# Assignment 2: Text Processing

## Introduction

Text is a very important part of the data that computers deal with. In this assignment you are going to learn two things: how to process text (the first step in understanding language) and how to write code in small, bite-size portions: functions!

## Functionality

Once again you will write your programs in increasingly more powerful pieces. But unlike the first assignment you will only produce two main programs and a number of functions. The number of functions that you write correctly will determine your grade.

Remember the **wcount.c** program in the textbook and in class. You will find a related program in the file **textRW.c**. In the original program, the text was read in a character at a time. The program **textRW.c** reads in text, one *line* at a time. The program assumes that no line contains more than 500 characters. Your code will not be tested on text that violates this constraint.

To test out **textRW** do the following:
```
$ make textRW
$ cat verneTest.txt | ./textRW
```

Produce the following functions that will work on a single line of text:

```
int chop ( char *line )
```
- Remove the '\n' from the end of the line.
- Returns 0 if successful and -1 on failure.

```
int convertLowerCase ( char *line )
```
- Convert all alphabetic characters to lower case (do not use the `tolower()` function in the string library for C or your grade will be zero for this function).
- Returns the number of characters converted from upper case to lower case.

```
int replaceDigits ( char *line )
```
- Replace the all numbers/digits (0-9) with blanks.
- Returns the number of digits replaced with blanks.

```
int replacePunc ( char *line )
```
- Replace all punctuation and symbols with blanks. A list of all punctuation and symbols appears in the **Information** section of this assignment.
- Returns the number of punctuation and symbols replaced with blanks.

```
int reduceSpace ( char *line )
```
- Reduce all white space (blank characters and tabs to a single blank space).
- Returns the number of white spaces removed.

```
int trim ( char *line )
```
- Trim all white space from the beginning and from the end of each line.
- Returns 0 on success and -1 on failure.

Now compile the program **findWords.c** with your functions using **Makefile** and run this against the test file called **verneTest.txt**.

```
$ make findWords
$ cat verneTest.txt | ./findWords
```

Your output should look like the output in file **resultsVerne.txt**.

The **findWords** program has found all the words in the file that are greater than **5 characters** in length.

Now run the following pipeline:

```
$ cat verneTest.txt | ./findWords | sort
```

The output should look like the output in file **sortedResultsVerne.txt**.

Your last task is to write a program called **wordBag.c** that prints out all the unique words produced by the previous pipeline with their counts (number of times that they occur). The output will look like the output in file **countsResultsVerne.txt**.

Now run the following pipelines:

```
$ cat verneTest.txt | ./findWords | sort | ./wordBag
$ cat verneTest.txt | ./findWords | sort | ./wordBag | sort -
gk1,1r -gk2,2
```

Now you have a set of programs and functions that show you how popular specific words are in a text.

