

INFORMATION RETRIEVAL AND RECOMMENDER SYSTEMS

---

## Exercise 3: Implementation and indexing

---

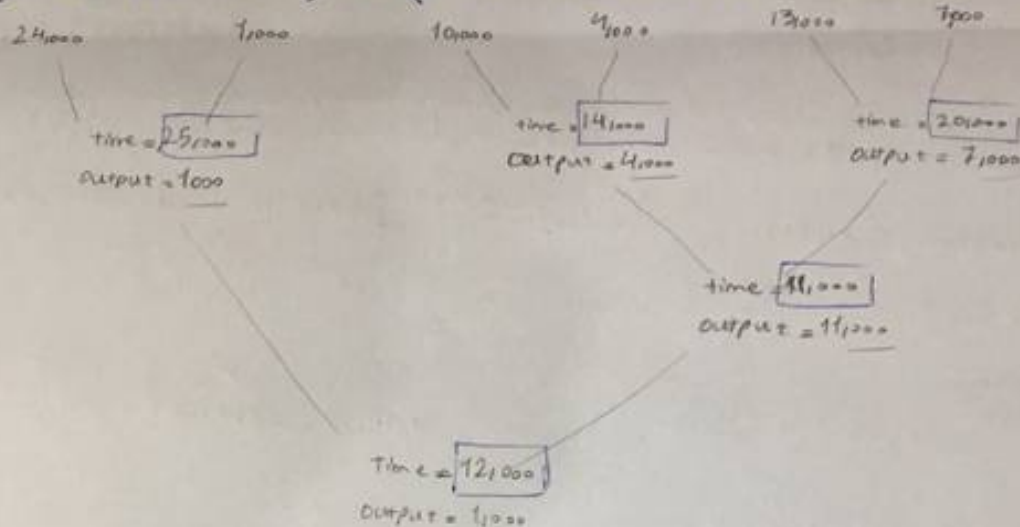
Ashwin Kumar Gururajan, Mohana Fathollahi

13th October, 2022

# Exercise 4:

Main:

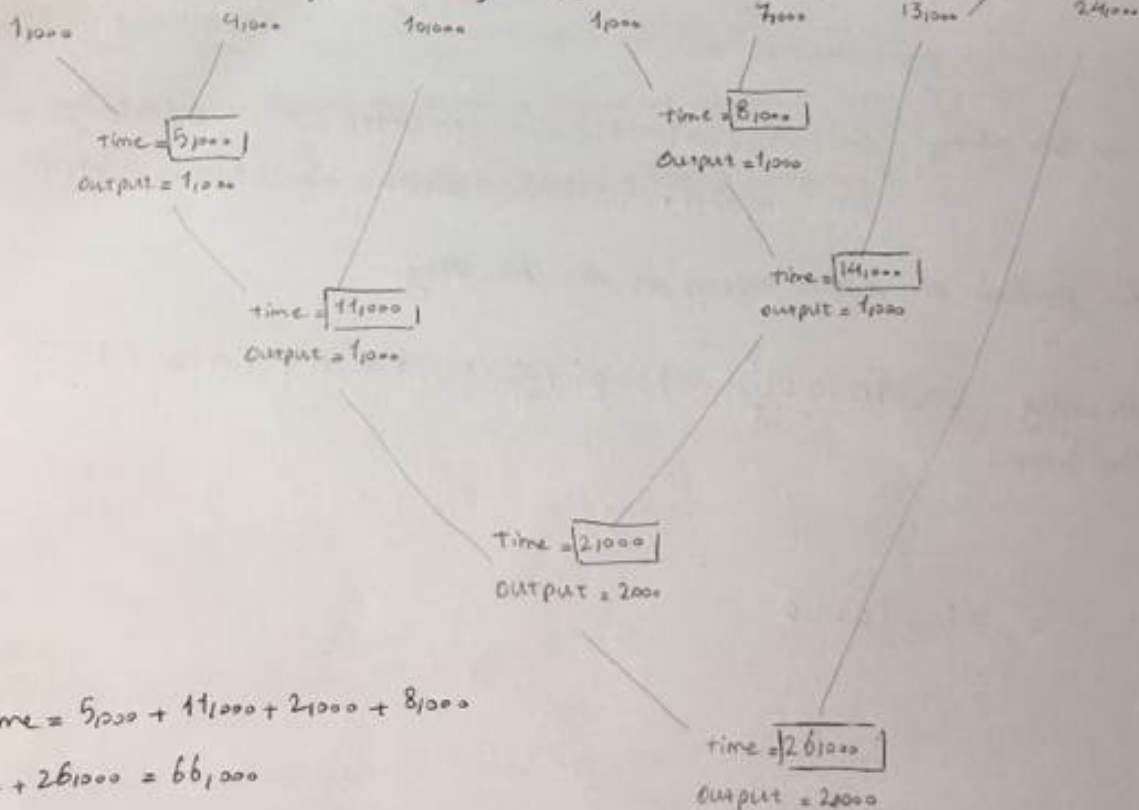
(Charles AND Dicken) AND ((Leon AND Tolstoi) OR (Anton AND chejov))



$$\text{Total Time} = 25,000 + 14,000 + 20,000 + 11,000 + 12,000 = 82,000$$

suggestion:

