

Assignment 10: Bloom Filters

The Internet is very large, very fast, and filled with degeneracy. The proliferation of foreign persons who not only do not speak English, they use it in corrupt and deviant ways. The untermenschen spend their days sending each other cat videos. Your hero, E.B. decided that in his country of Pacifica that a more neutral *newspeak* was required to keep the people content, pure, and from thinking too much. But how do you process so many words as they flow into your little country at 10Gbits/second? The answer that comes to your brilliant and pure mind is that you use a *Bloom filter*.

A Bloom filter is a hash table, where the entries in the hash table are simply single bits. Consider $h(\text{text}) = k$, then if $B_k = 0$ then the entry is definitely missing; if $B_k = 1$ then the entry may be present. The latter is called a false positive, and the false positive rate depends on the size of the Bloom filter and the number of texts that hash to the same position (hash collisions).

You assign your programming minions to take every *potentially offensive* word that they can find in an English dictionary and hash them into the Bloom filter. That is, for each word the corresponding bit in the filter is set to 1.

You will make *two* such Bloom filters. Why? to reduce the chance of a *false positive*. Each Bloom filter should be relatively large (thousands of entries), and use a different hash function. This is easily accomplished by giving our hash function a different input block.

When presented with a stream of text, it is first passed to the first layer Bloom filter. If the Bloom filter rejects any words then the person responsible is innocent of a *thoughtcrime*. But if the word has a corresponding bit set in the Bloom filter, it is likely that they are *guilty* of a thoughtcrime, and that is very ungood indeed. To make certain, we must consult the *hash table* and if the word is there as a *proscribed* word then they will be placed into the care of *Miniluv* and is sent off to *joycamp*. If the word passed both Bloom filters, and is not a forbidden word then the hash table will provide a translation that will replace offensive, insensitive, and otherwise dangerous words with new approved words. The advantage is that your government can augment this list at any time via the *Minitrue*.

Minitrue can also add words from the Bloom filter that you, in your wisdom, have determined are not wholesome for your people to use. There are three cases to consider:

1. Words that are approved, these will not appear in the Bloom filter.
2. Words that should be replaced, which will have a mapping from the old word to the new approved word, and
3. Words where no mapping to new approved words means that it's off to *joycamp*.

You will use the *bit vector* data structure that you developed for counting bit project to implement your Bloom filter.

```

1 # ifndef NIL
2 # define NIL (void *) 0
3 # endif
4 # ifndef _BF_H
5 # define _BF_H
6 # include <stdint.h>
7 # include <stdlib.h>
8 # include <stdio.h>
9
10 typedef struct bloomF {
11     uint8_t *v; // Vector
12     uint32_t l; // Length
13     uint32_t s[4]; // Salt
14 } bloomF;
15
16 // Each function has its own hash function, determined by the salt.
17
18 uint32_t hashBF(bloomF *, char *);
19
20 // Create a new Bloom Filter of a given length and hash function.
21
22 static inline bloomF *newBF(uint32_t l, uint32_t b[])
23 {
24     // Code
25 }
26
27 // Delete a Bloom filter
28
29 static inline void delBF(bloomF *v)
30 {
31     // Code
32 }
33
34 // Return the value of position k in the Bloom filter
35
36 static inline uint32_t valBF(bloomF *x, uint32_t k)

```

```

37 {
38     // Code
39 }
40
41 static inline uint32_t lenBF(bloomF *x) { return x->l; }
42
43 // Count bits in the Bloom filter
44
45 static inline uint32_t countBF(bloomF *b)
46 {
47     // Code
48 }
49
50 // Set an entry in the Bloom filter
51
52 static inline void setBF(bloomF *x, char * key)
53 {
54     // Code
55 }
56
57 // Clear an entry in the Bloom filter
58
59 static inline void clrBF(bloomF *x, char *key)
60 {
61     // Code
62 }
63
64 // Check membership in the Bloom filter
65
66 static inline uint32_t memBF(bloomF *x, char *key)
67 {
68     // Code
69 }
70
71
72 static inline void printBF(bloomF *x)
73 {
74     // Code
75 }
76
77 # endif

```

bf.h

bf.c is downloadable from the moodle web-site. You may change it to **bf.h** so that it can be included in other programs.

Sample text file to test bloom filter “**gettsburg.txt**” should be used in this project. Please ignore all punctuation marks.