



Swansea University  
Prifysgol Abertawe



## Cronfa - Swansea University Open Access Repository

---

This is an author produced version of a paper published in:

*Applied Linguistics*

Cronfa URL for this paper:

<http://cronfa.swan.ac.uk/Record/cronfa40792>

---

### **Paper:**

Fitzpatrick, T., Playfoot, D., Wray, A. & Wright, M. (2015). Establishing the Reliability of Word Association Data for Investigating Individual and Group Differences. *Applied Linguistics*, 36(1), 23-50.

<http://dx.doi.org/10.1093/applin/amt020>

---

This item is brought to you by Swansea University. Any person downloading material is agreeing to abide by the terms of the repository licence. Copies of full text items may be used or reproduced in any format or medium, without prior permission for personal research or study, educational or non-commercial purposes only. The copyright for any work remains with the original author unless otherwise specified. The full-text must not be sold in any format or medium without the formal permission of the copyright holder.

Permission for multiple reproductions should be obtained from the original author.

Authors are personally responsible for adhering to copyright and publisher restrictions when uploading content to the repository.

<http://www.swansea.ac.uk/library/researchsupport/ris-support/>



**Establishing the reliability of word association data for investigating individual and group differences**

Journal:	<i>Applied Linguistics</i>
Manuscript ID:	APPLING-12-08-198.R2
Manuscript Type:	Article
Keyword:	cognitive linguistics, language and aging, psycholinguistics, vocabulary, individual differences, Second Language Acquisition

SCHOLARONE™  
Manuscripts

View

## Establishing the reliability of word association data for investigating individual and group differences

### Abstract

This paper argues that, across different psychological contexts, the methods of data collection, treatment and analysis in word association tests have hitherto been inconsistent. We demonstrate that this inconsistency has resulted from inadequate control, in previous studies, of certain important variables including the basis of norm comparisons, and we present a principled method for collecting, scoring and analysing association responses, to address these issues. The method is evaluated using test and retest datasets from 16-year-old and over-65-year-old twins ( $n=636$ ), which enable us to (a) compare samples matched for key environmental variables, (b) assess the transferability of norming information between age cohorts, and (c) evaluate the reliability of the scoring protocols. We find systematic differences in the association behaviour of the two age cohorts, indicating the importance of evaluating data only against norms lists which are matched to the target population. Individual association behaviour is found to be consistent across test times, both in terms of response stereotypy and response type.

### Introduction

For over a century word association (WA) tasks have been used to investigate the content and organisation of words and concepts in the mind. In early studies the focus was conceptual, with responses interpreted as indicators of general behaviours (e.g. Galton 1879; Jung 1910) and, by extension, being used to diagnose psychological abnormality (e.g. Sommer 1901; Kent and Rosanoff 1910). More recently WA studies have adopted a lexical focus, and have investigated the development and organisation of the mental lexicon and the influence of specific variables on lexical access. In applied linguistics, interest has most often been on the integration of L2 items into the lexicon, and the ways in which WA responses might reflect the development of L2 proficiency (e.g. Kruse, Pankhurst and Sharwood Smith 1987; Wolter 2002; Henriksen 2008 and, for an overview, Meara 2009). However, the findings of these L2 studies have been inconsistent and inconclusive, and in this paper we propose that this is on account of an assumption about the nature of WA patterns that increasingly appears to be unsafe. It is an assumption that also pervades the L1 WA research context.

Most studies of WA in the L2 have evaluated learners' responses against 'native speaker norms'. The rationale is one of demonstrating that as proficiency increases, WA behaviour becomes more like that of an adult native speaker. However, recent investigations (e.g. Author 2007; Zareva and Wolter 2012) have questioned the validity of assuming there is a coherent norm behaviour in native speakers, with Author finding that 'not only do [native speakers] vary in the actual words they produce, they also seem to vary in the types of association they make' (2007, p327). On the other hand, consistency was found in the WA behaviour of individuals, both diachronically in the L1 and also synchronically across two languages (Author 2007; 2009).

A review of studies from outside mainstream applied linguistics, specifically from psychology, reveals that the idea of a 'normal' WA behaviour also anchors research and practice there. WA methods have been used (with informants operating

1  
2  
3 in their L1) to investigate the effects on association behaviour of age, personality,  
4 psychosis and cognitive function. While this indicates a recognition that there are  
5 individual differences in the L1 population, the focus has not been on capturing a  
6 range of normal behaviours so much as on interpreting the behaviour of an individual  
7 in relation to assumed normal responses. Specifically, norms lists are used here, just  
8 as they are in L2 research, as the core point of reference. We propose that it is perhaps  
9 for this reason that these L1 studies also present equivocal findings.

10  
11 The methodology we present in this paper was developed in order to maximise  
12 the opportunity to capture the nature of variation within L1 populations, and thus  
13 reveal the extent and nature of 'normal' WA behaviour as a reference point for  
14 research in both the L1 and L2 domains. The methodology was informed by theories  
15 of the mental lexicon and by previous WA research, and drew on a large sample of  
16 respondents (n=636). We evaluated the approach by exploiting several distinct  
17 features of our data set. Firstly, the informants were pairs of twins, making it possible  
18 to build two matched subsets of data. Secondly, a sub-group of informants completed  
19 the WA task at two separate test times, enabling us to assess reliability of response  
20 behaviour. Thirdly, the informants fell into two distinct age categories: 16 year olds  
21 and over 65s. This enabled us to examine the capacity of the methodology to capture  
22 differences between sub-populations that might inform future assumptions about  
23 reference norms.  
24

25 In addition, we had data for the informants regarding their zygoty (i.e. whether  
26 they were identical or non-identical twins) and their performance on a range of  
27 cognitive tests. However, these elements are not discussed in this paper, since they are  
28 not relevant to the methodology itself.  
29

30 In sum, our aim is to resolve the problem highlighted by Schmitt: 'It is clear that  
31 association data provides insights in the organization of the mental lexicon.....and it  
32 seems that this approach is still waiting for a breakthrough in methodology which can  
33 unlock its undoubted potential' (2010: 248). In the remainder of this section, we  
34 review the extent of variation in the management and analysis of WA data in a  
35 number of influential studies. The next two sections describe our dataset and analytic  
36 procedures. After this we present and evaluate our method for measuring WA  
37 responses by stereotypy, and we demonstrate evidence that norms lists must be  
38 selected appropriately for the test population. Finally we address the inherent  
39 complexities of categorising responses by type. Both the norms and categorisation  
40 measures are tested for reliability, using matched samples and longitudinal retests.  
41  
42  
43

#### 44 *A review of approaches to the management and analysis of WA data*

45  
46 WA protocols are attractive to the researcher for a number of reasons. They offer a  
47 relatively quick and straightforward method for gathering rich language data. The data  
48 they elicit are freely produced, but consist of discrete lexical items, or word pairs  
49 (cue→response), which lend themselves to quantitative analysis more readily than do  
50 discursive language data. They are also congruent with well-established  
51 psycholinguistic and applied linguistic theories, such as Connectionism and Latent  
52 Semantic Analysis (see Ellis 1998), the Bilingual Interaction Activation model (e.g.  
53 Dijkstra and van Heuven 1998), and other models of word knowledge and lexical  
54 storage and retrieval (e.g. Marslen-Wilson 1987; Nation 2001). Tracking changes in  
55 WA responses can inform the study of a dynamic, growing lexicon, in which links are  
56 being created and strengthened, and this is reflected in the amount of WA literature  
57  
58  
59  
60

1  
2  
3 published since the 1950s relating to the development of L1 (Ervin 1961; Entwisle  
4 1966; Nelson 1977) and L2 (Meara 2009). Furthermore, since the 1980s attention has  
5 been increasingly paid to the application of WA protocols to the study of lexical  
6 attrition (Gewirth, Shindler and Hier 1984; Gollan, Salmon and Paxton 2006).  
7