**Information** ASCII Table – all punctuation and symbols are indicated in the blue boxes.

http://www.asciitable.com/

| Dec | Hx | Oct | Char |  | Dec | Hx | Oct | Html | Chr |  | Dec | Hx | Oct | Html | Chr |  | Dec | Hx | Oct | Html | Chr |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 000 | NUL | (null) | 32 | 20 | 040 | &#32; | Space | | 64 | 40 | 100 | &#64; | @ | | 96 | 60 | 140 | &#96; | ` |
| 1 | 1 | 001 | SOH | (start of heading) | 33 | 21 | 041 | &#33; | ! | | 65 | 41 | 101 | &#65; | A | | 97 | 61 | 141 | &#97; | a |
| 2 | 2 | 002 | STX | (start of text) | 34 | 22 | 042 | &#34; | " | | 66 | 42 | 102 | &#66; | B | | 98 | 62 | 142 | &#98; | b |
| 3 | 3 | 003 | ETX | (end of text) | 35 | 23 | 043 | &#35; | # | | 67 | 43 | 103 | &#67; | C | | 99 | 63 | 143 | &#99; | c |
| 4 | 4 | 004 | EOT | (end of transmission) | 36 | 24 | 044 | &#36; | $ | | 68 | 44 | 104 | &#68; | D | | 100 | 64 | 144 | &#100; | d |
| 5 | 5 | 005 | ENQ | (enquiry) | 37 | 25 | 045 | &#37; | % | | 69 | 45 | 105 | &#69; | E | | 101 | 65 | 145 | &#101; | e |
| 6 | 6 | 006 | ACK | (acknowledge) | 38 | 26 | 046 | &#38; | & | | 70 | 46 | 106 | &#70; | F | | 102 | 66 | 146 | &#102; | f |
| 7 | 7 | 007 | BEL | (bell) | 39 | 27 | 047 | &#39; | ' | | 71 | 47 | 107 | &#71; | G | | 103 | 67 | 147 | &#103; | g |
| 8 | 8 | 010 | BS | (backspace) | 40 | 28 | 050 | &#40; | ( | | 72 | 48 | 110 | &#72; | H | | 104 | 68 | 150 | &#104; | h |
| 9 | 9 | 011 | TAB | (horizontal tab) | 41 | 29 | 051 | &#41; | ) | | 73 | 49 | 111 | &#73; | I | | 105 | 69 | 151 | &#105; | i |
| 10 | A | 012 | LF | (NL line feed, new line) | 42 | 2A | 052 | &#42; | * | | 74 | 4A | 112 | &#74; | J | | 106 | 6A | 152 | &#106; | j |
| 11 | B | 013 | VT | (vertical tab) | 43 | 2B | 053 | &#43; | + | | 75 | 4B | 113 | &#75; | K | | 107 | 6B | 153 | &#107; | k |
| 12 | C | 014 | FF | (NP form feed, new page) | 44 | 2C | 054 | &#44; | , | | 76 | 4C | 114 | &#76; | L | | 108 | 6C | 154 | &#108; | l |
| 13 | D | 015 | CR | (carriage return) | 45 | 2D | 055 | &#45; | - | | 77 | 4D | 115 | &#77; | M | | 109 | 6D | 155 | &#109; | m |
| 14 | E | 016 | SO | (shift out) | 46 | 2E | 056 | &#46; | . | | 78 | 4E | 116 | &#78; | N | | 110 | 6E | 156 | &#110; | n |
| 15 | F | 017 | SI | (shift in) | 47 | 2F | 057 | &#47; | / | | 79 | 4F | 117 | &#79; | O | | 111 | 6F | 157 | &#111; | o |
| 16 | 10 | 020 | DLE | (data link escape) | 48 | 30 | 060 | &#48; | 0 | | 80 | 50 | 120 | &#80; | P | | 112 | 70 | 160 | &#112; | p |
| 17 | 11 | 021 | DC1 | (device control 1) | 49 | 31 | 061 | &#49; | 1 | | 81 | 51 | 121 | &#81; | Q | | 113 | 71 | 161 | &#113; | q |
| 18 | 12 | 022 | DC2 | (device control 2) | 50 | 32 | 062 | &#50; | 2 | | 82 | 52 | 122 | &#82; | R | | 114 | 72 | 162 | &#114; | r |
| 19 | 13 | 023 | DC3 | (device control 3) | 51 | 33 | 063 | &#51; | 3 | | 83 | 53 | 123 | &#83; | S | | 115 | 73 | 163 | &#115; | s |
| 20 | 14 | 024 | DC4 | (device control 4) | 52 | 34 | 064 | &#52; | 4 | | 84 | 54 | 124 | &#84; | T | | 116 | 74 | 164 | &#116; | t |
| 21 | 15 | 025 | NAK | (negative acknowledge) | 53 | 35 | 065 | &#53; | 5 | | 85 | 55 | 125 | &#85; | U | | 117 | 75 | 165 | &#117; | u |
| 22 | 16 | 026 | SYN | (synchronous idle) | 54 | 36 | 066 | &#54; | 6 | | 86 | 56 | 126 | &#86; | V | | 118 | 76 | 166 | &#118; | v |
| 23 | 17 | 027 | ETB | (end of trans. block) | 55 | 37 | 067 | &#55; | 7 | | 87 | 57 | 127 | &#87; | W | | 119 | 77 | 167 | &#119; | w |
| 24 | 18 | 030 | CAN | (cancel) | 56 | 38 | 070 | &#56; | 8 | | 88 | 58 | 130 | &#88; | X | | 120 | 78 | 170 | &#120; | x |
| 25 | 19 | 031 | EM | (end of medium) | 57 | 39 | 071 | &#57; | 9 | | 89 | 59 | 131 | &#89; | Y | | 121 | 79 | 171 | &#121; | y |
| 26 | 1A | 032 | SUB | (substitute) | 58 | 3A | 072 | &#58; | : | | 90 | 5A | 132 | &#90; | Z | | 122 | 7A | 172 | &#122; | z |
| 27 | 1B | 033 | ESC | (escape) | 59 | 3B | 073 | &#59; | ; | | 91 | 5B | 133 | &#91; | [ | | 123 | 7B | 173 | &#123; | { |
| 28 | 1C | 034 | FS | (file separator) | 60 | 3C | 074 | &#60; | < | | 92 | 5C | 134 | &#92; | \ | | 124 | 7C | 174 | &#124; | | |
| 29 | 1D | 035 | GS | (group separator) | 61 | 3D | 075 | &#61; | = | | 93 | 5D | 135 | &#93; | ] | | 125 | 7D | 175 | &#125; | } |
| 30 | 1E | 036 | RS | (record separator) | 62 | 3E | 076 | &#62; | > | | 94 | 5E | 136 | &#94; | ^ | | 126 | 7E | 176 | &#126; | ~ |
| 31 | 1F | 037 | US | (unit separator) | 63 | 3F | 077 | &#63; | ? | | 95 | 5F | 137 | &#95; | _ | | 127 | 7F | 177 | &#127; | DEL |

