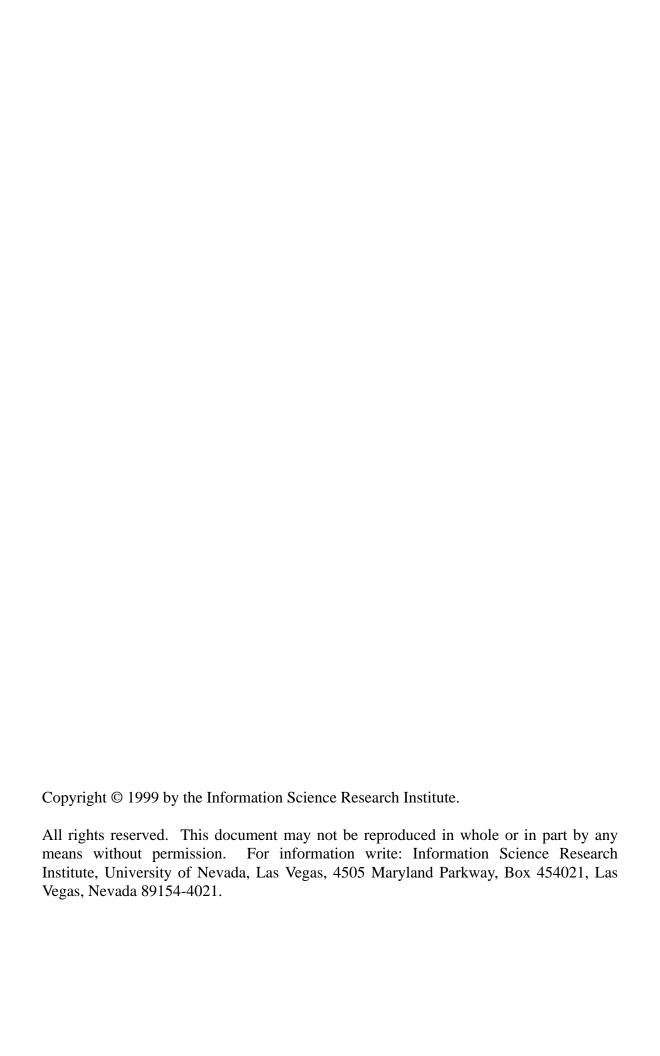
# The OCR Frontiers Toolkit

## Version 1.0

Andrew D. Bagdanov, Stephen V. Rice, and Thomas A. Nartker

**Information Science Research Institute** 

October 1999



# **Contents**

1	Introduction	1
	1.1 Operating Environment	1
	1.2 Special Characters	
	1.3 Zoning	
2	Character Accuracy	3
	2.1 The accuracy Program	
	2.2 The <b>synctext</b> Program	
	2.3 The <b>accsum</b> Program	
	2.4 The <b>groupacc</b> Program	
	2.5 The accci Program	
	2.6 The <b>accdist</b> Program	
	2.7 The <b>ngram</b> Program	
3	Word Accuracy	16
	3.1 The wordacc Program	
	3.2 The wordaccsum Program	
	3.3 The <b>nonstopacc</b> Program	
	3.4 The wordaccci Program	
	3.5 The wordaccdist Program	22
	3.6 The wordfreq Program	
R	References	23

#### 1 Introduction

The Information Science Research Institute (ISRI) at the University of Nevada, Las Vegas (UNLV) tested optical character recognition (OCR) systems on an annual basis from 1992 through 1996. The OCR systems that were evaluated are known as "page readers." These systems take as input a bitmapped image of any document page and attempt to locate and identify the machine-printed characters on the page. Each annual test is described in a technical report (see references [1] through [5]). Synopses of the test results have been published in *Inform* magazine (references [6] through [9]).

Since 1991, ISRI has been actively developing performance measures for pagereading systems. These measures have been used in the annual test and enable a comprehensive evaluation of these systems. The measures include character accuracy, marked character efficiency, word accuracy, non-stopword accuracy, phrase accuracy, and the cost of correcting automatic zoning errors. With the exception of the automatic zoning measure, which is described in [10], a formal specification of these measures, including algorithms for computing them, is presented in Steve Rice's doctoral dissertation [11].

This toolkit contains a subset of the tools developed at ISRI for conducting large-scale, automated tests of page readers. Each program is described in this document one by one with examples.

#### 1.1 Operating Environment

The Analytic Tools are command line programs that are suitable for use in batch-oriented shell scripts (which are essential to large-scale, automated testing). Each program is given command-line parameters and performs its task in a non-interactive manner. If processing is successful, a zero exit status is returned; otherwise, a non-zero status is returned and an error message is written to the standard error output.

The usage of a program can be displayed by entering the name of the program with no arguments. Specifying "-h" or "-help" as an argument will also display the usage. A Unix *man* page is available for each program.

#### 1.2 Special Characters

The user must supply OCR-generated text files and correct ("ground truth") text files. The Analytic Tools include programs to compare an OCR-generated file with the correct file and obtain measures of performance.

A tilde (~) in an OCR-generated text file is treated as a reject character. A circumflex (^) is interpreted as a suspect marker and serves to mark the following character as suspect. For example, in Ne^vada, the v is marked as suspect. The value of these special characters is assessed when computing marked character efficiency.