8 Typically, these studies have used one of two broad analytical approaches to the  
9 measurement of data. One entails examining the stereotypy of responses, that is, how  
10 similar an individual's response is to those in a reference set. The other approach  
11 examines the nature of the relationship between the cue and the response. Some  
12 studies combine the two approaches. The choice of analytic approach depends on the  
13 research question being addressed and the theoretical assumptions underlying the  
14 research. For instance, stereotypy approaches, which rely heavily on the similarity  
15 between a respondent's responses and 'normal responses', have been used in the  
16 context of cognitive and psychiatric disorders. Approaches categorising the type of  
17 link between cue and response tend to be used to map patterns of variation in normal  
18 populations.  
19

20 Research findings are of course dependent on the research questions and choice  
21 of analytic approach. However, a number of other factors also potentially impact  
22 heavily on the interpretation of data, so that different data-gathering procedures and  
23 materials may compromise the meaningfulness of cross-study comparisons. In  
24 addition to sample size, which influences the robustness of any quantitative empirical  
25 study, potential methodological variables to consider include:  
26

- 27 • Mode of elicitation: Cues may be read or heard, and responses spoken, written  
28 or typed.
- 29
- 30 • Cue choice: The number of cues in the WA task contributes to validity in the  
31 same way as population sample size. Less easy to quantify, but possibly even  
32 more important, is the way in which cue items are selected. Possible  
33 contributors to uncontrolled variation are word frequency, word class,  
34 imageability and the age at which the word was acquired. In addition,  
35 adequate attention has to be paid to the tendency for certain words to  
36 consistently cue a particular response, such as a highly probable collocate (e.g.  
37 bread→butter).  
38
- 39 • Norms lists: Studies using stereotypy measures depend on norms lists against  
40 which to score the responses of the target population. While some studies  
41 compile norms lists from the study participants themselves or create bespoke  
42 norms lists (e.g. Miller and Chapman 1983; Hirsh and Tree 2001), most use  
43 existing lists such as the Postman-Keppel lists (Postman and Keppel 1970) or  
44 the South Florida Association Norms (Nelson, McEvoy and Schreiber 1998).  
45 This second approach may not always allow for the possibility that responses  
46 are influenced by cohort characteristics such as generational differences,  
47 geographical location, and so on.  
48
- 49 • Treatment of responses: Researchers vary in their treatment of response items.  
50 Some correct spelling, some lemmatize responses, and problematic responses  
51 such as non-words, multi-word responses and blanks are dealt with in different  
52 ways.  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 Thus, although it would seem reasonable, when deciding on a specific  
4 methodology for a WA study, to replicate the protocols most commonly used in  
5 previous research so as to maximise opportunities for cross-study comparability, a  
6 brief review of studies that have used WA methods reveals very little commonality of  
7 approach. The studies listed in Table 1 have been selected to represent the main  
8 variables investigated through WA data: age, cognitive function, personality and  
9 psychosis. The studies with the highest number of citations have been selected for  
10 each variable, using the Publish or Perish database (Harzing 2007). As the table  
11 shows, there is considerable between-study variation in the selection of cues and  
12 norms lists, and in the treatment and analysis of responses, affording little  
13 methodological guidance to the researcher. This is exacerbated by the fact that many  
14 of these papers report strikingly little methodological detail. Most offer no  
15 justification for methodological or procedural decisions, and little or no reference to  
16 the way data has been collected, treated and analysed relative to other, comparable  
17 studies. There are exceptions to this of course, notably in the early studies of first  
18 language development (Ervin 1961; Entwisle, Forsyth and Muuss 1964). Even when  
19 studies addressing the same research question and using the same theoretical  
20 assumptions are compared, there is little consistency of approach, as seen in Table 2,  
21 which lists the most cited experimental studies using the production of WA responses  
22 to investigate L2 proficiency.  
23

24  
25 The methodology reported in the following sections of this paper is able to shed  
26 light on the potential impact of some of the previously uncontrolled variables listed  
27 above. We held constant the variables of mode of elicitation and cue choice, in order  
28 to explore the impact of norms sets and categorisation. Future research will be able to  
29 focus on the first two variables, using the findings from this study to anchor the latter  
30 two.  
31

32  
33 [TABLES 1 and 2 NEAR HERE]  
34

### 35 **The dataset**

36  
37 The opportunity to use WA data from twins arose in the context of our collaboration,  
38 since 2007, with a research team engaged in two large-scale twin studies: the Genes  
39 for Cognition Study and the Older Australian Twins Study (Author 2004; Sachdev et  
40 al. 2009; see <http://genepi.qimr.edu.au/> for further details)<sup>1</sup>. WA tasks were included  
41 in a battery of cognitive performance tests with the ultimate aim of exploring the roles  
42 of genes and environment in the relationships between different measures of linguistic  
43 and non-linguistic performance. For the norms lists and stereotypy analyses, the data  
44 are from 192 participants: 48 twin pairs aged 16 years and 48 twin pairs aged over 65.  
45 The categorisation of association types used the responses of 540 of the 16 year-old  
46 twins. Responses from a subset of the younger participant group (n=36), who  
47 performed the task twice, were used to assess the reliability of both the stereotypy and  
48 the categorisation methods. All participants in all analyses were native English  
49 speakers. The older twins were recruited through the Australian Twin Registry or  
50 publicity, and the 16 year olds through schools and word of mouth. The studies were  
51 subject to the strict ethics procedures of medical research. Participants completed the  
52 WA task as part of a suite of physical and cognitive tests during either a half (16 year-  
53 olds) or one day-long visit to the research unit, located in a hospital.  
54

55  
56 The WA task consisted of 100 cue words<sup>2</sup>, controlled for the impact of  
57 frequency by randomly selecting them from the 2k and 3k bands of the British  
58  
59  
60



1  
2  
3 National Corpus, <http://www.natcorp.ox.ac.uk>, (thus representing the second and third  
4 thousand most frequent words in English usage). Words from the first thousand band  
5 were not included, because previous research shows that frequently encountered  
6 words tend to produce strong dominant responses (Meara 1983) and a proliferation of  
7 predictable responses would mask potential differences between participants. On the  
8 other hand, restricting cue selection to the 2k and 3k bands (50 cues from each)  
9 ensured that cue items were familiar enough for the respondents to offer an  
10 association to them. The cues and their dominant responses are listed at [journal url].  
11 Although we did not explicitly control imageability or age of acquisition in the cues  
12 selected (see earlier note that these might affect responses), regression analyses  
13 indicated that these characteristics of the cue did not predict stereotypy or response  
14 category.  
15

16 The cues were presented in two columns of 25, on two pages. Next to each cue  
17 was a space for the participant to write a response<sup>3</sup>. Participants were instructed to  
18 write down the first word they thought of when reading each cue, and were told that  
19 there were no right or wrong answers. An excerpt from a completed task is shown in  
20 Figure 1. Participants were allowed up to 10 minutes to complete the task, and all  
21 participants finished it within this time.  
22  
23

24 [FIGURE 1 NEAR HERE]  
25

### 26 **Preparing the data for analysis**

27

28 The data were presented to the analysts with only identity codes that did not indicate  
29 gender or twin pairings. The hand-written responses were transcribed into an Excel  
30 file. In order to enable automatic searches, spelling was corrected, but only where the  
31 intention was clear (e.g. *controll* and *controle* were corrected to *control*). However,  
32 instances of possible spelling mistakes were not corrected if the response was a real  
33 word. For example, one participant wrote *backed* for the cue word *bean*. Although it  
34 is extremely likely in this particular case that the intended response was *baked*, many  
35 other cases rendered much less clear relationships between what was actually written  
36 and what might have been intended (e.g. both *council* and *counsel* are plausible as  
37 associates for the cue *session*). So, to avoid a kind of second guessing that would have  
38 imposed the analysts' own WA preferences, a blanket policy was adopted of treating  
39 real word responses at face value.  
40  
41

42 While the majority of responses (>95%) were single words, participants  
43 occasionally wrote two or more words or a short phrase. Where phrases could be  
44 construed as formulaic sequences with a single coherent meaning (Author 2002), they  
45 were transcribed as written. When multi-word responses did not represent strings in  
46 this way<sup>4</sup>, two procedures were employed to shorten them. The first, appropriate  
47 where two separate one-word responses had been offered, was to truncate responses at  
48 punctuation (comma, slash, etc). Thus, *bomb/explosion* was transcribed as *bomb*. The  
49 second entailed deleting function words, particularly conjunctions (*and*, *or*, *with*),  
50 pronouns (usually *I*), and infinitive *to*.  
51  
52

### 53 **Norms lists and stereotypy measures**

54

#### 55 *Use of norms lists*

56  
57  
58  
59  
60

1  
2  
3 Stereotypy determines how similar a participant's responses are to those of a  
4 comparison group and thus entails the use of a normative response corpus. As can be  
5 seen in Tables 1 and 2, many previous studies have used published norms lists.

