



SSDI 0887-6177(95)00033-X

Benton Controlled Oral Word Association Test: Reliability and Updated Norms

R. M. Ruff, R. H. Light, and S. B. Parker

University of California, San Diego

H. S. Levin

University of Maryland, Baltimore

The aim of this paper is to update the over 20-year-old normative data for the Benton Controlled Word Association (COWA) Test. In a sample of 360 normal volunteers, the age ranged between 16–70 years, and the educational level ranged from 7–22 years. Care was taken to ensure that the population was heterogenous, yet the two stratifications of gender, four age, and three educational groups led to 24 cells with 15 individuals in each. Test-retest reliability was established by testing 30% of the sample after a 6-month delay, which represents a typical follow up duration between testings in a clinical setting. The two forms of the COWA revealed significant test-retest reliability. Generally, our updated values fall above the original normative values, which were derived from a less well-educated and rural sample. No major gender or age trends were noted, but the COWA test performances were influenced by education, i.e., as the level of education increased, the performance on the COWA increased. The only gender differences that were found were for the women in the highest educational group (> 16 years), who performed significantly better than men in the highest educational group. An error analysis of repetitions or perseverations is provided, with cut-off scores according to age levels. Finally, the updated COWA norms are compared to the original norms as well as to other measures of word fluency.

Arthur Benton was stimulated by Brenda Milner's report (1964) of her application of the Thurston Word Fluency Test to patients with focal brain lesions. Thereafter, in his clinical setting, he administered this test to about 20 brain-diseased patients and some controls. He found the procedure to be unsuitable on a number of counts. It was not applicable to paresis of the writing hand, thus precluding its use in many patients with left hemisphere disease.

Dr. Arthur Benton was kind enough to review our manuscript, and he agreed that an update of the normative data is of clinical value. Dr. Benton also provided the history behind developing the COWA test, which is contained in the introduction. His editorial comments were appreciated for both expansion of key issues and his exceptional conceptual thinking. We were also impressed by the time and care the ACN reviewers devoted to their task. We are indebted to their contribution.

Address correspondence to: Ronald M. Ruff, PhD, Neurobehavioral Rehabilitation, St. Mary's Hospital and Medical Center, 450 Stanyan Street, San Francisco, CA 94117.

Patients with limited educational background found the test to be too difficult. Some older patients, perhaps because of minor arthritic disability, found the writing task to be effortful and unpleasant. Finally, Benton opined that the Thurston Test, which takes 9 minutes, was unduly lengthy, tedious, and fatiguing for evaluating patients who are often in poor physical condition. Therefore, Benton devised an oral version of the procedure, the FAS, which was shortened to 3 minutes and circumvented the above-stated weaknesses (Bechtoldt, Benton, & Fogel, 1962; Fogel, 1962). The FAS test also became part of the Neurosensory Center Examination for Aphasia (Benton, 1967; Spreen & Benton, 1969).

Some time later, Benton conceived the idea of developing a relatively brief aphasia test battery with equivalent versions (Benton, 1969). At that time, Norman Geschwind pointed out that the term "word fluency" was at risk of being confused with the "fluency/nonfluent" dimension of aphasic speech. Benton agreed with this point, and coined the term "Controlled Oral Word Association" to designate the procedure, which in fact, more described what was demanding it. Under this new name, the test became part of the Multilingual Aphasia Examination. In developing this examination, Benton used new sets of letters (CFL and PRW) as stimuli. These letters were not chosen at random (as had been the case with FAS), but on the basis of their difficulty as defined by the number of words beginning with each letter to be found in standard dictionaries of the English language. As a consequence, the two forms of the COWA have been determined to be of equivalent difficulty (Benton, Hamsher, & Sivan, 1994). Note that the same procedure was employed in selecting letter stimuli for the Spanish version of the test battery, with the result that different sets of letters more appropriate for Spanish-speaking subjects were adopted (Rey & Benton, 1991).

The aim of this study is to update the normative data and to analyze specific components of the Benton Controlled Oral Word Association (COWA) tests. The first reason for the update is that over 20 years have passed since the introduction of the COWA test (Bechtoldt, Benton, & Fogel, 1962). Secondly, Spreen and Strauss (1991) concluded that the original Spreen-Benton (1969) normative data were based on a rural sample with limited educational background. Because education may play a role (Read, 1987), and because gender differences have been reported (with girls performing at a superior rate to age-matched boys; Gaddes & Crockett, 1975), a normative sample stratified according to age, education, and gender is called for.