((Dickens AND Tolstoi) AND Leon) OR ((Dicken AND chejov) AND Anton) AND Charles



$$\begin{aligned} \text{Total time} &= 5,000 + 11,000 + 2,000 + 8,000 \\ &+ 14,000 + 26,000 = 66,000 \end{aligned}$$

To have smaller size of output, it is better to consider terms that have lower frequency inside of boolean operations and put terms that have higher frequency outside of boolean operations. In this case Leon/Anton/Charles have highest frequency.

# Exercise 5: -1.

[10, 1, 15, 3, 22, 2, 23, 4, 34, 1, 44, 1, 50, 2, 58, 8, 90, 1, 101, 1, 112, 2]

1. Gap Compression  $[(id_1, f_1), \dots, (id_k, f_k)] \xrightarrow{\text{compress}} [(id_1, f_1), \dots, (id_k - id_{k-1}, f_k)]$

[10, 1, 5, 3, 7, 2, 1, 4, 11, 1, 10, 1, 6, 2, 8, 8, 32, 1, 11, 1, 11, 2]

2. Elias' Gamma

[0001010, 1, 00101, 3, 00111, 2, 1, 4, 001011, 1, 001010, 1, 00110, 2]

[0001000, 8, 00000100000, 7, 0001011, 1, 0001011, 2]

3. Compress frequency, using unary self-delimiting

[0001010, 0, 00101, 110, 00111, 10, 1, 1110, 0001011, 0, 0001010, 0,

00110, 10, 0001000, 1111110, 00000100000, 0, 0001011, 1, 0001011, 10]

Final bit string: 00010100001011100011110111100001011000010100001101010  
00010001111110000001000000000010111000101110

2- perform the inverse process on the bit string

1. Decoding Elias' Gamma

$n=4$  00001010110 |  $n=2$  001000 | 1010 | 001000 | 10 | 00110110 |

[21, 10] [4, 0] [1, 0] [1, 0] [4, 0] [6, 110]

0100 | 0110

$n=1$  0100 | 0110

[2, 0] [3, 0]

2.

Decoding gap compression

Result of step 1:  $[21, 10, 4, 0, 1, 0, 1, 0, 4, 0, 1, 0, 6, 110, 2, 0, 3, 0]$

$[21, 40, 25, 0, 26, 0, 27, 0, 31, 0, 32, 0, 38, 110, 40, 0, 43, 0]$

3. Decoding unary self-delimiting

$[21, 2, 25, 1, 26, 1, 27, 1, 31, 1, 32, 1, 38, 3, 40, 1, 43, 1]$

Exercise 6:

$((A \text{ and } B) \text{ and } C) \text{ and } D) \text{ and } E$

10,000      20,000      40,000      80,000      120,000

$$10,000 \times \log_2 20,000 + 10,000 \times \log_2 40,000 + 10,000 \times \log_2 80,000 + 10,000 \times \log_2 120,000$$

$$= 10,000 \times (14,29 + 15,29 + 16,29 + 16,87) = 627,400$$

$A \text{ and } (B \text{ and } (C \text{ and } (D \text{ and } E)))$

20,000      40,000      80,000      120,000

$$80,000 \log_2 120,000 + 40,000 \log_2 80,000 + 20,000 \log_2 40,000 + 10,000 \log_2 20,000$$

$$= 80,000 \times 16,87 + 40,000 \times 16,29 + 20,000 \times 15,29 + 10,000 \times 14,29$$

$$= 2,449,900$$

$((A \text{ and } B) \text{ or } (C \text{ and } D)) \text{ or } (E \text{ and } F)$

10,000      40,000      120,000

$$10,000 \log_2 20,000 + 40,000 \log_2 80,000 + 150,000 \log_2 120,000 + (10,000 + 40,000) \times$$

$$(150,000 + 120,000) = 10,000 \times 14,29 + 40,000 \times 16,29 + 150,000 \times 16,87 + 50,000 + 170,000$$

$$= 3,545,100$$

$(A \text{ and } B) \text{ or } ((C \text{ and } D) \text{ or } (E \text{ and } F))$