6  
7 Selecting a norms list which has already been created, published and used in  
8 other studies can be a useful shortcut in stereotypy analysis. However, a norms list  
9 will only be reliable as a point of reference if it is able to transcend the impact of  
10 variables characterising sub-populations. Until more is known about how different  
11 variables affect WA behaviour, researchers should be cautious about using  
12 independently gathered norm data as the reference point. The best way to address this  
13 issue is to create a norms list specifically for the study at hand, reliably to reflect the  
14 maximum possible number of characteristics of the study population. In this way it  
15 will be possible to develop an understanding of the differences in such norms across  
16 populations and the contribution that those differences make in the interpretation of  
17 data. Accordingly, as outlined below, in this study separate norms lists were compiled  
18 for the two populations under investigation—16 year olds and over 65s, and it was  
19 these lists that were used to calculate stereotypy scores (see below)<sup>5</sup>.

20  
21 Each norms list represented the associations of 96 participants in the respective  
22 age group. The lists were created by compiling a full list of the responses for each cue  
23 word, and counting up how many times each response was given. In order to do this,  
24 it was necessary to determine a definition of 'word'. For example, some scholars count  
25 every different word form as a different response (so that *walk* is different to *walked*  
26 or *walking* or *walker*), while others group such responses together as versions of the  
27 same lemma. The decision we took here was to lemmatize inflectional variants but not  
28 derivational ones. Specifically, words which corresponded to level 2 of Bauer and  
29 Nation's (1993) description of word families were considered the same. In practice  
30 that meant affixes producing plural nouns or verb participles were ignored, so that *cat*  
31 was considered the same as *cats*, *think* the same as *thinking*, and *walk* the same as  
32 *walked*. Derivational affixes, though, were retained, so that *health* and *healthy* were  
33 considered different responses, as were *teach* and *teacher*. The justification for this  
34 decision was that while any kind of lemmatising potentially impacts on gaining a full  
35 understanding of collocational behaviour (compare *attack* and *attacked* as responses  
36 to *heart*), the impact of not lemmatising is arguably greater, because it considerably  
37 reduces the incidence of common responses across the population. The key  
38 consideration is consistency and transparency, so that the way is clear for future  
39 empirical interrogations of the potential impact of the decisions taken.

40  
41 The norms lists were finalized by ordering the responses according to their  
42 frequency for that cue word, along with a record of those frequencies.

### 43 44 45 **Scoring for stereotypy**

46  
47 Previous studies have scored stereotypy in different ways (see Tables 1 and 2, last  
48 column), variously awarding 'stereotypy' points

- 49 a. for any response in the top 3 (or 5) in the norms list
- 50 b. for each percentage point of the norming population giving the response
- 51 c. according to percentage bands of the norming population giving the  
52 response
- 53 d. according to the ranking of the response on the norms list
- 54 e. for any response that appears anywhere in the norms list
- 55 f. for a response which is the dominant response on the norms list
- 56  
57  
58  
59  
60



1  
2  
3 In this paper we focus on a method using procedure (f), as this represents the  
4 measure most commonly used in the studies cited in Tables 1 and 2.<sup>6</sup> It should be  
5 noted, though, that the decision about which stereotypy measure to use will be  
6 dependent on the context of that particular study. In L2 research, for example, where  
7 participants typically have limited lexical resources, method (e) above might be more  
8 appropriate. Using scoring method (f), a response was considered 'stereotypical' if it  
9 was the most frequently recorded response on the norms list for the participant's age  
10 cohort. Participants scored 1 point for every stereotypical response, and all their other  
11 responses scored zero. For cues where two (or more) responses were equally popular,  
12 a point could be scored for either response.

13  
14 The data used in this analysis were from participants who had provided  
15 responses to more than 90% of the cues. In studies like this one, which use relatively  
16 frequent cue words from the participants' L1, and where participants are adults with  
17 no cognitive impairment, blank responses are rare. However, in other contexts a  
18 proliferation of blank responses might affect the analysis of some data sets, and  
19 appropriate methodological adjustments (typically the exclusion of data sets with  
20 more than *n* blank responses, or scores calculated on proportional rather than raw  
21 counts) have to be implemented.

### 22 23 24 *Assessing the validity of the norms list approach*

25  
26 In order to assess the effect of norms list characteristics on the profiling of the data,  
27 age was used as a variable. The 192 participants were split on the basis of age and  
28 twin birth order (1 or 2) to create four groups (young twin 1, young twin 2, older twin  
29 1, older twin 2<sup>7</sup>). A separate norms list was created for each group following the  
30 procedures described above, with each norms list therefore representing the responses  
31 of 48 participants. The prediction here was that differences between groups matched  
32 for age would be smaller than those not so matched.

33  
34 Using the four separate norms lists as the reference, four stereotypy scores were  
35 calculated for each participant, according to the procedure described above. The first  
36 score was calculated from the norms lists to which the participant had contributed (i.e.  
37 a young twin 1 was given a point for every response which was a dominant response  
38 on the norms list compiled from all young twin 1 participants). The second stereotypy  
39 score was calculated from the norms list of responses from the group of the same age,  
40 different twin number (i.e. young twin 1 was given a point for every response that was  
41 a dominant one on the young twin 2 norms list). The third and fourth stereotypy  
42 scores were calculated from the norms list of twin 1 in the other age group, and twin 2  
43 in the other age group. The four stereotypy scores therefore represent the similarity to  
44 'own list', 'same age, other twin', 'twin 1, other age group' and 'twin 2, other age group'  
45 norms. Group mean stereotypy scores and standard deviations are presented in Table  
46 3.  
47  
48

49  
50 [TABLE 3 NEAR HERE]

51  
52 Three patterns are apparent. First, twin 1s and twin 2s have similar mean scores  
53 irrespective of the norms list. This is consistent with the assumption that there would  
54 be no material differences between first- and second-born twins in the context of  
55 stereotypy score. Second, the levels of stereotypy for any given condition of  
56 comparison (i.e. the figures in each column) are similar, which indicates that the four  
57 groups' responses are related to each other in a consistent way. Third, all participants'  
58  
59  
60

1  
2  
3 responses are more typical of their own age group than of the other age group, as  
4 shown by the lower mean stereotypy scores when using the norms derived from the  
5 other age twin lists.  
6

7 To test the significance of the observations derived from these descriptive  
8 statistics, stereotypy data were entered into age (2) by twin (2) by norms list (4)  
9 repeated measures ANOVA analyses by subjects and by items. Age and twin were  
10 entered as between subject variables in the analysis by subjects, and as within subject  
11 variables in the analysis by items. 'Norms list' was treated as a within subjects variable  
12 in both analyses. Greenhouse-Geisser corrections were applied to all analyses  
13 including the norms list factor as it violated the assumption of sphericity. The analysis  
14 was conducted to establish whether a) the choice of norms list for comparison had a  
15 significant effect on the stereotypy scores of the participants and b) whether there  
16 were overall differences in stereotypy levels between age groups or twin pairs once  
17 norms list factors were taken into account. The main effect of norms list was  
18 significant by subjects and by items [ $F_1(3, 564) = 136.948$ ,  $MSe = 25.730$ ,  $p < .001$ ,  
19  $\eta^2 = .421$ ;  $F_2(3, 297) = 69.319$ ,  $MSe = 22.181$ ,  $p < .001$ ,  $\eta^2 = .412$ ]. Bonferroni  
20 corrected follow up t-tests ( $\alpha/6 = .0083$ ) revealed that mean 'own list' and 'same age  
21 group' stereotypy scores (27.18 and 25.56) were both significantly higher than those  
22 calculated from the other age group norms lists (19.23 and 18.96). The mean 'own list'  
23 stereotypy score (27.18) was significantly higher than stereotypy on the other norms  
24 list from the same age group (25.56), as is predictable given that participants'  
25 responses by definition all appear on their own norms list, and thus potentially  
26 contributed to the dominance of that response. The small difference in mean  
27 stereotypy in relation to other age twin 1 and other age twin 2 lists was not significant.  
28 The main effects of age and twin number did not reach significance, and no  
29 interactions were significant.  
30  
31