Two versions of the COWA test are available (Version A, using the stimulus letters C, F, and L; version B, using P, R, and W), and the test-retest reliability was examined in our study. Within the test, three letters are chosen with an increasing level of difficulty for each successive letter. The rate of fluency for each of these letters was analyzed and compared, providing aspects of consistency within our sample. A further analysis addressed the errors made, which principally include (at least in normal samples) repetitions or perseverations of words. No normative values for potential cutoff scores for these perseverative errors have been available, and cutoff scores which are clinically relevant were analyzed.

A further aim of this study was to compare our updated normative values for the COWA test to the earlier norms. Normative values of other fluency measures are also discussed in the context of trends for age and gender.

METHOD

Subjects

The total sample of 360 normal volunteers ranged in age between 16 and 70 years and in education from 7 to 22 years. All participants were native English-speaking individuals.

About 65% of the participants resided in California, 30% resided in Michigan, and the rest resided on the eastern seaboard. The majority of our sample resided in urban or suburban regions, and only a small minority (approximately 5–10%) resided in rural areas. Care was taken to ensure that the population was heterogeneous with respect to age and education; there were four age and three education groups (see Table 1). To assess test-retest reliability as well as stability, five or more randomly selected subjects from each of the 12 cells (see Table 1) were retested after a 6-month delay (i.e., 30% of sample was retested). This interval was chosen because it typically represents the interval often selected for follow-up testing within a clinical setting for evaluating gains or deterioration of neuropsychological functioning. All participants were screened to exclude those with a positive history of psychiatric hospitalization, chronic poly-drug abuse, or neurological disorders.

Procedure

The examiner's instructions for the Benton Controlled Oral Word Association (COWA) Test are as follows: "I am going to say a letter of the alphabet, and I want you to say as quickly as you can all the words you can think of which begin with that letter. You may say any word at all except proper names, such as names of people or places. So you would not say "Rochester" or "Robert." Also, do not use the same words again with a different ending, such as "eat" and "eating."

"For example, if I say "S," you could say, "sun," "sit," "shoe," or "slow." Can you think of other words beginning with the letter "S?" Wait for the subject to give a word, indicate if the word is correct, and ask the subject to give another word beginning with the letter "S." Once two appropriate words beginning with the demonstration letter are given, say, "That is fine. Now I'm going to give you another letter, and again say all the words beginning with that letter that you can think of. Remember, no names of people or places, just ordinary words. Also, if you should draw a blank, I want you to keep on trying until the time limit is up. You will have a minute for each one." The first letter is C, and 1 minute is allowed, and the same applies for the letters F and L of Version A. In Version B, the same procedure is used with the alternate letters PRW.

The record sheet provides numbered lines on which the subject's responses can be entered. If the speed of production is too fast to permit verbatim recording, then a + sign should be entered. However, all incorrect responses should be entered verbatim. Many words have two or more meanings, and a repetition of the word is accepted only in those cases where the subject definitely indicated an alternate meaning. If the patient produces one or more questionable responses (e.g., "frank," which could represent a proper name), the association is simply recorded, and the subject is not interrupted. However, at the end of the 1-minute period of association the subject should be asked what was meant by this

TABLE 1
Sample of Participants Arranged by Age and Education

Age	Education					
	12 years or less		13–15 years		16 years or more	
	Men	Women	Men	Women	Men	Women
16–24	15	15	15	15	15	15
25–39	15	15	15	15	15	15
40–54	15	15	15	15	15	15
55–70	15	15	15	15	15	15

word. Slang terms are generally admissible, as well as certain foreign words (e.g., "lasagna"), as long as these words are listed as standard English.

A perseverative error is defined as a word that is repeated. In the case where a word has two or more meanings, a perseveration is scored if the subject does not indicate the different meanings, for example, "four" (the number) and "for" spelled f-o-r.

In addition to COWA Test, we administered a comprehensive neuropsychological battery, and for the present study, we coanalyzed the IQ scores from the WAIS-R (Wechsler, 1981) intrusion errors from the Selective Reminding Tests (Buschke & Fuld, 1974; Ruff, Light & Quayhagen, 1988), and also the perseverative errors from the Ruff Figural Fluency Test (Ruff, 1988; Ruff, Light, & Evans 1987). All measures were administered by specifically trained psychometrists.

