

DESIGN

Purpose: The purpose of assignment 7 or "The Great Firewall" is to take in a list of badspeak words and newspeak words and determine whether any given user input violates those of outline transgressions. To accomplish this system our directory makes use of a bloom filter, hash table, and linked lists of nodes.

Layout / Structure

Bloom Filter • A DS to determine if a word by the user is definitely not a bad word or probably is a bad word. It uses a bitvector as its "filter". If a word has 3 set bits in the three indexes assigned to it then it probably is a bad word. If it is missing one then it definitely isn't. This is bc some bad words can share indices, thus "false positives" occur.

Its attributes are 3 salts and a BitVector type filter.

BF insert() • Function that generates 3 indices & sets a bit at all three indices. (for a word)
index1, index2, index3 = getindex(), getindex(), getindex()
Setbit(index1, index2, index3).

PSuedocode

bf-probe(7)

• Function to check if the appropriate indices generated from ~~getindex~~ the word

Pseudocode

```
index1, index2, index3 = getindex(7), getindex, getindex(7)
if getbit(index1, index2, index3) == true:
    return true
else:
    return false
```

• To return true all indices must have a bit set.

bf_count(7)

• Returns number of bits set.

```
Count = 0
for bit in bits:
    if b == 1:
        Count += 1
return Count
```

Bit Vector

• Has components length (of bits) and
a vector which is an array of
unsigned integers that has memory allocated
precisely according to * the byte length.

• We calculate byte length by

Byte length = $x / 8$ (if x is divisible
by 8)

Byte length = $(x / 8) + 1$ (if x is not
divisible by 8)

For example $x = 17$ bits

Byte length = $(\frac{17}{8}) + 1 = 3$ bytes

We must waste a little memory since we must
store a single bit in one whole byte.

Set-Bit

- we set a bit by finding the index of the bit vector we'd like to set and then shifting it w/ a 1 left-shifted to position:

→ If we want to set 3rd bit of Cl000 then:

1000 0101

$$\text{OR } 0001\ 0000 = 1001\ 0101$$



Note: OR operation preserves all other bits b/c OR w/ anything returns identity.

Clr-bit

Similarly to clear a bit we use AND operation w/ not and a left shifted 1

$$\begin{array}{r} \text{AND } 1000\ 0101 \\ 0111\ 1111 \end{array}$$

Cl000 0101
↑

Note: All zeroes will stay zero when AND'd w/ 1

and all the 1's we wish to clear will turn to zero

Get-bit

• We right shift our vector by the index AND'd w/ 1.

Xor-bit

• We locate index and add the bit we wish to xor w/

• We Modulo 2 to ensure that if it were two 1's to be zeroed

100
100
000

Capacity: 7

IP: 4

7 - 4 = 3

Hash Table

- ADT to check if a word is certain a bad word
- Uses linked lists ~~and an array is stored~~ inside of a linked list

ht_lookup()

Psuedocode

- We get index ~~MAX~~ generated w/ the word and look for the word appropriate list

```
index = get_index(word)
if (list[index] == NULL):
    return NULL
ll_lookup(list[index], word)
Node = ll_lookup()
return Node
```

- If we don't have a list we return NULL. Ultimately we look through the list to find our appropriate node

ht_insert()

Psuedocode

- Inserts a node ~~at the end of the list~~ for specified ~~MAX~~ word

```
index = get_index(word)
if (list[index] == NULL):
    list[index] = newlist()
```

```
ll_insert(list[index])
```

- If a list does not exist we create one

- Ultimately we use ll_insert() to insert a new node into the linked list

~~MAX~~

int count()

- We wish to return the amount of non-NULL lists ~~non-empty~~, in our list.

Linked List

- ADT is used by hash table & it uses nodes to link together a linked lists,
- It has two sentinel nodes; Head & tail

||-lookup()

- Wish to check if a given bad word is in the linked list

Psudocode

```
for node in nodes:
```

```
    if node.word == word:
```

```
        return true
```

```
return false
```

||-insert()

- Wish to insert a new word in half form at the head of list, In front of head.

Psudocode

```
node.next = head.next
```

```
head.next = head
```

```
node.prev = head
```

```
node.next, prev = node
```

- Special case if it's just Head & tail

```
head.next = node
```

```
tail.prev = node
```

```
node.prev = head
```

```
node.next = tail
```

Node

- ADT used by linked list ADT
- Can tell us what is next, previous, and word/new punctuation

Node.create()

Procedure

node.word = word

node.next = new

Node

Banhammer.c • This file is the main file that uses all modules and ADT's to determine transgressions & whether a user violated them.

