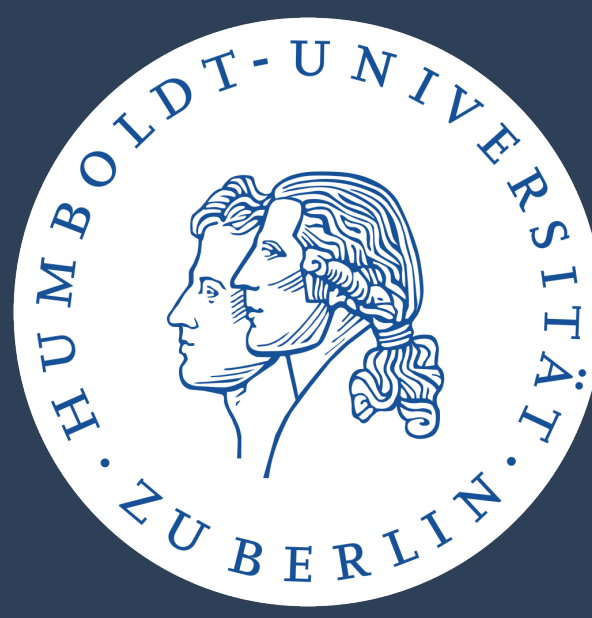


# SuSe: Summary Selection for Regular Expression Subsequence Aggregation over Streams



Steven Purtzel and Matthias Weidlich

Humboldt-Universität zu Berlin

purtzesc@hu-berlin.de

## Problem Setting

### Given ... :

- stream  $s = \langle e_1, e_2, \dots \rangle$
- summary size  $n$
- Regex query  $(\gamma, \tau, A)$ 
  - $\gamma = AB^*C$
  - $\tau = 10$  time units
  - $A = \text{COUNT}$
- (unknown) evaluation time points  $\mathcal{E}$

### ... select a substream ...

Stream  $s$ :  $A_1 \ A_2 \ B_3 \ C_4 \ B_5 \ C_6 \ A_7 \ B_8 \ C_9 \ A_{10} \ \dots$

Summary Size: 5; RegEx:  $AB^*C$ ; Eval. Time Point: 10; Matches: 28

Random  $S$ :  $B_3 \ B_5 \ C_6 \ C_9 \ A_{10}$  Matches: 0

FIFO  $S$ :  $C_6 \ A_7 \ B_8 \ C_9 \ A_{10}$  Matches: 2  
 $\{(A_7, C_6), (A_7, B_8, C_9)\}$

SuSe  $S$ :  $A_1 \ A_2 \ B_3 \ B_5 \ C_6$  Matches: 8  
 $\{(A_1, C_6), (A_2, C_6), (A_1, B_3, C_6), (A_2, B_3, C_6), (A_1, B_5, C_6), (A_2, B_5, C_6), (A_1, B_3, B_5, C_6), (A_2, B_3, B_5, C_6)\}$

### ... such that:

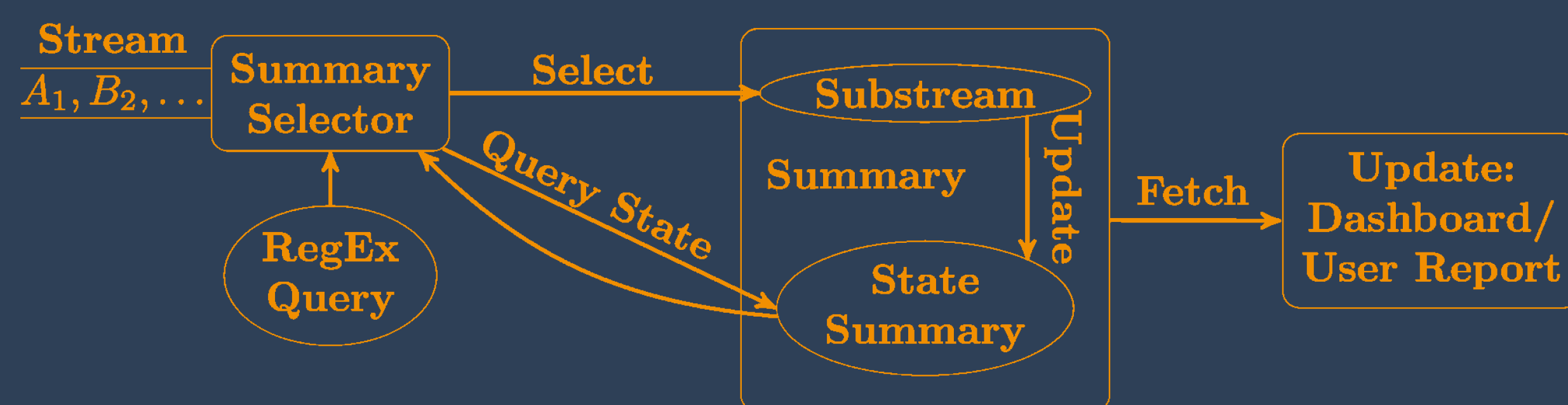
- at each evaluation point, the number of subsequence matches is maximized to minimize aggregation loss