RESULTS

Reliability

Reliability of the COWA was assessed in two ways. First, a coefficient alpha was computed by taking the total number of words generated for each letter separately as three individual items, and the summation of these scores as the COWA total test score; this provides a measure of internal consistency. The coefficient alpha of $R = .83$ was acceptably high to indicate that even though the test is relatively short (i.e., three item letters), the items (number of words generated for each letter) had a high enough average intercorrelation ($R = .61$) to guarantee accurate measurement and high test homogeneity.

The second procedure for evaluating reliability examined test score stability over time. A randomly selected sample of 120 subjects was retested on the alternate version of the COWA 6 months following the initial testing. As is the clinical practice, version one was always preceded by version two. Table 2 compares demographic and baseline verbal fluency and Intelligence variables between the selected subsample of 120 subjects vs. those 240 subjects not retested. None of the differences were significant. The correlation between scores from the first and the second testing (test-retest reliability) was significant ($R = .74$,

TABLE 2
Comparison of Subsamples Used for Reliability Testing

Variable	Subsample Retested ($n = 120$)	Subsample Not Retested ($n = 240$)
Demographic		
Age	40.5	40.4
Education	14.0	14.2
Gender m/f	60/60	120/120
Controlled Oral Word Association— Baseline		
Sum total	39.7	40.3
Perseveratives	0.82	0.70
Percentile	64.7	66.6
Intelligence		
Verbal IQ	110.3	111.3
Performance IQ	107.3	107.9
Full Scale IQ	110.0	111.0

Note. None of the differences in the above variables were found to be statistically significant according to a *t*-test for independent sample at the $p > .05$ level of confidence; gender was not statistically compared. IQ was measured according to the WAIS-R (Wechsler, 1981).

$p < .001$). However, the overall mean from the first testing for the same 120 subjects was 39.7 ($SD = 10.48$) and on the second testing 42.5 ($SD = 9.9$). This gain of approximately three words proved to be significant [$t(1, 119) = 4.19, p < .0001$], and may indicate a practice effect. Nonetheless, the two reliability measures indicate that the COWA is a reliable instrument that is reasonably stable over time.

Effect of Age, Gender, and Education

A three-way analysis of variance (age \times gender \times education) was carried out on the mean combined words for all three letters. Age had no significant effect. However, gender moderated the effect of education on the number of words produced [$F(2, 336) = 4.33, p = .014$]. This interaction accounted for 2% of the total variance in COWA scores. Although the effect of the interaction of gender and education was reliable, this interaction was ordinal, in that the main effect of level of education positively predicted COWA performance for both men and women. Thus, the effect of education was significant [$F(2, 336) = 16.21, p < .0001$], and differences due to education alone accounted for a greater proportion of total variance (8%).

To determine the sources of the significance of the interaction, the data were further analyzed. A Fisher protected least significant difference (Cohen & Cohen, 1975) was computed and relevant means were compared; overall p -values were maintained at the .05 level. On average, men with 12 years of education or less produced significantly fewer words than their cohorts with some college education or those with college degrees or more education. There were no significant mean differences between groups of men with some college education and the group of men with 16 years of education or more. With respect to the women, there were significant mean group differences across all three educational categories. Furthermore, in the first two educational groups, there were no significant mean differences between men and women; however, in the highest educational group, women produced significantly more words on average than did their male counterparts. Table 3 presents the means and standard deviations of measurement for total COWA scores segregated by gender and education.

Consistency Ratings Among Letters

The letters were originally chosen so that an increasing level of difficulty was encountered for each successive letter. A set of paired sample t -tests confirmed a significant increase in

TABLE 3
Mean Values for the Controlled Oral Word Association Test, Separately by Gender and Education

Education	Men $n = 180$		Women $n = 180$		Combined Sex $n = 360$	
	Mean	SD	Mean	SD	Mean	SD
12 years or less	36.9 ^{a,b}	9.8	35.9 ^{c,d}	9.6	36.5	9.9
13–15 years	40.5 ^a	9.4	39.4 ^{c,e}	10.1	40.0	9.7
16 years and up	41.0 ^{b,f}	9.3	46.5 ^{d,e,f}	11.2	43.8	10.6
All education levels	39.5	9.8	40.6	11.2	40.1	10.5

* $n = 180$ represents a combination of the three educational subgroups.