32 This analysis demonstrates the importance of using age-appropriate norms lists  
33 in the study of WA stereotypy. Participants gained an advantage of more than 6  
34 stereotypy points (average 25.56 versus average 19.1) when scored against age-  
35 appropriate lists. There are several possible reasons for an age-related difference in  
36 the norms lists. One is that certain changes in WA selection strategies occur as a  
37 function of ageing. A second is that each generation has its own preferred set of  
38 vocabulary and/or associations. The first explanation predicts that the 16 year olds'  
39 responses would, over time, come to resemble more closely the norms of the 65+  
40 group. This means that the appropriacy of norms for new experimental groups could  
41 be calculated as a gradation on the basis of age. The second explanation predicts that  
42 the 16 year olds would, in 50 years time, display norms rather similar to those they  
43 produced in teenage, but that a new cohort of 16 year olds at that time would produce  
44 new norms. A third possibility is that age and generation interact, such that as one gets  
45 older one attends to different concepts and words in the environment, as a function of  
46 one's changing interests and common activities, themselves influenced by prevailing  
47 generational cultural preferences. This more complex explanation, if correct, would  
48 predict that neither of the norms lists developed in this study would be a good match  
49 for the 16 year olds when they got to 65+. Common to all three explanations is the  
50 caution about using as a reference point any norms list that is not derived directly  
51 from the target population.  
52  
53  
54

### 55 *Assessing the reliability of the stereotypy measure*

56  
57  
58  
59  
60

1  
2  
3 For a measure to be considered reliable, it should produce comparable results at two  
4 test events using the same participants, always assuming participant performance is a  
5 stable factor. Key reasons why participant performance might not be replicable are  
6 practice effects including memory for the previous iteration (if the test events are very  
7 close in time) and developmental or attritional changes in the participant's underlying  
8 organisation of response options (if the test events are temporally very distant). The  
9 interval between test events here was approximately 3 months, which was considered  
10 large enough to minimize practice effects without reflecting substantial inherent  
11 changes in lexical knowledge or organisation.

12  
13 Thirty-six of the younger participants provided the data for this analysis, having  
14 completed the WA task on two separate occasions. Following the finding reported  
15 above, age appropriate norms lists were used to score participants' responses for  
16 stereotypy. Table 4 presents descriptive statistics for stereotypy test and retest scores.  
17

18 [TABLE 4 NEAR HERE]  
19

20  
21 Mean scores were broadly similar across test times, with a significant, positive  
22 test retest correlation indicating consistency in WA behaviour over time. A  
23 calculation of repeated responses revealed that this consistency in scoring is not  
24 explained by participants producing the same responses to the same cues at each test  
25 time: on average identical responses were only produced for 25.5 of the 100 cues (see  
26 Table 5).  
27

28 [TABLE 5 NEAR HERE]  
29

### 30 **Word association response type measures**

31  
32  
33 Word association behaviour has also conventionally been assessed in terms of the  
34 types of link between the cue and the response. In early studies of this nature, analyses  
35 of the links were based on the Saussurian definitions of syntagmatic and paradigmatic  
36 relationships. A distinction was made between pairs of words which co-occur in text  
37 (syntagmatic, e.g. *van-drive*) and pairs of words which can be substituted for one  
38 another without changing the grammaticality of the sentence (paradigmatic, e.g. *van-*  
39 *train*). A third category, known as 'clang', was later added to this framework to  
40 represent responses based on the form of the cue, typically phonological (e.g. *van-*  
41 *fan*). Of the studies summarized in Tables 1 and 2, some (e.g. Ervin 1961; Gewirth,  
42 Shindler and Hier 1984) use variations of this framework and terminology, and  
43 there has more recently been a partial shift towards a change in terms to increase  
44 transparency, e.g. 'collocational', 'semantic' and 'phonological'. Developments in  
45 cognitive linguistics relating to the categorisation of sense relations (e.g. Croft and  
46 Cruse 2004), insights from natural language processing research (e.g. latent semantic  
47 analysis, Landauer, Foltz and Laham 1998) and the development of large-scale lexical  
48 databases such as WordNet (Miller, 1995) have some potential to challenge and  
49 inform WA categorisation systems, especially in the case of semantic (paradigmatic)  
50 connections. However, the recurrence in WA data of syntactic (usage-based) and  
51 orthographic/phonological associations has endorsed the continued inclusion of  
52 categories which accommodate these, such as the syntagmatic and clang categories in  
53 the conventional classification system.  
54

55  
56 These broad categories have revealed some qualitative differences in the  
57 response behaviours of children and adults (see Nelson 1977). However, category  
58  
59  
60

1  
2  
3 comparisons between responses of other participant groups have been less conclusive,  
4 with studies sometimes producing contradictory findings (see Meara 2009 for a  
5 summary of these in relation to L2 investigations). Author (2006), also focussing on  
6 L2 WA processes, proposes a categorisation based on a word knowledge framework  
7 (Nation 2001), which specifies subtypes of association response within each main  
8 category. She argues that this fine-grained approach provides greater insight into how  
9 learners of English engage with words. Her studies of distributions across these sub-  
10 categories reveal differences between WA behaviour of L1 and L2 users of English,  
11 and between L2 users of different proficiency levels, which had hitherto been masked  
12 by the broad category approach (Author 2006, 2009).  
13  
14

### 15 *Categorisation of responses*

16  
17 The system of categorisation used in the present analysis was based on Author (2006),  
18 and informed by the findings of subsequent studies (Author 2007; 2009;  
19 Higginbotham 2010, Author 2011). Key features of the revised system are, first, a  
20 rationalisation of the number of sub-categories, so as to ensure definitions are clear  
21 and the number of responses for each type is large enough for formal analysis.  
22 Second, the framework allows for responses to be coded as a potential combination of  
23 multiple links. For example, *knife* is commonly followed by *fork* in general usage (a  
24 collocation), but they are also items from the same lexical set (cutlery). In previous  
25 WA categorisation systems the researcher would be forced to make a choice as to  
26 which of these reasons was more likely. Here, the response can be classified as being  
27 both *lexical set* and *cue-response collocation*. It is advantageous to be able to  
28 recognize this level of complexity in light of the finding that participants are  
29 particularly quick to respond when the cue and the response are linked in more than  
30 one aspect (Author 2011).  
31  
32

33 The new framework comprises 14 sub-category headings in total, and is  
34 summarized in Table 6, with examples drawn from data in the present study.  
35

36 [TABLE 6 NEAR HERE]  
37

### 38 *Scoring word association responses using categories*

39  
40 The rationale when devising a categorisation framework is to sustain a balance  
41 between consistency and common sense, while adequately accommodating all the  
42 responses. This is not an easy task, nor an exact science, because the analyst's belief  
43 that a participant probably had a reason for giving a particular response is not always  
44 enough to create a warrantable assumption about the link. In order to avoid second-  
45 guessing, the balance of power must lie with consistency. In this study two specific  
46 procedures were employed to maximize such consistency. First, to ensure that the  
47 raters were not influenced by the respondent's previous behaviour patterns, or by the  
48 popularity of a particular response across the sample, the categorisation was done by  
49 cue not by participant. Thus, the complete list of responses to each cue was compiled  
50 into a single list, and duplicate answers were deleted, so that each response was listed  
51 only once per cue word. The relationship between cue and response was thereby  
52 neutralized, meaning that when raters were assigning responses to categories, they  
53 were not tempted to think 'this person has given a lot of collocations already so this is  
54 probably one too', or 'only one person said this so it's likely to be an erratic response'.  
55  
56  
57  
58  
59  
60

