

# Btsort - a sort program using a binary tree

## Introduction

Why write a sort program? Well if you need to do a linear search through UTF-8 data which includes ordinary ASCII data, the sort command available at the shell prompt by default sorts into dictionary order based on the current locale. That means that if we are using eg a form of English then the default ordering will be like 'Agnes' 'agriculture' ... 'Bronwyn' 'brown' and so on. Moreover words containing apostrophes are required to collate after similar words without. Now this ordering is perfectly fine for many purposes but when doing a linear search through data sorted in such an order the search will fail because the natural bitwise order of ASCII characters is A..Z then a..z. The situation gets far worse when you need to search through strings of Thai characters, or for that matter Burmese, Lao or Cambodian. In these languages some of the vowel symbols are placed before their attached consonants as well as over, under or after. Compound vowel sounds may be formed from symbols in combination also. It's not so weird; look at 'hat' and 'hate'. In the second word the vowel sound is formed by symbols wrapped around the terminal consonant. In the SE Asian languages the wrapping happens around the initial consonant never the terminal. So I wrote this program so as to be able to sort bitwise left to right ascending or optionally descending regardless of any locale setting.

Why a binary tree? Binary trees are a fascinating data structure, at least to me. Insertion is of the order of  $\log(n)$  where  $n$  is the number of items being sorted. They work very well on random ordered data but do degenerate to  $n^2$  behaviour if the data arrives already in order or in reverse order. To overcome that problem you can use instead a self balancing tree, such as the AVL tree, or another rather well balanced tree called a red-black tree. This program uses an AVL tree. Google on AVL Tree and follow the links to Wikipedia; this will unearth more than enough information on these particular data structures. A binary tree sort is easily made to produce a stable sort, ie one where equal keys are sorted in order of arrival. Btsort is almost as fast as qsort when it is unfettered from any requirement to produce a dictionary order. Of course you can force qsort to make a stable sort by appending an ASCII formatted record number but once you do that any speed advantage of qsort is well and truly lost.

## What the program does.

The program sorts a list of lines eg the dictionary used to solve Jumbles, and sorts it in character set order. The program reads from a text file and writes to stdout. I use a dictionary called 'mydict' which is derived from some source that used Webster's spelling. That is fine for puzzles like Jumble which originates in USA. But I also run a program called 'xword' to cheat on crossword puzzles which are mostly based on Oxford or Macquarie spelling. Consequently I need to add words to 'mydict' from time to time. I do it like this:

```
user> cp /usr/local/etc/mydict .
user> echo new_word >> mydict
user> btsort mydict > newdict
user> sudo mv newdict /usr/local/etc/mydict
```

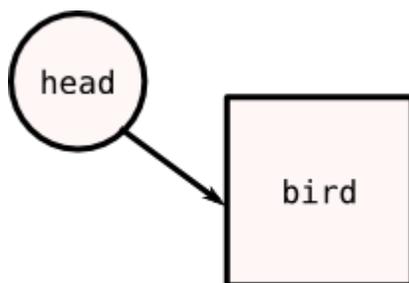
The program has one option, '-d' for a reverse character set ordered sort.

## How the program works.

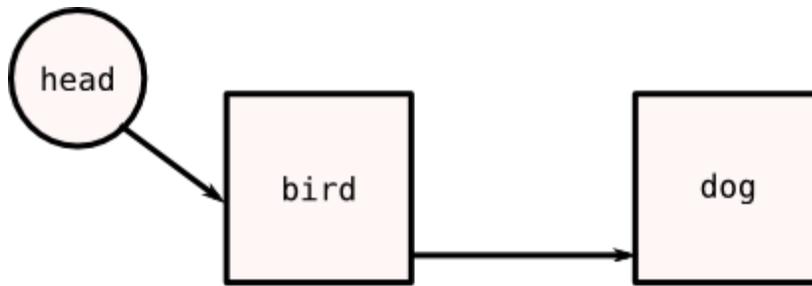
Before getting down to describing binary trees I will start by examining what they are not, the simple linked list. For each node of a linked list we have a structure like this:

```
typedef struct ln {
    char *data;
    char struct ln *next;
} LNODE;
```