^{a,b}For men, the educational subgroup comparisons were significantly lower for those with up to 12 years of education versus both higher educational groups ($p < .05$).

^{c,d,e}For women, all three educational groups were significantly different ($p < .05$) with a higher rate with increasing years of education.

^fWomen achieved a significantly higher fluency rate as compared to men only in the educational subgroup with the highest years of education ($p < .05$).

difficultly from C to F to L. Specifically, the mean number of words produced for the letter C, Mean = 14.1 ($SD = 4.15$), was higher than the mean number for the letter F, Mean = 13.3 ($SD = 4.10$) [$t(1, 359) = 3.79, p < .0001$]. Furthermore, the mean number of words produced for the letter F was significantly greater than the mean number for the letter L, Mean = 12.7 ($SD = 4.0$) [$t(1, 359) = 3.67, p < .0001$].

With respect to the consistency of this trend across our sample, 62% gave as many or more words to the letter C than F. Similarly, 68% of the sample produced more words to the letter C than L, and 61% generated more words to the letter F than L.

Error Analysis

In analyzing the frequency of repetitions or perseverations on the COWA, a majority (56%) of subjects produced no perseverations at all. The distribution of perseveration scores decreased logarithmically, with the highest number of perseverations being six, which occurred only once in the sample. Thus, the perseveration measure was dichotomized into those who did not perseverate versus those who did. No significant differences were noted in the rate of perseveration due to gender or level of education. However, age moderated the perseveration rate; $\chi^2(3, n = 360) = 8.01, p < .05$. No significant differences in perseveration rate were found between the age groups between 25 to 70; however, those aged 16–24 perseverated at a significantly lower rate (32%) than did the combined age groups from 25 to 70 (49%); $\chi^2(1, n = 360) = 7.26, p < .01$. Table 4 presents the perseveration rates across age levels.

Are the perseverative errors on the COWA related to other measures of perseveration? There was no significant association between COWA perseveration and perseverations on the Ruff Figural Fluency Test. However, the COWA perseverations were significantly positively associated with the intrusion errors on the Selective Reminding Test ($R = .20, p < .001$); intrusion errors are words provided by the subject that are not part of the original 12-word list.

Current Normative Tables and Correction Factors

In order to obtain a percentile and T -score ranking for a given individual, an education adjustment must first be made, which differs between men and women. The correction factors for total number of words generated on the COWA were computed for each cell in a gender by education group matrix. These correction factors are presented in Table 5. After correction of scores, the percentile ranks and normalized T -scores from the 360 subjects are presented in Table 6. With respect to COWA perseveration, suggested clinical cutoffs are presented in Table 7.

DISCUSSION

The aim of the present study was to establish normative values. Performance on the COWA was influenced to the largest degree by education; as the level of education increased,

TABLE 4
Rate of Perseveration on the Controlled Oral Word
Association Test, Separately by Age

Age	n^a	Percent
16–24 years	29	32
25–39 years	42	47
40–54 years	47	52
55–70 years	42	47

^aNumber exhibiting at least one perseveration.

TABLE 5
Correction Factors for the Controlled Oral Word
Association Test, by Gender and Education

Education	Men	Women
12 years or less	+3	+4
13–15 years	–1	+1
16 years or more	–1	–7

the performance on the COWA increased. However, there was an ordinal interaction, because gender accounted for the determining effect of education. Specifically, men in our sample with some college education (13 years and up) performed at a significantly better

TABLE 6
Percentile Ranks, Normalized T-Scores, and Interpretation of Corrected Scores
for the Controlled Oral Word Association Test