1  
2  
3 The complete set of responses to all the cues was categorized by two raters  
4 separately, according to the definitions above. Once the categorisation had been  
5 completed by both raters, the scoring of responses was compared, revealing that  
6 76.9% of response items had been assigned to the same category in the initial coding.  
7 A further 22.8% of the classifications were agreed after a short discussion and close  
8 reference to the definitions. The non-alignments in the initial categorisation of these  
9 responses were usually attributable to one rater missing a possible sense of the cue  
10 word. For example, one rater had missed the fact that *routine* could mean 'dull, boring  
11 and monotonous', while the other missed the meaning of *establish* as 'to prove'. This  
12 highlights the necessity for multiple raters, particularly given the demands on raters to  
13 pay close attention to such large amounts of data. Agreement about the categorisation  
14 of a very small number of responses (0.3%) could not be reached even after  
15 discussion. In these cases, a third party was consulted, and the link identified by the  
16 third party was used to arbitrate between the two options. During the categorisation  
17 process, two cue words were found to be problematic, in that participants commonly  
18 mistook them for a (near-) homophone. *Miner* was mistaken for *minor*, and responded  
19 to as such, and *instance* was responded to as *instant*. These cues and the responses  
20 they elicited were excluded from the categorisation analysis.

21  
22  
23 Using a spreadsheet, the responses were allocated their category type, and the  
24 instances of each category were summed to create individual response profiles.

### 25 26 ***Assessing the reliability of the categorisation system***

27  
28 Having categorised participants' responses according to the process described above,  
29 an assessment of the reliability of this method was undertaken. The aim was to  
30 establish whether, irrespective of specific items in responses, the distributional  
31 patterns of response types were replicable—these patterns are the basis on which  
32 observations might be made about differences in participant profiles. Data from the  
33 thirty-six test-retest participants were used. Responses were categorized according to  
34 the framework in Table 6, and profiles were produced for all participants at time one  
35 and time two. The mean number of responses in each subcategory is presented in  
36 Table 7, along with test-retest correlation coefficients (categories represented by, on  
37 average, less than one response per participant are not listed). Of the six main  
38 subcategories, significant positive correlations were observed for all but the erratic  
39 response category. High scorers on a given category in the initial test were likely to be  
40 high scorers on the same category in the retest.

41  
42 As observed in connection with the stereotype analyses reported above, this  
43 consistency cannot be attributed to participants providing identical response items at  
44 each test time (see Table 5); the consistency here is in the type of response given, not  
45 the item itself.

46  
47  
48 [TABLE 7 NEAR HERE]

### 49 50 ***Assessing the validity of the category clusters: a principal components analysis***

51  
52 As mentioned previously, a common analytic approach to WA data is to cluster  
53 responses into semantic, collocational and form-based groups, and indeed the  
54 subcategories proposed by Author were originally presented as subdivisions of these  
55 three groups. While there are theoretical grounds for making these distinctions,  
56 whether responses actually cluster in this way is an empirical question, which can be  
57  
58  
59  
60



1  
2  
3 explored by submitting WA profile data (i.e. category scores) to a principle  
4 components analysis.

5 Principal components analysis is a technique designed to organize large  
6 numbers of inter-correlated variables into clusters such that the information can be  
7 described using only a small number of 'components'. This has advantages in terms of  
8 statistical power, and avoids multi-collinearity problems when using regression  
9 analyses. For example, imagine you have a bowl containing 100 sweets and you ask a  
10 child to pick five. There are a large number of possible combinations of five sweets  
11 that the child could choose. When asked, the child tells you that he decided which  
12 sweets to take on the basis of their colour, picking only red ones. A second child  
13 chooses his five sweets from the bowl, and also takes only red sweets, but this child  
14 tells you that his decision was based on flavour. As there is a strong correlation  
15 between the colour and flavour of sweets, the identical selections of these two  
16 children, in the context of a larger set of children choosing on other grounds, could  
17 not be explained reliably using either of these variables, since both are possible  
18 explanations for their choice. A principal components analysis identifies patterns like  
19 this in the data set, and suggests a single 'colour-flavour' factor instead. Another child  
20 chooses his five sweets from the bowl, but his strategy is to take the sweets closest to  
21 the surface. His selection has nothing to do with the 'colour-flavour' factor, and the  
22 variance in sweet picking is instead explained by proximity.

23  
24  
25 This analysis takes the total variance in the WA behaviour and attempts to  
26 partition it into linear components. The procedure results in clusters of variables (in  
27 this case, WA categories) which explain a proportion of the variance not explained by  
28 anything else. If the three major conventional categories are valid, they should  
29 manifest as clusters. Our initial categorisation matrix contained 14 possible  
30 classifications for a response. Response data from 540 participants (all aged 16), in  
31 the form of response profiles were entered into a principal components analysis. The  
32 sample size was determined to be adequate using the Keyser-Meyer Olkin measure  
33 (KMO = .51). The data met the sphericity assumption as determined by a significant  
34 Bartlett's test statistic [ $\chi^2(78) = 1069.056, p < .001$ ]. The principal components  
35 analysis extracted five factors (rotated using the varimax procedure with Kaiser  
36 normalization) to explain the data. The rotated component matrix is presented in  
37 Table 8. The component labels in the table represent our interpretation of the  
38 component clusters; the analysis merely identifies them as discrete components.

39  
40  
41 [TABLE 8 NEAR HERE]

42  
43  
44 Table 8 lists components from left to right, in order of the proportion of  
45 variance in the data they account for, with the largest proportion being attributed to  
46 the first rows. The first component identified comprises synonym, lexical set and  
47 other conceptual link categories. This can be described as a meaning-based (semantic)  
48 component, as a conceptual link between cue and response underlies each of these  
49 subcategories. A second component includes both cue-response and response-cue  
50 collocations. This can be described as a position-based (collocational) component, as  
51 the link is determined by the close occurrence of the two items in language use. The  
52 third component comprises form-only, two-step, affix manipulation and erratic  
53 responses. It is suggested that this is a form-based component. In Author's original  
54 system only two of these sub-categories, *form only* and *affix* constituted the broad  
55 category *form*. The components analysis suggests that two additional subcategories  
56 may belong in this group, and a closer analysis of these subcategories provides a  
57  
58  
59  
60



1  
2  
3 principled explanation for this. First, in *two-step* associations, one step is nearly  
4 always form-based. This is illustrated by examples such as *bean* → *stork*. Here there  
5 has been an intermediate association involving the collocation *stalk*, a homophone  
6 (similar in form only) of the response *stork*. Second, the *erratic* response category  
7 encompasses potential spelling mistakes (i.e. form errors). The fact that these two  
8 categories load on the same component supports the notion that the *bean*→ *stalk/stork*  
9 response type might indeed be caused by erratic spelling (similarly, the *bean*→*backed*  
10 example cited earlier in this paper). Component 4 includes only the cue-response-  
11 response-cue collocations; note that these did not load with the other position-based  
12 categories, though given that very few of these responses were produced (less than  
13 0.5%), it is unwise to speculate about the reason for this.

14  
15 The final component includes dual-link associations: synonym plus cue-  
16 response collocations and lexical set plus response-cue collocations. The separation of  
17 these associations from the main groups supports Author's (2011) finding that dual-  
18 link associations are particularly strong and quick to retrieve, and do not behave in the  
19 same way as either semantically or position-based responses. The last two  
20 components contribute an extremely small proportion of the total variance, and indeed  
21 items with these double links were uncommon in the data.

