognitive Science Society*, Sapporo, Japan.
- 897 Frassinelli, D. and Schulte im Walde, S. (2019). Distributional Interaction of Concreteness and Abstractness  
898 in Verb–Noun Subcategorisation. In *Proceedings of the 13th International Conference on Computational*  
899 *Semantics*, pages 38–43, Gothenburg, Sweden.
- 900 Frassinelli, D., Naumann, D., Utt, J., and Schulte im Walde, S. (2017). Contextual Characteristics of  
901 Concrete and Abstract Words. In *Proceedings of the 12th International Conference on Computational*  
902 *Semantics*, Montpellier, France.
- 903 Glenberg, A. M. and Robertson, D. A. (2000). Symbol Grounding and Meaning: A Comparison of  
904 High-Dimensional and Embodied Theories of Meaning. *Journal of Memory and Language*, **43**, 379–401.
- 905 Gorman, A. M. (1961). Recognition Memory for Nouns as a Function of Abstractness and Frequency.  
906 *Journal of Experimental Psychology*, **61**, 23–29.
- 907 Gross, D. and Miller, K. J. (1990). Adjectives in wordnet. *International Journal of Lexicography*, **3**(4),  
908 265–277.
- 909 Hare, M., Jones, M., Thomson, C., Kelly, S., and McRae, K. (2009). Activating Event Knowledge.  
910 *Cognition*, **111**, 151–167.
- 911 Harris, Z. (1954). Distributional structure. *Word*, **10**(23), 146–162.
- 912 Hearst, M. (1992). Automatic Acquisition of Hyponyms from Large Text Corpora. In *Proceedings of the*  
913 *14th International Conference on Computational Linguistics*, pages 539–545, Nantes, France.
- 914 Hearst, M. (1998). Automated Discovery of WordNet Relations. In Fellbaum (1998b).
- 915 Hill, F., Korhonen, A., and Bentz, C. (2014). A Quantitative Empirical Analysis of the Abstract/Concrete  
916 Distinction. *Cognitive Science*, **38**, 162–177.
- 917 Hoffman, P. and Woollams, A. M. (2015). Opposing Effects of Semantic Diversity in Lexical and Semantic  
918 Relatedness Decisions. *Journal of Experimental Psychology: Human Perception and Performance*,  
919 **41**(2), 385–402.
- 920 Hoffman, P., Lambon Ralph, M. A., and Rogers, T. T. (2013). Semantic Diversity: A Measure of Semantic  
921 Ambiguity based on Variability in the Contextual Usage of Words. *Behavior Research Methods*, **45**,  
922 718–730.
- 923 Kanske, P. and Kotz, S. A. (2010). Leipzig Affective Norms for German: A Reliability Study. *Behavior*  
924 *Research Methods*, **42**(4), 987–991.

- 925 Köper, M. and Schulte im Walde, S. (2016). Automatically Generated Affective Norms of Abstractness,  
926 Arousal, Imageability and Valence for 350 000 German Lemmas. In *Proceedings of the 10th International*  
927 *Conference on Language Resources and Evaluation*, pages 2595–2598, Portoroz, Slovenia.
- 928 Köper, M. and Schulte im Walde, S. (2017). Improving Verb Metaphor Detection by Propagating  
929 Abstractness to Words, Phrases and Individual Senses. In *Proceedings of the 1st Workshop on Sense,*  
930 *Concept and Entity Representations and their Applications*, pages 24–30, Valencia, Spain.
- 931 Köper, M., Schulte im Walde, S., Kisselew, M., and Padó, S. (2016). Improving Zero-Shot-Learning for  
932 German Particle Verbs by using Training-Space Restrictions and Local Scaling. In *Proceedings of the*  
933 *5th Joint Conference on Lexical and Computational Semantics*, pages 91–96, Berlin, Germany.
- 934 Kotlerman, L., Dagan, I., Szpektor, I., and Zhitomirsky-Geffet, M. (2010). Directional Distributional  
935 Similarity for Lexical Inference. *Natural Language Engineering*, **16**(4), 359–389.
- 936 Kousta, S.-T., Vigliocco, G., Vinson, D. P., Andrews, M., and Del Campo, E. (2011). The Representation of  
937 Abstract Words: Why Emotion Matters. *Journal of Experimental Psychology: General*, **140**(1), 14–34.
- 938 Lahl, O., Göritz, A. S., Pietrowsky, R., and Rosenberg, J. (2009). Using the World-Wide Web to obtain  
939 Large-Scale Word Norms: 190,212 Ratings on a Set of 2,654 German Nouns. *Behavior Research*  
940 *Methods*, **41**(1), 13–19.
- 941 Lenci, A. and Benotto, G. (2012). Identifying Hypernyms in Distributional Semantic Spaces. In *Proceedings*  
942 *of the 1st Joint Conference on Lexical and Computational Semantics*, pages 75–79, Montréal, Canada.
- 943 Lenci, A., Lebani, G. E., and Passaro, L. C. (2018). The Emotions of Abstract Words: A Distributional  
944 Semantic Analysis. *Topics in Cognitive Science*, **10**, 550–572.
- 945 Lindeman, R. H., Merenda, P., and Gold, R. (1980). *Introduction to Bivariate and Multivariate Analysis*,  
946 volume 119. Glenview, IL.
- 947 Lynott, D. and Connell, L. (2009). Modality Exclusivity Norms for 423 Object Properties. *Behavior*  
948 *Research Methods*, **41**(2), 558–564.
- 949 Lynott, D. and Connell, L. (2013). Modality Exclusivity Norms for 400 Nouns: The Relationship between  
950 Perceptual Experience and Surface Word Form. *Behavior Research Methods*, **45**, 516–526.
- 951 Lynott, D., Connell, L., Brysbaert, M., Brand, J., and Carney, J. (2020). The Lancaster Sensorimotor Norms:  
952 Multidimensional Measures of Perceptual and Action Strength for 40,000 English Words. *Behavior*  
953 *Research Methods*, **52**, 1–21.
- 954 Mandera, P., Keuleers, E., and Brysbaert, M. (2015). How useful are Corpus-based Methods for  
955 Extrapolating Psycholinguistic Variables? *The Quarterly Journal of Experimental Psychology*, **68**(8),  
956 1623–1642.
- 957 McDonald, S. A. and Shillcock, R. C. (2001). Rethinking the Word Frequency Effect: The Neglected Role  
958 of Distributional Information in Lexical Processing. *Language and Speech*, **44**(3), 295–323.
- 959 Miller, G. A. and Fellbaum, C. (1991). Semantic Networks of English. *Cognition*, **41**, 197–229.
- 960 Miller, G. A., Beckwith, R., Fellbaum, C., Gross, D., and Miller, K. J. (1990). Introduction to Wordnet: An  
961 On-line Lexical Database. *International Journal of Lexicography*, **3**(4), 235–244.
- 962 Mohammad, S. M. (2018). Obtaining Reliable Human Ratings of Valence, Arousal, and Dominance for  
963 20,000 English Words. In *Proceedings of the 56th Annual Meeting of the Association for Computational*  
964 *Linguistics*, Melbourne, Australia.
- 965 Murphy, M. L. (2003). *Semantic Relations and the Lexicon*. Cambridge University Press.
- 966 Naumann, D., Frassinelli, D., and Schulte im Walde, S. (2018). Quantitative Semantic Variation in the  
967 Contexts of Concrete and Abstract Words. In *Proceedings of the 7th Joint Conference on Lexical and*  
968 *Computational Semantics*, pages 76–85, New Orleans, LA, USA.