When inserting into a linked list, the first node is appended to a head node and then after that any insert goes after the nodes that collate earlier and before those that collate later if any. If there is no bigger item the new node is appended to the list. In diagram form it is like the following:

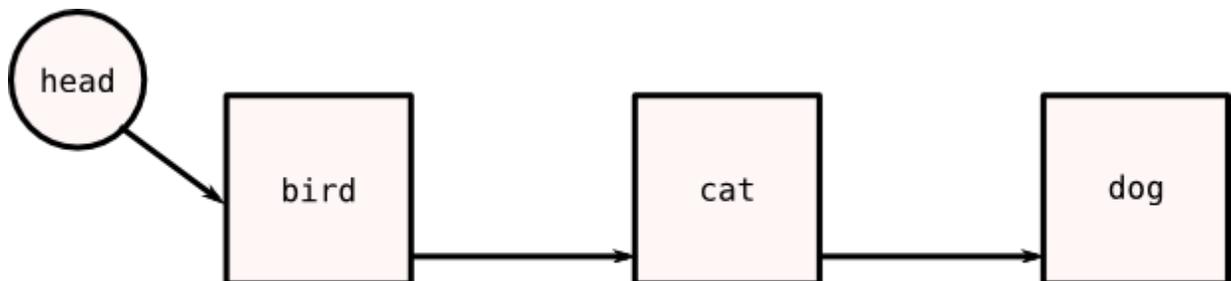


Linked List With One Item



Linked List With Two Items