Each tilde ( $\sim$ ) in a correct text file is treated as a wildcard and allows zero or one arbitrary character to be generated for it without an error being charged. For example, suppose the page contains a character that a page reader is not expected to identify, such as a degree symbol ( $^{\circ}$ ) or Greek letter ( $\delta$ ). Then the character should be represented in the correct file by a tilde so that a page reader may generate any single character for it, or no character at all, without penalty. For more information on the use of wildcards in ground truth, see reference [12].

Extraneous spacing characters are ignored in text files and have no effect on the computation of performance measures. Specifically, blank lines are disregarded, as well as leading and trailing blanks on a line. Multiple consecutive blanks within a line are treated as a single blank. The "newline" character at the end of each line is not ignored and appears as "<\n>" within the accuracy reports. Other whitespace characters, such as tabs and formfeeds, are treated as blanks. Page-reading systems and human ground-truth preparers may freely utilize spacing characters to format their text.

#### 1.3 Zoning

We assume that an ordered set of zones has been defined for a page image, and that the coordinates of these zones are communicated to the page-reading system when it processes the image. The resulting OCR-generated text file contains the recognized text for zone 1, followed by the text for zone 2, and so on. The correct text file for this page must match this sequence, with the ground truth for zone 1 appearing first, followed by the ground truth for zone 2, and so on. Character and word accuracy are computed from these files as described in Sections 2 and 3.

### 2 Character Accuracy

#### 2.1 The **accuracy** Program

accuracy correctfile generatedfile [ accuracy\_report ]

The **accuracy** program compares the correct text found in *correctfile* with the OCR-generated text found in *generatedfile*. A character accuracy report is written to *accuracy\_report* if specified; otherwise, it is written to "stdout," the standard error output file.

We will illustrate this program with an example. The following is a small page image containing two columns of text. Assume that zone 1 contains the left-hand column and zone 2 contains the right-hand column.

crushed under vacuum in stainless steel tubes. Liberated water was extracted at 200°C and converted, using uranium, into hydrogen for D/H analyses. The deuterium content is expressed in parts per thousand difference (per mil) relative to standard mean ocean water (SMOW) [normalized to the V-SMOW/SLAP scale (7)]. The  $\delta D$  values are plotted against age in Fig. 2.

We cannot attribute the changes in deuterium to water-mineral exchange because the water-bearing fractures in the regional carbonate aquifer, feeding the modern (and fossil) flow system, are typically coated with calcite or dolomite (8). This coating precludes the exchange of hydrogen between water and clay minerals during flow from recharge to discharge areas. In fact, the difference in

The following is the correct text for this page. Notice the two wildcards.

crushed under vacuum in stainless steel tubes. Liberated water was extracted at 200~C and converted, using uranium, into hydrogen for D/H analyses. The deuterium content is expressed in parts per thousand difference (per mil) relative to standard mean ocean water (SMOW) [normalized to the V-SMOW/SLAP scale (7)]. The ~D values are plotted against age in Fig. 2. We cannot attribute the changes in deuterium to water-mineral exchange because the water-bearing fractures in the regional carbonate aquifer, feeding the modern (and fossil) flow system, are typically coated with calcite or dolomite (8). This coating precludes the exchange of hydrogen between water and clay minerals during flow from recharge to discharge areas. In fact, the difference in Here is OCR-generated text for this page. Notice the reject characters and suspect markers.

```
crushed under vacuum in stainless steel
tubes. Liberated water was extracted at
200"C and converted. using uranium,
into hydrogen for D/H analyses. The
deuterium content is expressed in parts
per thousand difference (per mil) relative
to standard mean ocean water (SMOW)
[normalized to the V-SMOWISLAP
scale (7)1. The 6D values are plotted
against age in Fig. 3.
We cannot attribute the changes in dcu-
terium to water-mineral exchange be-
cause the water-bearing ir .acturss in the
regional carbonaie aquif^er, feeding the
modern (and fo~sil) now ~vstem. are
typically coated ~-^.ith calci~s or dolomite
(6). This coating pr-ecludes the exchanjie
of hyd^l-ogen bet~^.etrn water and clay
minerals during now from 'I'.echarge to
discharge areas. in f,rct, the di~ference in
```

On the next page is the character accuracy report produced by the **accuracy** program when given these correct and generated files as inputs. A character accuracy report consists of six sections. The first section indicates the number of characters in the ground truth (756), the number of errors made by the page reader (39), and the character accuracy (94.84%). The second section gives the number of reject characters (6), suspect markers (7), and false marks (1), followed by information relating to marked character efficiency: if a user examines the marked characters (1.72% of the text) and corrects the marked errors, the character accuracy will increase (to 96.96%).

Errors are actually edit operations (character insertions, substitutions, and deletions) that are needed to correct the OCR-generated text. The third section of the character accuracy report gives a breakdown of marked errors, unmarked errors, and total errors by edit operation.

The fourth section shows the accuracy by character class. Here the ground-truth characters are divided into classes and the percentage of characters recognized in each class is reported. The total number of ground-truth characters missed (29) is always equal to the number of insertions plus substitutions.

The fifth section lists the "confusions" sorted by the number of errors charged for each. In this example, four errors were charged because n was generated for fl. Since this confusion requires only two edit operations to correct (one insertion and one substitution), then this confusion must have occurred twice to account for a total of four errors.

