

## Task 5 Report: Spelling checking & correction Algorithms

### 1. terminology

$D$  = dict size

$L$  = dict average word len

$Q$  = docu size

$N$  = docu average word len

$M$  = docu average word edit dist

$R_{|R=k}$  = ratio of words with  $\text{dist}=k \in \{0, 1, 2, 3, 3+\}$  for docu

$T$  = duplicate rate of docu words

Possible(len(word), dist)  $(53 * (n + \text{dist}) + 26)$

= estimate possible  $\text{dist}=1$  #words

HashTable = separating chain based hash table

$P_{\text{larger}|\text{dist}=k}$  = possible combination of word(len=n) with  $\text{dist}=k > D$

### 2. Assumption

- (\*) xor hash can perfectly distribute  $n$  keys for  $n$  bucket  $\rightarrow$  insert and compare the bucket for separating chain  $\in O(1)$  i.e.  $E(\text{collision}=k) \in O(1)$
- (\*) assume  $L \equiv N$  for easier stricmp analysis
- assume malloc() free()  $\in O(1)$

### 3. HashTable Operation

- hash by xor\_hashing(key)  $\in O(\text{len}(\text{key}))$
- HashTable has  $\in O(\text{strcmp}(\text{key}, \text{query}) + \text{hash})$  as (\*)
- HashTable put  $\in O(\text{strcmp}(\text{key}, \text{query}) + \text{hash})$  as (\*)
- HashTable get  $\in O(\text{strcmp}(\text{key}, \text{query}) + \text{hash})$  as (\*)
- make HashTable(List) by inserting each word in List  
 $\in O(\text{len}(\text{List}) * \text{HashTable put}(\text{List} \rightarrow \text{word}))$

### 4. Search Operation

- strcmp(dict word, docu word)  $\in O(\text{MIN}(L, N))$
- cal edit dist(key, query)  $\in O(L * N)$   
 $\in O(N^2)$  by (\*)
- linear match word(dict, query) by comparing query with each word in dict  
 $\in O(D * \text{strcmp}(\text{dic} \rightarrow \text{word}, \text{query}))$
- linear search word(dict, query) by calculate edit dist(dict  $\rightarrow$  word, query) and chose one with min edit dist and return first found when multiple choices faces  $\in O(D * N^2)$
- generated dist1(query) find “deleted, replaced, inserted” possible combinaitons  
 $\in O(\text{Possible}(\text{query}))$
- generated dist>1(query, dist=k) based on previous

generated combination

$$\in O(\prod_{k=1}^{k=\text{dist}} \text{Possible}(\text{query}, k - 1))$$

- linear search word for dist(dict, query, dist=k) by comparing query with each word in dict and return when dist=k matched found  $\in O(D * N^2)$

### 5. Small tricks

- when linear search  $|\text{diff}(\text{key}, \text{query})| > 3$  implies edit dist  $> 3$
- when linear search cal edit dist(key, query) for a given dist\_bound, terminate calculation early when cur all elements in calculating row and column(from left top to right bottom layer by layer)  $> \text{dist\_bound}$ .
- store search result to prevent extra time spent on duplicate queries
- assign pointer to HashTable key but not strcpy() to save malloc() and free() time.

### 6. Task 3 analysis:

Approaches:

#### 1) linear search look up:

step1: create HashTable(hist) for storing search history  
 $\in O(1)$

step2: for each word in docu

**Loop( $O(Q)$ )**

(1) if query was not searched before  $\in O(N + \text{MIN}(L, N))$

(2) linear match word  $\in O(D * \text{MIN}(L, N)|_{\text{search once}}^1)$

(3) if this is a new search store {query:match}  
 $\in O(N + \text{MIN}(L, N)|_{\text{search once}})$

**$\therefore$  Total(approach1 average/worst)**

$$\in O(Q * (N + \text{MIN}(L, N) + \text{MIN}(L, N) * D|_{\text{search once}}))$$

$$\in O(Q * D * N) \text{ as } (*)$$

#### 2) hash table look up:

step1: make HashTable(dict) for List(dictionary)

$$\in O(D * (N + \text{MIN}(L, N)))$$

step2: for each word in docu

**Loop( $O(Q)$ )**

(1) HashTable(dict) has word  $\in O(N + \text{MIN}(L, N))$

**$\therefore$  Total(approach2 average/worst)**

$$\in O(D(N + \text{MIN}(L, N)) + Q(N + \text{MIN}(L, N)))$$

$$\in O((D + Q) * N) \text{ as } (*)$$

### 7. Task 3 space analysis

$$\begin{aligned} \text{Space}(\text{approach1}) &= \text{HashTable}(\text{hist}) + \text{List}(\text{dict}) + \\ &\quad \text{List}(\text{docu}) \quad \in O(Q + D) \end{aligned}$$

---

1 search once for duplicate query

$$\text{Space}(\text{approach2}) = \text{HashTable}(\text{dist}) + \text{List}(\text{dist}) + \text{List}(\text{docu}) \in O(Q + D)$$

### 8. Task 3 Comparisons

For space both approaches are the same but approach 2 requires less time complexity.