So when inserting the item "cat", the link from "bird" to "dog" must be broken so that "bird" points to "cat" which in turn points to "dog". The list looks like this:



Linked List After Insertion

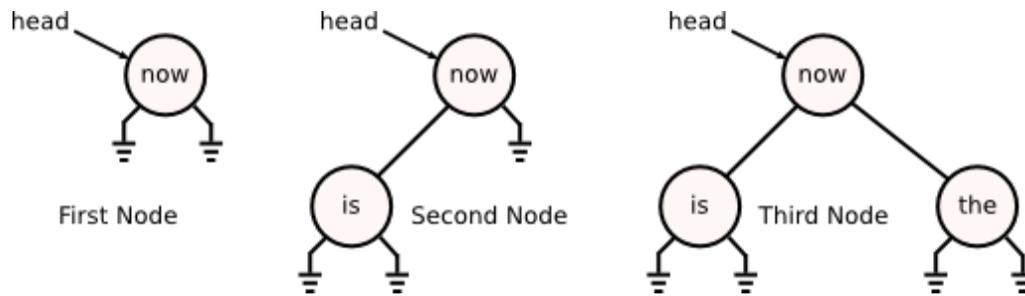
Well a binary tree is quite unlike the above. First the the data structure:

```

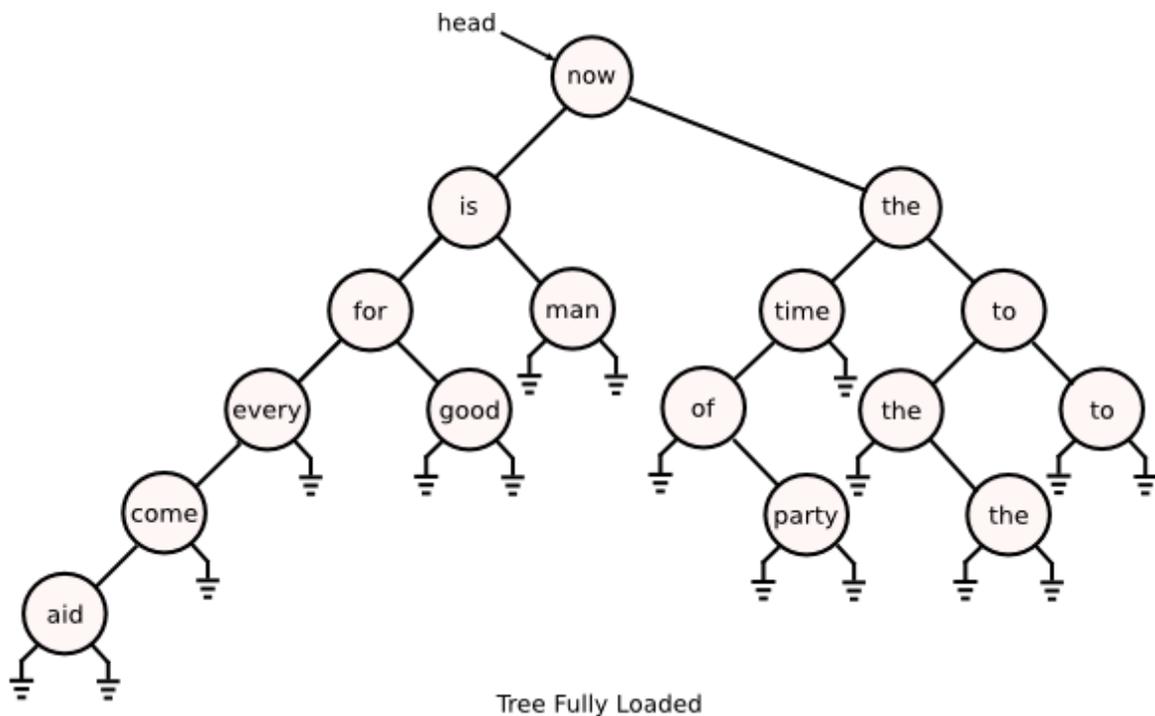
typedef struct tn {
    char *key;
    struct tn *left;
    struct tn *right;
}TNODE;
  