The sixth and last section of a character accuracy report provides a complete enumeration of the ground-truth characters.

```
_____
   756 Characters
    39 Errors
  94.84% Accuracy
     6 Reject Characters
     7 Suspect Markers
     1 False Marks
   1.72% Characters Marked
  96.96% Accuracy After Correction
   Ins
        Subst
                 Del Errors
                 6 16 Marked
     0
         10
                   4
                           23
     2
            17
                              Unmarked
     2
            27
                   10
                          39
                              Total
  Count Missed %Right
         0 100.00 ASCII Spacing Characters
   117
            4 87.10 ASCII Special Symbols
    31
    6
            2 66.67 ASCII Digits
    24
            1 95.83 ASCII Uppercase Letters
    578
          22 96.19 ASCII Lowercase Letters
    756
          29 96.16 Total
       Marked Correct-Generated
0 {f1}-{n}
3 {w}-{~-.}
 Errors
     4
     3
            2 {r}-{I.}
     2
     2
            2 \{r\} - \{1-\}
     2
            2 {sy}-{~v}
     2
            2 {te}-{~s}
     2
            2 {w}-{~.}
            0 {,}-{.}
0 {a}-{,r}
0 {e}-{c}
     2
     2
     2
           0 (e)-{tr}
     2
           0 	 {g}-{ji}
     2
            1 {f}-{~}
     1
     1 1 1
            1 {s}-{~}
            1 {}-{.}
            0 {/}-{I}
            0 {2}-{3}
     1
            0
               {8}-{6}
{I}-{i}
     1
     1
            0
            0 {1}-{1}
     1
            0 {e}-{s}
     1
            0 {f}-{i}
     1
     1
            0 \{t\} - \{i\}
     1
            0 {}-{-}
  Count Missed %Right
               100.00
                       {<\n>}
        0
    20
    97
            0
                100.00
                       { }
            0 100.00
                       {(}
     5
            0 100.00 {)}
     5
            2 60.00 {,}
     5
            0 100.00 {-}
     7
            0 100.00 {.}
     2
            1 50.00 {/}
            0 100.00 {0}
     2
```

UNLV-ISRI OCR Accuracy Report Version 5.1

```
1
               50.00
                        {2}
              100.00
                        {7}
                0.00
1
          1
                        {8}
              100.00
1
          Ω
                        {A}
                        {C}
          0
              100.00
1
2
              100.00
          0
                        {D}
1
          0
              100.00
                        {F}
1
          0
              100.00
                        {H}
               0.00
1
          1
                        {I}
2
          0
              100.00
                        {L}
2
              100.00
          0
                        { M }
2
          0
              100.00
                        {0}
1
          Ω
              100.00
                        {P}
              100.00
3
          0
                        \{S\}
              100.00
3
          0
                        \{T\}
1
          0
              100.00
                        {V}
3
          0
              100.00
                        { W }
              100.00
1
          0
                        {[]}
1
          1
               0.00
                        { ] }
               98.21
56
          1
                        {a}
7
              100.00
                        {b}
26
          0
              100.00
                        {c}
27
          0
              100.00
                        {d}
88
          5
               94.32
                        {e}
14
               71.43
                        {f}
16
          1
               93.75
                        {g}
              100.00
20
          0
                        {h}
              100.00
37
          Ω
                        {i}
21
          2
              90.48
                        {1}
13
          0
             100.00
                        {m}
44
          0
              100.00
              100.00
28
          0
                        {o}
7
          0
              100.00
                        {p}
1
              100.00
          0
                        \{q\}
45
          2
               95.56
                        {r}
               93.55
31
          2
                        \{s\}
51
          2
               96.08
                        {t}
             100.00
20
          0
                        {u}
              100.00
          0
                        {v}
10
          2
               80.00
                        {w}
4
          0
              100.00
                       \{x\}
7
              85.71
          1
                        {y}
              100.00
                        \{z\}
```

### 2.2 The **synctext** Program

```
synctext [-H][-i][-s][-T] textfile1 textfile2 ... > resultfile
```

The **synctext** program can be used to show the alignment of two or more text files. If exactly two input files are specified (one correct file and one OCR-generated file), then the alignment that is computed is the same as the one computed by the **accuracy** program. This allows the user to see where errors occurred within the text. The algorithm used to compute this alignment is described in reference [11].

If more than two input files are specified, or if the "-H" option is given, then the algorithm described in reference [13] is utilized to perform the alignment. This algorithm is also used by the **vote** program to align multiple text files (see Section 2.8).

The "-T" option selects yet another alignment algorithm which can find transposed matches between two input files (see reference [10]). This algorithm is used by the **editop** program for automatic zoning evaluation (see Section 4.1). The "-i" option specifies that the alignment is to be performed on a case insensitive basis (the default is case sensitive), and the "-s" option displays suspect markers in the output.

