

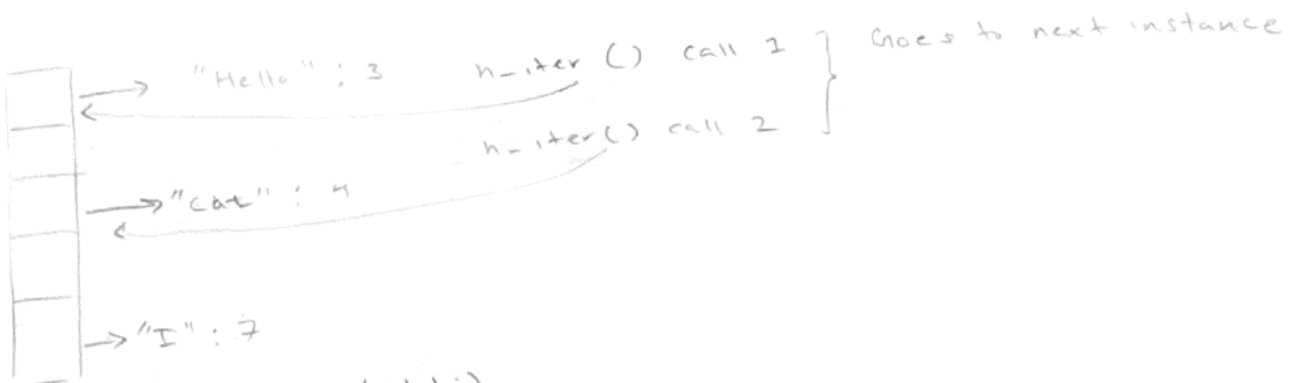
## Assignment 7:

In some 3D space, which of known texts are closest to new value.



Hash table size is  $1.3 \rightarrow 1.5 \times$  the expected unique words

ht-size (\*ht)  
return (top); every add of unique word increments to top.



Node \*ht-iter(\*hti)

do { • do atleast once  
slot ++; move to next slot

} while (ht[slot % ht-size] → frequency == 0 or slot > ht-size)

• If the slot goes past max exit % to Loop  
• go to next non-empty instance

return &ht[slot % ht-size];

$\Omega$  ← lower bound ← best case  $\Omega(1)$

$\Theta$  ← Average case

$O$  ← upper bound ← worst case  $O(n)$

## SPECK Cipher

-  $\text{hash}(\text{salt } 1, \text{"cat"}) \rightarrow 6$   
-  $\text{hash}(\text{salt } 2, \text{"cat"}) \rightarrow 100$  } Different salts result in different values

Salt is an array of 2 uint64-ts

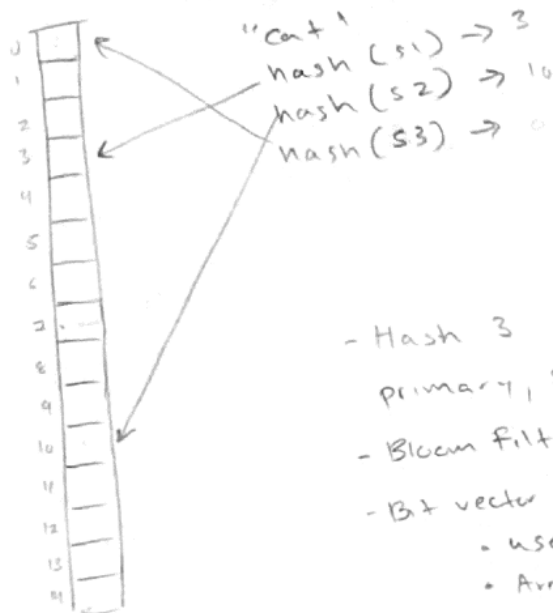
- we use 128 bit salts

Hash table uses 2 salt

## Bloom Filters

- A bit vector, reports if an entry has not been seen before. Can result in false positives because probabilistic data structure. Salt [2]

{ WENT OVER IN LECTURE  
REVIEW }



- Hash 3 times for Bloom Filter  
primary, secondary, tertiary

- Bloom filter is constructed from a bit vector

- Bit vector is a dynamically size of bits

- use malloc, and realloc
- Array of uint64s

bv\_create (n bits)

Allocate  $\lceil n/8 \rceil$  - if  $n/8$  (+1) bytes

bv\_length

return bv  $\rightarrow$  length

bv\_clrbit

$\text{bv}[n/8] \gg n\%8 \& 0$

bv\_setbit

$\text{bv}[n/8] \gg n\%8 \mid 1$

## Distance functions

$\vec{u}$  = Anon.txt (infile)

$\vec{v}$  = known.txt (refile)