22  
23 Specific research questions and hypotheses can demand a focus on particular  
24 subcategories (for example, Author 2006 found that synonyms make a much larger  
25 contribution to the semantic category in L1 responses than in L2). However, it is often  
26 advantageous, for reasons of statistical analysis, to group data into larger categories,  
27 and this principal components analysis has identified a convincing framework for  
28 doing so.  
29

## 30 31 32 **Conclusions**

33  
34 We have demonstrated that norms lists differ between age cohorts, and we  
35 strengthened the evidence by using two uniquely matched participant groups, enabling  
36 within-group comparisons to constitute a point of reference. The implications of this  
37 for stereotypy-based measures of association behaviour are clear: norms lists must be  
38 selected, or compiled, to reflect the demographic profile of the target population. In  
39 this study we have found an age, or generational, difference, and this has direct  
40 relevance, for example, to the way WA tasks have been used in SLA research to  
41 assess L2 proficiency: often the experiment group has a somewhat restricted age  
42 profile (they are typically university undergraduates), which differs considerably from  
43 that of the norming group (see Meara (1978) and Kruse et al. (1987) in Table 2). It is  
44 possible that other factors such as educational background or gender might also affect  
45 response norms.  
46

47 Using the age-appropriate norms lists we produced stereotypy scores for all  
48 participants, reflecting the number of primary dominant responses (i.e. those at the top  
49 of the norms lists) they produced. Large individual differences in stereotypy proved  
50 consistent, with a significant test-retest correlation of .855. In terms of response  
51 category analysis, a principal components analysis indicated a slightly different  
52 grouping of subcategories from that used in previous studies. Again, a test-retest  
53 analysis produced significant positive correlations in all main categories.  
54

55 Taken together, the evidence presented in this study moves the field of WA  
56 research forward in a number of ways. Firstly, the test-retest data, the establishment of  
57 norming criteria and the confirmation of category clusters all contribute towards an  
58  
59  
60

1  
2  
3 argument for the construct validity and the reliability of this method of investigation.  
4 Secondly, it proposes a principled protocol for the analysis of WA data, facilitating  
5 comparison of data sets and making transparent the assumptions and procedures that  
6 underpin the methodology and analytic framework. As we have acknowledged  
7 throughout, specific research questions may motivate changes to the way association  
8 data is measured. For example, measures of idiosyncrasy will complement stereotypy  
9 scores, and particular subcategories of association type will be salient to the study of  
10 certain variables. The studies summarized in Tables 1 and 2 of this paper are evidence  
11 that researchers in diverse fields, for well over half a century, have seen the potential  
12 of WA protocols to investigate lexical behaviour in conditions of development,  
13 decline and impairment. By understanding the implications of methodological  
14 decisions, and by basing further studies on a consistent approach, it will be possible to  
15 maximize both the mutually informative nature of inter-study comparisons, and the  
16 degree to which findings can be interpreted in a meaningful way.  
17  
18

## 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60

### References

- Bauer, L. M., and I. S. P. Nation.** 1993. 'Word families.' *International Journal of Lexicography*, 6, 253-279.
- Coxhead, A.** 2000. 'A new Academic Word List', *TESOL Quarterly*, 34(2), 213-238.
- Croft, W. and Cruse, D. A.** 2004. *Cognitive Linguistics*. Cambridge: Cambridge University Press.
- Dijkstra, T., and W. J. B. van Heuven.** 1998. 'The BIA model and bilingual word recognition', in J. Grainger and A. M. Jacobs (Eds.), *Localist Connectionist Approaches to Human Cognition* (pp. 189-225). New Jersey: Lawrence Erlbaum.
- Ellis, N. C.** 1998. 'Emergentism, connectionism and language learning.' *Language Learning*, 48(4), 631-664.
- Entwisle, D. R.** 1966. *The Word Associations of Young Children*. Baltimore: John Hopkins University Press.
- Entwisle, D. R., D. F. Forsyth, and R. Muuss.** 1964. 'The syntagmatic-paradigmatic shift in children's word associations.' *Journal of Verbal Learning and Verbal Behaviour*, 3, 19-29.
- Ervin, S.** 1961. 'Changes with age in the verbal determinants of word association.' *American Journal of Psychology*, 74, 361-372.
- Author** 2006.
- Author** 2007.
- Author** 2009.
- Author** 2011.
- Galton, F.** 1879. 'Psychometric experiments.' *Brain* 2, 149-162.
- Gewirth, L. R., A. G. Shindler, and D. B. Hier.** 1984. 'Altered patterns of word associations in dementia and aphasia.' *Brain and Language*, 21(2), 307-317.
- Gollan, T. H., D. P. Salmon, and J. L. Paxton.** 2006. 'Word association in early Alzheimer's disease.' *Brain and Language*, 99, 289-303.
- Gough, H. G.** 1976. 'Studying creativity by means of word association tests.' *Journal of Applied Psychology*, 61(3), 348-353.
- Harzing, A.W.** 2007. *Publish or Perish*, from <http://www.harzing.com/pop.htm>.
- Henriksen, B.** 2008. 'Declarative lexical knowledge', in D. Albrechtsen, K. Hastrup and B. Henriksen, *Vocabulary and Writing in a First and Second Language: processes and development* (pp. 22-66). Basingstoke: Palgrave Macmillan.

- 1  
2  
3 **Higginbotham, G. M.** 2010. 'Individual learner profiles from word association tests:  
4 The effect of word frequency.' *System*, 38(3), 379-390.
- 5 **Hirsh, K. W., and J. T. Tree.** 2001. 'Word association norms for two cohorts of  
6 British adults.' *Journal of Neurolinguistics*, 14(1), 1-44.
- 7 **Jung, C. G.** 1910. 'The association method.' *The American Journal of Psychology*,  
8 21(2), 219-269.
- 9  
10 **Kent, G. H., and A. J. Rosanoff.** 1910. 'A study of association in insanity.' *American*  
11 *Journal of Insanity*, 67, 37-96, 317-390.
- 12 **Kiss, G.R., C. Armstrong, and R. Milroy.** 1973. *An Associative Thesaurus of*  
13 *English*. EP Microfilms, Wakefield.
- 14 **Kruse, H., J. Pankhurst, and M. Sharwood Smith.** 1987. 'A multiple word  
15 association probe in second language acquisition research.' *Studies in Second*  
16 *Language Acquisition*, 9(2), 141-154.
- 17 **Landauer, T. K., Foltz, P. W. and Laham, D.** 1998. 'An Introduction to Latent  
18 Semantic Analysis.' *Discourse Processes*, 25(2&3), 259-284
- 19 **Marslen-Wilson, W. D.** 1987. 'Functional parallelism in spoken word-recognition.'  
20 *Cognition*, 25(1-2), 71-102.
- 21 **Meara, P.** 1978. 'Learners' word associations in French.' *Interlanguage Studies*  
22 *Bulletin*, 3, 192-211.
- 23  
24 **Meara, P.** 1983. 'Word associations in a foreign language: a report on the Birkbeck  
25 vocabulary project.' *Nottingham Linguistic Circular*, 11(2), 29-38.
- 26 **Meara, P.** 2009. *Connected Words*. Amsterdam: John Benjamins.
- 27 **Merten, T.** 1992. 'Wortassoziation und Schizophrenie - eine empirische Studie.'  
28 *Nervenarzt*, 63, 401-408.
- 29 **Merten, T.** 1993. 'Word association responses and psychoticism.' *Personality and*  
30 *Individual Differences*, 14, 837-839.
- 31  
32 **Merten, T., and I. Fischer.** 1999. 'Creativity, personality and word association  
33 responses: associative behaviour in forty supposedly creative persons.' *Personality*  
34 *and Individual Differences*, 27(5), 933-942.
- 35 **Miller, E. N., and L. J. Chapman.** 1983. 'Continued word association in  
36 hypothetically psychosis-prone college students.' *Journal of Abnormal Psychology*,  
37 92(4), 468-478.
- 38 **Miller, G. A.** 1995. 'WordNet: A Lexical Database for English.' *Communications of*  
39 *the ACM*, 38(11), 39-41.
- 40  
41 **Author.** 2011.
- 42 **Author.** 2010.
- 43 **Namei, S.** 2004. Bilingual lexical development: a Persian-Swedish word association  
44 study. *International Journal of Applied Linguistics*, 14(3) 363-388.
- 45 **Nation, I. S. P.** 2001. *Learning Vocabulary in Another Language*. Cambridge:  
46 Cambridge University Press.
- 47 **Nelson, D. L., C. L. McEvoy, and T. A. Schreiber.** (1998). The University of South  
48 Florida word association, rhyme, and word fragment norms.
- 49 **Nelson, K.** 1977. 'The syntagmatic-paradigmatic shift revisited: A review of research  
50 and theory.' *Psychological Bulletin*, 84, 93-116.
- 51 **Palermo, D. S., and J. J. Jenkins.** 1964. *Word Association Norms: Grade School*  
52 *through College*. Minneapolis: University of Minnesota Press.
- 53 **Postman, L. J., and G. Keppel (Eds.).** 1970. *Norms of Word Association* New York:  
54 Academic Press.
- 55  
56 **Rosenzweig, M. R.** 1970. 'International Kent-Rosanoff word association norms  
57 emphasizing those of French male and female students and French workmen', in L. J.  
58  
59  
60

1  
2  
3 Postman and G. Keppel (Eds.), *Norms of Word Association* (pp. 95-176). New York:  
4 Academic Press.

