

I117 (10) テキストファイル (その3)

知念

北陸先端科学技術大学院大学 情報科学研究科
School of Information Science,
Japan Advanced Institute of Science and Technology

文字列頻度計上

ファイル中の文字列を登場頻度を数える
以下のようなファイルを対象とする。

```
Aikawa  
Aikawa  
Aoyama  
Asaka  
Asaba
```

先日登場した辞書プログラムを応用できる
メモリ上の配列ではなくファイル内容を格納する

文字列頻度計上 (*cont.*)

辞書領域確保

```
#define DICTLEN      (10)
dict_t *dict;
int      dictlen, dictuse;
...
dictlen = DICTLEN;
dict = (dict_t*)malloc(
        sizeof(dict_t)*dictlen);
dictuse = 0;
```

必要数が分からないので、ひとまず固定

読み込んだ行をそのまま格納、頻度計上

```
while(fgets(line, BUFSIZ, stdin)) {
    ref.value = line;
    ppos = bsearch(&ref, dict, dictuse,
                  sizeof(dict[0]), dictcmp);
    if(!ppos) { /* not found */
        dict[dictuse].value = strdup(ref.value);
        dict[dictuse].count = 1; dictuse++;
        qsort(dict, dictuse, sizeof(dict[0]),
              dictcmp);
    }
    else { ppos->count++; }
}
```

```
for(i=0;i<dictuse;i++){  
    printf("%2d %d %s\n", i, dict[i].count,  
           dict[i].value);  
}
```

結果 — 行区切りが含まれている

0 2 Aikawa

A	i	k	a	w	a	\n
---	---	---	---	---	---	----

1 1 Aoyama

A	o	y	a	m	a	\n
---	---	---	---	---	---	----

2 1 Asaba

A	s	a	b	a	\n
---	---	---	---	---	----

3 1 Asaka

A	s	a	k	a	\n
---	---	---	---	---	----

文字列頻度計上 (*cont.*)

大筋動作する

行区切りを除去

- 読み込む時点で行区切りを除去する
区切り文字を '\0' にする関数を用意

```
int chomp(char *src) {
    char *p=src;
    while(*p) {
        if(*p=='\n' || *p=='\r') {
            *p = '\0';    return 0;    }
        p++;
    }
    return 0;    }
```

行区切りを除去 (*cont.*)

行を読み込んだ後にこの関数を呼び出す

```
while(fgets(line, BUFSIZ, stdin)) {  
    chomp(line);  
    ref.value = line;  
}
```

結果 — 行区切りがなくなる

```
0 2 Aikawa  
1 1 Aoyama  
2 1 Asaba  
3 1 Asaka
```


格納容量の調整 – マクロの利用

固定では使いづらいため調整する方法を紹介する
未定義ならプログラム内で定義（デフォルト指定）

```
#ifndef DICTLEN
#define DICTLEN      (10)
#endif
```

足りなくなるとエラーを表示

```
if(!ppos) { if(dictuse>=dictlen) {
    fprintf(stderr, "no more memory\n");
    continue; } }
```

格納容量の調整 – マクロの利用 (*cont.*)

コンパイル時に特定の数を与えて調整可能

```
% cc -DDICTLEN=4 wfc01c.c      % cc -DDICTLEN=3 wfc01c.c
% cat x0 | ./a.out             % cat x0 | ./a.out
0 2 Aikawa                     no more memory
1 1 Aoyama                     0 2 Aikawa
2 1 Asaba                      1 1 Aoyama
3 1 Asaka                      2 1 Asaka
```