10,000      40,000      120,000

$$10,000 \log_2 20,000 + 40,000 \log_2 80,000 + 120,000 \log_2 150,000 + (40,000 + 120,000) \times$$

$$10,000) = 3,320,000$$



<sup>10,000</sup>  
(A and E) or <sup>20,000</sup>  
(B and E)

$$10,000 \log_2 120,000 + 20,000 \log_2 120,000 + (10,000 + 20,000) = 540,000$$

(A and B) or E

$$10,000 \log_2 20,000 + (10,000 + 120,000) = 280,000$$

2. Assuming mutual independence between occurrences

For this part we used formula for independence probability.

$$P(A \cap B) = P(A) \cdot P(B)$$

And to get expected length of the list, we multiply the above probability by the total number of words.

1

2

1

1)  $((A \text{ and } B) \text{ and } C) \text{ and } D) \text{ and } E$

10000 \ / 20000

# = 30000

output = 200 =  $10000 * (20000 / 1000000)$

|

40000

# 40200

output = 8

|

# 80008

output = 1

|

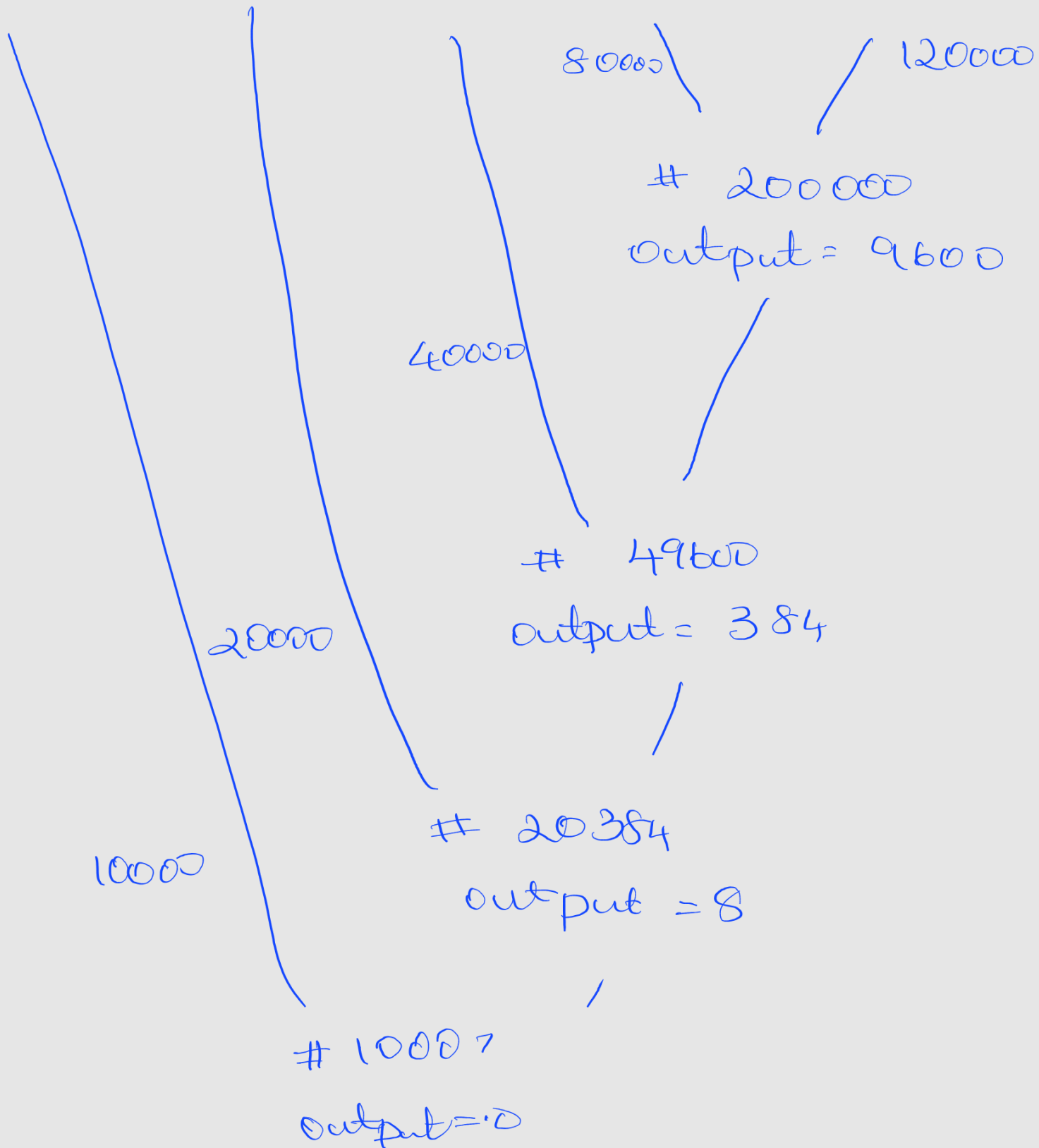
120000

# 120001

output = 0

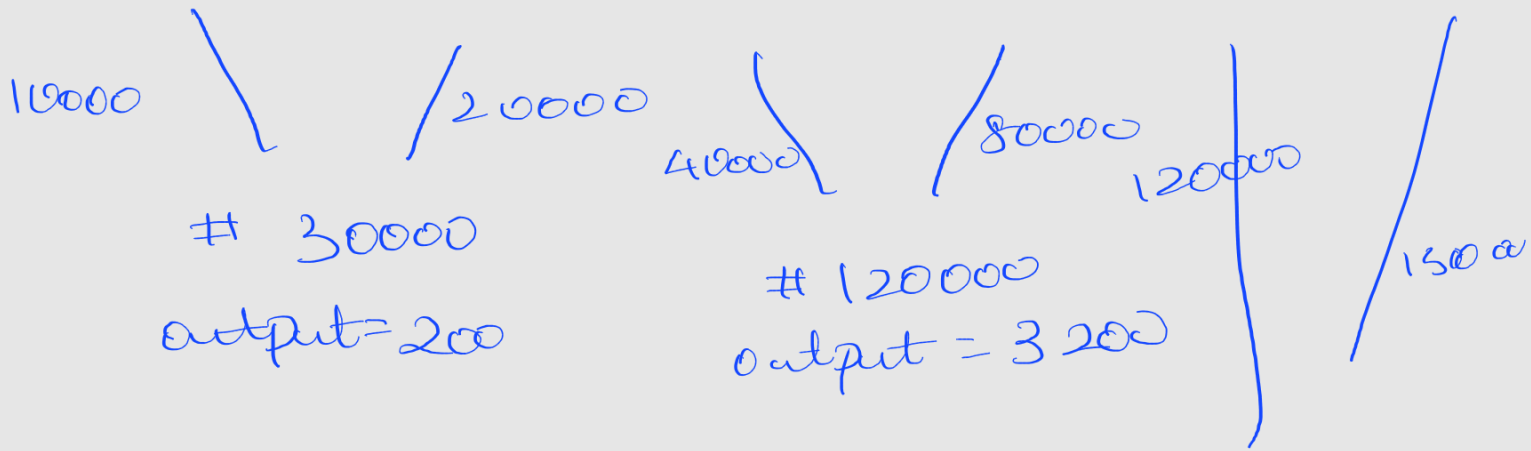
# of comparisons =  $30000 + 40200 + 80008 + 120001 = 270209$

2) A and (B and (C and (D and E)))