$\vec{u} = \langle 1, 0, 1 \rangle$

$\vec{v} = \langle 0, 2, 1 \rangle$

Anon.txt "Hello World"

known.txt "Goodbye, goodbye, world"

$\langle \text{"hello"}, \text{"goodbye"}, \text{"world"} \rangle$

• use regex to make lowercase

$u[\text{unique}], v[\text{unique}]$  // dimension is # of unique words in both infile and refile (maybe calloc)

use a loop for  $\left[ \begin{array}{l} \text{in\_wc} = 0; // \text{infile wordcount} \\ \text{ref\_wc} = 0; // \text{refile wordcount} \end{array} \right.$  unique = ht-size(2\*ht)  
? is there a ht for each file?

for ( $i=0$ ;  $i < \text{unique}$ ;  $i++$ )

word = ht\_iter(ht)  $\rightarrow$  word; // Node  $\rightarrow$  word

if (word == NULL) break;

$u[i] = \text{text\_frequency}(\text{infile}, \text{word})$  //  $\vec{u}_i = \frac{\vec{u}_i}{|\vec{u}_i|}$

$u[i] /= \text{in\_wc}$

$v[i] = \text{text\_frequency}(\text{refile}, \text{word})$  //  $\vec{v}_i = \frac{\vec{v}_i}{|\vec{v}_i|}$

$v[i] /= \text{ref\_wc}$

float dist = 0 // dist of 1 of the 3 dist methods

Manhattan

for ( $i < \text{unique}$ )

dist +=  $|u[i] - v[i]|$

return dist;

Euclidean

for ( $i < \text{unique}$ )

dist +=  $\sqrt{(u[i] - v[i])^2}$

return dist;

Cosine

for ( $i < \text{unique}$ )

dist +=  $u[i] \cdot v[i]$

return  $1 - \text{dist}$

file (anonymous)

"Kim likes cats"

noise  
"the and of ..."

```
Text {
  word-count
  bf
  ht
}
```



 "Kim" wc = 1  
 "likes" wc = 2  
 "cats" wc = 3

lib.db

```
2
Author A
texts/a.txt
Author B
texts/b.txt
```

→ "Alice likes cats"  
 → "Bob likes dogs"

Manhattan dist  
 $\vec{u}, \vec{v} \rightarrow \sum_{i=1}^n |\vec{u}_i - \vec{v}_i|$   
 Anon A

Anon ht

"Kim": 1  
 "likes": 2  
 "cats": 2

A ht

"Alice": 1  
 "likes": 1  
 "cats": 1

for (key, value in anon.txt) {  
 Anon-freq = text-freq  
 A-freq = text-freq

Anon.txt

"Kim"	Anon-freq	$\frac{1}{3}$	A-freq	$\frac{0}{3}$	→	$\frac{1}{3}$
"likes"		$\frac{1}{2}$		$\frac{1}{3}$		0
"cats"		$\frac{1}{2}$		$\frac{1}{3}$		0

$\sum |\vec{u}_i - \vec{v}_i|$

$\in a.txt \rightarrow \in anon.txt$

(word-contains function)

"Alice"		$\frac{0}{2}$		$\frac{1}{3}$		$\frac{1}{3}$
						<hr/>
						$\frac{2}{3}$

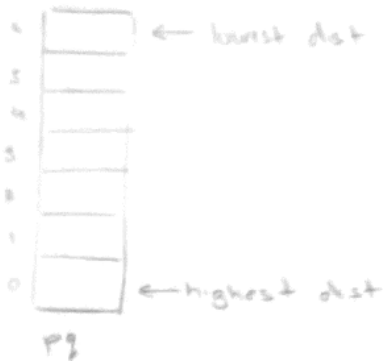
dist to A

# Priority Queue

- getting the  $k^{th}$  nearest neighbors  
→ have the lowest value distance (text-dist(...))

pq stores capacity elements

? what exactly is capacity based on?



- remove from end first.
- lowest distance has highest priority, min queue
- use heap or quick sort because pq can become very large

distance away from author's text compared to anonyms

bool enqueue(\*q, \*author, dist)

if (pq-full(q))  
return False

q[top] = (author, dist)

top++ // move up top  
return True

create a new instance of AuthorDist  
ad-create(author, dist)

? what is an author, dist pair?

• ? A struct?

• ? 2 element array?