DICTLEN が 3 以下になると容量が足りなくなる
容量を変更できていることが確認できた

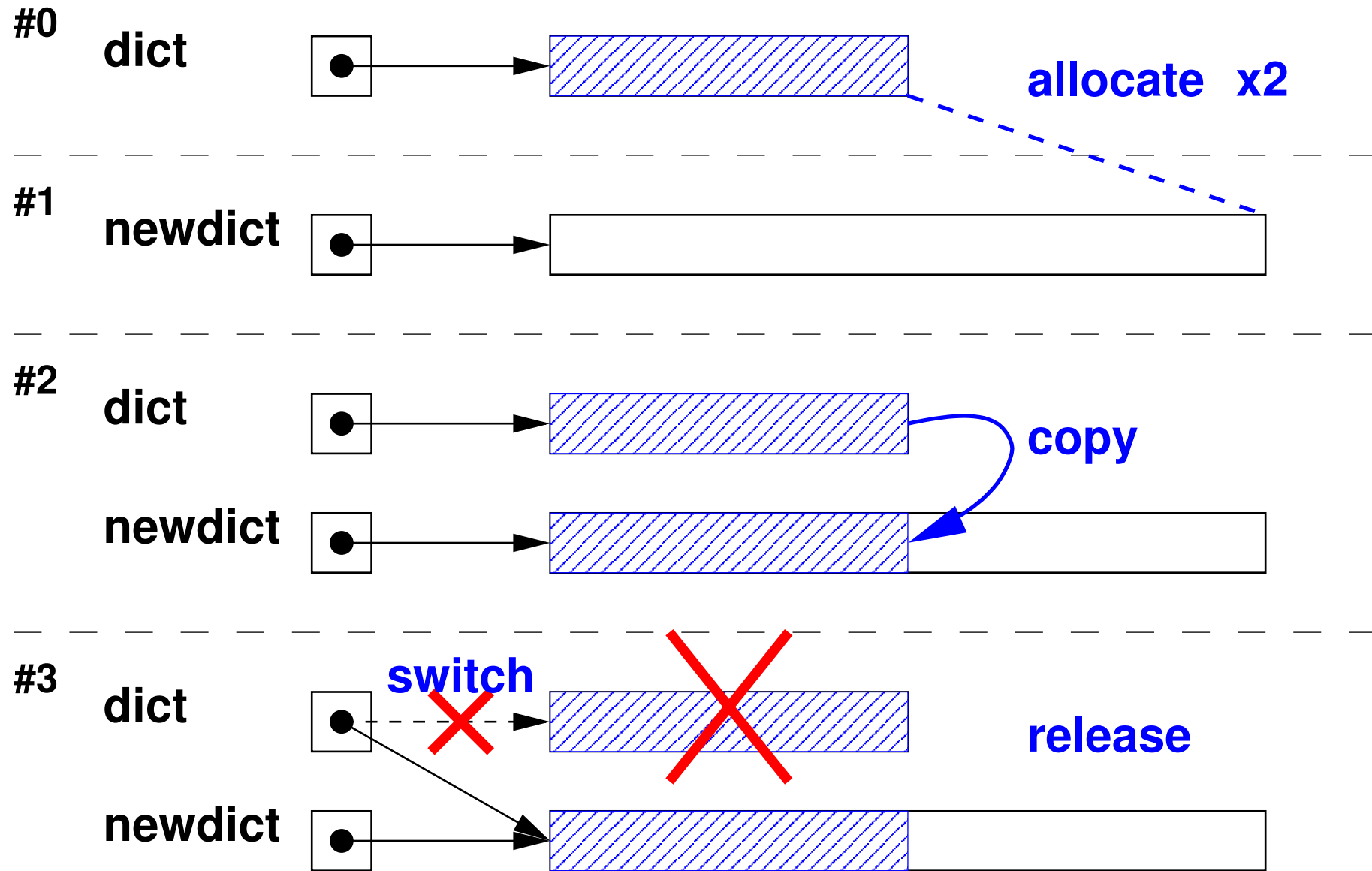
格納容量の調整 – 不足時に再確保

領域拡大 (2倍) の関数を用意

```
int expand() {
    dict_t *newdict;  int newlen, i;
    newlen = dictlen*2;
    newdict = (dict_t*)malloc(
        sizeof(dict_t)*newlen);
    for(i=0;i<dictuse;i++){ newdict[i] = dict[i]; }
    free(dict);  dict = newdict;
    fprintf(stderr, "#expand %d -> %d\n",
        dictlen, newlen);
    dictlen = newlen; }
}
```