- 969 Navigli, R. and Ponzetto, S. (2012). BabelNet: The Automatic Construction, Evaluation and Application  
970 of a Wide-Coverage Multilingual Semantic Network. *Artificial Intelligence*, **193**, 217–250.
- 971 Nguyen, K. A., Köper, M., Schulte im Walde, S., and Vu, N. T. (2017). Hierarchical Embeddings for  
972 Hypernymy Detection and Directionality. In *Proceedings of the Conference on Empirical Methods in*  
973 *Natural Language Processing*, pages 233–243, Copenhagen, Denmark.
- 974 Paivio, A. (1971). Imagery and Language. In S. J. Segal, editor, *Imagery: Current Cognitive Approaches*,  
975 pages 7–32. Academic Press, New York and London.
- 976 Paivio, A. and Begg, I. (1971). Imagery and Comprehension Latencies as a Function of Sentence  
977 Concreteness and Structure. *Perception and Psychophysics*, **10**(6), 408–412.
- 978 Paivio, A., Yuille, J. C., and Madigan, S. A. (1968). Concreteness, Imagery, and Meaningfulness Values  
979 for 925 Nouns. *Journal of Experimental Psychology (Monograph Supplement)*, **76**(1/2), 1–25.
- 980 Pantel, P. and Pennacchiotti, M. (2006). Espresso: Leveraging Generic Patterns for Automatically  
981 Harvesting Semantic Relations. In *Proceedings of the 21st International Conference on Computational*  
982 *Linguistics and the 44th Annual Meeting of the Association for Computational Linguistics*, pages  
983 113–120, Sydney, Australia.
- 984 Pecher, D., Boot, I., and Van Dantzig, S. (2011). Abstract Concepts. Sensory-Motor Grounding, Metaphors,  
985 and Beyond. *Psychology of Learning and Motivation – Advances in Research and Theory*, **54**, 217–248.
- 986 Pollock, L. (2018). Statistical and Methodological Problems with Concreteness and other Semantic  
987 Variables: A List Memory Experiment Case Study. *Behavior Research Methods*, **50**, 1198–1216.
- 988 Rabinovich, E., Sznajder, B., Spector, A., Shnayderman, I., Aharonov, R., Konopnicki, D., and Slonim, N.  
989 (2018). Learning Concept Abstractness Using Weak Supervision. arXiv:1809.01285.
- 990 Radovanović, M., Nanopoulos, A., and Ivanović, M. (2010). Hubs in Space: Popular Nearest Neighbors in  
991 High-Dimensional Data. *Journal of Machine Learning Research*, **11**, 2487–2531.
- 992 Recchia, G. and Jones, M. N. (2012). The Semantic Richness of Abstract Concepts. *Frontiers in Human*  
993 *Neuroscience*, **6**(315).
- 994 Recchia, G. and Louwerse, M. M. (2015). Reproducing Affective Norms with Lexical Co-Occurrence  
995 Statistics: Predicting Valence, Arousal, and Dominance. *The Quarterly Journal of Experimental*  
996 *Psychology*, **68**(8), 1584–1598.
- 997 Reilly, M. and Desai, R. H. (2017). Effects of Semantic Neighborhood Density in Abstract and Concrete  
998 Words. *Cognition*, **169**, 46–53.
- 999 Richens, T. (2008). Anomalies in the WordNet Verb Hierarchy. In *Proceedings of the 22nd International*  
1000 *Conference on Computational Linguistics*, pages 729–736, Manchester, UK.
- 1001 Rimell, L. (2014). Distributional Lexical Entailment by Topic Coherence. In *Proceedings of the 14th*  
1002 *Conference of the European Chapter of the Association for Computational Linguistics*, pages 511–519,  
1003 Gothenburg, Sweden.
- 1004 Roth, M. and Schulte im Walde, S. (2014). Combining Word Patterns and Discourse Markers for  
1005 Paradigmatic Relation Classification. In *Proceedings of the 52nd Annual Meeting of the Association for*  
1006 *Computational Linguistics*, pages 524–530, Baltimore, MD, USA.
- 1007 Sagi, E., Kaufmann, S., and Clark, B. (2009). Semantic Density Analysis: Comparing Word Meaning  
1008 across Time and Phonetic Space. In *Proceedings of the EMNLP Workshop on Geometrical Models for*  
1009 *Natural Language Semantics*, pages 104–111, Athens, Greece.
- 1010 Salton, G., Wong, A., and Yang, C.-S. (1975). A Vector Space Model for Automatic Indexing.  
1011 *Communications of the ACM*, **18**(11), 613–620.