```

The binary tree node has, in the simplest form at least, a minimum of two pointers, one I'll call left, the other right. Now the usual convention is if an item is less than any existing item the program examines the path pointed to by left, otherwise it follows the path to the right. But unlike the linked list above, no linkage is ever broken and reassigned, the program simply follows the sorting rule until it finds an empty path, ie a NULL pointer, and assigns that pointer to the address of the new node. In other words all new insertions are leaf nodes.

Lets examine what happens when we insert that favourite sentence of the tty guys of yore: "now is the time for every good man to come to the aid of the party". The diagram below shows the tree with the first words inserted:



And here is the tree complete:



Since the purpose of the tree in this program to sort data in character set order the next process is an 'inorder' traversal of the tree.

## Inorder traversal

The inorder traversal requires a visit to the leftmost node of the tree and then as we return up the tree.

1. Output this node.
2. Up to the parent node. Output this node.
3. Traverse the right subtree following the same rules as above.

There is also the 'post order' traversal of the tree wherein we travel as far right as possible in a mirror of the inorder traversal and thus obtain a

descending order sort.

## How it works (The program listing).

```
1 /* btsort.c - binary tree sort program for lines of chars eg as in a
2 * dictionary. Sorting is done using an AVL tree so it's well behaved
3 * on pre sorted data. The sort is stable, ie identical elements are
4 * output in order of receipt. The program reads from argv[1] or
5 * argv[2] and writes to stdout.
6 */
7 #include<stdio.h>
8 #include<string.h>
9 #include<stdlib.h>
10 #include<ctype.h>
11
12 // #define DEBUG 1
13 #ifndef DEBUG
14 typedef struct tn {
15     char *key;
16     struct tn *left;
17     struct tn *right;
18     int balance;
19     int number;
20 }TNODE;
21 #endif
22 #ifndef DEBUG
23 typedef struct tn {
24     char *key;
25     struct tn *left;
26     struct tn *right;
27     int balance;
28 }TNODE;
29 #endif
30 TNODE *newnode(char *line);
31 char *getline(FILE *fp); // so named because getline() exists;
32 // it returns a '\n' I don't want.
33 TNODE *tinsert(TNODE *parent, TNODE *node, char *line, int dir);
34 void tprint(TNODE *node);
35 void rprint(TNODE *node);
36 void inprint(TNODE *node);
37 #ifndef DEBUG
38 int node_count = 0;
39 void node_print(TNODE *parent, TNODE *node, char *text);
```

```
40 #endif
41 void do_error(char *msg);
42
43 enum{Left = -1, Right = 1, None = 0};
44
45 int main (int argc, char** argv) {
46     TNODE *head = NULL;
47     char *line;
48     int sortdir = 0;
49     FILE *fpi;
50     char *infile;
51 #ifndef DEBUG
52     if (system("ls track > /dev/null") == 0)
53         system("rm track");
54 #endif
55     if (argc == 3){
56         if (strcmp("-d", argv[1]) == 0)
57             sortdir = 1;
58         else
59             do_error("invalid option");
60         infile = argv[2];
61     } else {
62         infile = argv[1];
63     }
64     fpi = fopen(infile, "r");
65     if (!(fpi)) {
66         fprintf(stderr, "Failed to open %s\n", infile);
67         exit(1);
68     }
69     while ((line = getlin(fpi))) {
70         head = tinsert(NULL, head, line, None);
71     } // while()
72     if (sortdir == 1)
73         rprint(head);
74     else
75         tprint(head);
76     //puts("");
77 /* inprint(head);
78 puts("");*/
79     return 0;
80 } // main()
81 TNODE *newnode(char *line){
82     char *p;
83     TNODE *tmp = (TNODE *)malloc(sizeof(TNODE));
```

```

84  if (tmp && (p = strdup(line))) {
85      tmp->left = NULL;
86      tmp->right = NULL;
87      tmp->key = p;
88      tmp->balance = 0;
89  #ifdef DEBUG
90      node_count++;
91      tmp->number = node_count;
92  #endif
93  } else {
94      fprintf(stderr, "Could not get memory\n");
95      exit(1);
96  }
97  return tmp;
98 }// newnode()
99 TNODE *tinsert(TNODE *parent, TNODE *node, char *line, int dir){
100  int result;
101  if (node == NULL) {
102      node = newnode(line);
103  } else if ((result = strcmp(line, node->key)) >= 0) {
104      node->balance++;
105      node->right = tinsert(node, node->right, line, Right);
106  } else {
107      node->balance--;
108      node->left = tinsert(node, node->left, line, Left);
109  }
110  if (node->balance == -2 ){ /* rotate right
111      have to make the left child the parent of
112      the node we are looking at */
113      TNODE *np, *op, *ll; // new parent, old parent, left link
114  #ifdef DEBUG
115      node_print(parent, node, "Before right rotation\n");
116  #endif
117  /* Terminology:
118      Old parent, the node we are looking at
119      New parent, the left child of the old parent
120  What changes:
121  1. Left link of old parent to become the right link
122      of the new parent.
123  2. Right link of new parent -> old parent.
124  3. Balance of both new and old becomes 0
125  What stays the same:
126  1. Left link of new parent remains as is.
127  2. Right link of old parent remains as is.