格納容量の調整 – 不足時に再確保 (*cont.*)

- 新しいサイズ newlen を dictlen の2倍とする
 - 新しい領域 newdict を malloc で確保
 - newdict に dict の内容をコピー
 - dict を開放して、newdict を代入する
 - newdict に newlen を代入する
- ※ 経過が分かるように stderr に長さの変化を出力



格納容量の調整 – 不足時に再確保 (*cont.*)

メモリ領域不足時にこの領域拡大関数を呼び出す

```
if (!ppos) { /* not found */
    if (dictuse >= dictlen) {
        expand();
    }
}
```

こうすると、可能な限り格納する
上限は OS やシェルなどの設定による

格納容量の調整 – 不足時に再確保 (*cont.*)

```
% cc -DDICTLEN=3 wfc01d.c      % cc -DDICTLEN=2 wfc01d.c
% cat x0 | ./a.out            % cat x0 | ./a.out
#expand 3 -> 6                #expand 2 -> 4
 0 2 Aikawa                    0 2 Aikawa
 1 1 Aoyama                    1 1 Aoyama
 2 1 Asaba                     2 1 Asaba
 3 1 Asaka                     3 1 Asaka
```

長さの初期値が 3 や 2 でも拡張して格納している

出力の工夫 – 頻度順整列

データが多いと、頻度順の出力が欲しくなる

```
...  
185 1 text  
186 1 text.  
187 11 the  
188 2 these  
190 3 to  
191 1 to,  
...
```


出力の工夫 – 頻度順整列 (*cont.*)

頻度の比較関数を書いて qsort を適用

```
int freqcmp(const void *ap,
            const void *bp) {
    int a, b;
    a = ((dict_t*)ap)->count;
    b = ((dict_t*)bp)->count;
    return a-b; }

...
qsort(dict, dictuse,
      sizeof(dict[0]), freqcmp);
```

出力の工夫 – 頻度順整列 (*cont.*)

結果 — 昇順、最後の 177 件は空行

```
... (略)
197 4 |
198 6 If
199 7 man
200 7 SunOS
201 7 User
202 11 The
203 11 the
204 177
```

整形

プログラムの出力を整形することも多い

- テキスト
 - ◇ pr 等
- 整形（組版）システムへ送る
 - ◇ HTML
 - ◇ (La)TeX, roff
- 印刷システムへ送る
 - ◇ PostScript

pr

印刷向けにテキストを整形するプログラム
ヘッダ・フッタなしで 38文字に複数段を整形すると

```
% ls *-?.tex | pr -t -w 38 -2
```

```
I117-0.tex           I117-6.tex
```

```
I117-1.tex           I117-7.tex
```

```
I117-2.tex           I117-8.tex
```

(略)

```
% ls *-?.tex | pr -t -w 38 -3
```

```
I117-0.tex          I117-4.tex          I117-8.tex
```

```
I117-1.tex          I117-5.tex          I117-9.tex
```

(略)

横幅が足りない際には切り捨てられる点に注意

```
% ls *-?.tex | pr -t -w 38 -4
I117-0.t I117-3.t I117-6.t I117-9.t
I117-1.t I117-4.t I117-7.t I117-p.t
I117-2.t I117-5.t I117-8.t
```

作成プログラムの出力を pr に与えて整形

```
% ./a.out |pr -t -w 38 -2
197 4 |                201 7 User
198 6 If                202 11 The
199 7 man                203 11 the
200 7 SunOS             204 177
```