5 **Russell, W. A., and J. J. Jenkins.** 1970. 'The complete Minnesota norms for  
6 responses to 100 words from the Kent-Rosanoff Word Association Test', in *Norms of*  
7 *Word Association*. New York: Academic Press.

8 **Author.** 2009.

9 **Schmitt, N.** 2010. *Researching Vocabulary*. Basingstoke: Palgrave Macmillan.

10 **Söderman, T.** 1993. 'Word associations of foreign language learners and native  
11 speakers – different response types and their relevance to lexical development', in B.  
12 Hammarberg (Ed.) *Problems, Process and Product in Language Learning*. Abo:  
13 AfinLA.

14 **Sommer, R.** 1901. *Diagnostik der Geisteskrankheiten*. Berlin/Wien: Urban und  
15 Schwarzenberg.

16 **Wolter, B.** 2002. 'Assessing proficiency through word associations: is there still  
17 hope?' *System*, 30, 315-329.

18 **Author.** 2002.

19 **Author.** 2008.

20 **Author.** 2003.

21 **Author.** 2004.

22 **Zareva, A., and B. Wolter.** 2012. 'The 'promise' of three methods of word association  
23 analysis to L2 lexical research.' *Second Language Research*, 28(1), 41-67.

24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

---

<sup>1</sup> Previous outputs from this collaboration include Author 2010; Author 2011.

<sup>2</sup> Two of these cue words, and the responses they elicited, were subsequently excluded from analyses

<sup>3</sup> The WA task was presented in written rather than spoken mode for three reasons. Firstly, it was not feasible to collect both written and spoken responses from the same informants, unless in the same short timeslot of the same day, when fatigue and/or repetition effects would confound the results. The data were collected as part of a larger study, with very little scope to manipulate the order of presentation or to extend the overall time taken for the WA element. Given this constraint, the main consideration was which mode to prefer. The written mode was preferable because, secondly, a team of research assistants was involved in data collection, and it would not be possible to guarantee consistency of delivery of spoken cues. And thirdly, the majority of WA studies in applied linguistics use written data, and employing that same elicitation method maximised the relevance of our study to others. Clearly the mode of delivery is a significant variable, and future research needs to extend to a methodical comparison of the responses from participants under both conditions.

<sup>4</sup> For a practical approach to justifying the identification of wordstrings as formulaic sequences, see Author 2003, Author 2008 chapter 9.

<sup>5</sup> Subsequently, for the purposes of validity evaluation, the norming groups were further divided to enable both within- and between- age group analyses.

<sup>6</sup> We also calculated 'weighted stereotypy' and 'idiosyncrasy' scores for some other aspects of our study. In the former, respondents gained a score derived from the number of norms list contributors providing the same response; in the latter, respondents gained a score for every response they gave that no-one else has produced.

<sup>7</sup> The assignment to 'twin 1' or 'twin 2' was random: on the advice of the geneticists in the team, birth order was not considered a variable.

Age: \_\_\_\_\_ Date: \_\_\_\_\_ Total time taken: \_\_\_\_\_

Please write down the first word you think of when you read each of the words listed below.  
There are no right or wrong answers.

abuse	child	joint	bones
agenda	Plan	landlord	bossy
annoy	nuisance	loss	like
attack	injure	mathematics	confusing
bean	green	miner	dirty
blame	fault	nail	varnish
bread	salt	nurse	care
candidate	perfect	owe	payback
cheese	crackers	permit	allow
cloud	rain	plug	drain
concentrate	focus	prevent	accident
cope	like	pudding	yummy
cupboard	dresses	reflect	think
delay	buses	repair	mend
diet	healthy	rock	concert
domestic	housework	sand	grainy
effort	try	session	class
establish	facts	sin	bad
extension	house	source	food
fence	white	store	items
fraction	layers	swear	oath
gold	jewellery	thick	cream
heaven	god	tour	guide
ideal	best	variety	park
instance	moment	weak	link

1

Figure 1: Excerpt from data set

**Table 1: Subjects, cues, norms lists, response treatment and measures used in the most cited WA studies investigating function, personality and psychosis.**

STUDY and VARIABLE		SUBJECTS	CUES	NORMS LIST	TREATMENT OF RESPONSES	MEASURES
age	Entwisle, Forsyth & Muuss (1964) <b>The syntactic-paradigmatic shift in children's WAs</b>	500 x children aged 5-11	24 high-frequency words: 8 nouns; 8 adjectives; 8 verbs	n/a	grammatical analysis; subjective judgement made of 'transitional probabilities'	1) syntactic/non-syntactic (class) 2) homogeneous/heterogeneous 3) form class of responses
	Ervin (1961) <b>Changes with age in the verbal determinants of WA</b>	23 x kindergarten 10 x 1 <sup>st</sup> grade 52 x 3 <sup>rd</sup> grade 99 x 6 <sup>th</sup> grade	46 cues in vocabulary range of youngest children, 39 of which elicit antonyms or coordinates	n/a	principled classification according to grammatical class, sequential analysis	paradigmatic (strict grammatical interpretation)/syntactic and text-informed interpretation
	Hirsh & Tree (2001) <b>WA norms for two cohorts of British adults</b>	45 x young adults 45 x older adults	90 concrete nouns and items likely to elicit concrete nouns	compiled from participant responses	plurals lemmatized	1) dominant/unique/syntactic 2) response variation 3) propositional-relational/hierarchical
cognitive function	Gewirth Shindler & Hier (1984) <b>Altered patterns of WA in dementia and aphasia</b>	38 x demented 17 x aphasic 22 x normal	16 cues from Palermo & Jenkins (1964): 4 nouns; 4 verbs; 5 adjectives; 3 adverbs	Palermo & Jenkins (1964)	no information given	1) popular/unpopular (list) 2) paradigmatic/syntactic identity (identical or not)
	Gollan Salmon & Paxton (2006) <b>WA in early Alzheimer's disease</b>	18 x probable AD 18 x elderly normals	52 cues from Nelson, McEvoy & Schreiber (1998): 26 eliciting strong and 26 eliciting weak associations	Nelson, McEvoy & Schreiber (1998)	in multi-word responses, most strongly associated word is scored; responses lemmatized to strongest association	1) 'mean response strength' (according to % of non-same responses) 2) semantic/form/both word/unrelated/non-verb



personality	Gough (1976) <b>Studying creativity by means of WA tests</b>	45 x research scientists 66 x engineering students	100 Kent & Rosanoff (1910) cues	Russell & Jenkins (1970)	no information	close/remote associations following percentages: >50%; 25-50%; 10-25%
	Merten & Fischer (1999) <b>Creativity, personality and WA responses</b>	40 x 'artistic' professionals 40 x schizophrenics 40 x normals	25 common nouns	normative sample from Merten (1992)	no information	For each of 3 conditions (common and individual response) 1) number 'common' response) 2) number of 'individual' norms list)
psychosis	Merten (1993) <b>WA responses and psychoticism</b>	46 psychiatric hospital staff and non-medical professionals	25 common nouns	normative sample from Merten (1992)	no information	For each of 3 conditions (common and individual response) 1) number 'common' response) 2) number of 'individual' norms list)
	Miller & Chapman (1983) <b>Continued WA in hypothetically psychosis-prone college students</b>	60 x probable psychosis-prone 21 x controls	32 cues from Kent & Rosanoff (1910) lists with >10 consensual associations	norms compiled from 120 male students	no information	1) popular (>25% in norms) one occurrence in norms) 2) of idiosyncratic responses (deviant) /unusual (deviant)

