
Overlap Set Similarity Join

Dong Deng

Problem Definition

- *Input:*
 - a collection of sets R
 - a constant integer threshold c
 - *Output:*
 - all pairs $(X, Y) \in R \times R$ s.t. $|X \cap Y| \geq c$
-

Example

- *Input:*

- R

- $c = 2$

- *Output:*

- all set pairs with overlap size no smaller than c

$$|R_1 \cap R_3| = 2, \quad |R_1 \cap R_2| = 2, \quad |R_2 \cap R_3| = 3, \quad |R_4 \cap R_5| = 4, \quad |R_6 \cap R_7| = 8$$

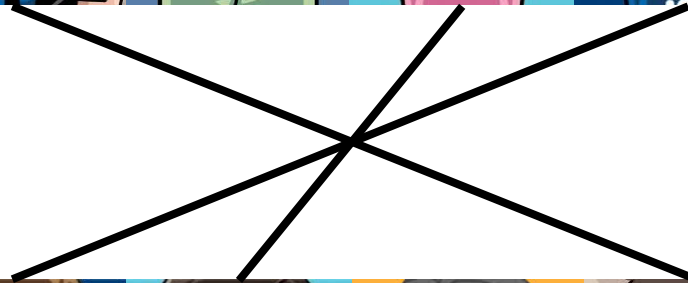
<i>id</i>	set
R_1	$\{e_1, e_2, e_3\}$
R_2	$\{e_1, e_3, e_4, e_7\}$
R_3	$\{e_1, e_3, e_5, e_7\}$
R_4	$\{e_2, e_4, e_5, e_6\}$
R_5	$\{e_2, e_4, e_5, e_6, e_8, e_9, e_{10}, e_{11}\}$
R_6	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}\}$
R_7	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}, e_{19}\}$

Application - Friend Recommendation

User 1



User 2



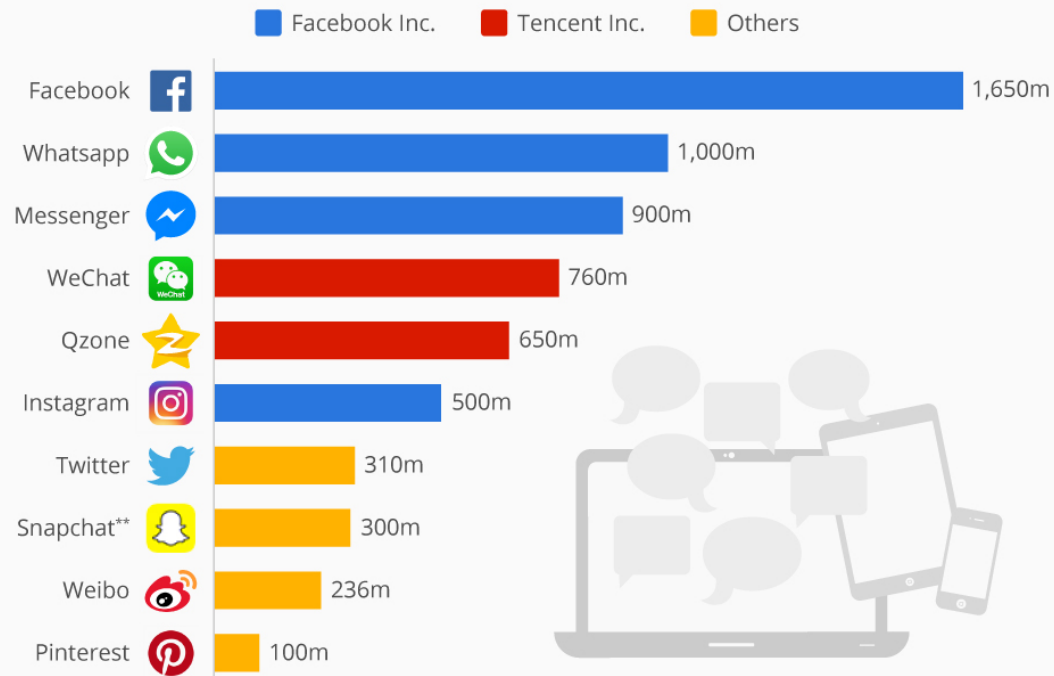
More Applications

- Data Management
 - Data Integration and Cleaning
 - Keyword Subscription
 - Data Mining
 - Frequent Pattern Mining
 - Recommendation
 - Computer Vision
 - Scene Reconstruction
 - Machine Learning
 - Large Entries Retrieval in Matrix Productions
 - Non-negative Matrix Factorization
 - Singular Value Decomposition
-

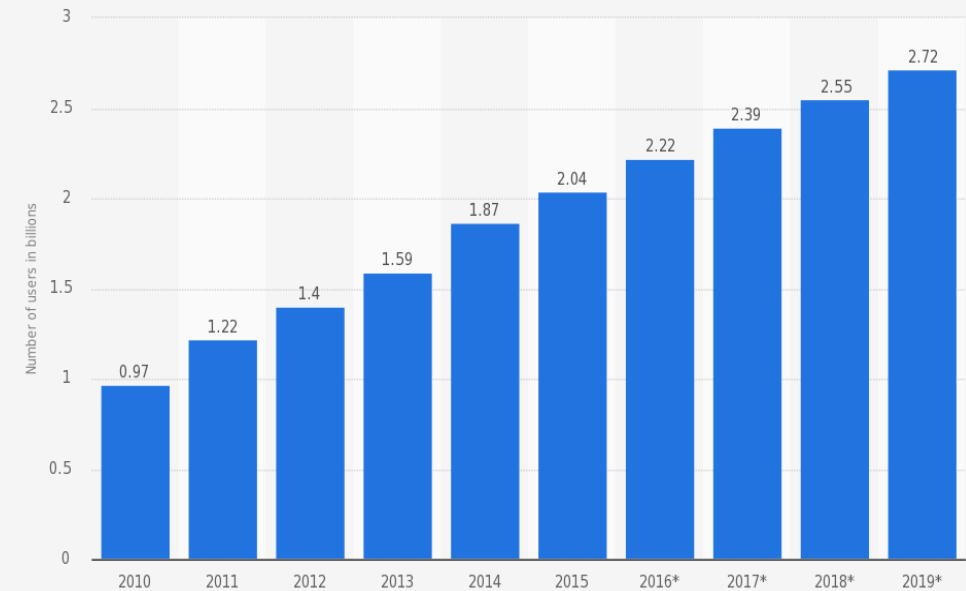
Challenge

Facebook Inc. Dominates the Social Media Landscape

Monthly active users of selected social networks and messaging services*



Number of social network users worldwide from 2010 to 2019 (in billions)



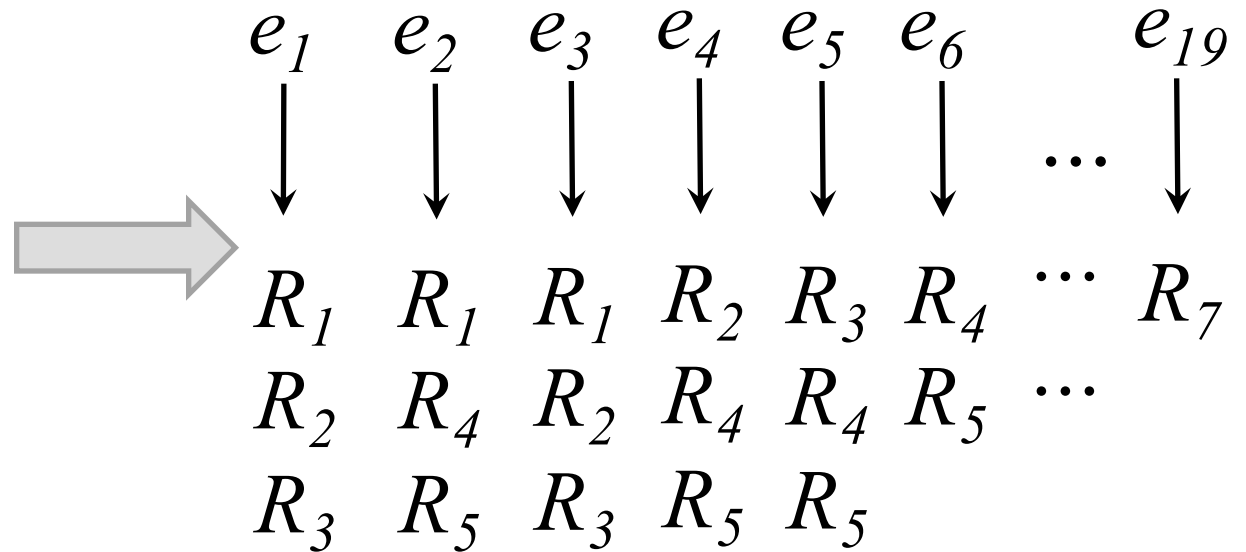
Source:
eMarketer
© Statista 2016

Additional Information:
Worldwide, eMarketer, 2010 to 2015

Solution?