- We begin by opening badpeak.txt & newpeak.txt and add them to our Bloom Filter & Hash Table
- ~~We~~ We add them using the API functions bf_insert() & ht_insert()
- Now we parse through a user's message using our parcon.c module. Our regular expression is
`^([a-zA-Z]+[-\s]?[a-zA-Z]+[\s-\s]+)*$`
this insures that we accept letters, dashes, and apostrophes as a word.
- We initialized two new linked lists in which we will save the words that the user typed that is a transgression of bad words or/and bad words w/ translations
- Depending on our lists we will present what transgressions the user committed

Pseudocode

Open(badSpeak.txt)

Open(newSpeak.txt)

bf = BloomFilter()

ht = Hash-table()

for word in badSpeak.txt:

 bf.insert(word)

 ht.insert(word)

for word in newSpeak.txt:

 bf.insert(word)

 ht.insert(word)

LinkedList thought, correction

for word in userInput.txt:

 if (word == bf_probe):

 if (bf_probe(word) == true):

 if (word != newSpeak):

 correction.insert(word)

 else:

 thought.insert(word)

 if (thought != 0 & correction != 0):

 print(message)

 print(thought)

 print(correction)

 if (correction != 0):

 print(message)

 print(correction)

 else:

 print(message)

 print(thought)

Bad Speak.txt

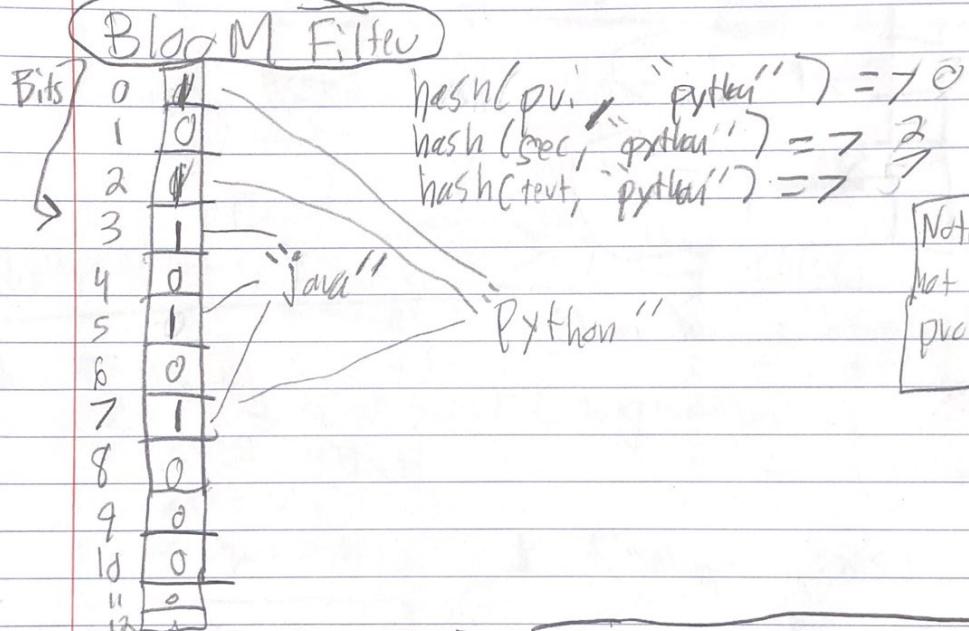
Java
Python
bash
VScode

NewSpeak.txt

covid → plague
C++ → C--
banana → berry

1.7 Make Bloom Filter

- An array of bits w/ 3 hash funcs



Note: 0 means definitely not there, 1 means probably

Probe(bf, "python")
hash(pui, "python") = 0
hash(sec, "python") = 2
hash(text, "python") = 7

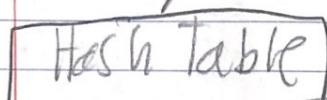
Checking all indices of "python"

is index 0, 2, 7 set? We have most likely seen python before

Probe(bf, "dog")
hash(pui) = 3
hash(sec) = 7
hash(text) = 11

Since index 11 isn't set, we can definitely say "dog" is not in bloom filter.

Also must add to our hash table... but why?
• Array of linked lists...



• We only add once w/ a hash table
salt ~~whh~~

$$\text{hash}(\text{ht} \rightarrow \text{salt}, \text{"Java"}) = 8$$

$$\text{hash}(\text{ht} \rightarrow \text{salt}, \text{"python"}) = 4$$

$$\text{hash}(\text{ht} \rightarrow \text{salt}, \text{"bash"}) = 0$$

$$\text{hash}(\text{ht} \rightarrow \text{salt}, \text{"vscale"}) = 4$$

$$\text{hash}(\text{ht} \rightarrow \text{salt}, \text{"banana"}) = 6$$

Is "Computer" a bad word?

$$\text{hash('C')} = 2$$

$$\text{hash('o')} = 4$$

$$\text{hash('m')} = 9$$

It is in all 3 indices
of B.F., so we check
that.

$$\text{h}(\text{ht} \rightarrow \text{salt}, \text{"Computer"}) = 4$$

Is "Computer" the dispeak word of word "vscale"? Node "python"
? No, it is not. Computer is definitely not a bad word.