$\therefore$  prefer approach 2 than 1 as  $O(Q+D) \leq O(Q * D)$

### 9. Task 4 analysis:

#### 1) linear search look up:

step1: create HashTable(hist) for storing search history  
 $\in O(1)$

step2: for each word in docu

**Loop( $O(Q)$ )**

(1) if query was not searched before

$$\in O(N + \text{MIN}(L, N))$$

(2) linear search word

$$\in O(D * N^2_{|\text{search once}|})$$

(3) if this is a new search store {query:match}

$$\in O(N + \text{MIN}(L, N)_{|\text{search once}|})$$

$\therefore$  **Total(approach1)**

$$\in O(Q * (N + \text{MIN}(L, N) + D * N^2_{|\text{search once}|}))$$

$$\in O(Q * D * N^2) \text{ as } (*)$$

$$\text{space} \in O(Q+D)$$

#### 2) linear combined with hash table look up:

step1: make HashTable(dict) for List (dictionary)

$$\in O(D * (N + \text{MIN}(L, N)))$$

step2: create HashTable(hist) for storing search history  
 $\in O(1)$

step3: for each word in docu

**Loop( $O(Q)$ )**

(1) if searched before then **output** and next<sup>2</sup>

$$\in T * O(N + \text{MIN}(L, N))$$

(2) if HashTable(dict) has query then **output** and next

$$\in (1-T) * R_{/R=0} * O(N + \text{MIN}(L, N))$$

(3) if HashTable(dict) has generated dist1(query) then **output** and next

$$\in (1-T) * R_{/R=1} * O(N * (N + \text{MIN}(L, N)))$$

(4) if Possible combination for dist=2 < D<sup>3</sup>

if HashTable(dict) has generated dist>1(query, 2) then

**output** and next

$$\in (1-T) * R_{/R=2} * O(N^2 * (N + \text{MIN}(L, N)))$$

else

linear search word for dist(dict, query, 2) then **output**

$$\text{and next} \in (1-T) * R_{/R=2} * O(D * N^2)$$

(5) repeat previous with dist=3

generate and search:  $\in (1-T) * R_{/R=3} * O(N^3 * (N + \text{MIN}(L, N)))$

$$\text{linear search:} \in (1-T) * R_{/R=3} * O(D * N^2)$$

(5) if this is a new search store {query:match}

$$\in O(N + \text{MIN}(L, N)_{|\text{search once}|})$$

$\therefore$  **Total(approach2)**

$$\in O(D * N + Q * (T * N +$$

$$(1-T) * (R_{/R=0} * N +$$

$$R_{/R=1} * N^2 +$$

$$\text{MIN}((1-P_{\text{larger}/\text{dist}=2}) * R_{/R=2} * N^3, P_{\text{larger}/\text{dist}=2} * D * N^2) +$$

$$\text{MIN}((1-P_{\text{larger}/\text{dist}=3}) * (N^3 + R_{/R=3} * N^4),$$

$$P_{\text{larger}/\text{dist}=3} * (N^3 + D * N^2)))$$

*average: with constants T, P dropped*

$$\in O(D * N + R_0 * Q * N + R_1 * Q * N^2 + R_2 * \text{MIN}(Q * N^3 + R_3 * \text{MIN}(Q * N^3 + Q * N^4, Q * N^3 + D * Q * N^2), D * Q * N^2))$$

$$\text{worst} \in O(D * N + Q * N^3 + D * Q * N^2)$$

*that is dist=2 generated and looked up but no matched and then linear search*

$$\text{space} \in O(Q+D)$$

### 10. Task 4 analysis

If we assume  $O(\text{string}) \in O(1)$  then  $O(\text{approach1}) \in O(QD)$ ,  $O(\text{approach2}) \in O(Q+D)$ . From the above analysis, **Approach2 is better than Approach1 by experimental result and prevent linear search if generating and hash table lookup can save time, given R has a larger proportion for  $k \in \{0, 1, 2\}$  than  $k=3$  or greater than 3.**

jabberwocky.txt	Approach1	Approach2
100K	4s	0.12s
250K	6s	0.3s
1M	8s	1.6s
5M	12s	4.0s

by dimefox user+sys

<sup>2</sup> P(searched before) = T/Q print search result and "continue;" for next docu word

<sup>3</sup>  $(1-P_{\text{larger}/\text{dist}=2})$