Naïve Method: ScanCount

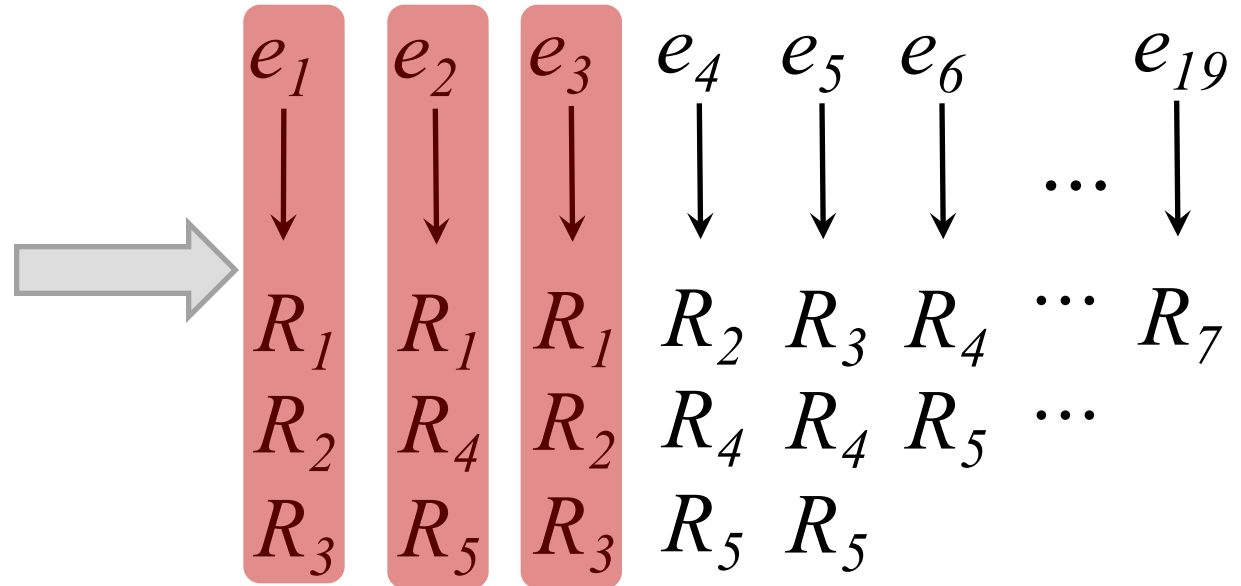
<i>id</i>	set
R_1	$\{e_1, e_2, e_3\}$
R_2	$\{e_1, e_3, e_4, e_7\}$
R_3	$\{e_1, e_3, e_5, e_7\}$
R_4	$\{e_2, e_4, e_5, e_6\}$
R_5	$\{e_2, e_4, e_5, e_6, e_8, e_9, e_{10}, e_{11}\}$
R_6	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}\}$
R_7	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}, e_{19}\}$



Step 1: build an inverted index

Naïve Method: ScanCount

id	set
R_1	$\{e_1, e_2, e_3\}$
R_2	$\{e_1, e_3, e_4, e_7\}$
R_3	$\{e_1, e_3, e_5, e_7\}$
R_4	$\{e_2, e_4, e_5, e_6\}$
R_5	$\{e_2, e_4, e_5, e_6, e_8, e_9, e_{10}, e_{11}\}$
R_6	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}\}$
R_7	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}, e_{19}\}$



Step 1: build an inverted index

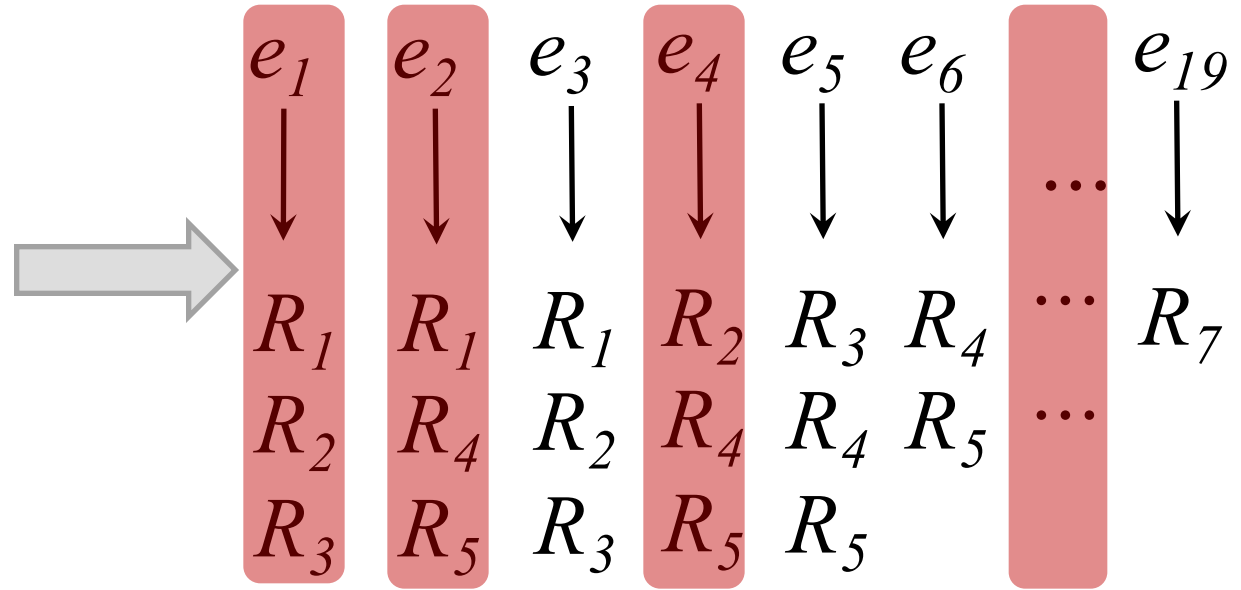
Step 2: scan each set and count

Count: $R_2: 2$ ✓ $R_3: 2$ ✓ $R_4: 1$ ✗ $R_5: 1$ ✗

Results: (R_1, R_2) (R_1, R_3)

Naïve Method: ScanCount

id	set
R_1	$\{e_1, e_2, e_3\}$
R_2	$\{e_1, e_3, e_4, e_7\}$
R_3	$\{e_1, e_3, e_5, e_7\}$
R_4	$\{e_2, e_4, e_5, e_6\}$
R_5	$\{e_2, e_4, e_5, e_6, e_8, e_9, e_{10}, e_{11}\}$
R_6	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}\}$
R_7	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}, e_{19}\}$



Step 1: build an inverted index

Step 2: scan each set and count

Count: $R_3: 2$ ✓ $R_4: 2$ ✗ $R_5: 1$ ✗

Results: (R_2, R_3)

Analysis

How many times is this list scanned?

$$e \longrightarrow R_1 R_2 R_3 \dots R_m$$

is scanned m times, each time takes $O(m)$; in total $O(m^2)$

Analysis

*Some very frequent elements yield excessive long inverted lists
the long inverted lists are scanned many times.*

$$e \longrightarrow R_1 R_2 R_3 \dots R_m$$

for $m=1$ million, ScanCount takes ~ 1 trillion operations

Naïve Method: Subset Enumeration

<i>id</i>	set
R_1	$\{e_1, e_2, e_3\}$
R_2	$\{e_1, e_3, e_4, e_7\}$
R_3	$\{e_1, e_3, e_5, e_7\}$
R_4	$\{e_2, e_4, e_5, e_6\}$
R_5	$\{e_2, e_4, e_5, e_6, e_8, e_9, e_{10}, e_{11}\}$
R_6	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}\}$
R_7	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}, e_{19}\}$

threshold $c = 2$

e_1e_2 e_1e_3 e_2e_3
 e_1e_3 e_1e_4 e_1e_7 e_3e_4 e_3e_7 e_4e_7
 \vdots
 \vdots
 \vdots

(R_1, R_2)

Step 1. for each set, enumerate all subsets of size c (c -subset for short)

Step 2. output all set pairs sharing a common c -subset

Analysis

how many c-subsets are generated from each set?

$$R = \{e_1 \ e_2 \ e_3 \ \dots \ e_m\}$$

$\binom{m}{c}$ *c-subsets.*

Analysis

some very large sets yield a huge number of c-subsets.

$$R = \{e_1 \ e_2 \ e_3 \ \dots\dots e_m \}$$

for $m=1000$ and $c=3$, it generates 166 million!

DossJoin: Combination of two methods

Small Sets

R_1	$\{e_1, e_2, e_3\}$
R_2	$\{e_1, e_3, e_4, e_7\}$
R_3	$\{e_1, e_3, e_5, e_7\}$
R_4	$\{e_2, e_4, e_5, e_6\}$



R_1	$\{e_1, e_2, e_3\}$
R_2	$\{e_1, e_3, e_4, e_7\}$
R_3	$\{e_1, e_3, e_5, e_7\}$
R_4	$\{e_2, e_4, e_5, e_6\}$

Naïve Algorithm 1 – Subset Enumeration

size boundary: x - - - - -

Large Sets

R_5	$\{e_2, e_4, e_5, e_6, e_8, e_9, e_{10}, e_{11}\}$
R_6	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}\}$
R_7	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}, e_{19}\}$



R_1	$\{e_1, e_2, e_3\}$
R_2	$\{e_1, e_3, e_4, e_7\}$
R_3	$\{e_1, e_3, e_5, e_7\}$
R_4	$\{e_2, e_4, e_5, e_6\}$
R_5	$\{e_2, e_4, e_5, e_6, e_8, e_9, e_{10}, e_{11}\}$
R_6	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}\}$
R_7	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}, e_{19}\}$

Naïve Algorithm 2 – ScanCount

DossJoin: Combination of two methods

Small Sets

R_1	$\{e_1, e_2, e_3\}$
R_2	$\{e_1, e_3, e_4, e_7\}$
R_3	$\{e_1, e_3, e_5, e_7\}$
R_4	$\{e_2, e_4, e_5, e_6\}$

bounded by x

size boundary: x

Large Sets

R_5	$\{e_2, e_4, e_5, e_6, e_8, e_9, e_{10}, e_{11}\}$
R_6	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}\}$
R_7	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}, e_{19}\}$

bounded by $\frac{n}{x}$ Why?