- 1012 Santus, E., Lenci, A., Lu, Q., and Schulte im Walde, S. (2014). Chasing Hypernyms in Vector Spaces  
1013 with Entropy. In *Proceedings of the 14th Conference of the European Chapter of the Association for*  
1014 *Computational Linguistics*, pages 38–42, Gothenburg, Sweden.
- 1015 Santus, E., Lenci, A., Chiu, T.-S., Lu, Q., and Huang, C.-R. (2016). Unsupervised Measure of Word  
1016 Similarity: How to Outperform Cooccurrence and Vector Cosine in VSMs. In *Proceedings of the 13th*  
1017 *AAAI Conference on Artificial Intelligence*, pages 4260–4261, Phoenix, Arizona, USA.
- 1018 Schäfer, R. (2015). Processing and Querying Large Web Corpora with the COW14 Architecture. In  
1019 *Proceedings of the 3rd Workshop on Challenges in the Management of Large Corpora*, pages 28–34,  
1020 Mannheim, Germany.
- 1021 Schäfer, R. and Bildhauer, F. (2012). Building Large Corpora from the Web Using a New Efficient Tool  
1022 Chain. In *Proceedings of the 8th International Conference on Language Resources and Evaluation*,  
1023 pages 486–493, Istanbul, Turkey.
- 1024 Schlechtweg, D., Eckmann, S., Santus, E., Schulte im Walde, S., and Hole, D. (2017). German in  
1025 Flux: Detecting Metaphoric Change via Word Entropy. In *Proceedings of the 21st Conference on*  
1026 *Computational Natural Language Learning*, pages 354–367, Vancouver, Canada.
- 1027 Schulte im Walde, S. (2020). Distinguishing between Paradigmatic Semantic Relations across Word  
1028 Classes: Human Ratings and Distributional Similarity. *Journal of Language Modelling*, **8**(1), 53–101.
- 1029 Schulte im Walde, S. and Köper, M. (2013). Pattern-based Distinction of Paradigmatic Relations for  
1030 German Nouns, Verbs, Adjectives. In I. Gurevych, C. Biemann, and T. Zesch, editors, *Language*  
1031 *Processing and Knowledge in the Web. Proceedings of the 25th International Conference of the German*  
1032 *Society for Computational Linguistics and Language Technology*, volume 8105 of *Lecture Notes in*  
1033 *Computer Science*, pages 184–198. Springer.
- 1034 Schwanenflugel, P. J. and Shoben, E. J. (1983). Differential Context Effects in the Comprehension of  
1035 Abstract and Concrete Verbal Materials. *Journal of Experimental Psychology: Learning, Memory, and*  
1036 *Cognition*, **9**(1), 82–102.
- 1037 Shwartz, V., Santus, E., and Schlechtweg, D. (2017). Hypernyms under Siege: Linguistically-motivated  
1038 Artillery for Hypernymy Detection. In *Proceedings of the 15th Conference of the European Chapter of*  
1039 *the Association for Computational Linguistics*, pages 65–75, Valencia, Spain.
- 1040 Siegel, S. and Castellan, N. J. (1988). *Nonparametric Statistics for the Behavioral Sciences*. McGraw-Hill,  
1041 Boston, MA, USA.
- 1042 Snow, R., Jurafsky, D., and Ng, A. Y. (2006). Semantic Taxonomy Induction from Heterogenous Evidence.  
1043 In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics*, pages  
1044 801–808, Sydney, Australia.
- 1045 Spreen, O. and Schulz, R. W. (1966). Parameters of Abstraction, Meaningfulness, and Pronunciability for  
1046 329 Nouns. *Journal of Verbal Learning Behavior*, **5**, 459–468.
- 1047 Theijssen, D., van Halteren, H., Boves, L., and Oostdijk, N. (2011). On the Difficulty of making  
1048 Concreteness Concrete. *Computational Linguistics in the Netherlands Journal*, **1**, 61–77.
- 1049 Troche, J., Crutch, S., and Reilly, J. (2014). Clustering, Hierarchical Organization, and the Topography of  
1050 Abstract and Concrete Nouns. *Frontiers in Psychology*, **5**(360).
- 1051 Turney, P. D. and Pantel, P. (2010). From Frequency to Meaning: Vector Space Models of Semantics.  
1052 *Journal of Artificial Intelligence Research*, **37**, 141–188.
- 1053 Vigliocco, G., Meteyard, L., Andrews, M., and Kousta, T. (2009). Toward a Theory of Semantic  
1054 Representation. *Language and Cognition*, **1**(2), 219–247.