Corrected Score	Percentile	T-Score	Interpretation
17 or less	1	26.7	
20	2	29.5	Seriously deficient
21	3	31.2	Deficient
23	4	32.5	
25	5	33.5	Deficient
26	8	35.8	Borderline
27	9	36.6	
28	10	37.2	Borderline
29	13	38.7	Low average
30	16	40.2	
31	19	41.2	
32	21	41.9	
33	27	43.9	Low average
34	30	44.7	Average
35	34	45.9	
36	38	46.9	
37	43	48.2	
38	47	49.2	
39	51	50.3	
40	58	52.0	
41	61	52.8	
42	64	53.6	
43	67	54.4	
44	69	55.0	
45	72	55.8	Average
46	76	57.0	High average
47	78	57.7	
48	80	58.5	
49	82	59.1	
50	85	60.4	
51	87	61.3	
52	89	62.3	High average
53	91	63.4	Superior
54	92	64.1	
55	94	65.6	
56	95	66.5	
58	97	68.9	Superior
60	98	70.6	Very superior
64 and up	99	73.3	

TABLE 7
Suggested Cutting Scores for Controlled Oral Word
Association Perseverations

Perseverations	Percent ^a	Interpretation
0	56	Intact
1	26	Low average
2	11	Borderline
3	5	Deficient
4 and up	2	Seriously deficient

^aPercent of total sample, total $n = 360$.

rate than those with 12 years of education or less. Furthermore, there was no difference between those men with some college (13 years to 15 years) and those that had completed degrees (16 years and up). The results were somewhat different for the women in our sample, as they demonstrated significant differences in performances across all three educational groups. When contrasting the female with the male subjects directly, no differences were found for the first two educational groups, while women in the highest educational group performed significantly better than men in the highest educational group.

An attempt was made to analyze, in addition to the fluency rate, the specific performances on the three letters as well as the perseveration rate. Approximately 60% of our subjects generated more responses for the letter C than F, C, than L, and F than L. This percentage represents a trend yet clinically provides limited utility. However, the analysis of perseverative errors reveals a potentially greater clinical application. Our data indicate that three or more perseverative errors corresponds with a deficient performance.

Assessment of reliability provides further evidence of the value of the COWA as a clinical instrument. The COWA was reliable according to both alpha coefficients and test-retest correlations. Note that with our entire sample, we administered Version A in the initial testing, and Version B in the second testing, and overall a three-word gain was noted, which represented a significant improvement. This may, indeed, be due to a learning effect; however, it is not possible to rule out that the two letter sets may not be totally equivalent. However, for this to be evaluated, the two sets of letters would need to be administered to half the normative sample in the inverse order. This may be the subject of a future study.

Comparison With Earlier COWA Test Norms

The updated values for the three-letter combined scores are consistently above the original norms of Benton, Hamsher, and Sivan (1994). For example, a "seriously deficient" score in the original version was ≤ 16 words, and in the updated norms, this cutoff is increased to ≤ 20 words. At the other extreme, the 99th percentile used to correspond with 58+ words, whereas the updated score increased to 64+ words for the same percentile ranking.

In the middle ranges the increases are slightly less, but consistently one to three words above the original values. With respect to the median adjusted score, an increase between 37.5 to 39 words is noted. This general increase in the updated norms is not necessarily unexpected. Indeed, as Benton (1981) has pointed out, particularly for verbal tests, norms should not be considered stable over time. Changes can be due to educational practices and cultural factors (including the influence of television). The first edition of the WAIS was published in 1955, and when this was compared with the revised WAIS-R (published in 1981), an IQ score of 100 corresponded with an IQ score of 108 on the old version. This represents an increase of approximately 0.5 standard deviation (Wechsler, 1981).

With respect to age and education, no strict comparisons are possible between the original and updated COWA scores because different age breakdowns are utilized. The key difference is that in the original data, a minor but steady correction for age was indicated, whereas in our updated norms, no age dependent corrections were indicated. In part, this could be due to sampling differences. Even when tests are carefully standardized with due regard to the influence of age, a relatively wide range of scores within each age cell can be represented. It is also possible that the greater range in education negated this "age" effect.

Comparison of Normative Trends Using Other Measures of Word Fluency

Cauthen (1978) analyzed the age factor by comparing 51 normal volunteers (39 women and 12 men) between the ages of 20–59, with an older sample of 64 participants (36 women and 28 men) between the ages of 64–94. Each subject was given a 60-second time limit for each of the eight selected letters. No significant differences emerged for the younger sample across decades of age or for gender. However, for the older group, verbal fluency was lower, except for the brighter participants with IQs above 119. Note that in the younger sample, no similar IQ effect was observed.