1. each large set has at least x elements
 2. there are n elements in total
-

Time Complexity Analysis for Large Sets

size boundary: x -----

Large Sets

R_5	$\{e_2, e_4, e_5, e_6, e_8, e_9, e_{10}, e_{11}\}$
R_6	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}\}$
R_7	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}, e_{19}\}$



R_1	$\{e_1, e_2, e_3\}$
R_2	$\{e_1, e_3, e_4, e_7\}$
R_3	$\{e_1, e_3, e_5, e_7\}$
R_4	$\{e_2, e_4, e_5, e_6\}$
R_5	$\{e_2, e_4, e_5, e_6, e_8, e_9, e_{10}, e_{11}\}$
R_6	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}\}$
R_7	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}, e_{19}\}$

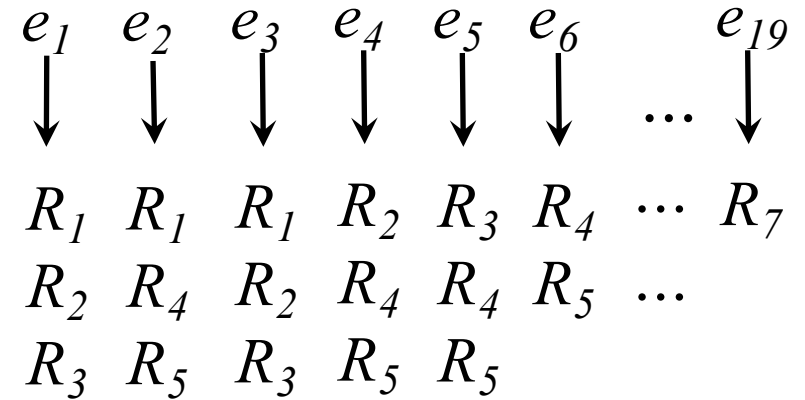
Naïve Algorithm 2 – ScanCount

Time Complexity Analysis for Large Sets

size boundary: x -----

Large Sets

R_5	$\{e_2, e_4, e_5, e_6, e_8, e_9, e_{10}, e_{11}\}$
R_6	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}\}$
R_7	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}, e_{19}\}$



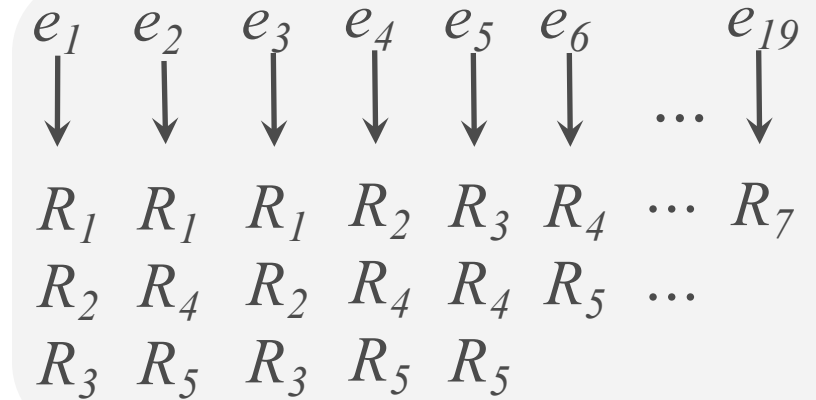
inverted index on all sets

Time Complexity Analysis for Large Sets

size boundary: x -----

Large Sets

R_5	$\{e_2, e_4, e_5, e_6, e_8, e_9, e_{10}, e_{11}\}$
R_6	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}\}$
R_7	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}, e_{19}\}$



each set takes $O(n)$ time

Why?

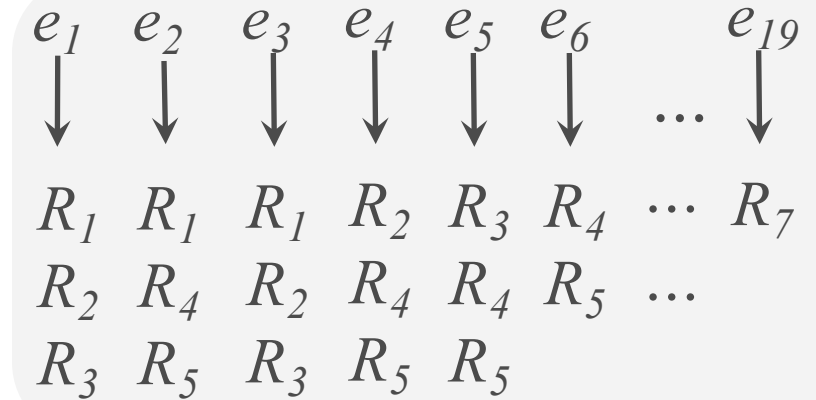
Time Complexity Analysis for Large Sets

size boundary: x

Large Sets

R_5	$\{e_2, e_4, e_5, e_6, e_8, e_9, e_{10}, e_{11}\}$
R_6	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}\}$
R_7	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}, e_{19}\}$

at most $\frac{n}{x}$ large sets as discussed



each set takes $O(n)$ time

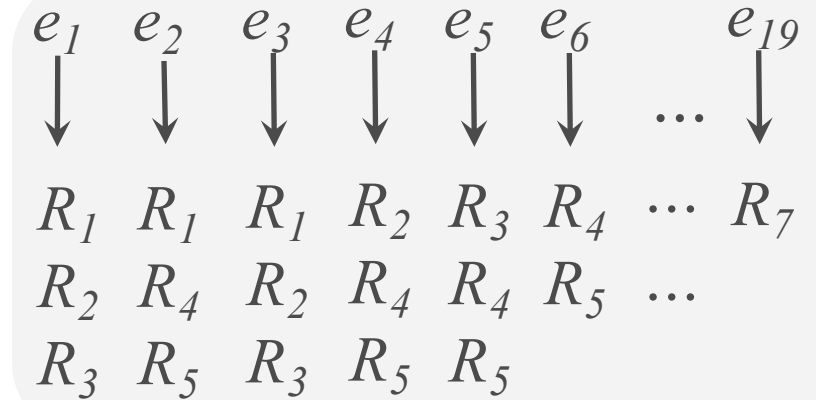
Time Complexity Analysis for Large Sets

size boundary: x

Large Sets

R_5	$\{e_2, e_4, e_5, e_6, e_8, e_9, e_{10}, e_{11}\}$
R_6	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}\}$
R_7	$\{e_{11}, e_{12}, e_{13}, e_{14}, e_{15}, e_{16}, e_{17}, e_{18}, e_{19}\}$

at most $\frac{n}{x}$ large sets as discussed



each set takes $O(n)$ time

Time Complexity: $O\left(\frac{n^2}{x}\right)$

Time Complexity Analysis for Small Sets

Step 1: enumerate c-subsets, takes $O(\text{total \# of c-subsets})$

R_1	$\{e_1, e_2, e_3\}$	$\longrightarrow \leq \binom{ R_1 }{c}$
R_2	$\{e_1, e_3, e_4, e_7\}$	
	\dots	
R_3	$\{e_1, e_3, e_5, e_7\}$	
R_4	$\{e_2, e_4, e_5, e_6\}$	$\longrightarrow \leq \binom{ R_4 }{c}$

Time Complexity Analysis for Small Sets

Step 1: enumerate c-subsets, takes $O(\text{total \# of c-subsets})$

R_1	$\{e_1, e_2, e_3\}$	\longrightarrow	$\leq \binom{ R_1 }{c} \leq R_1 ^c$
R_2	$\{e_1, e_3, e_4, e_7\}$		
		\dots	\dots
R_3	$\{e_1, e_3, e_5, e_7\}$		
R_4	$\{e_2, e_4, e_5, e_6\}$	\longrightarrow	$\leq \binom{ R_4 }{c} \leq R_4 ^c$

Time Complexity Analysis for Small Sets

Step 1: enumerate c-subsets, takes $O(\text{total \# of c-subsets})$

R_1	$\{e_1, e_2, e_3\}$	\longrightarrow	$\leq \binom{ R_1 }{c} \leq R_1 ^c \leq R_1 * x^{c-1}$
R_2	$\{e_1, e_3, e_4, e_7\}$		
	\dots	\dots	
R_3	$\{e_1, e_3, e_5, e_7\}$		
R_4	$\{e_2, e_4, e_5, e_6\}$	\longrightarrow	$\leq \binom{ R_4 }{c} \leq R_4 ^c \leq R_4 * x^{c-1}$

$$\Sigma |R_{\text{small}}| \leq n$$

Time Complexity Analysis for Small Sets

Step 1: enumerate c-subsets, takes $O(\text{total \# of c-subsets})$

R_1	$\{e_1, e_2, e_3\}$
R_2	$\{e_1, e_3, e_4, e_7\}$
R_3	$\{e_1, e_3, e_5, e_7\}$
R_4	$\{e_2, e_4, e_5, e_6\}$

$$\begin{aligned}
 &\longrightarrow \leq \binom{|R_1|}{c} \leq |R_1|^c \leq |R_1| * x^{c-1} \\
 &\qquad \qquad \qquad \dots \qquad \dots \\
 &\longrightarrow \leq \binom{|R_4|}{c} \leq |R_4|^c \leq |R_4| * x^{c-1}
 \end{aligned}
 \left. \vphantom{\begin{aligned} &\longrightarrow \leq \binom{|R_1|}{c} \leq |R_1|^c \leq |R_1| * x^{c-1} \\ &\qquad \qquad \qquad \dots \qquad \dots \\ &\longrightarrow \leq \binom{|R_4|}{c} \leq |R_4|^c \leq |R_4| * x^{c-1} \end{aligned}} \right\} \leq nx^{c-1}$$

$$\Sigma |R_{\text{small}}| \leq n$$

$$\text{total \# of c-subsets} \leq nx^{c-1}$$

Time Complexity Analysis for Small Sets

Step 2: output all set pairs sharing a c-subset

R_1	$\{e_1, e_2, e_3\}$
R_2	$\{e_1, e_3, e_4, e_7\}$
R_3	$\{e_1, e_3, e_5, e_7\}$
R_4	$\{e_2, e_4, e_5, e_6\}$



build an inverted index for all c-subsets

for each inverted list, output all pairs of sets

	e_1e_2	e_1e_3	e_1e_4	e_1e_5	e_1e_7	e_2e_3	e_2e_4	e_2e_5	e_2e_6	e_3e_4	e_3e_5	e_3e_7	e_4e_5	e_4e_6	e_4e_7	e_5e_6	e_5e_7
R_1	■	■				■											
R_2		■	■		■					■		■			■		
R_3		■		■	■						■	■					■
R_4							■	■	■				■	■		■	