The following is the output of the **synctext** program when given the example correct and OCR-generated files as inputs. The characters upon which the input files agree are shown first, and everywhere there is a difference, a number in braces identifies a footnote below that indicates what the difference is.

crushed under vacuum in stainless steel tubes. Liberated water was extracted at 200{1}C and converted{2} using uranium, into hydrogen for D/H analyses. The deuterium content is expressed in parts per thousand difference (per mil) relative to standard mean ocean water (SMOW) [normalized to the V-SMOW{3}SLAP scale  $(7)\{4\}$ . The  $\{5\}D$  values are plotted against age in Fig. {6}. We cannot attribute the changes in  $d\{7\}u$ terium to water-mineral exchange because the water-bearing  $\{8\}r\{9\}$ actur $\{10\}$ s in the regional carbona [11] e aquifer, feeding the modern (and fo $\{12\}$ sil)  $\{13\}$ ow  $\{14\}$ stem $\{15\}$  are typically coated {16}ith calci{17} or dolomite ( $\{18\}$ ). This coating  $pr\{19\}$ eclud $\{20\}$ s the exchan $\{21\}$ e of hyd{22}ogen bet{23}e{24}n water and clay minerals during {25}ow from {26}echarge to discharge areas. {27}n f{28}ct, the di{29}ference in {1} Correct {~} Generated {"} \_\_\_\_\_\_ {2} Generated {.} \_\_\_\_\_\_ {3} Generated {I} \_\_\_\_\_\_ {4} Correct {1} Generated {1} \_\_\_\_\_\_ Correct {~}

Generated {6}

{6} Correct Generated	{2} {3}
{7}	
{8}	{f}
{9} Correct Generated	
{10} Correct Generated	
{11} Correct Generated	{t}
{12} Correct Generated	{~}
Generated	<pre>{fl} {n}</pre>
{14} Correct Generated	{sy}
{15} Correct Generated	<b>{.}</b>
Generated	{w} {~}
Generated	{te} {~s}
{18} Correct Generated	
{19} Correct Generated	
Generated	{e} {c}

```
{21}
Correct {g}
Generated {ji}
______
{22}
Correct
   {r}
Generated {1-}
_____
{23}
Correct {w}
Generated {~.}
______
{24}
Correct
   {e}
Generated {tr}
Correct {fl}
Generated {n}
Correct {r}
Generated {I.}
______
Correct {I}
Generated {i}
______
Correct {a}
Generated {,r}
_____
Correct
Generated {~}
```

### 2.3 The **accsum** Program

```
accsum accuracy_report1 accuracy_report2 ... > accuracy_report
```

The **accsum** program combines two or more character accuracy reports to produce an aggregate report. While the results for a single page are interesting, tabulating the results for a set of pages yields important insights into the page-reading process.

On the next page is an aggregate character accuracy report that was produced by combining 175 individual reports using the **accsum** program. We see that there is a total of 293,493 characters on these 175 pages, and that the page reader made 5,916 errors for an overall character accuracy of 97.98%.

Given this large amount of data, some interesting observations can be made. For example, only 94.53% of the digits were recognized correctly versus 99.07% of the lowercase letters. Only 91.98% of the occurrences of the numeral 1 were correctly

#### identified.

The list of confusions is very long so it has been truncated here. The most common confusion for this set of pages, contributing 236 errors, is the generation of a space where there should be none, which causes a word to be split (e.g., Nev ada). There are 43 errors attributed to the opposite case, i.e., no space is generated where there should be one, which causes two words to be joined (e.g., LasVegas).

The numeral zero (0) was generated for the letter O a total of 92 times on these pages; it is difficult also for humans to distinguish these symbols. The confusion charged 39 errors was a single sequence of 39 characters that was omitted by the page reader. An ellipsis (...) in a confusion indicates that a long sequence of characters has been truncated.

```
UNLV-ISRI OCR Accuracy Report Version 5.1
_____
  293493 Characters
   5916 Errors
   97.98% Accuracy
    700 Reject Characters
      5 Suspect Markers
      5 False Marks
    0.24% Characters Marked
   98.37% Accuracy After Correction
                   Del Errors
          Subst
     Ins
    71 650 415 1136 Market
807 2699 1274 4780 Unmark
878 3349 1689 5916 Total
                                      Marked
                                      Unmarked
   Count Missed %Right
   47850 234 99.51 ASCII Spacing Characters
            603 95.78 ASCII Special Symbols
   14306
  14341 784 94.53 ASCII Digits
         768 95.90 ASCII Uppercase Letters
1838 99.07 ASCII Lowercase Letters
4227 98.56 Total
  18737
  198259
  293493
 Errors Marked Correct-Generated
     236
          0 {}-{}
              0 {0}-{0}
            0 {,}-{.}
      90
      61
             0 {}-{'}
      52
               0 {1}-{1}
      52
               0
                   {}-{.}
      43
                    { }-{}
                   { }-{}
{tion of<\n>temperature a...}-{}
      39
               0
               0 {.}-{,}
      38
              0 {,}-{}
      37
              0 {1}-{I}
      31
      25
              0 {-}-{}
      24
              0 {0}-{o}
      24
              0 {2}-{z}

\begin{array}{ccc}
 & \{z\} - \{z\} \\
0 & \{in\} - \{m\} \\
0 & \vdots
\end{array}

      2.4
              0 {1}-{1}
21 {}-{~}
      24
      21 21 \{\}-\{\sim\}
21 0 \{1\}-\{I\}
```