Yeudall, Fromm, Reddon, and Stefanyk (1986) provided normative values for 225 volunteers ranging from 15–40 years of age using the letters F, A, and S. No significant gender differences emerged for both the oral and written rates of production. Moreover, no marked trends indicated age differences. The sum score for all three letters was at or above 40, which was considerably above the Spreen-Benton (1969) mean score of 33. Note that our updated values fell also at about 40 total words for those with education of more than 13 years. The Spreen-Benton sample was less well educated, and our volunteers of 12 and fewer years of education produced a mean of 36 total words. Thus, our updated normative values are closer to Yeudall et al.'s norms on the FAS. Despite the fact that the FAS and COWA are two versions of the same procedures, and the raw scores on the two versions are not comparable. Only percentile or standard scores are comparable. Finally, it should be pointed out that the written version of the Thurston Word Fluency and the COWA are different procedures, and are not comparable.

REFERENCES

- Bechtoldt, H. P., Benton, A. L., & Fogel, M. L. (1962). An application of factor analysis in neuropsychology. *Psychological Record*, *12*, 147–156.
- Benton, A. L. (1967). Problems of test construction in the field of aphasia. *Cortex*, *3*, 32–58.
- Benton, A. L. (1969). Development of a multilingual aphasia battery: Progress and problems. *Journal of the Neurological Sciences*, *9*, 39–48.
- Benton, A. L. (1981). Basic approaches to neuropsychological assessment. In S. R. Steinhauer, J. H. Gruzelier, & J. Zubin (Eds.), *Handbook of schizophrenia* (Vol. 5). New York: Elsevier Science Publishers.
- Benton, A. L., Hamsher, K. de S., & Sivan, A. B. (1994). *Multilingual Aphasia Examination*. Iowa City: AJA Associates.
- Buschke, H., & Fuld, P. A. (1974). Evaluating storage, retention, and retrieval in disordered memory and learning. *Neurology*, *11*, 1019–1025.
- Cauthen, N. R. (1978). Verbal fluency: Normative data. *Journal of Clinical Psychology*, *34*, 126–129.
- Cohen, J., & Cohen, P. (1975). *Applied multiple regression/correlation analysis for the behavioral sciences*. Hillsdale, NJ: Lawrence Erlbaum Associates.
- Fogel, M. L. (1962). The Gerstmann syndrome and the parietal symptom-complex. *Psychological Record*, *12*, 85–90.
- Gaddes, W. H., & Crockett, D. J. (1975). Spreen-Benton aphasia tests: Normative data as a measure of normal language development. *Brain and Language*, *2*, 3, Jul, 257–280.
- Milner, B. (1964). Some effects of frontal lobectomy in man. In J. M. Warren & K. Akert (Eds.), *The frontal granular cortex and behavior*. New York: McGraw Hill.

- Rey, G. J. & Benton, A. L. (1991). *Examen de Afasia Multilingue*. Iowa City, IA: AJA Associates.
- Read, D. E. (1987). *Neuropsychological assessment of memory in early dementia: Normative data for a new battery of memory tests*. Unpublished manuscript. British Columbia: University of Victoria.
- Ruff, R. M. (1988). *Ruff figural fluency test. Administrative manual*. San Francisco: Neuropsychological Associates.
- Ruff, R. M., Light, R. H., & Evans, R. W. (1987). The Ruff Figural Fluency Test: A normative study with adults. *Developmental Neuropsychology*, 3, 37–51.
- Ruff, R. M., Light, R. H., & Quayhagen, M. (1989). Selective reminding tests: A normative study of verbal learning in adults. *Journal of Clinical and Experimental Psychology*, 11, 539–550.
- Spreen, O., & Benton, A. L. (1969). *Neurosensory Center Comprehensive Examination for Aphasia: Manual of directions*. Victoria, BC: Neuropsychology Laboratory, University of Victoria.
- Spreen, O., & Strauss, E. (1991). *A compendium of neuropsychological tests: Administration, norms, and commentary*. New York: Oxford University Press.
- Wechsler, D. A. (1981). *Manual for the Wechsler Adult Intelligence Scale — Revised*. New York: Psychological Corporation.
- Yeudall, L. T., Fromm, D., Reddon, J. R., & Stefanyk, W. O. (1986). Normative data stratified by age and sex for 12 neuropsychological tests. *Journal of Clinical Psychology*, 42, 918–946.