Time Complexity Analysis for Small Sets

Step 2: output all set pairs sharing a c -subset

	e_1e_2	e_1e_3	e_1e_4	e_1e_5	e_1e_7	e_2e_3	e_2e_4	e_2e_5	e_2e_6	e_3e_4	e_3e_5	e_3e_7	e_4e_5	e_4e_6	e_4e_7	e_5e_6	e_5e_7
R_1	■	■				■											
R_2		■	■		■					■		■			■		
R_3		■		■	■						■	■					■
R_4							■	■	■				■	■		■	

what's the total length of all inverted lists?

exactly the total # of c -subsets $\leq nx^{c-1}$

Time Complexity Analysis for Small Sets

Step 2: output all set pairs sharing a c-subset

	e_1e_2	e_1e_3	e_1e_4	e_1e_5	e_1e_7	e_2e_3	e_2e_4	e_2e_5	e_2e_6	e_3e_4	e_3e_5	e_3e_7	e_4e_5	e_4e_6	e_4e_7	e_5e_6	e_5e_7
R_1	■	■				■											
R_2		■	■		■					■		■			■		
R_3		■		■	■						■	■					■
R_4							■	■	■				■	■		■	

} $\leq \sqrt{k}$

what's the maximum length of any inverted list?

$$\text{number of results} = L^2 \leq k \qquad L \leq \sqrt{k}$$

Time Complexity Analysis for Small Sets

Step 2: output all set pairs sharing a c-subset

	e_1e_2	e_1e_3	e_1e_4	e_1e_5	e_1e_7	e_2e_3	e_2e_4	e_2e_5	e_2e_6	e_3e_4	e_3e_5	e_3e_7	e_4e_5	e_4e_6	e_4e_7	e_5e_6	e_5e_7
R_1	■	■				■											
R_2		■	■		■					■		■			■		
R_3		■		■	■						■	■					■
R_4							■	■	■				■	■		■	

L_1	...	L_5	L_{17}
↓		↓				↓
$\leq L_1 * L_1 $		$\leq L_5 * L_5 $		maximum length of any inverted list		$\leq L_{17} * L_{17} $
$\leq L_1 * \sqrt{k}$		$\leq L_5 * \sqrt{k}$		total length of all inverted lists		$\leq L_{17} * \sqrt{k}$

$$\leq nx^{c-1}\sqrt{k}$$

Time Complexity Analysis

overall time complexity: $\mathcal{O}\left(\frac{n^2}{x} + x^{c-1}n\sqrt{k}\right)$

let $x = (n/\sqrt{k})^{1/c}$

the time complexity is $\mathcal{O}\left(n^{2-\frac{1}{c}}k^{\frac{1}{2c}}\right) = o(n^2) + O(k)$

why?

case 1: $k = o(n^2)$

case 2: $k = O(n^2)$

Two Practical Problems

- Needs to enumerate all c -subsets, whose number can be huge in practice.
 - The size boundary x set by the theoretical analysis is not practical, as it overestimates the cost for small sets a lot.
-

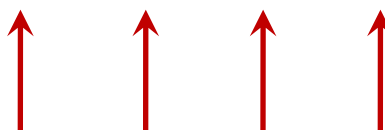
Skip Unnecessary c-subsets

	e_1e_2	e_1e_3	e_1e_4	e_1e_5	e_1e_7	e_2e_3	e_2e_4	e_2e_5	e_2e_6	e_3e_4	e_3e_5	e_3e_7	e_4e_5	e_4e_6	e_4e_7	e_5e_6	e_5e_7
R_1	■	■				■											
R_2		■	■		■					■		■			■		
R_3		■		■	■						■	■					■
R_4							■	■	■				■	■		■	

Observation 1: *Unique c-subsets* cannot generate any result and we can skip them

Skip Unnecessary c-subsets

	e_1e_2	e_1e_3	e_1e_4	e_1e_5	e_1e_7	e_2e_3	e_2e_4	e_2e_5	e_2e_6	e_3e_4	e_3e_5	e_3e_7	e_4e_5	e_4e_6	e_4e_7	e_5e_6	e_5e_7
R_1	■	■				■											
R_2		■	■		■					■		■			■		
R_3		■		■	■						■	■					■
R_4			⏟				■	■	■				■	■		■	



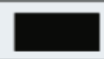

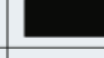

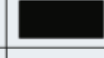
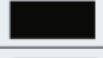
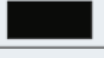

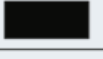
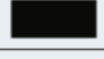
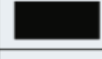

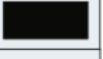








Observation 2: *Redundant c-subsets* only generate duplicate results and we can skip them

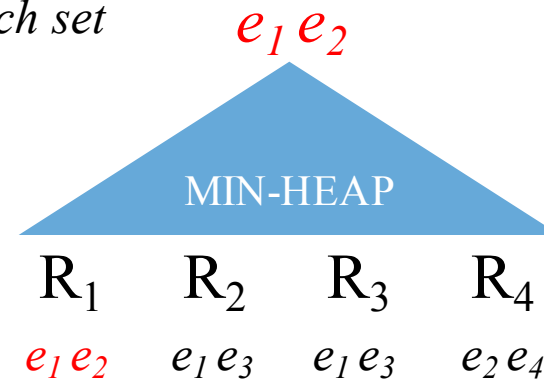
How to skip the unnecessary c-subsets?

Heap-based Method

→ Global Order

	e_1e_2	e_1e_3	e_1e_4	e_1e_5	e_1e_7	e_2e_3	e_2e_4	e_2e_5	e_2e_6	e_3e_4	e_3e_5	e_3e_7	e_4e_5	e_4e_6	e_4e_7	e_5e_6	e_5e_7
R_1																	
R_2																	
R_3																	
R_4																	

 *min-subset: the smallest unvisited c-subset in each set*



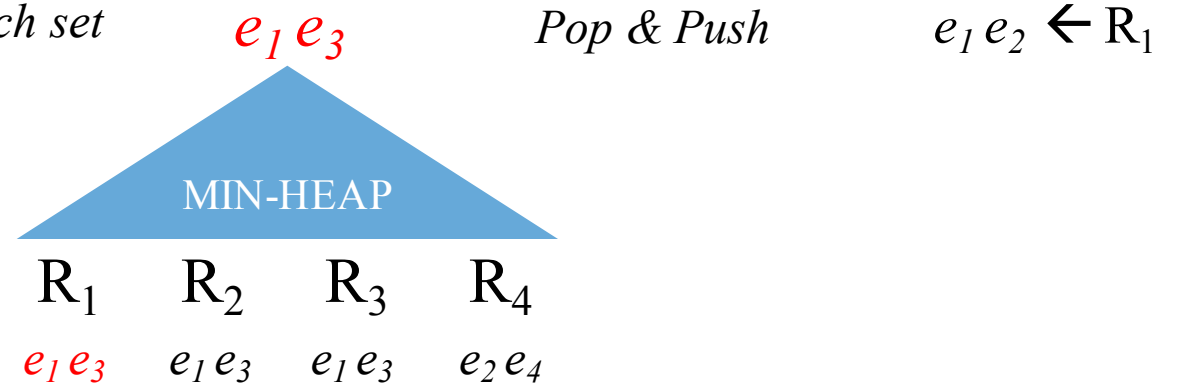
Goal: generate all the c-subsets in the global order, one by one

Heap-based Method

Global Order

	e_1e_2	e_1e_3	e_1e_4	e_1e_5	e_1e_7	e_2e_3	e_2e_4	e_2e_5	e_2e_6	e_3e_4	e_3e_5	e_3e_7	e_4e_5	e_4e_6	e_4e_7	e_5e_6	e_5e_7
R_1	■	■				■											
R_2		■	■		■					■		■			■		
R_3		■		■	■						■	■					■
R_4							■	■	■				■	■		■	

min-subset: the smallest unvisited c-subset in each set

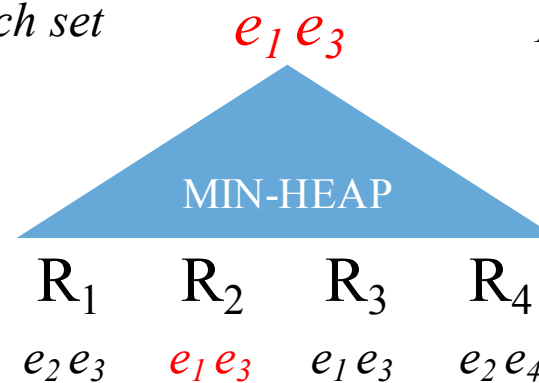


Heap-based Method

Global Order

	e_1e_2	e_1e_3	e_1e_4	e_1e_5	e_1e_7	e_2e_3	e_2e_4	e_2e_5	e_2e_6	e_3e_4	e_3e_5	e_3e_7	e_4e_5	e_4e_6	e_4e_7	e_5e_6	e_5e_7
R_1	■	■				■											
R_2		■	■		■					■		■			■		
R_3		■		■	■						■	■					■
R_4						■	■	■					■	■		■	

min-subset: the smallest unvisited c-subset in each set



Pop & Push

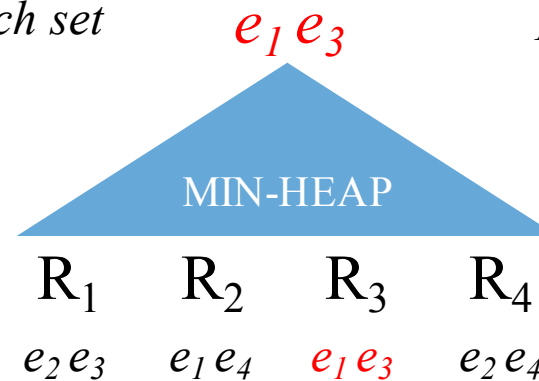
$e_1e_2 \leftarrow R_1$
 $e_1e_3 \leftarrow R_1$