# Extended ASCII Codes

| 128 | Ç | 144 | É | 160 | á | 176 | ░ | 192 | └ | 208 | ╨ | 224 | α | 240 | ≡ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 129 | ü | 145 | æ | 161 | í | 177 | ▒ | 193 | ┴ | 209 | ╤ | 225 | ß | 241 | ± |
| 130 | é | 146 | Æ | 162 | ó | 178 | ▓ | 194 | ┬ | 210 | ╥ | 226 | Γ | 242 | ≥ |
| 131 | â | 147 | ô | 163 | ú | 179 | │ | 195 | ├ | 211 | ╙ | 227 | π | 243 | ≤ |
| 132 | ä | 148 | ö | 164 | ñ | 180 | ┤ | 196 | ─ | 212 | ╘ | 228 | Σ | 244 | ⌠ |
| 133 | à | 149 | ò | 165 | Ñ | 181 | ╡ | 197 | ┼ | 213 | ╒ | 229 | σ | 245 | ⌡ |
| 134 | å | 150 | û | 166 | ª | 182 | ╢ | 198 | ╞ | 214 | ╓ | 230 | μ | 246 | ÷ |
| 135 | ç | 151 | ù | 167 | º | 183 | ╖ | 199 | ╟ | 215 | ╫ | 231 | τ | 247 | ≈ |
| 136 | ê | 152 | ÿ | 168 | ¿ | 184 | ╕ | 200 | ╚ | 216 | ╪ | 232 | Φ | 248 | ° |
| 137 | ë | 153 | Ö | 169 | ⌐ | 185 | ╣ | 201 | ╔ | 217 | ┘ | 233 | Θ | 249 | ∙ |
| 138 | è | 154 | Ü | 170 | ¬ | 186 | ║ | 202 | ╩ | 218 | ┌ | 234 | Ω | 250 | · |
| 139 | ï | 155 | ¢ | 171 | ½ | 187 | ╗ | 203 | ╦ | 219 | █ | 235 | δ | 251 | √ |
| 140 | î | 156 | £ | 172 | ¼ | 188 | ╝ | 204 | ╠ | 220 | ▄ | 236 | ∞ | 252 | ⁿ |
| 141 | ì | 157 | ¥ | 173 | ¡ | 189 | ╜ | 205 | = | 221 | ▌ | 237 | φ | 253 | ² |
| 142 | Ä | 158 | ₧ | 174 | « | 190 | ╛ | 206 | ╬ | 222 | ▐ | 238 | ε | 254 | ■ |
| 143 | Å | 159 | ƒ | 175 | » | 191 | ┐ | 207 | ╧ | 223 | ▀ | 239 | ∩ | 255 |  |