20 20	0	{96}-{%} {o}-{a}	
20	U	(U)-(a)	
Count	Missed	%Right	
6064	23	99.62	$\{<\n>\}$
41786	211	99.50	{ }
159	3	98.11	{ " }
13	6	53.85	{#}
3 56	0 13	100.00 76.79	{\$}
6	1	83.33	{ % } { & }
70	3	95.71	{ , }
883	46	94.79	{(}
881	44	95.01	{ ) }
75	6	92.00	{ * }
54	13	75.93	{+}
3192	150	95.30	{,}
2353 5620	113 88	95.20 98.43	{-} {.}
224	17	92.41	{/}
2513	110	95.62	{0}
2806	225	91.98	{1}
1797	122	93.21	{2}
1328	41	96.91	{3}
1109	57	94.86	{4}
1131	54	95.23	{5}
1021 858	65 21	93.63	{6}
812	31 35	96.39 95.69	{7} {8}
966	44	95.45	{9}
234	4	98.29	{:}
157	4	97.45	<pre>; ;</pre>
25	0	100.00	{<}
170	38	77.65	{ = }
6	2	66.67	{>}
1067	2	93.33	{?}
1267 512	43 22	96.61 95.70	{A} {B}
1393	31	97.77	{C}
686	29	95.77	{D}
1494	39	97.39	(E)
668	14	97.90	$\{F\}$
500	28	94.40	{G}
542	19	96.49	{H}
1097 198	38	96.54 98.48	{I} {J}
177	7	96.05	{K}
710	13	98.17	{L}
774	55	92.89	(M)
1135	45	96.04	{N}
902	147	83.70	{0}
811	38	95.31	{P}
1107	7	91.46	{Q}
1187 1546	41 41	96.55 97.35	{R} {S}
1656	31	98.13	(T)
428	38	91.12	{U}
279	11	96.06	{V}
452	20	95.58	{ W }
67	2	97.01	{x}
127	2	98.43	{Y}

```
47
          4 91.49
                      { Z }
         10 69.70
  33
                      {[]}
         12 63.64
  33
                      { ] }
         2
  3
               33.33
                      {`}
             99.35
17028
         111
                      {a}
               98.10
2638
         50
                      {b}
7269
         75
               98.97
                      {c}
         32
7500
               99.57
                      {d}
         135
              99.47
25274
                      {e}
4819
         93 98.07
                      {f}
        50 98.54
3434
                      {a}
             99.22
7603
         59
15387
         167
               98.91
                      {i}
        12
               91.11
 135
                      {j}
 851
         16
               98.12
                      {k}
9304
         156
               98.32
                      {1}
               97.99
4966
         100
                      { m }
               99.18
14430
         118
                      {n}
15151
         107
             99.29
                      {0}
         40 99.08
4371
                      {q}
         6 98.35
                      \{q\}
13574
         107
             99.21
         96
             99.25
12743
                      \{s\}
17702
         103
               99.42
                      {t}
5627
         58
               98.97
                      {u}
2023
         22
               98.91
                      {v}
         39
               98.24
2216
                      { w }
         30
             94.92
 590
                      \{x\}
2898
         49 98.31
                      {y}
             98.07
                      \{z\}
  3
         3
             0.00 {{}
          20
  20
               0.00
                     { | }
                0.00
   3
          3
                      {}}
```

### 2.4 The **groupacc** Program

```
groupacc groupfile accuracy_report [ groupacc_report ]
```

In accuracy by character class, ground-truth characters are grouped according to predefined classes, such as the ASCII digits and the ASCII lowercase letters. The **groupacc** program allows user-defined groupings. This program summarizes the accuracy data from *accuracy\_report* for the characters specified in *groupfile*. The results are written to *groupacc\_report* if specified; otherwise, they are written to stdout.

For example, suppose we wish to focus on lowercase letters with descenders. We would create a group file containing the desired letters: gjpqy. Given this group file and the aggregate character accuracy report shown above, the **groupacc** program produces the following display:

```
Count Missed %Right
3434 50 98.54 {g}
135 12 91.11 {j}
4371 40 99.08 {p}
```

```
364 6 98.35 {q}
2898 49 98.31 {y}
11202 157 98.60 Total
```

#### 2.5 The **accci** Program

```
accci accuracy_report1 accuracy_report2 ... > resultfile
```

Given a set of character accuracy reports as input, the **accci** program computes an approximate 95% confidence interval for character accuracy using the method of jackknife estimation (which is described in reference [14]). Each input report is treated as one observation. For best results, at least 30 observations are needed.

The following is the output from the **accci** program when given 175 character accuracy reports as input:

```
175 Observations
293493 Characters
5916 Errors
97.98% Accuracy
97.75%, 98.22% Approximate 95% Confidence Interval for Accuracy
```

#### 2.6 The **accdist** Program

```
accdist accuracy_report1 accuracy_report2 ... > xyfile
```

The **accdist** program computes the distribution of the character accuracies found in a set of character accuracy reports. The results are written in "xy format" to stdout. In this format, each line contains the x- and y-coordinates of a data point separated by spaces. A file in this format can be easily imported into most graphing programs.