Heap-based Method

Global Order

	e_1e_2	e_1e_3	e_1e_4	e_1e_5	e_1e_7	e_2e_3	e_2e_4	e_2e_5	e_2e_6	e_3e_4	e_3e_5	e_3e_7	e_4e_5	e_4e_6	e_4e_7	e_5e_6	e_5e_7
R_1	■	■				■											
R_2		■	■		■					■		■			■		
R_3		■		■	■						■	■					■
R_4							■	■	■				■	■		■	

min-subset: the smallest unvisited c-subset in each set



Pop & Push

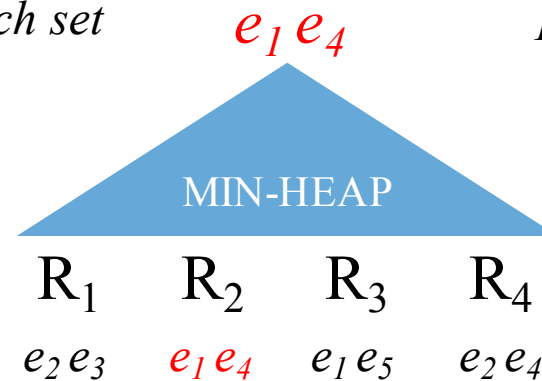
$e_1e_2 \leftarrow R_1$
 $e_1e_3 \leftarrow R_1$
 $e_1e_3 \leftarrow R_2$

Heap-based Method

→ *Global Order*

	e_1e_2	e_1e_3	e_1e_4	e_1e_5	e_1e_7	e_2e_3	e_2e_4	e_2e_5	e_2e_6	e_3e_4	e_3e_5	e_3e_7	e_4e_5	e_4e_6	e_4e_7	e_5e_6	e_5e_7
R_1	■	■				■											
R_2		■	■		■					■		■			■		
R_3		■		■	■						■	■					■
R_4							■	■	■				■	■		■	

○ *min-subset: the smallest unvisited c-subset in each set*



Pop & Push

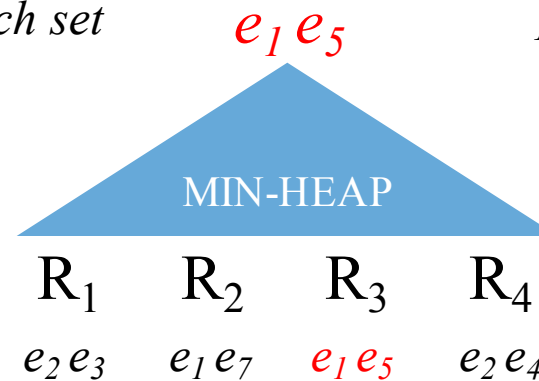
$e_1e_2 \leftarrow R_1$
 $e_1e_3 \leftarrow R_1$
 $e_1e_3 \leftarrow R_2$
 $e_1e_3 \leftarrow R_3$

Heap-based Method

→ *Global Order*

	e_1e_2	e_1e_3	e_1e_4	e_1e_5	e_1e_7	e_2e_3	e_2e_4	e_2e_5	e_2e_6	e_3e_4	e_3e_5	e_3e_7	e_4e_5	e_4e_6	e_4e_7	e_5e_6	e_5e_7
R_1	■	■				■											
R_2		■	■		■					■		■			■		
R_3		■		■	■						■	■					■
R_4							■	■	■				■	■		■	

○ *min-subset: the smallest unvisited c-subset in each set*

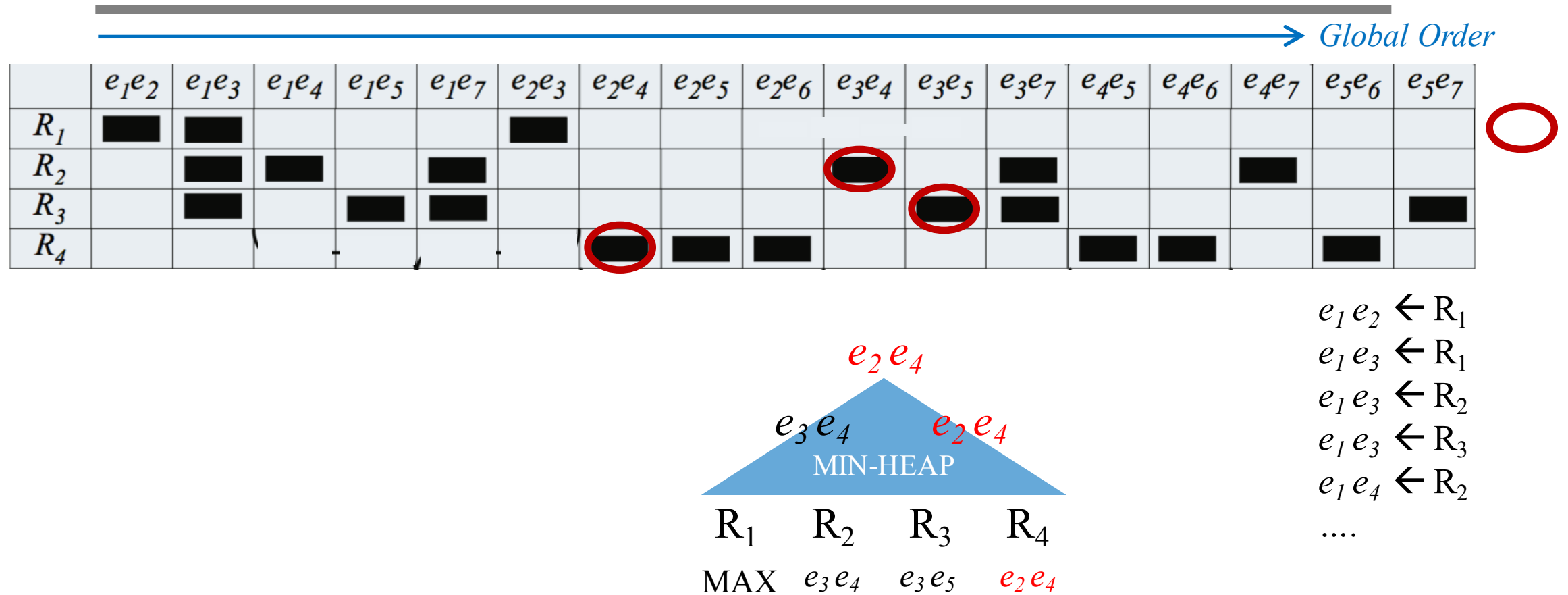


Pop & Push

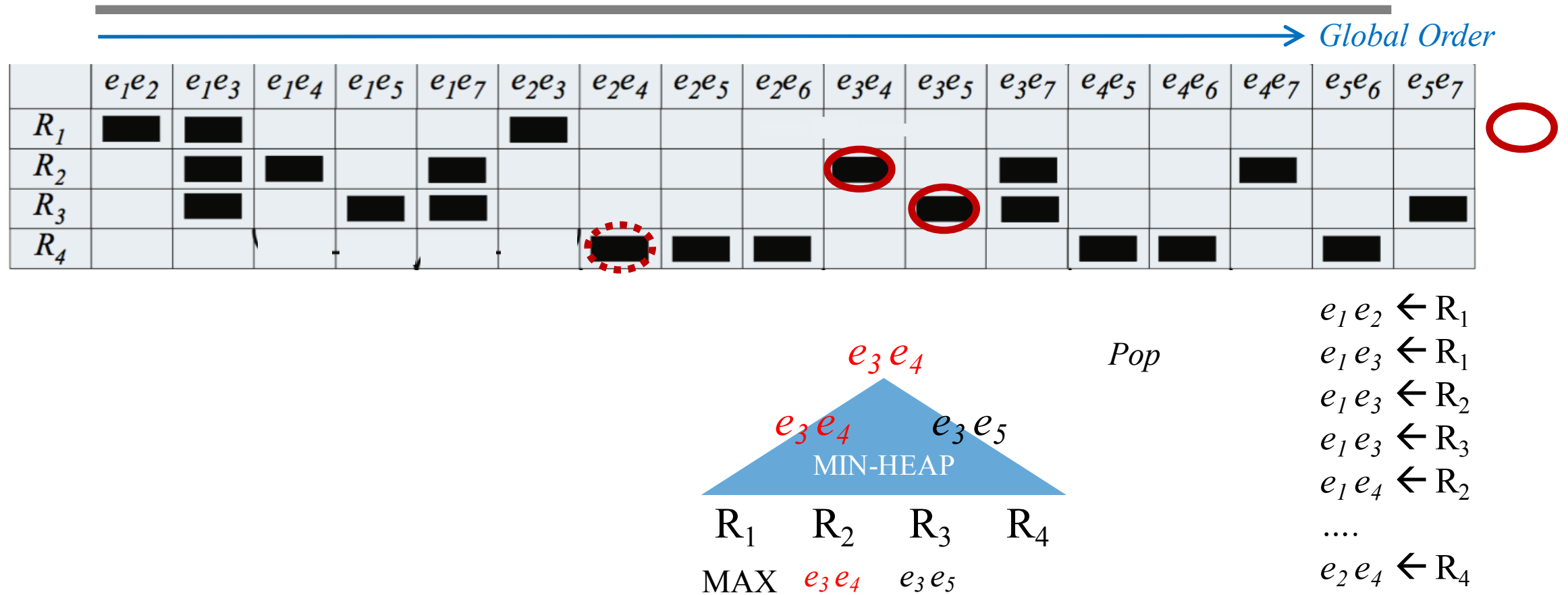
$e_1e_2 \leftarrow R_1$
 $e_1e_3 \leftarrow R_1$
 $e_1e_3 \leftarrow R_2$
 $e_1e_3 \leftarrow R_3$
 $e_1e_4 \leftarrow R_2$