整形 HTML table タグ

WWW で使われているテキスト形式

ID	freq	string
198	6	If
199	7	man
200	7	SunOS

```
<table>
<tr><th>ID</th><th>freq</th><th>string</th></tr>
<tr><td align="right">198</td>
  <td align="right">6</td><td>If</td></tr>
<tr><td align="right">199</td>
  <td align="right">7</td><td>man</td></tr>
<tr><td align="right">200</td>
  <td align="right">7</td><td>SunOS</td></tr>
</table>
```

整形 HTML table タグ

いくらか省略可能

ID	freq	string
198	6	If
199	7	man
200	7	SunOS

```
<table>
<tr><th>ID<th>freq<th>string
<tr><td align="right">198<td align="right">6<td>If
<tr><td align="right">199<td align="right">7<td>man
<tr><td align="right">200<td align="right">7<td>SunOS
</table>
```

整形 LaTeX tabular 環境

```
\begin{tabular}{rrl}
\multicolumn{1}{c}{ID} & & 
\multicolumn{1}{c}{freq} & & 
\multicolumn{1}{c}{string} \\
198 & 6 & If \\
199 & 7 & man \\
200 & 7 & SunOS
\end{tabular}
```

ID	freq	string
198	6	If
199	7	man
200	7	SunOS

整形 roff tbl

	ID	freq	string
	198	6	If

.TS

199	7	man
-----	---	-----

ccc

200	7	SunOS
-----	---	-------

rri.

ID → freq → string

198 → 6 → If

199 → 7 → man

200 → 7 → SunOS

.TE

タブ (0x09) を → で表現

整形 PostScript

- プリンタ向けページ記述言語
- グラヒック機能が多彩
- 表作成の専用機能はない
- 文字の大きさや座標を考えて整形する

ID	freq	string
198	6	lf
199	7	man
200	7	SunOS

整形 PostScript (cont.)

フォントを選択して、関数群を定義

```
%!PS-Adobe-2.0 EPSF-2.0
%%BoundingBox: 0 0 250 160
/Helvetica findfont 32 scalefont setfont
/rjshow { /s exch def s stringwidth pop sub
  exch /y exch def add y moveto s show } def
/cjshow { /s exch def s stringwidth pop sub 2
  div exch /y exch def add y moveto s show } def
/ljshow { /s exch def pop moveto s show } def
```

rjshow 右詰め、cjshow センタリング、ljshow 左詰め

```
10 130 50 (ID) cjshow
70 130 60 (freq) cjshow
140 130 100 (string) cjshow
10 90 50 (198) rjshow
70 90 60 (6) rjshow
140 90 100 (If) ljshow
10 50 50 (199) rjshow
70 50 60 (7) rjshow
140 50 100 (man) ljshow
10 10 50 (200) rjshow
70 10 60 (7) rjshow
140 10 100 (SunOS) ljshow
```

スタック言語、慣れると読めるようになる

整形 一般化

一般化すると単純

前処理

```
for (i=0; i<dictuse; i++) {  
    dict[i] 出力  
}
```

後処理

後段プログラムに合わせて前・後処理を記述
さらに前後に処理が必要な場合もある

整形 一般化 (*cont.*)

LaTeX向けの出力ルーチン例

```
printf( "\\begin{tabular}{rrl}\n" );  
printf( "\\multicolumn{1}{c}{ID} &\n" );  
printf( "\\multicolumn{1}{c}{freq} & \n" );  
printf( "\\multicolumn{1}{c}{string} \\\\ \n" );  
for(i=0;i<dictuse;i++){  
    printf( "%d & %d & %s \\\\ \n",  
           i, dict[i].count, dict[i].value );  
}  
printf( "\\end{tabular}\n" );
```

実行例

```
% ./a.out < x0
\begin{tabular}{rrl}
\multicolumn{1}{c}{ID} & & 
\multicolumn{1}{c}{freq} & & 
\multicolumn{1}{c}{string} \\
0 & 1 & Aoyama \\
1 & 1 & Asaba \\
2 & 1 & Asaka \\
3 & 2 & Aikawa \\
\end{tabular}
```