```

```
128  */
129  /* preserve existing states before we destroy any
130     existing linkage
131  */
132  np  = node->left;
133  op  = node;
134  ll  = np->right; // What changes 1.
135  np->right = op; // What changes 2.
136  op->left  = ll; // What changes 1.
137  np->balance = op->balance = 0; // What changes 3.
138  node  = np; // New parent
139  #ifdef DEBUG
140  node_print(parent, node, "After right rotation\n");
141  #endif
142  } // if (node->bal...
143  if (node->balance == 2) { /* rotate left
144     have to make the right child the parent of
145     the node we are looking at */
146     TNODE *np, *op, *rl; // new parent, old parent, right link
147  #ifdef DEBUG
148  node_print(parent, node, "Before left rotation\n");
149  #endif
150  /* Terminology:
151     Old parent, the node we are looking at
152     New parent, the right child of the old parent
153  What changes:
154  1. Right link of old parent to become the left link
155     of the new parent.
156  2. Left link of new parent -> old parent.
157  3. Balance of both new and old becomes 0
158  What stays the same:
159  1. Right link of new parent remains as is.
160  2. Left link of old parent remains as is.
161  */
162  /* preserve existing states before we destroy any
163     existing linkage
164  */
165  np  = node->right;
166  op  = node;
167  rl  = np->left; // What changes 1.
168  np->left  = op; // What changes 2.
169  op->right = rl; // What changes 1.
170  np->balance = op->balance = 0; // What changes 3.
171  node  = np; // New parent
```

```
172 #ifdef DEBUG
173     node_print(parent, node, "After left rotation\n");
174 #endif
175     } // if (node->bal...
176     return node;
177 } // tinsert()
178 #define MAX 1000
179 char *getlin(FILE *fp){
180     static char buf[MAX];
181     int ch, count;
182     count = 0;
183     while ((ch = fgetc(fp)) != EOF && (ch != '\n'))
184         buf[count++] = ch;
185
186     buf[count] = '\0';
187     if (ch == EOF)
188         return NULL;
189     else
190         return buf;
191 } // getlin()
192 void tprint(TNODE *node){
193     // pre-order traversal
194     if (node->left)
195         tprint(node->left);
196     printf("%s\n", node->key);
197     if (node->right)
198         tprint(node->right);
199     return;
200 } // tprint()
201 void rprint(TNODE *node){
202     // post-order traversal
203     if (node->right)
204         rprint(node->right);
205     printf("%s\n", node->key);
206     if (node->left)
207         rprint(node->left);
208     return;
209 } // rprint()
210 void inprint(TNODE *node){
211     // inorder traversal
212     printf("%3d %s\n", node->balance, node->key);
213     if (node->left)
214         inprint(node->left);
215     if (node->right)
```

```
216     inprint(node->right);
217     return;
218 }// inprint()
219 #ifdef DEBUG
220 void node_print(TNODE *parent, TNODE *node, char *text) {
221     FILE *fp;
222     fp = fopen("track", "a");
223     fputs(text, fp);
224     fprintf(fp, "node count %d\n", node_count);
225     if(parent)
226         fprintf(fp, "parent->key %s ..->number %d\n", parent->key,
227             parent->number);
228     fprintf(fp, "node->key %s ..->number %d\n", node->key,
node->number);
229     if (node->left) {
230         fprintf(fp, "node->left->key %s ..->number %d\n",
231             node->left->key,
232             node->left->number);
233         if (node->left->left)
234             fprintf(fp, "node->left->left->key %s ..->number %d\n",
235                 node->left->left->key,
236                 node->left->left->number);
237         if (node->left->right)
238             fprintf(fp, "node->left->right->key %s ..->number %d\n",
239                 node->left->right->key,
240                 node->left->right->number);
241     }// if (node->left)
242     if (node->right) {
243         fprintf(fp, "node->right->key %s ..->number %d\n",
244             node->right->key,
245             node->right->number);
246         if (node->right->right)
247             fprintf(fp, "node->right->right->key %s ..->number %d\n",
248                 ,node->right->right->key,
249                 node->right->right->number);
250         if (node->right->left)
251             fprintf(fp, "node->right->left->key %s ..->number %d\n",
252                 node->right->left->key,
253                 node->right->left->number);
254     }// if (node->right)
255     fputs("\n", fp);
256     fflush(fp);
257     fclose(fp);
258     return;
```

```
259 }// node_print()
260 #endif
261 void do_error(char *msg) {
262     fputs(msg, stderr);
263     fputs("\n", stderr);
264     exit(1);
265 } // do_error()
```

## Afterwords

Some improvements are strongly needed:

1. It needs a help function so the initial processing should be replaced with standard options processing, maybe allowing the output file to be specified but default to standard out.
2. The sorted file ends up entirely in memory along with the necessary tnode structures which will number slightly more than 50% of the input line count. Each line read is strdup() on read. Possibly I can gain some speed advantage by reading the entire file into memory after opening and then preallocate the space for the tnodes. There would be some waste of memory because I'd allocate one tnode per line.
3. Implementing the above would allow me to have memory allocation failures followed by sorting the file in two halves, four quarters etc and then merge the smaller parts of the file. Whether I'd go that far depends on the how much use this program gets in the wild.
4. That leads me to the next necessity. Put it up somewhere! Sourceforge or Ubuntu One. Suggestions are welcome.