$$\begin{aligned} \text{no. of comparisons} &= 200000 + 49600 + 20384 + 10007 \\ &= 279991 \end{aligned}$$

3) ((A and B) or (C and D)) or (E and F)



$$\begin{aligned} \# & 3400 \\ \text{output} &= 200 + 3200 - 1 \\ &= 3399 \end{aligned}$$

$$\# 273399$$

$$\text{output} = 3399 + 18000$$

$$\begin{aligned} \# &= 30000 + 120000 + \\ & 270000 + 3400 + \\ & 273399 = 696799 \end{aligned}$$

--(U but it'll be small)

$$21399$$



4. (A and B) or (C and D) or (E and F))

10000 \	/ 20000	40K \	/ 80K	120K \	/ 150K
# 30000		# 120000		# 270000	
output = 200		output = 3200		output = 18000	

# 21200  
 output = 3200 + 18000  
 - 58  
 = 21142

# = 21342

output = 200 + 21142  
 - 5  
 = 21337

# = 30000 + 120000 + 270000  
 + 21200 + 21342  
 = 462542 //

5. (A and E) or (B and E)

10K \      / 120K      20K \      / 120K

# 130K

output = 1200

# 40K

output = 2400

\      /

# = 3600

output = 3600 - 3

= 3797

# = 130000 + 140000 + 3600

= 273600 //

6) CA and B) or E

10K

20K

# = 30K

output = 200

120K

# 120200

$$\# = 30000 + 120200$$

$$= 150200 //$$