- 1055 Vigliocco, G., Kousta, S.-T., Anthony Della Rosa, P., Vinson, D. P., Tettamanti, M., Devlin, J. T., and  
1056 Cappa, S. F. (2014). The Neural Representation of Abstract Words: The Role of Emotion. *Cerebral*  
1057 *Cortex*, **24**, 1767–1777.
- 1058 Warriner, A. B., Kuperman, V., and Brysbaert, M. (2013). Norms of Valence, Arousal, and Dominance for  
1059 13,915 English Lemmas. *Behavior Research Methods*, **45**(4), 1191–1207.
- 1060 Weeds, J. and Weir, D. (2005). A Flexible Framework for Lexical Distributional Similarity. *Computational*  
1061 *Linguistics*, **31**(4), 439–476.
- 1062 Weeds, J., Weir, D., and McCarthy, D. (2004). Characterising Measures of Lexical Distributional Similarity.  
1063 In *Proceedings of the 20th International Conference of Computational Linguistics*, pages 1015–1021,  
1064 Geneva, Switzerland.
- 1065 Weeds, J., Clarke, D., Reffin, J., Weir, D., and Keller, B. (2014). Learning to Distinguish Hypernyms  
1066 and Co-Hyponyms. In *Proceedings of the 25th International Conference of Computational Linguistics*,  
1067 pages 2249–2259, Dublin, Ireland.
- 1068 Wiemer-Hastings, K. and Xu, X. (2005). Content Differences for Abstract and Concrete Concepts.  
1069 *Cognitive Science*, **29**, 719–736.
- 1070 Yap, W. and Baldwin, T. (2009). Experiments on Pattern-based Relation Learning. In *Proceedings of the*  
1071 *18th ACM Conference on Information and Knowledge Management*, pages 1657—1660, Hong Kong.

## 1 EXAMPLES: CONTEXT AND NEIGHBOUR WORDS

1072 Table 7 shows the five strongest context and neighbour words for a small subset of noun and verb targets,  
 1073 in order to get an impression of conceptual differences between context and neighbour words. Note that  
 1074 strongest noun context words as used in the density-CC and density-TC variants have been selected based  
 1075 on target–context plmi scores, and that strongest nearest neighbours as used in the density-NN and density-  
 1076 TN variants have been selected based on target–neighbour cosine scores, here showing the respective noun  
 1077 neighbours for noun targets and verb neighbours for verb targets.

targets		mean ratings	strongest contexts			strongest neighbours	
			word	plmi	cosine	word	cosine
N concrete	wine	4.79	bottle	174,811	0.75	vino	0.83
			glass	158,713	0.63	demijohn	0.81
			beer	92,498	0.58	rosé	0.81
			grape	69,048	0.65	sommelier	0.79
			food	55,781	0.18	tasting	0.79
	trout	4.72	fishing	45,436	0.55	grayling	0.81
			salmon	31,941	0.70	steelhead	0.77
			rainbow	28,065	0.62	salmon	0.70
			fish	19,793	0.38	whitefish	0.68
			lake	14,159	0.41	kokanee	0.65
N abstract	wisdom	1.53	knowledge	33,767	0.29	fount	0.51
			word	21,322	0.20	foolishness	0.47
			man	14,678	0.11	prajna	0.46
			love	12,914	0.20	sagacity	0.43
			power	10,539	0.12	folly	0.41
	sensibility	1.52	sense	6,559	0.12	aesthetic	0.43
			film	2,520	0.19	humor	0.42
			pop	2,347	0.19	expansiveness	0.38
			sensuality	2,277	0.32	rootlessness	0.38
			art	1,589	0.21	purposefulness	0.38
V concrete	sit	4.80	room	152,949	0.46	seat	0.68
			table	144,106	0.57	plop	0.61
			chair	134,806	0.63	scoot	0.61
			front	75,121	0.42	slouch	0.51
			seat	71,815	0.30	plonk	0.50
	breathe	4.07	air	64,932	0.51	humidify	0.59
			sigh	38,937	0.42	dehumidify	0.58
			life	27,472	0.25	condition	0.56
			relief	24,369	0.29	rarefy	0.55
			breath	23,892	0.41	gasp	0.54
V abstract	expect	1.89	result	29,932	0.22	anticipate	0.60
			price	29,408	0.35	forecast	0.43
			week	28,254	0.26	come	0.38
			level	28,192	0.28	rise	0.37
			month	27,556	0.33	disappoint	0.36
	overrate	1.86	player	1,529	0.33	underrate	0.73
			film	616	0.14	cogitate	0.70
			opinion	521	0.12	crystalize	0.70
			game	448	0.18	mistake	0.61
			bit	394	0.19	delude	0.59