[can be anything  
struct is probably best]

bool dequeue(\*q, \*\*author, \*dist)

if (pq-empty(q)): return False

top-- // move down top

~~\*author, \*dist = q[top]~~ // get values of author and dist from queue

\*author = q[top] → author  
\*dist = q[top] → distance ] // maybe have unpack function

ad-delete(&q[top]) // AD node no longer needed, can be freed  
return True

struct AuthorDist {

char author

float distance

}

## Hash Table

speck  $\rightarrow$  takes in a word outputs a number  
use mod to make it fall into the array

if collision add to linked list (can do binary search tree but seems like a lot of work)

0	$\rightarrow$ "a" : 2
1	
2	$\rightarrow$ "b" : 3
3	$\downarrow$ "c" : 4
4	
5	$\rightarrow$ "d" : 5

Node \*ht\_insert(ht, \*word)

Msize = ht  $\rightarrow$  Max-Size // total spaces in hash table

key; // stores key

hash(ht  $\rightarrow$  salt, &key) // function of speck takes

key % Msize // keep within bounds

Node \*link\_node = ht[key]

while (link\_node != 0 OR != NULL)

if (link\_node  $\rightarrow$  next == NULL):

link\_node  $\rightarrow$  next = create\_node(word);

link\_node = link\_node  $\rightarrow$  next

break;

ht  $\rightarrow$  size ++;

return link\_node;

Node \*ht\_lookup(ht, \*word)

key; hash(ht  $\rightarrow$  salt, &key);

key % ht  $\rightarrow$  size;

Node \*look\_node = ht[key]

while (look\_node != Null or 0)

if (look\_node  $\rightarrow$  word == word)

return look\_node;

if (look\_node  $\rightarrow$  next != Null)

look\_node = look\_node  $\rightarrow$  next

return NULL; // Nothing

Bloom Filter if False stop  
if true check

ht\_create (:size)

HashTable \*ht = (HashTable \*) malloc (sizeof (HashTable))

ht → size = size

ht → n-items = 0

ht → slots = (Node \*\*) calloc (sizeof (Node), size)

return ht ;

ht\_delete (\*\*ht)

free \*ht → slots

free (\*ht)

ht\_size (\*ht) // USED WRONG IN OTHER AREAS

return ht → size ;

## Assignment 7: Author Identification

purpose of the program is to take an anonymous txt file and find  $k$  authors and how different the styles are and maybe find a match

### General Outline

- 1) get 2 files Anon.txt, Known.txt
- 2) Create 2 Text ADT, each with
  - Hashtable  $ht\_create \rightarrow ht$
  - Bloom Filter  $bf\_create \rightarrow bf$
  - word-count; 1A for loop
- 3) Find Distances with either Man, Euclid, or Cos
  - Fill the  $ht$  of both, exclude noise  $\rightarrow$  also find the # of unique words in both files together  $u\_wc$
  - Create 2 arrays  $All$  both with frequency of a corresponding word. (might need to make 3rd Array for words) and divide by their own individual  $wc$  from their  $ht$
  - push onto priority-queue
- 4) Sort Priority Queue  $0^{th}$  index should be greatest distance, a min queue with heap or quick sort
- 5) Output take the first  $k$  off the queue and output it.



## Text

### Struct Text

Hash table ht

BloomFilter bf

int wordcount ← total number of words in text not unique

text\_create (inFile, noise)

noise "a", "the", "it" ← not indicative of author's dictation  
"him", "her", "there"

• create a text file of only noise words

→ text\_create (file, NULL) // create noise text

If current word being looked at in noise table do not add it to table

text\_dist (text1, text2, metric) ← 1 of the 3 measurement methods

text\_frequency (text, word)

return (double) f / (# unique words)

→ f is frequency of word

text\_contains (text, word)

bf → if (No) return false

bf → if (YES) check ht

$$\frac{12}{160} = 0.075$$

$$12 \times 16$$

$$x = 0.6L$$

$$1.6L \quad 3:5$$

Priority Queue

use Last Assignment

$$0.6L : 1L \text{ Rum}$$

$$0.4L : 0.5L \text{ GF}$$

$$0L : 1.5L \text{ whiskey}$$

$$1L : 3L$$

$$\frac{L(A)}{L(A) + L(W)} \quad (Pr.A) = w\_prior$$

parsing words at a time

- word character should contain letters from a-z, A-Z and accept contractions, and hyphenations

msg = "Hello world"

[a-z]+: ello world

[a-zA-Z]+: Hello world

+ means 1 or more \* means 0 or more ? 0 or 1

[.] character set

Ex 0b101101101 "0b(0|1)+"

WORD "[a-zA-Z]+"

while (nextword (file, WORD) != NULL)

← going through words

- for apostrophe    bob's ⇒ word "bob's"  
                          bobs' ⇒ word "bobs"

- for hyphen        bo-be ⇒ "bo-be"  
                          bo-be- ⇒ "bo-be"

' and - must have a letter before and after it  
if that is met can have any number of them