Is "banana" a bad word?

$$\text{hash('b')} = .$$

$$\text{hash('a')} = .$$

"banana" matches dispeak and has a new speak "berry"

Phase 1

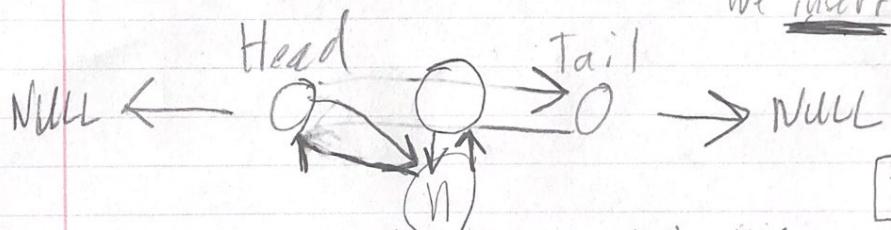
- Add words from badwords.txt to Bloom Filter & Hash Table

- Add words from Newspeak.txt to B.F. & H.T.

Phase 2

- Firewall part... Read user input, for each word query file B.F.
 - If in B.F., check H.T.
 - Else not a bad word

Linked Lists



Say, we add node

"banana"
"berry"

We insert at head.

Insert

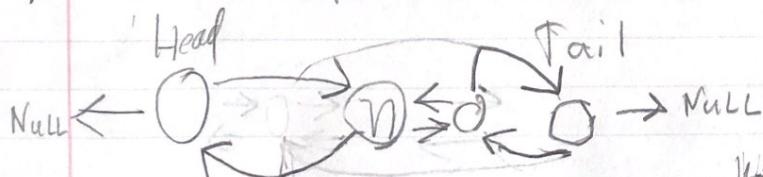
$$\begin{aligned} n \rightarrow \text{next} &= \text{Head} \rightarrow \text{next} \text{ (originally)} \\ n \rightarrow \text{prev} &= \text{Head} \\ \text{Head} \rightarrow \text{next} &= n \\ n \rightarrow \text{next} \rightarrow \text{prev} &= n \end{aligned}$$

Why move to front?

"the"
null

is common and it should be in front for time efficiency

- If move to front is enabled



$$\begin{aligned} n \rightarrow \text{prev} \rightarrow \text{next} &= n \rightarrow \text{next} \\ n \rightarrow \text{next} \rightarrow \text{prev} &= n \rightarrow \text{prev} \end{aligned}$$

$$\begin{aligned} n \rightarrow \text{next} &= \text{head} \rightarrow \text{next} \\ n \rightarrow \text{prev} &= \text{head} \\ \text{head} \rightarrow \text{next} &= n \\ n \rightarrow \text{next} \rightarrow \text{prev} &= n \end{aligned}$$

Basic
expressions
for regex

Regular Expressions	Pattern matching
$+$: one or more
$*$: zero or more
$?$: zero or one (optional)
$ $: OR
$\cdot \cdot \cdot$: for precedence

a^* : matches strings w/ one or more a's...

a^* : matches strings w/ zero or more a's...

"aaaaaa" matched by "a^{*}"

"a^{*}" matched or accepted by "a^{*}"
"a^{*}" does not match/accepted by "a^{*}"

(0|1)⁺; binary (a string of 0's & 1's)

"0|1|0|0|1|1|0|1" accepted by "(0|1)⁺"

[a-z]: matches any symbol b/w a-z inclusive

[a-z]⁺: matches one or more symbols b/w a-z.

"(a-z)⁺ - (a-z)⁺"
hyphens
matches hyphenated words

"(a-z)⁺" (a-z)⁺
contractions
don't
can't

Purpose of Regex is to parse out input

FIVE STAR

FIVE STAR

FIVE STAR

FIVE STAR

Input.txt

Dear class
Today is my penultimate section.

- All different words... we need regular expression to parse these words out
- We use parsing module to do this
- " " would not be a word

"[a-zA-Z]+-[a-zA-Z]+" will match...
(match) pseudo-code
(not match) pseudo

Whatever we read from input.txt needs to be converted to lowercase

Deleting L.L

Head ↓ ↑ ↓ ↓ Tail
∅ - ∅ = ∅ = ∅ = ∅ = ∅

Note Save Note;
for each node:

Save pointer to next
delete-node
go to saved pointer

Begin assignment w/...

1.7 Nodl

2.7 Linked List

3.7 Hash Table (try adding to hash table)

4.7 Bloom Filter

5.7 Add parser

6.7 Messages

Note!

hash table ~~by~~ default size = 10000

B.F. default size is 2^{20}

Note!

If hash table is smaller, we will have
longer ~~int~~ L.L's bc more collisions

If B.F. is smaller we will have more
false-positives bc we have less bits to use which means
~~more~~ words will share bits more often.