**Table 7.** Strongest context and neighbour words for a selection of target nouns and verbs.

## 2 FULL TABLES OF RESULTS

1078 Tables 8–11 provide the full results for pair-wise distinctions between degrees of abstraction in  
 1079 terms of concreteness and hypernymy, both for nouns and for verbs. We applied symmetric co-  
 1080 occurrence windows of  $\pm 2$  and  $\pm 20$  words; vector spaces including only co-occurring nouns (space: N)  
 1081 vs. nouns/verbs/adjectives (space: N-V-A); and density variants taking only nouns or verbs or  
 1082 nouns/verbs/adjectives (all) as context/neighbour words into account. The results show precision scores in  
 1083 combination with the number of pairs for which the distinctions were made. The best result per package is  
 1084 highlighted.

	window 2		window 20	
	space: N	space: N-V-A	space: N	space: N-V-A
baseline: frequency	0.4574 (250,000)			
weeds-token	0.3797 (166,457)	0.4263 (245,173)	0.3642 (250,000)	0.4157 (250,000)
weeds-type	0.4243 (166,457)	0.4330 (245,173)	0.4673 (250,000)	0.4758 (250,000)
entropy	0.4451 (249,000)	0.4355 (250,000)	0.5255 (250,000)	0.5230 (250,000)
density-CC-5 (nouns)	0.6833 (247,000)	0.6663 (247,000)	0.6965 (250,000)	0.7087 (250,000)
density-CC-5 (all)	0.6513 (250,000)	0.6567 (250,000)	0.6798 (250,000)	0.7044 (250,000)
density-CC-10 (nouns)	0.6863 (247,000)	0.6707 (247,000)	0.7142 (250,000)	0.7272 (250,000)
density-CC-10 (all)	0.6524 (250,000)	0.6554 (250,000)	0.6900 (250,000)	0.7150 (250,000)
density-CC-20 (nouns)	0.6878 (247,000)	0.6505 (247,000)	0.7088 (250,000)	0.7244 (250,000)
density-CC-20 (all)	0.6257 (250,000)	0.6417 (250,000)	0.6648 (250,000)	0.6979 (250,000)
density-CC-50 (nouns)	0.6479 (247,000)	0.5673 (247,000)	0.6395 (250,000)	0.6547 (250,000)
density-CC-50 (all)	0.5647 (250,000)	0.5823 (250,000)	0.5784 (250,000)	0.6233 (250,000)
density-TC-5 (nouns)	0.5713 (248,000)	0.5882 (249,000)	0.7068 (250,000)	0.7799 (250,000)
density-TC-5 (all)	0.6037 (248,500)	0.6475 (250,000)	0.7357 (250,000)	0.7740 (250,000)
density-TC-10 (nouns)	0.5834 (248,000)	0.6066 (249,000)	0.7235 (250,000)	0.7930 (250,000)
density-TC-10 (all)	0.6171 (249,000)	0.6572 (250,000)	0.7391 (250,000)	0.7777 (250,000)
density-TC-20 (nouns)	0.5904 (248,000)	0.6108 (249,000)	0.7200 (250,000)	0.7870 (250,000)
density-TC-20 (all)	0.6144 (249,000)	0.6647 (250,000)	0.7147 (250,000)	0.7690 (250,000)
density-TC-50 (nouns)	0.5874 (248,000)	0.6002 (249,000)	0.6962 (250,000)	0.7613 (250,000)
density-TC-50 (all)	0.6019 (249,000)	0.6520 (250,000)	0.6698 (250,000)	0.7318 (250,000)
density-NN-5 (nouns)	0.5160 (249,000)	0.4931 (249,000)	0.6541 (250,000)	0.6296 (250,000)
density-NN-5 (all)	0.5028 (250,000)	0.5002 (250,000)	0.6311 (250,000)	0.6249 (250,000)
density-NN-10 (nouns)	0.5053 (249,000)	0.4804 (249,000)	0.6608 (250,000)	0.6380 (250,000)
density-NN-10 (all)	0.4944 (250,000)	0.4888 (250,000)	0.6229 (250,000)	0.6185 (250,000)
density-NN-20 (nouns)	0.4779 (249,000)	0.4501 (249,000)	0.6453 (250,000)	0.6397 (250,000)
density-NN-20 (all)	0.4795 (250,000)	0.4684 (250,000)	0.6185 (250,000)	0.6181 (250,000)
density-NN-50 (nouns)	0.4683 (249,000)	0.4247 (249,000)	0.6188 (250,000)	0.6276 (250,000)
density-NN-50 (all)	0.4480 (250,000)	0.4409 (250,000)	0.5815 (250,000)	0.6015 (250,000)
density-TN-5 (nouns)	0.4995 (249,000)	0.4898 (249,000)	0.7325 (250,000)	0.7350 (250,000)
density-TN-5 (all)	0.4921 (250,000)	0.4930 (250,000)	0.7031 (250,000)	0.7224 (250,000)
density-TN-10 (nouns)	0.5005 (249,000)	0.4818 (249,000)	0.7228 (250,000)	0.7246 (250,000)
density-TN-10 (all)	0.4916 (250,000)	0.4885 (250,000)	0.6892 (250,000)	0.7090 (250,000)
density-TN-20 (nouns)	0.4910 (249,000)	0.4655 (249,000)	0.7065 (250,000)	0.7102 (250,000)
density-TN-20 (all)	0.4824 (250,000)	0.4764 (250,000)	0.6685 (250,000)	0.6913 (250,000)
density-TN-50 (nouns)	0.4749 (249,000)	0.4418 (249,000)	0.6665 (250,000)	0.6780 (250,000)
density-TN-50 (all)	0.4641 (250,000)	0.4539 (250,000)	0.6266 (250,000)	0.6595 (250,000)