The **accdist** program produces one data point for each value of x from 0 to 100. The y value is the percentage of characters recognized with at least x% character accuracy. Below is a portion of the output for the 175-page sample. It shows, for example, that 98.25% of the sample was recognized with at least 95% character accuracy. It is widely believed that it costs more to correct OCR-generated text that is less than 95% accurate than it costs to type the entire text from scratch. Thus, pages recognized with less than 95% accuracy should be sent to a manual data-entry operation. In this example, pages containing 1.75% of the characters should be entered manually.

```
0 100.00
1 100.00
...
94 98.78
```

```
95 98.25
96 95.09
97 88.51
98 63.68
99 16.99
100 0.53
```

#### 2.7 The **ngram** Program

```
ngram [-n \ 1|2|3] textfile1 textfile2 ... > resultfile
```

The **ngram** program computes n-gram statistics for one or more text files. If the "-n" option is omitted, or if "-n 1" is specified, this program displays the frequency of each character found in the input files ("unigrams"). If "-n 2" is chosen, the frequency of every distinct character pair is shown ("bigrams"), and if "-n 3" is selected, the frequency of each unique triple of characters (i.e., three consecutive characters) is displayed ("trigrams").

The *n*-gram statistics are displayed twice: first in the collating order of the characters, and then in order by decreasing frequency. Below is abbreviated unigram output for 175 correct text files. The blank character occurs most often, with 41,786 occurrences found in the input files. The letter e is the next most common, with 25,274 occurrences. In this example, none of the occurrences have been marked as suspect.

```
Count Suspect
 6064
           0
                \{<\n>\}
              { }
41786
            0
           0 {"}
  159
           0 {#}
   13
   3
            0 {$}
   56
            0 {%}
    6
            0
              { & }
  590
                \{x\}
 2898
            0
                {y}
  362
            0
                \{z\}
   3
            Ω
               { { }
   20
           0
              {|}
   3
            0 {}}
 5453
           0
              {~}
298946
               Total
Count Suspect
                { }
41786
25274
            0
               {e}
           0 (t)
17702
           0 {a}
17028
15387
          0 {i}
          0 {o}
15151
14430
           0 {n}
```

14

### 3 Word Accuracy

#### 3.1 The **wordacc** Program

**wordacc** [-S stopwordfile] correctfile generatedfile [wordacc\_report]

The **wordacc** program compares the correct text found in *correctfile* with the OCR-generated text found in *generatedfile*. A word accuracy report is written to *wordacc\_report* if specified; otherwise, it is written to stdout.

Only the words found in *stopwordfile* are considered to be stopwords. (If the "-S" option is omitted, the default set of 110 stopwords from the BASISplus text retrieval system is used. See reference [15] for information on this system.) If stopwords are placed in the file in order by decreasing frequency of usage (as determined from some large corpus), then the file can also be used with the **nonstopacc** program (see Section 3.3). Here is an example of such a file containing 200 English stopwords in decreasing order of frequency:

the of and to a in that is was he for it with as his on be at by i this had not are but from or have an they which one you were her all she there would their we him been has when who will more no if out so said what up its about into than them can only other new some could these two may then do first any my now such like our over man me even most made after also did many before must through back years where much your way well down should because each just those mr how too state good very make still see men work long get here between both being under never same another know while last might us great old year off come since against go came right used take three states himself few use during without again place around however small mrs thought went say part once general high upon every does got number until always away something fact though less put think almost enough far took yet better nothing end why find going asked later knew point next give group toward young let room side given

On the next page is the word accuracy report produced by the **wordacc** program when given this stopword file and the example correct and generated files. A word accuracy report consists of seven sections. The first section indicates the number of words in the ground truth (119), the number of words misrecognized by the page reader (18), and the word accuracy (84.87%). The second and third sections show the stopword accuracy (92.86%) and the non-stopword accuracy (80.52%), respectively, with a breakdown by word length for each.

The fourth section presents the results for distinct non-stopword accuracy, which differs from non-stopword accuracy. See reference [4] for a description of this performance measure.

The fifth section gives the phrase accuracy for phrases of lengths 1 through 8. The sixth and seventh sections provide a complete enumeration of the stopwords and non-stopwords, respectively. Word accuracy is determined on a case-insensitive basis, so the

## words are displayed here in lowercase only.

UNLV-ISRI		Accuracy	Report Versio	n 5.1
119	Words			
18	Misreco	gnized		
84.87%	Accurac	У		
Stopwords				
Count	Missed	%Right	Length	
17	0	100.00	2	
16	0	100.00	3	
5	2	60.00	4	
1	0	100.00	5	
1	0	100.00	6	
2	1	50.00	7	
42	3	92.86	Total	
Non-stopw	ords			
Count	Missed	%Right	Length	
5	0	100.00	1	
6	1	83.33	3	
7	4	42.86	4	
13	0	100.00	5	
8	2	75.00	6	
11	1	90.91	7	
12	3	75.00	8	
12	3	75.00	9	
3	1	66.67	10	
77	15	80.52	Total	
Distinct 1	Non-stopw	ords		
Count	Missed		Occurs	
58	9	84.48	1	
7	1	85.71	2	
1	0	100.00	5	
66	10	84.85	Total	
Phrases				
Count	Missed	%Right	Length	
119	18	84.87	1	
118	31	73.73	2	
117	39	66.67	3	
116	47	59.48	4	
115	53	53.91	5	
114	57	50.00	6	
113	59	47.79	7	
112	61	45.54	8	
C+				
Stopwords	364	0 D d - d- L		
Count	Missed	%Right		
1	0	100.00	against	
3	0	100.00	and	
2	0	100.00	are	
1	0	100.00	at	
1	0	100.00	be	
1	1	0.00	between	
1	0	100.00	during	
1	1	0.00	fact	
1	0	100.00	for	
1	0	100.00	from	
7	0	100.00	in	