Skip the Unique c-subsets using heap

Heap-based Method



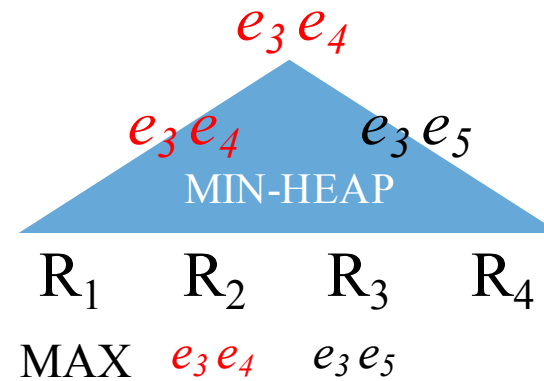
Heap-based Method



Heap-based Method

This is the c-subset on top of the heap

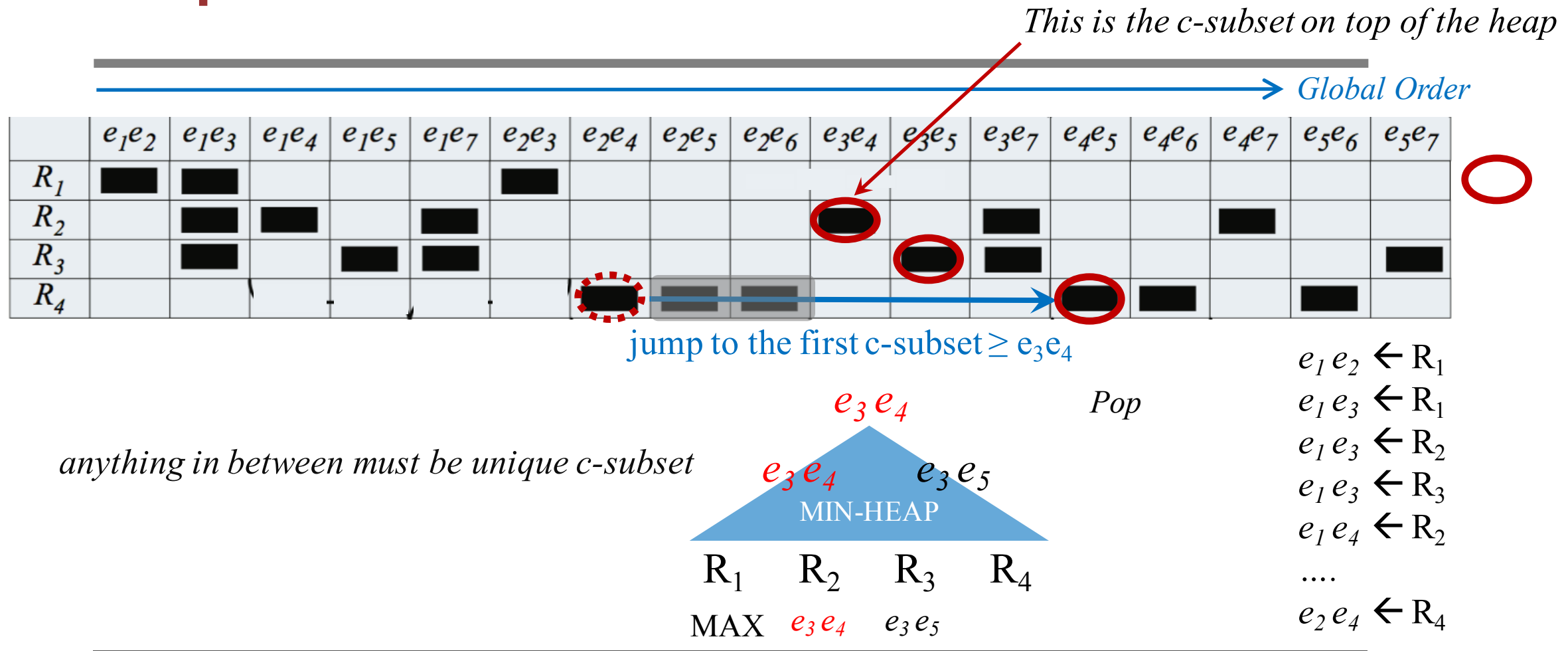
	e_1e_2	e_1e_3	e_1e_4	e_1e_5	e_1e_7	e_2e_3	e_2e_4	e_2e_5	e_2e_6	e_3e_4	e_3e_5	e_3e_7	e_4e_5	e_4e_6	e_4e_7	e_5e_6	e_5e_7
R_1	■	■				■											
R_2		■	■		■					■		■			■		
R_3		■		■	■					■	■	■					■
R_4							■	■	■				■	■		■	



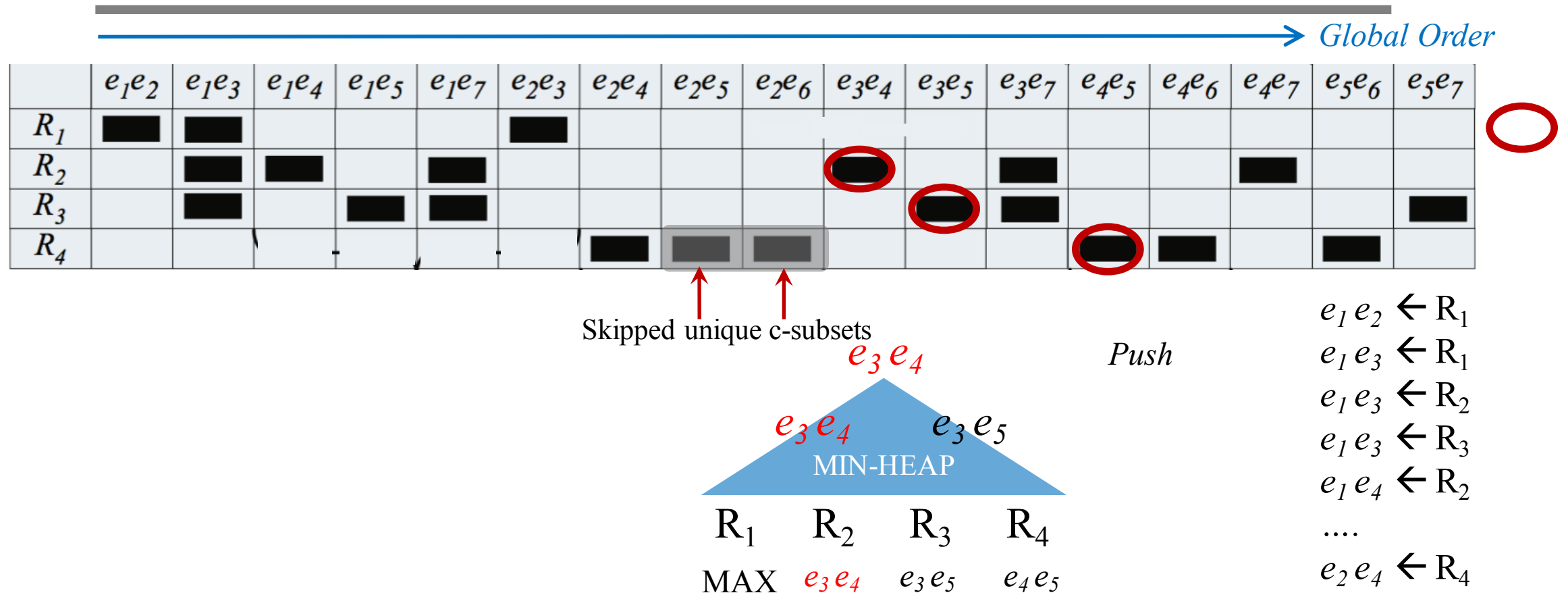
Pop

$e_1e_2 \leftarrow R_1$
 $e_1e_3 \leftarrow R_1$
 $e_1e_3 \leftarrow R_2$
 $e_1e_3 \leftarrow R_3$
 $e_1e_4 \leftarrow R_2$
 \dots
 $e_2e_4 \leftarrow R_4$