**Table 8.** Full results for pair-wise distinctions between degrees of concreteness: nouns.

	window 2				window 20			
	space: N		space: N-V-A		space: N		space: N-V-A	
baseline: frequency	0.5421 (40,000)							
weeds-token OLD	0.5176	(36,966)	0.5543	(38,956)	0.6108	(40,000)	0.6083	(40,000)
weeds-type OLD	0.4771	(36,966)	0.4797	(38,956)	0.5463	(40,000)	0.5723	(40,000)
weeds-token	0.5072	(36,966)	0.5084	(38,956)	0.5163	(40,000)	0.5373	(40,000)
weeds-type	0.5241	(36,966)	0.5270	(38,956)	0.5477	(40,000)	0.5501	(40,000)
entropy	0.5371	(40,000)	0.5280	(40,000)	0.5654	(40,000)	0.5646	(40,000)
density-CC-5 (nouns)	0.4731	(39,800)	0.4212	(39,800)	0.5322	(40,000)	0.5295	(40,000)
density-CC-5 (all)	0.4646	(40,000)	0.4316	(40,000)	0.5058	(40,000)	0.5202	(40,000)
density-CC-10 (nouns)	0.4506	(39,800)	0.3460	(39,800)	0.5115	(40,000)	0.4980	(40,000)
density-CC-10 (all)	0.4148	(40,000)	0.3546	(40,000)	0.4680	(40,000)	0.4883	(40,000)
density-CC-20 (nouns)	0.4059	(39,800)	0.2989	(39,800)	0.4806	(40,000)	0.4556	(40,000)
density-CC-20 (all)	0.3983	(40,000)	0.3212	(40,000)	0.4398	(40,000)	0.4556	(40,000)
density-CC-50 (nouns)	0.3891	(39,800)	0.2427	(39,800)	0.4324	(40,000)	0.3927	(40,000)
density-CC-50 (all)	0.3646	(40,000)	0.2840	(40,000)	0.3698	(40,000)	0.3899	(40,000)
density-TC-5 (nouns)	0.5142	(39,800)	0.5538	(40,000)	0.5697	(40,000)	0.6650	(40,000)
density-TC-5 (all)	0.5139	(40,000)	0.5591	(40,000)	0.6151	(40,000)	0.6475	(40,000)
density-TC-10 (nouns)	0.5142	(39,800)	0.5500	(40,000)	0.5454	(40,000)	0.6381	(40,000)
density-TC-10 (all)	0.5211	(40,000)	0.5613	(40,000)	0.5659	(40,000)	0.6257	(40,000)
density-TC-20 (nouns)	0.5389	(39,800)	0.5664	(40,000)	0.5141	(40,000)	0.6028	(40,000)
density-TC-20 (all)	0.5188	(40,000)	0.5658	(40,000)	0.5289	(40,000)	0.5938	(40,000)
density-TC-50 (nouns)	0.5492	(39,800)	0.5637	(40,000)	0.4870	(40,000)	0.5604	(40,000)
density-TC-50 (all)	0.4932	(40,000)	0.5378	(40,000)	0.4625	(40,000)	0.5292	(40,000)
density-NN-5 (verbs)	0.5925	(40,000)	0.5698	(40,000)	0.5789	(40,000)	0.5454	(40,000)
density-NN-5 (all)	0.5624	(40,000)	0.5756	(40,000)	0.6319	(40,000)	0.6035	(40,000)
density-NN-10 (verbs)	0.6020	(40,000)	0.5436	(40,000)	0.5695	(40,000)	0.5284	(40,000)
density-NN-10 (all)	0.5962	(40,000)	0.6049	(40,000)	0.6319	(40,000)	0.6186	(40,000)
density-NN-20 (verbs)	0.5861	(40,000)	0.5509	(40,000)	0.5353	(40,000)	0.5023	(40,000)
density-NN-20 (all)	0.6048	(40,000)	0.6043	(40,000)	0.6223	(40,000)	0.6075	(40,000)
density-NN-50 (verbs)	0.5832	(40,000)	0.5355	(40,000)	0.4829	(40,000)	0.4409	(40,000)
density-NN-50 (all)	0.6054	(40,000)	0.5813	(40,000)	0.6211	(40,000)	0.5976	(40,000)
density-TN-5 (verbs)	0.5081	(40,000)	0.4818	(40,000)	0.5275	(40,000)	0.5120	(40,000)
density-TN-5 (all)	0.4891	(40,000)	0.4656	(40,000)	0.5586	(40,000)	0.5499	(40,000)
density-TN-10 (verbs)	0.5241	(40,000)	0.4919	(40,000)	0.5206	(40,000)	0.5098	(40,000)
density-TN-10 (all)	0.5128	(40,000)	0.4932	(40,000)	0.5640	(40,000)	0.5605	(40,000)
density-TN-20 (verbs)	0.5260	(40,000)	0.4903	(40,000)	0.5057	(40,000)	0.4972	(40,000)
density-TN-20 (all)	0.5305	(40,000)	0.5167	(40,000)	0.5644	(40,000)	0.5638	(40,000)
density-TN-50 (verbs)	0.5087	(40,000)	0.4762	(40,000)	0.4608	(40,000)	0.4569	(40,000)
density-TN-50 (all)	0.5436	(40,000)	0.5288	(40,000)	0.5548	(40,000)	0.5529	(40,000)