1	U	100.00	ınto
1	0	100.00	is
1	0	100.00	of
1	0	100.00	or
9	0	100.00	the
1	0	100.00	this
4	0	100.00	to
1	0	100.00	under
1	0	100.00	was
1	0	100.00	we
1	1	0.00	with

Non-stopwo			
Count	Missed	%Right	
1	0	100.00	age
1	0	100.00	analyses
1	0	100.00	aquifer
1	0	100.00	areas
1	0	100.00	attribute
1	0	100.00	bearing
1	0	100.00	С
1	1	0.00	calcite
1	0	100.00	cannot
1	1	0.00	carbonate
1	0	100.00	cause
1	0	100.00	changes
1	0	100.00	clay
1	0	100.00	coated
1	0	100.00	coating
1	0	100.00	content
1	0	100.00	converted
1	0	100.00	crushed
2	0	100.00	d
1	1	0.00	deu
1	0	100.00	deuterium
2	1	50.00	difference
1	0	100.00	discharge
1	0	100.00	dolomite
2	1	50.00	exchange
1	0	100.00	expressed
1	0	100.00	extracted
1	0	100.00	feeding
1	0	100.00	fig
2	2	0.00	flow
1	1	0.00	fossil
1	1	0.00	fractures
1	0		h
2	1	100.00 50.00	n hydrogen
	0		
1 1	0	100.00	liberated
			mean
1 1	0	100.00	mil mineral
	0	100.00	
1	0	100.00	minerals
1	0	100.00	modern
1	0	100.00	normalized
1	0	100.00	ocean
1	0	100.00	parts
2	0	100.00	per
1	0	100.00	plotted
1	1	0.00	precludes
1	1	0.00	recharge

```
1 0 100.00 regional
           0 100.00 relative
1
           0 100.00 scale
1
           1 0.00 slap
2 1 1
           1
                   50.00 smow
           0 100.00 stainless
0 100.00 standard
1
           0 100.00 steel
1
                   0.00 system
1 1 0.00 system
1 0 100.00 terium
1 0 100.00 thousand
1 0 100.00 tubes
1 0 100.00 typically
1 0 100.00 uranium
1 0 100.00 using
1 0 100.00 vacuum
1 0 100.00 vacuum
1 0 100.00 values
5 0 100.00 water
1
          1
          0 100.00 water
```

#### 3.2 The **wordaccsum** Program

```
wordacc_report1 wordacc_report2 ... > wordacc_report
```

The **wordaccsum** program combines two or more word accuracy reports to produce an aggregate report. Below is the aggregate report that was produced by combining 175 individual reports. Of 43,928 words on these 175 pages, 2,211 were misrecognized for an overall word accuracy of 94.97%. It is usual to see a higher overall stopword accuracy (97.80%) than non-stopword accuracy (93.22%). The lists of stopwords and non-stopwords are very long and have been truncated here.

```
UNLV-ISRI OCR Word Accuracy Report Version 5.1
  43928 Words
   2211 Misrecognized
  94.97% Accuracy
Stopwords
  Count Missed %Right Length
         46 94.32
    810
                          1
   6069
           128 97.89
           92 98.48
61 97.50
   6065
   2438
            27 97.15
    949
    204
            10 95.10
            5
                   97.84
    232
                                7
              0 100.00
     1
                                9
           369
                           Total
   16768
                   97.80
Non-stopwords
  Count Missed %Right Length

      2275
      250
      89.01
      1

      1050
      249
      76.29
      2
```

1425	169	88.14	3	
3025	160	94.71	4	
3705 3327	186 172	94.98 94.83	5 6	
3216	139	95.68	7	
3132	166	94.70	8	
2345	125	94.67	9	
1560	96	93.85	10	
1032	57	94.48	11	
599	38	93.66	12	
260	17	93.46	13	
118	5	95.76	14	
58	11	81.03	15	
23	2	91.30	16	
3	0	100.00	17	
5	0	100.00	18	
2	0	100.00	19	
27160	1842	93.22	Total	
Distinct	Non-stopw	ords		
Count	Missed	%Right	0ccurs	
11889	808	93.20	1	
2437	54	97.78	2	
868	15	98.27	3	
431	7	98.38	4	
241	1	99.59	5	
155 107	4	97.42 96.26	6 7	
69	1	98.55	8	
62	1	98.39	9	
35	0	100.00	10	
112	1	99.11	>10	
16406	896	94.54	Total	
Phrases				
Count	Missed	%Right	Length	
43928	2211	94.97	1	
43753	3927	91.02	2	
43578	5446	87.50	3	
43403	6812	84.31	4	
43228	8082	81.30	5	
43053	9249	78.52	6	
42878	10334	75.90	7	
42704	11351	73.42	8	
Stopwords				
Count	Missed	%Right		
723	32	95.57	a	
69	1	98.55	about	
32	1	96.88	after	
9 10	1	88.89 100.00	again against	
10	U	100.00	against	
Non-stopwords				
Count	Missed	%Right	,	
2	2	0.00	ab	
3	0	100.00	abbreviations	
3 6	0	100.00	abernethy	
6	0	100.00	ability able	
0	0	100.00	anic	

. . .