Heap-based Method



Heap-based Method



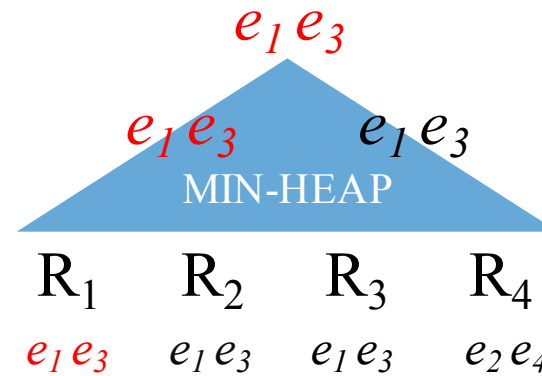
Skip the redundant c-subsets using heap

Heap-based Method

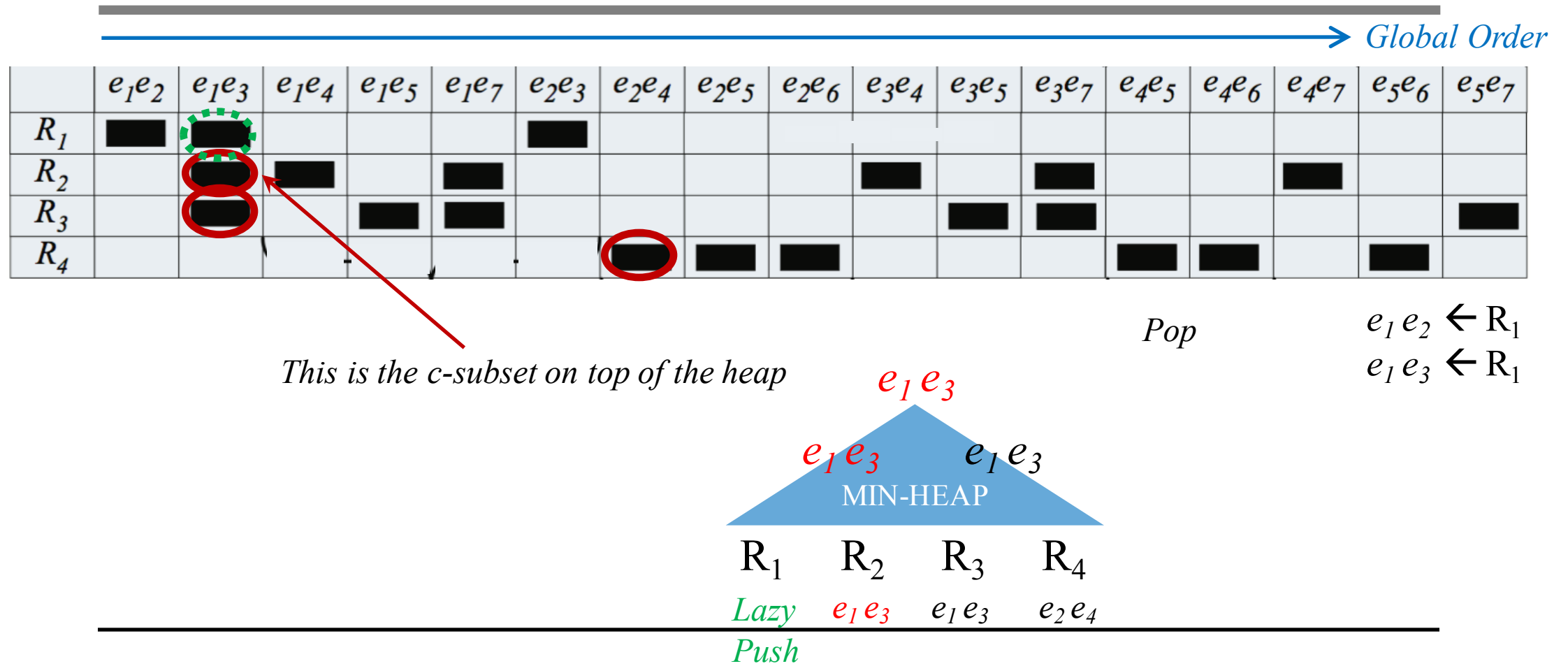
→ Global Order

	e_1e_2	e_1e_3	e_1e_4	e_1e_5	e_1e_7	e_2e_3	e_2e_4	e_2e_5	e_2e_6	e_3e_4	e_3e_5	e_3e_7	e_4e_5	e_4e_6	e_4e_7	e_5e_6	e_5e_7
R_1	■	■				■											
R_2		■	■		■					■		■			■		
R_3		■		■	■						■	■					■
R_4							■	■	■				■	■		■	