**Table 9.** Full results for pair-wise distinctions between degrees of concreteness: verbs.

	window 2		window 20	
	space: N	space: N-V-A	space: N	space: N-V-A
baseline: frequency	0.7276 (86,636)			
weeds-token OLD	0.5110 (38,890)	0.5424 (46,677)	0.5387 (58,382)	0.5382 (60,985)
weeds-type OLD	0.4246 (38,890)	0.4054 (46,677)	0.3845 (58,382)	0.3916 (60,985)
weeds-token	0.7064 (38,890)	0.7141 (46,677)	0.7220 (58,382)	0.7221 (60,985)
weeds-type	0.7167 (38,890)	0.7227 (46,677)	0.7279 (58,382)	0.7250 (60,985)
entropy	0.7068 (49,139)	0.7152 (53,735)	0.7241 (61,152)	0.7244 (62,882)
density-CC-5 (nouns)	0.4138 (42,371)	0.4342 (42,371)	0.4934 (57,062)	0.4895 (57,062)
density-CC-5 (all)	0.3904 (48,114)	0.4016 (48,114)	0.4572 (59,475)	0.4665 (59,475)
density-CC-10 (nouns)	0.4114 (42,371)	0.4293 (42,371)	0.4903 (57,062)	0.4862 (57,062)
density-CC-10 (all)	0.3637 (48,114)	0.3755 (48,114)	0.4487 (59,475)	0.4613 (59,475)
density-CC-20 (nouns)	0.4172 (42,371)	0.4313 (42,371)	0.4823 (57,062)	0.4797 (57,062)
density-CC-20 (all)	0.3612 (48,114)	0.3713 (48,114)	0.4451 (59,475)	0.4556 (59,475)
density-CC-50 (nouns)	0.4381 (42,371)	0.4556 (42,371)	0.4850 (57,062)	0.4844 (57,062)
density-CC-50 (all)	0.3695 (48,114)	0.3806 (48,114)	0.4396 (59,475)	0.4539 (59,475)
density-TC-5 (nouns)	0.4664 (46,866)	0.4569 (47,724)	0.5089 (61,006)	0.5020 (61,016)
density-TC-5 (all)	0.4691 (47,669)	0.4609 (50,526)	0.4671 (61,067)	0.4676 (62,775)
density-TC-10 (nouns)	0.4638 (46,866)	0.4498 (47,724)	0.4977 (61,006)	0.4903 (61,016)
density-TC-10 (all)	0.4588 (47,734)	0.4496 (50,526)	0.4449 (61,067)	0.4497 (62,775)
density-TC-20 (nouns)	0.4640 (46,866)	0.4473 (47,724)	0.4954 (61,006)	0.4836 (61,016)
density-TC-20 (all)	0.4534 (47,744)	0.4431 (50,526)	0.4346 (61,067)	0.4408 (62,775)
density-TC-50 (nouns)	0.4649 (46,866)	0.4447 (47,724)	0.4981 (61,006)	0.4846 (61,016)
density-TC-50 (all)	0.4439 (47,744)	0.4317 (50,526)	0.4245 (61,067)	0.4336 (62,775)
density-NN-5 (nouns)	0.4756 (48,770)	0.4934 (53,452)	0.5117 (61,037)	0.5172 (62,797)
density-NN-5 (all)	0.4640 (48,868)	0.4890 (53,517)	0.4857 (61,090)	0.4990 (62,813)
density-NN-10 (nouns)	0.4778 (48,770)	0.4785 (53,456)	0.5187 (61,037)	0.5149 (62,797)
density-NN-10 (all)	0.4753 (48,868)	0.4872 (53,517)	0.4933 (61,090)	0.5017 (62,813)
density-NN-20 (nouns)	0.4679 (48,770)	0.4717 (53,456)	0.5256 (61,037)	0.5141 (62,797)
density-NN-20 (all)	0.4691 (48,868)	0.4801 (53,517)	0.4965 (61,090)	0.4983 (62,813)
density-NN-50 (nouns)	0.4556 (48,770)	0.4576 (53,456)	0.5294 (61,037)	0.5129 (62,797)
density-NN-50 (all)	0.4569 (48,868)	0.4714 (53,517)	0.5021 (61,090)	0.4987 (62,813)
density-TN-5 (nouns)	0.5197 (48,894)	0.5211 (53,564)	0.4676 (61,055)	0.4789 (62,821)
density-TN-5 (all)	0.5283 (48,977)	0.5329 (53,597)	0.4476 (61,108)	0.4663 (62,821)
density-TN-10 (nouns)	0.5019 (48,894)	0.5015 (53,564)	0.4611 (61,055)	0.4708 (62,821)
density-TN-10 (all)	0.5083 (48,977)	0.5156 (53,597)	0.4397 (61,108)	0.4558 (62,821)
density-TN-20 (nouns)	0.4864 (48,894)	0.4810 (53,564)	0.4569 (61,055)	0.4587 (62,821)
density-TN-20 (all)	0.4913 (48,977)	0.4971 (53,597)	0.4340 (61,108)	0.4464 (62,821)
density-TN-50 (nouns)	0.4627 (48,894)	0.4521 (53,564)	0.4494 (61,055)	0.4414 (62,821)
density-TN-50 (all)	0.4677 (48,977)	0.4739 (53,597)	0.4255 (61,108)	0.4318 (62,821)