表の分はできたが、コンパイルにはまだ足りない

表の外側、文書の前処理と後処理を追加

```
printf( "\\documentclass[a4j]{jarticle}\\n" );  
printf( "\\begin{document}\\n" );  
  
...  
printf( "\\end{document}\\n" );
```

コンパイル成功

```
% ./a.out < x0 > x1.tex  
% platex x1.tex  
This is pTeX, Version 3.14159-p3.1.3 (euc) (Web2C 7.4.5)  
(./x1.tex  
pLaTeX2e <2001/09/04>+0 (based on LaTeX2e <2001/06/01> patch level 0)  
(/app/teTeX-2.0.2-pTeX-3.1.3/share/texmf/ptex/platex/base/jarticle.cls  
Document Class: jarticle 2001/10/04 v1.3 Standard pLaTeX class  
(/app/teTeX-2.0.2-pTeX-3.1.3/share/texmf/ptex/platex/base/jsize10.clo)  
(./x1.aux) [1] (./x1.aux) )  
Output written on x1.dvi (1 page, 384 bytes).  
Transcript written on x1.log.
```


出力内容

```
% ./a.out < x0
\documentclass[a4j]{jarticle}
\begin{document}
\begin{tabular}{rrl}
\multicolumn{1}{c}{ID} & & 
\multicolumn{1}{c}{freq} & & 
\multicolumn{1}{c}{string} \\
0 & 1 & Aoyama \\
1 & 1 & Asaba \\
2 & 1 & Asaka \\
3 & 2 & Aikawa 
\end{tabular}
\end{document}
```

演習

1) 実際に文字列頻度を調査せよ★

この講義のホームページに LIST01 というファイルをおいた

URLは以下のとおり

**[http://www2.jaist.ac.jp/is/private/
lecture/i117/2008/LIST01](http://www2.jaist.ac.jp/is/private/lecture/i117/2008/LIST01)**

文字列頻度を計上するプログラムを作成し、このファイル中の文字列とその頻度の一覧を作れ

演習 (cont.)

WWW 上のデータは wget で取得できる

```
% wget http://www2.jaist.ac.jp/is/private/lecture/i117/2008/LIST01
--16:25:19--  http://www2.jaist.ac.jp/is/private/lecture/i117/2008/LIST01
              => `LIST01'
Resolving www2.jaist.ac.jp... 150.65.38.80
Connecting to www2.jaist.ac.jp|150.65.38.80|:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 624 [text/plain]

100%[=====>] 624          ---K/s

16:25:20 (11.12 MB/s) - `LIST01' saved [624/624]
```

WWW ブラウザでも取得可能

演習 (cont.)

- 2) 先のプログラムをコメントや空行を無視するよう変更せよ★
- 3) 先のプログラムを頻度で整列し降順に出力するようプログラムを変更せよ★
- 4) 実際に文字列頻度を調査せよ★★
先と同様にファイル LIST02 を用意した
そのファイルに対して以下の処理のプログラムと
その結果を確認せよ
 - a) 昇順(頻度の多い順)に文字列を表示、頻度の多い上位 5つ

演習 (cont.)

b) 降順(頻度の少ない順)に文字列を表示、頻度の少ない上位5つ

ただし、以下の条件をつける

- 先頭に # が付いている行はコメントとして除く

5) 先の問題を C 言語以外のプログラムで解く方法を見付けよ★

- スクリプト言語でもよい
- データベースに格納して問い合わせ言語で解いてもよい

演習 (cont.)

- 6) 頻度で整列して出力するプログラムが HTML の表を出力するよう変更せよ★
 - WWW ブラウザなど、HTML を解釈するプログラムで表示して確認せよ
- 7) 先のプログラムに、枠を付ける等、装飾がつくよう変更せよ★