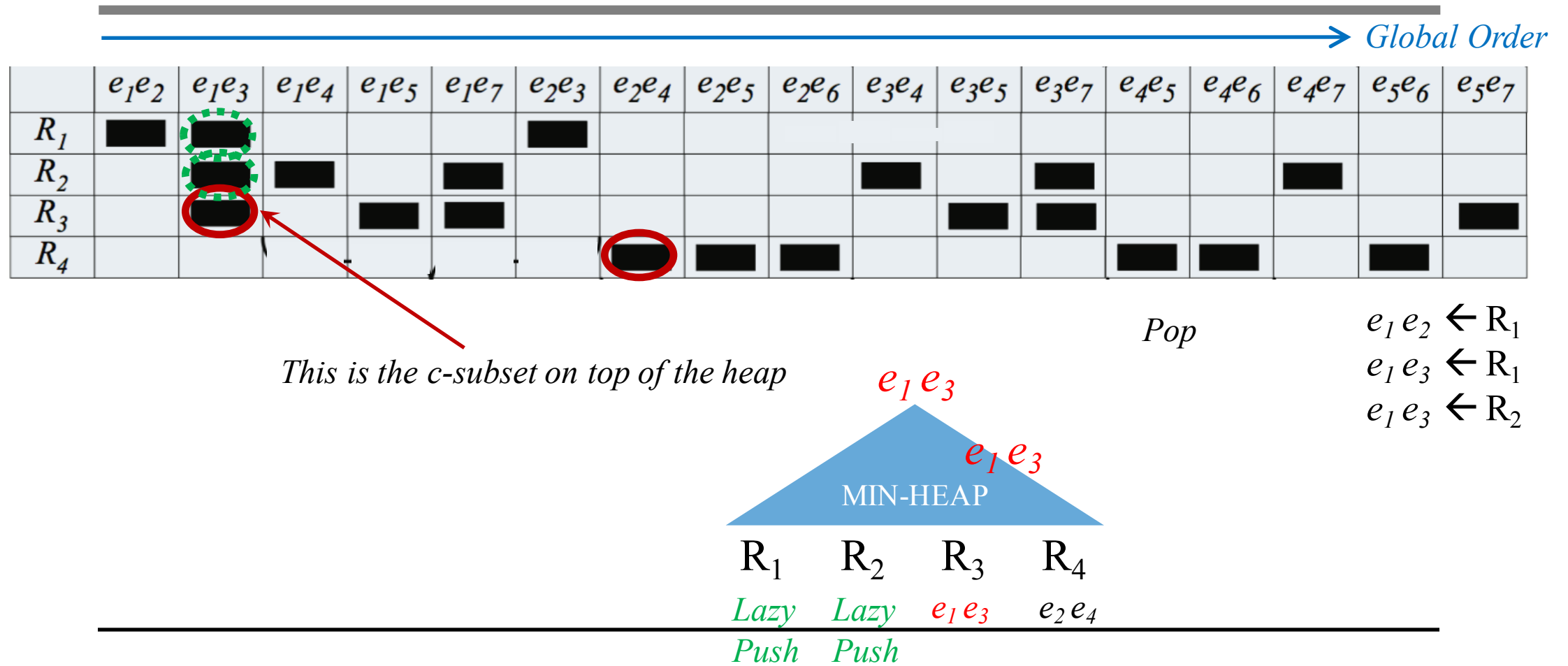
$e_1e_2 \leftarrow R_1$



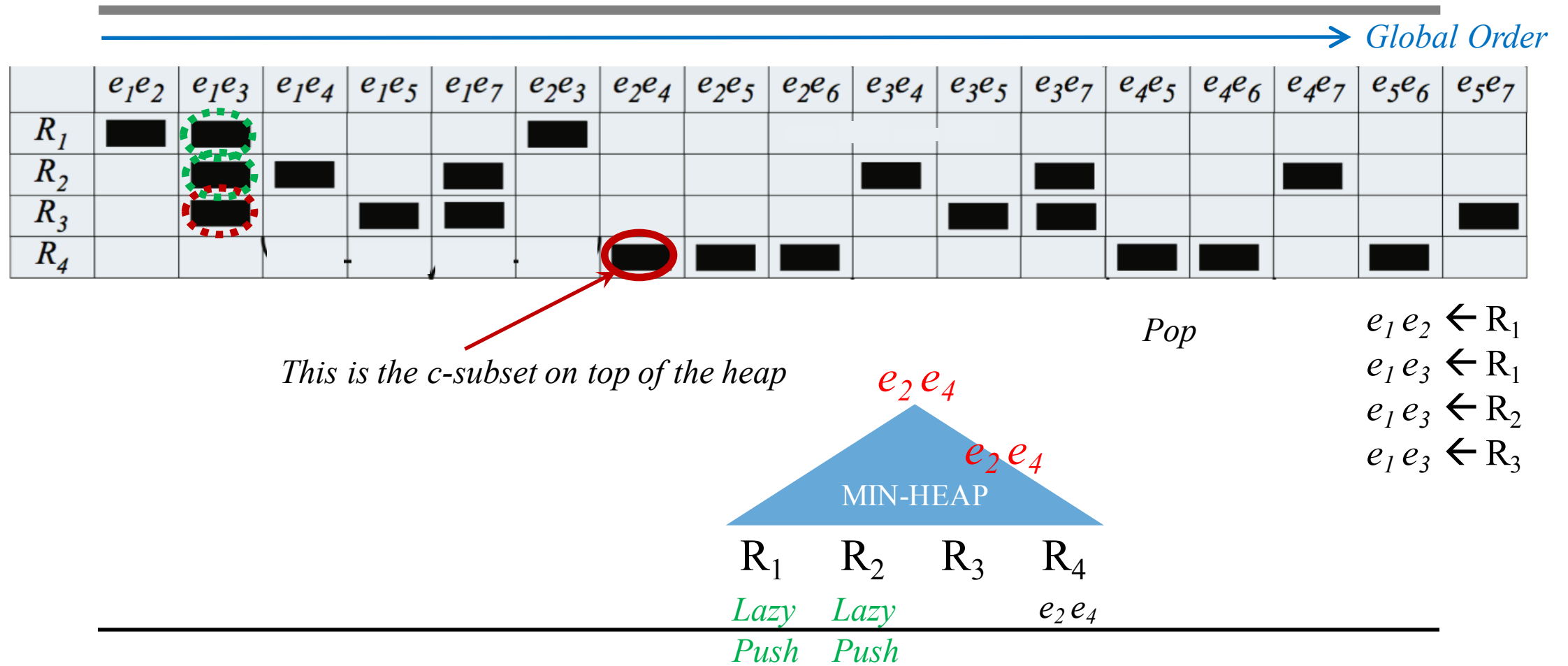
Heap-based Method



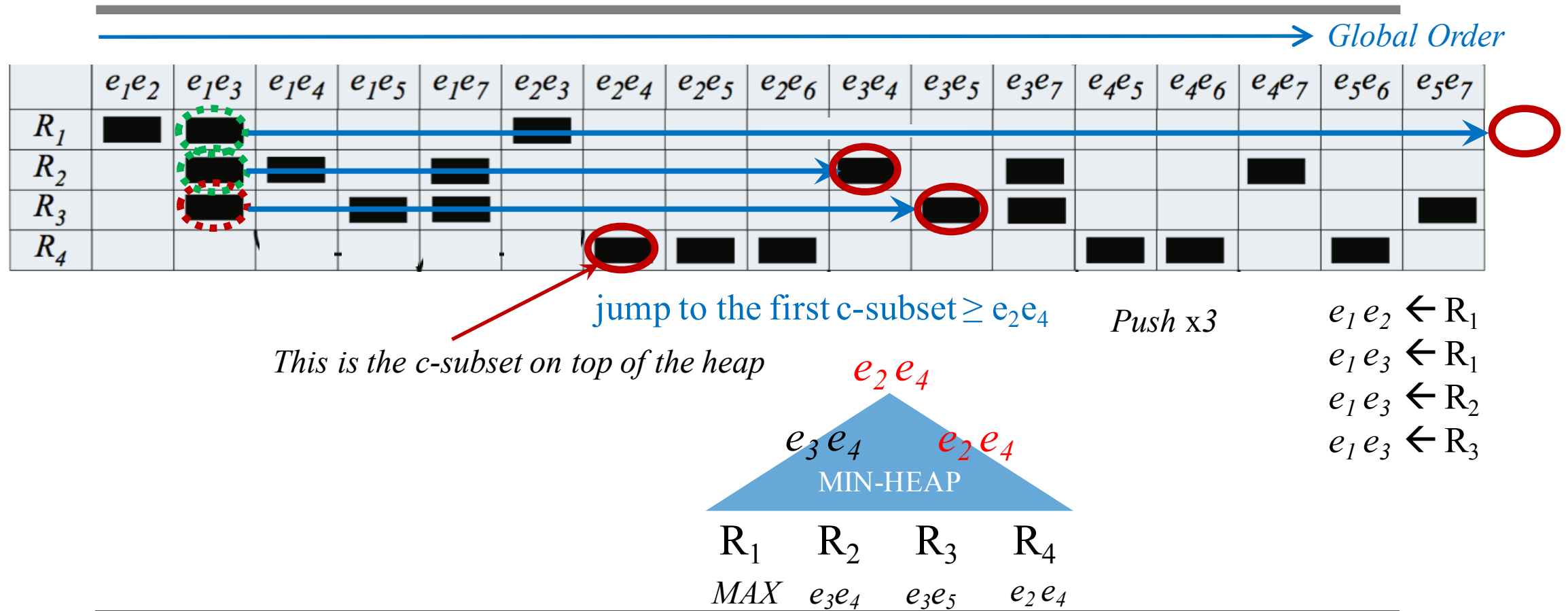
Heap-based Method



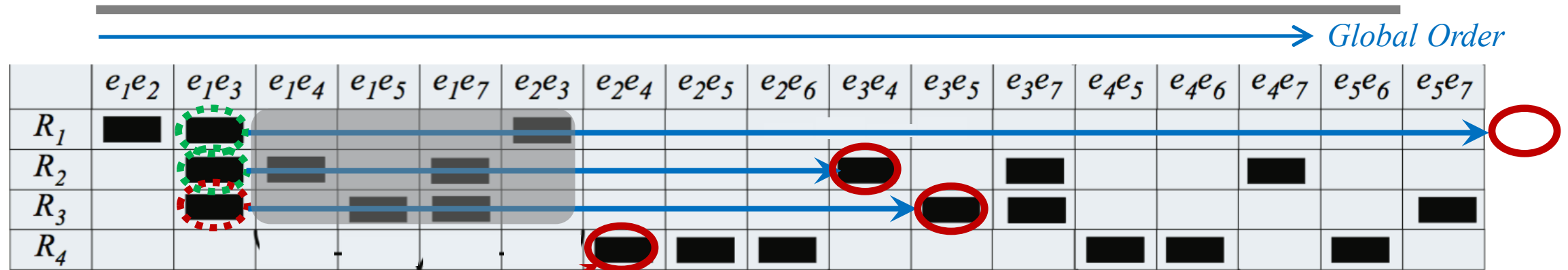
Heap-based Method



Heap-based Method



Heap-based Method



jump to the first c-subset $\geq e_2e_4$

Push x3

$e_1e_2 \leftarrow R_1$

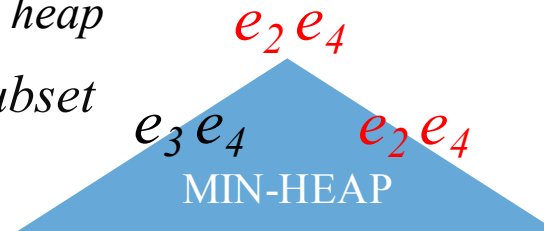
$e_1e_3 \leftarrow R_1$

$e_1e_3 \leftarrow R_2$

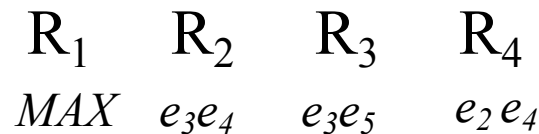
$e_1e_3 \leftarrow R_3$

This is the c-subset on top of the heap

anything in between must be redundant c-subset



Easy push when c-subset on top remains the same



How to reduce the heap-adjustment cost?

each heap adjustment takes $\log_2 m$ time, where m is the number of small c-subsets

Heap Blocking

block by the minimum element in the c -subsets

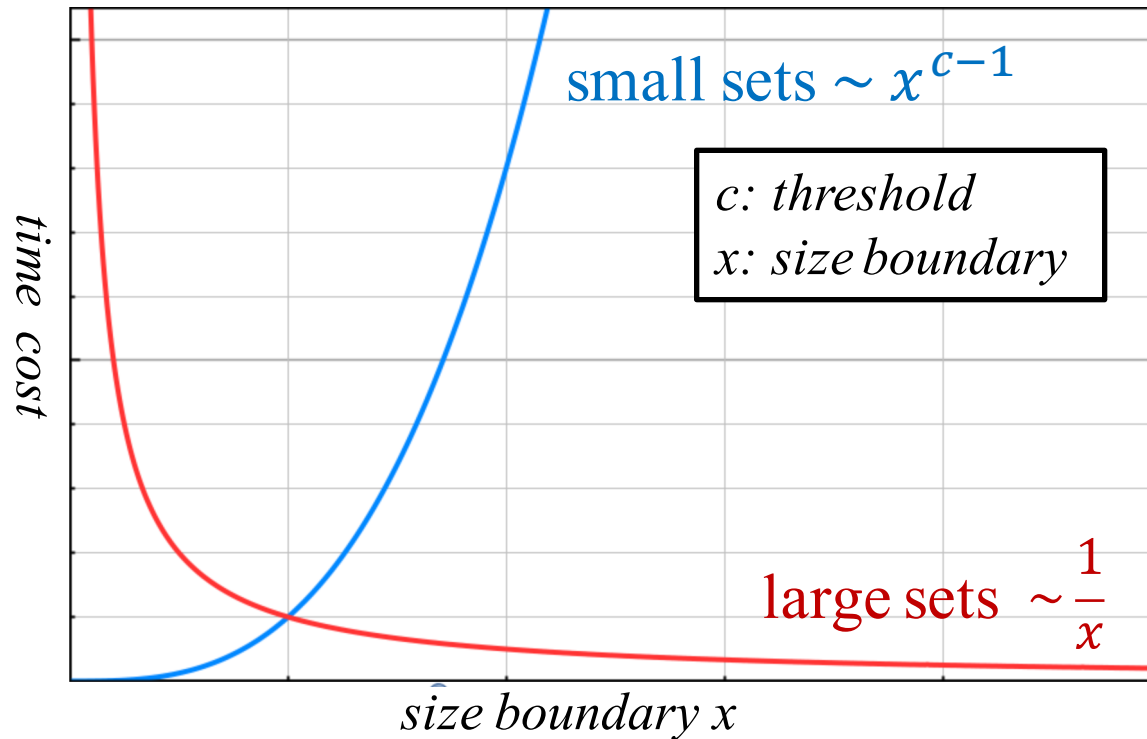
	<i>Block e_1</i>					<i>Block e_2</i>				<i>Block e_3</i>			<i>Block e_4</i>			<i>Block e_5</i>	
	e_1e_2	e_1e_3	e_1e_4	e_1e_5	e_1e_7	e_2e_3	e_2e_4	e_2e_5	e_2e_6	e_3e_4	e_3e_5	e_3e_7	e_4e_5	e_4e_6	e_4e_7	e_5e_6	e_5e_7
R_1	■	■				■											
R_2		■	■		■					■	■	■			■		
R_3		■		■	■					■	■	■					■
R_4							■	■	■				■	■		■	

invoke the heap-based method in each block with threshold $c - 1$.

the heap size will be much smaller than m

How to select a practical size boundary?

Size Boundary Selection



cost (slope) goes up smoothly first, and then rapidly

benefit (slope) goes down rapidly first, and then smoothly

Benefit: the reduction of time spend on large sets

Cost: the increase of time spend on small sets

increase x little by little, stop right after estimated benefit is less than estimated cost

Estimating large set processing cost

- For each large set, add up all its corresponding inverted list lengths
 - Use the summation as the cost
-

Estimating small set processing cost

- Heap adjustment cost and binary search costs.
 - Sampling some blocks and run the heap-based method to get the estimation
 - Result Generation Costs
 - $R = \{e_1 \ e_2 \ e_3 \ e_4\}$ and $S = \{e_1 \ e_2 \ e_3 \ e_5\}$
 - The pair $\langle R, S \rangle$ is generated 3 times by our method (Why?)
 - ***Thus the cost is proportional to the # of c-subsets shared by small sets***
 - Sample a small number of set pairs
 - For each set pair, compute their overlap as P
 - The # of shared c-subsets between them is $\binom{P}{c}$
-

Conclusion

- $o(n^2) + O(k)$ -- Sub-quadratic Algorithm whenever possible
 - Heap-based Methods for Small Sets
 - Size-boundary Selection Method
-