**Table 10.** Full results for pair-wise distinctions between degrees of specificity: nouns.

	window 2				window 20			
	space: N		space: N-V-A		space: N		space: N-V-A	
baseline: frequency	0.7110 (39,572)							
weeds-token OLD	0.5158 (27,094)	0.5310 (28,500)	0.5191 (32,686)	0.5273 (33,438)				
weeds-type OLD	0.4212 (27,094)	0.4259 (28,500)	0.4038 (32,686)	0.4146 (33,438)				
weeds-token	0.7054 (27,094)	0.7083 (28,500)	0.7104 (32,686)	0.7088 (33,438)				
weeds-type	0.7111 (27,094)	0.7107 (28,500)	0.7112 (32,686)	0.7095 (33,438)				
entropy	0.7039 (30,622)	0.7049 (31,529)	0.7072 (33,704)	0.7068 (34,255)				
density-CC-5 (nouns)	0.4888 (28,306)	0.4273 (28,372)	0.5149 (32,445)	0.4780 (32,445)				
density-CC-5 (all)	0.4972 (29,517)	0.4167 (29,572)	0.5001 (33,241)	0.4750 (33,241)				
density-CC-10 (nouns)	0.4813 (28,306)	0.4045 (28,372)	0.5143 (32,445)	0.4751 (32,445)				
density-CC-10 (all)	0.4869 (29,517)	0.4005 (29,572)	0.4971 (33,241)	0.4643 (33,241)				
density-CC-20 (nouns)	0.4803 (28,306)	0.4067 (28,372)	0.5164 (32,445)	0.4742 (32,445)				
density-CC-20 (all)	0.4776 (29,517)	0.4020 (29,572)	0.5027 (33,241)	0.4735 (33,241)				
density-CC-50 (nouns)	0.4938 (28,306)	0.4387 (28,372)	0.5213 (32,445)	0.4883 (32,445)				
density-CC-50 (all)	0.5017 (29,517)	0.4253 (29,572)	0.5158 (33,241)	0.4907 (33,241)				
density-TC-5 (nouns)	0.4500 (30,092)	0.4479 (30,292)	0.4941 (33,704)	0.4869 (33,704)				
density-TC-5 (all)	0.4555 (30,292)	0.4582 (30,652)	0.4831 (33,704)	0.4808 (34,251)				
density-TC-10 (nouns)	0.4498 (30,092)	0.4467 (30,292)	0.4807 (33,704)	0.4678 (33,704)				
density-TC-10 (all)	0.4491 (30,292)	0.4518 (30,652)	0.4615 (33,704)	0.4593 (34,251)				
density-TC-20 (nouns)	0.4509 (30,092)	0.4469 (30,292)	0.4789 (33,704)	0.4642 (33,704)				
density-TC-20 (all)	0.4428 (30,292)	0.4459 (30,652)	0.4499 (33,704)	0.4440 (34,251)				
density-TC-50 (nouns)	0.4506 (30,092)	0.4433 (30,292)	0.4801 (33,704)	0.4631 (33,704)				
density-TC-50 (all)	0.4377 (30,292)	0.4359 (30,652)	0.4420 (33,704)	0.4408 (34,251)				
density-NN-5 (verbs)	0.5191 (30,602)	0.5230 (31,494)	0.5265 (33,704)	0.5340 (34,251)				
density-NN-5 (all)	0.5307 (30,611)	0.5162 (31,508)	0.5562 (33,704)	0.5586 (34,251)				
density-NN-10 (verbs)	0.5123 (30,602)	0.5149 (31,494)	0.5166 (33,704)	0.5298 (34,251)				
density-NN-10 (all)	0.5288 (30,611)	0.5201 (31,508)	0.5552 (33,704)	0.5625 (34,251)				
density-NN-20 (verbs)	0.4941 (30,602)	0.5084 (31,494)	0.5012 (33,704)	0.5173 (34,251)				
density-NN-20 (all)	0.5132 (30,611)	0.5169 (31,508)	0.5455 (33,704)	0.5628 (34,251)				
density-NN-50 (verbs)	0.4867 (30,602)	0.4933 (31,494)	0.4754 (33,704)	0.4929 (34,251)				
density-NN-50 (all)	0.4975 (30,611)	0.5057 (31,508)	0.5315 (33,704)	0.5526 (34,251)				
density-TN-5 (verbs)	0.5194 (30,609)	0.5213 (31,508)	0.5047 (33,704)	0.5009 (34,251)				
density-TN-5 (all)	0.5731 (30,614)	0.5698 (31,511)	0.5616 (33,704)	0.5420 (34,251)				
density-TN-10 (verbs)	0.5056 (30,609)	0.5053 (31,508)	0.4875 (33,704)	0.4895 (34,251)				
density-TN-10 (all)	0.5634 (30,614)	0.5596 (31,511)	0.5509 (33,704)	0.5361 (34,251)				
density-TN-20 (verbs)	0.4908 (30,609)	0.4909 (31,508)	0.4667 (33,704)	0.4758 (34,251)				
density-TN-20 (all)	0.5472 (30,614)	0.5430 (31,511)	0.5363 (33,704)	0.5278 (34,251)				
density-TN-50 (verbs)	0.4654 (30,609)	0.4644 (31,508)	0.4356 (33,704)	0.4506 (34,251)				
density-TN-50 (all)	0.5222 (30,614)	0.5232 (31,511)	0.5103 (33,704)	0.5149 (34,251)				

**Table 11.** Full results for pair-wise distinctions between degrees of specificity: verbs.