#### 3.3 The **nonstopacc** Program

```
nonstopacc stopwordfile wordacc_report > xyfile
```

Given an ordered list of N stopwords in *stopwordfile* and a word accuracy report in *wordacc\_report*, the **nonstopacc** program writes the data for the non-stopword accuracy curve to stdout. This curve presents non-stopword accuracy as a function of the number of stopwords. Examples of these curves can be found in reference [5]. The data points are written in xy format for values of x ranging from 0 to x. Each y value is the non-stopword accuracy when using the x most frequently-occurring stopwords from the stopword file. Here is a portion of the output for the example stopword file and aggregate word accuracy report:

```
0 94.97

1 94.69

2 94.50

3 94.35

4 94.26

5 94.23

6 94.13

7 94.09

8 94.02

9 94.00

...

200 93.22
```

### 3.4 The wordaccci Program

```
wordacci wordacc_report1 wordacc_report2 ... > resultfile
```

Given a set of word accuracy reports as input, the **wordaccci** program computes an approximate 95% confidence interval for word accuracy using the method of jackknife estimation [14]. This program is analogous to the **accci** program for character accuracy (see Section 2.5). The following is the output from the **wordaccci** program when given 175 word accuracy reports as input:

```
175 Observations
43928 Words
2211 Misrecognized
94.97% Accuracy
94.40%, 95.54% Approximate 95% Confidence Interval for Accuracy
```

#### 3.5 The wordaccdist Program

```
wordaccdist wordacc_report1 wordacc_report2 ... > xyfile
```

The **wordaccdist** program computes the distribution of the word accuracies found in a set of word accuracy reports. The results are written in *xy* format to stdout. This program is analogous to the **accdist** program for character accuracy (see Section 2.6).

The **wordaccdist** program produces one data point for each value of x from 0 to 100. The y value is the percentage of words recognized with at least x% word accuracy.

#### 3.6 The **wordfreq** Program

```
wordfreq textfile1 textfile2 ... > resultfile
```

The **wordfreq** program determines the frequency of words in one or more text files. The frequency data are displayed twice: first in the collating order of the words, and then in order by decreasing frequency. Below is abbreviated **wordfreq** output for 175 correct text files. The word "the" occurs most often, with 2,981 occurrences found in the input files.

```
Count
 723
       а
   2
      ab
   3
      abbreviations
      abernethy
      ability
  28
      zone
   9
      zones
      zonplt
   2
   1
      zro
   3
       zubovic
43928
       Total
Count.
2981
      the
1907 of
1366 and
1048
      in
 723
      zirconates
      zirconium
      zirconolite
   1
   1 zonal
   1 zro
43928 Total
```

### References

- [1] S. V. Rice, J. Kanai, and T. A. Nartker. A report on the accuracy of OCR devices. Technical Report 92-02, Information Science Research Institute, University of Nevada, Las Vegas, March 1992.
- [2] S. V. Rice, J. Kanai, and T. A. Nartker. An evaluation of OCR accuracy. Technical Report 93-01, Information Science Research Institute, University of Nevada, Las Vegas, April 1993.
- [3] S. V. Rice, J. Kanai, and T. A. Nartker. The third annual test of OCR accuracy. Technical Report 94-03, Information Science Research Institute, University of Nevada, Las Vegas, April 1994.
- [4] S. V. Rice, F. R. Jenkins, and T. A. Nartker. The fourth annual test of OCR accuracy. Technical Report 95-04, Information Science Research Institute, University of Nevada, Las Vegas, April 1995.
- [5] S. V. Rice, F. R. Jenkins, and T. A. Nartker. The fifth annual test of OCR accuracy. Technical Report 96-01, Information Science Research Institute, University of Nevada, Las Vegas, April 1996.
- [6] T. A. Nartker, S. V. Rice, and J. Kanai. OCR accuracy: UNLV's second annual test. *Inform*, Association for Information and Image Management, 8(1):40+, January 1994.
- [7] T. A. Nartker and S. V. Rice. OCR accuracy: UNLV's third annual test. *Inform*, Association for Information and Image Management, 8(8):30+, September 1994.
- [8] T. A. Nartker, S. V. Rice, and F. R. Jenkins. OCR accuracy: UNLV's fourth annual test. *Inform*, Association for Information and Image Management, 9(7):38+, July 1995.
- [9] T. A. Nartker, S. V. Rice, and F. R. Jenkins. OCR accuracy: UNLV's fifth annual test. *Inform*, Association for Information and Image Management, to appear, 1996.
- [10] J. Kanai, S. V. Rice, T. A. Nartker, and G. Nagy. Automated evaluation of OCR zoning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 17(1):86-90, 1995.
- [11] S. V. Rice. *Measuring the Accuracy of Page-Reading Systems*. Ph.D. dissertation,

- University of Nevada, Las Vegas, 1996.
- [12] S. V. Rice, J. Kanai, and T. A. Nartker. Preparing OCR test data. Technical Report 93-08, Information Science Research Institute, University of Nevada, Las Vegas, June 1993.
- [13] S. V. Rice, J. Kanai, and T. A. Nartker. An algorithm for matching OCR-generated text strings. *International Journal of Pattern Recognition and Artificial Intelligence*, 8(5):1259-1268, 1994.
- [14] E. J. Dudewicz and S. N. Mishra. *Modern Mathematical Statistics*, pages 743-748. John Wiley & Sons, 1988.
- [15] Information Dimensions, Inc., Dublin, Ohio. *BASISplus Database Administration Reference, Release L*, June 1990.