Table 2: Subjects, cues, norms lists, response treatment and measures used in the most cited WA studies investigated

STUDY	SUBJECTS	CUES	NORMS LIST	TREATMENT OF RESPONSES	MEASURES
Fitzpatrick (2006) <b>Habits and rabbits: word associations and the L2 lexicon</b>	40 learners of English (mixed L1) 40 native speakers of English	60 cues selected from the Academic Word List (Coxhead 2000)	n/a	post-task interviews to confirm motivation for response	divided into 3 categories: position-, form-based and content-based subcategories
Kruse, Pankhurst & Sharwood Smith (1987) <b>A multiple WA probe in second language acquisition</b>	15 x Dutch learners of English 7 x native speakers of English	10 cues selected from Postman & Keppel (1970)	Postman & Keppel (1970)	no information	1) number of responses 2) weighted scores for each response according to whether response was correct 3) non-weighted scores for whether response was correct
Meara (1978) <b>Learners' WAs in French</b>	76 x female English learners of French	French translations of 100 Kent & Rosanoff (1910) cues	Rosenzweig (1970) (female list)	no information	1) primary response 2) primary response list 3) primary response list
Namei (2004) <b>Bilingual lexical development: a Persian-Swedish WA study</b>	100 x Persian-Swedish bilinguals aged 6-22: 50 Swedish L1 aged 6-18 50 Persian L1 aged 6-19	Persian and Swedish translations of 100 Kent & Rosanoff (1910) cues	n/a	responses 1) phonemically transcribed 2) translated into English	categorised as class paradigmatic / misparadigmatic
Söderman (1993) <b>Word associations of foreign language learners and native speakers</b>	112 x Finnish learners of English: 28 each from 7 <sup>th</sup> grade; Gymnasium; 1 <sup>st</sup> yr university; advanced learners Expt 2 only: 28 native speakers of English	Expt 1: 100 Kent & Rosanoff (1910) cues Expt 2: 64 cues (mostly adjectives): 32 frequent 32 infrequent	n/a	no information	categorised as class paradigmatic / other
Wolter (2002) <b>Assessing proficiency through word associations: is there still hope?</b>	30 x Japanese learners of English 42 x native speakers of English	20 verbs from Edinburgh Associative Thesaurus, excluding items eliciting dominant primary response or high number of idiosyncratic responses	Edinburgh Associative Thesaurus (Kiss et al. 1973)	1) multiword responses reduced to head word 2) responses lemmatised	1) non-weighted scores whether response was correct 2) weighted scores for native speakers whether response was correct

**Table 3: Mean scores (and standard deviations) for four measures of WA stereotypy**

		Comparison norms list			
		Own list	Same age other twin	Other age twin 1	Other age twin 2
<b>Participant group</b> (n=48 per group)	<b>Young Twin 1</b>	28.31 (6.68)	25.21 (6.39)	18.71 (7.18)	19.81 (7.01)
	<b>Young Twin 2</b>	27.31 (6.12)	27.33 (6.98)	18.54 (5.60)	19.38 (6.02)
	<b>Older Twin 1</b>	27.10 (10.36)	25.00 (9.69)	19.02 (7.35)	17.85 (7.32)
	<b>Older Twin 2</b>	26.00 (8.05)	24.71 (7.78)	20.65 (6.35)	18.81 (6.08)
<b>Overall mean</b>		27.18	25.56	19.23	18.96

For Peer Review

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

For Peer Review

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**Table 4: Test retest - stereotypy scores with correlation coefficient**

n=36	Test 1			Test 2			Correlation
	min	max	mean	min	max	mean	
<b>Stereotypy</b>	4	42	23.86 (8.371)	8	39	23.78 (7.388)	.855**

\*\* p < .01

For Peer Review

**Table 5: Response items repeated at test time two (maximum 100)**

n=36	min	max	mean	sd
<b>repeated items</b>	8	54	25.53	9.667

For Peer Review



Table 6: Sub-categories used to classify WA responses

Sub-category	Definition	Example
Synonym	Cue and response are synonymous in some situations	<i>delay</i> → <i>impede</i> <i>fraction</i> → <i>portion</i> <i>establish</i> → <i>build</i>
Lexical set	Cue and response share a hyponym, or one word in the pair is an example of the other; includes antonyms	<i>bean</i> → <i>pea</i> <i>bean</i> → <i>vegetable</i> <i>permit</i> → <i>deny</i>
Other conceptual	Cue and response are related in meaning, but are not synonyms or in the same lexical set	<i>fence</i> → <i>field</i> <i>sin</i> → <i>prayer</i> <i>nurse</i> → <i>illness</i>
Cue-response collocation	Cue is followed by the response in common usage; includes compound nouns	<i>fence</i> → <i>post</i> <i>rock</i> → <i>roll</i> <i>swear</i> → <i>word</i>
Response-cue collocation	Cue is preceded by the response in common usage; includes compound nouns	<i>fence</i> → <i>electric</i> <i>candidate</i> → <i>nominate</i> <i>plug</i> → <i>spark</i>
Cue-response and response-cue collocation	Cue could precede or follow the response in a common phrase(s)	<i>rock</i> → <i>hard</i> <i>dog</i> → <i>eat</i>
Affix manipulation	Cue is the response with the addition, deletion or changing of an affix	<i>irony</i> → <i>ironic</i> <i>abuse</i> → <i>abusive</i> <i>plug</i> → <i>unplug</i>
Similar in form only	Cue and response are similar in orthography and/or phonology but do not share meaning	<i>fence</i> → <i>hence</i> <i>weak</i> → <i>week</i>
Two step association	Cue and response appear linked only through another word	<i>weak</i> → <i>monday</i> (via <i>week</i> ) <i>owe</i> → <i>mine</i> (via <i>own</i> )
Erratic	The link between cue and response seems illogical. Includes repetition of the cue	<i>wolf</i> → <i>and</i> <i>heaven</i> → <i>heaven</i>
Lexical set <i>and</i> cue-response collocation		<i>bread</i> → <i>cheese</i> <i>gold</i> → <i>silver</i> <i>heaven</i> → <i>hell</i>
Lexical set <i>and</i> response-cue collocation		<i>cheese</i> → <i>bread</i> <i>nurse</i> → <i>doctor</i>
Synonym <i>and</i> cue-response collocation		<i>torch</i> → <i>light</i>
Synonym <i>and</i> response-cue collocation		<i>shove</i> → <i>push</i>

**Table 7: Test retest - mean category scores and correlation coefficients (categories represented by an average of <1 response per participant are not included)**

n=36	Test 1 (sd)	Test 2 (sd)	Correlation
Synonym	17.17 (8.062)	14.61 (6.478)	.721**
Lexical set	5.81 (2.877)	6.06 (3.189)	.521**
Other conceptual	51.42 (9.749)	52.28 (9.254)	.824**
Cue-response collocation	10.86 (6.095)	12.25 (5.406)	.724**
Response-cue collocation	6.47 (3.247)	6.97 (2.932)	.518**
Erratic	1.06 (1.548)	1.22 (1.606)	.259 ns

\*\* p<.001

**Table 8: Rotated Component Matrix (Factor loadings below 0.5 have been suppressed)**

	Component				
	Meaning	Position	Form	Multi-position	Position plus meaning
Other conceptual	-.822				
Synonym	.717				
Lexical set	.709				
Cue-Response		.816			
Response-Cue		.672			
Two step			.641		
Erratic			.617		
Affix			.548		
Form only			.535		
Cue-Response-Response-Cue				.788	
Lexical set plus Response-Cue					-.743
Synonym plus Cue-Response					.685