$$\text{minimize} \sum_{\varepsilon \in \mathcal{E}} (|CM_{s(\varepsilon)}| - |CM_{S(\varepsilon, n)}|)$$

## SuSe Architecture

### 1 Two Main Components:

1) **summary selector** & 2) **summary** (substream & StateSummary)



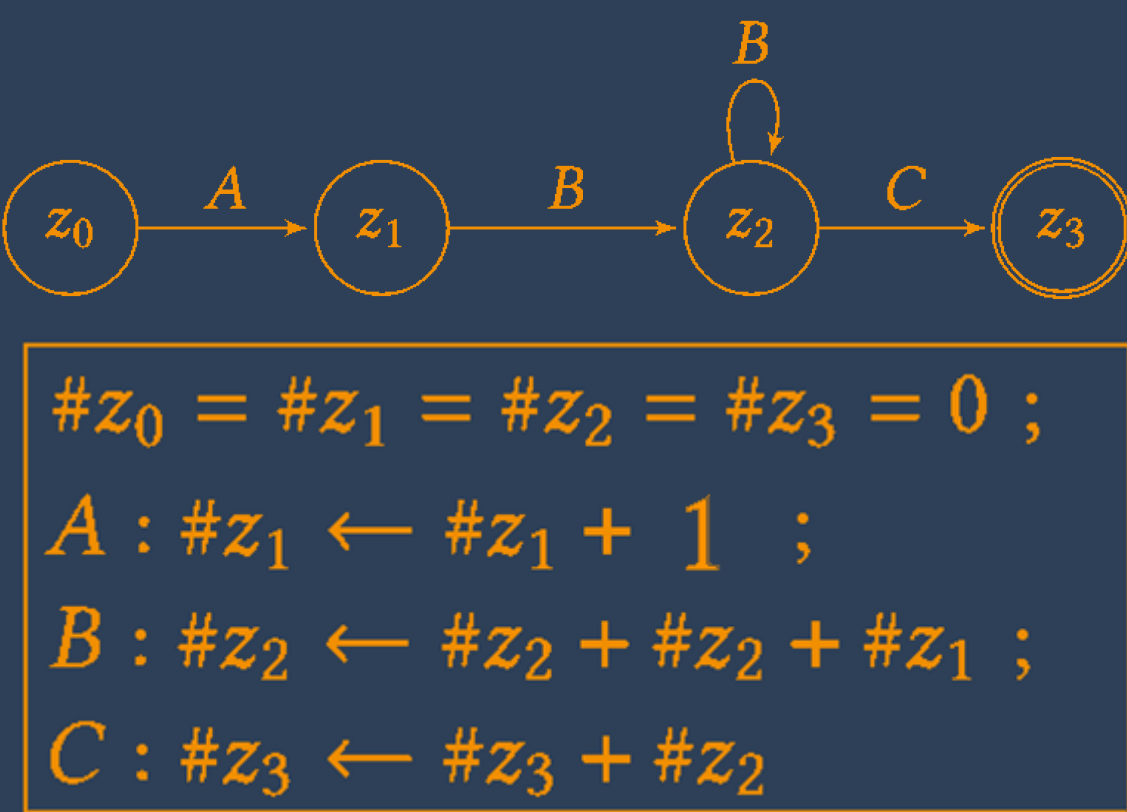
### summary

- consists of the selected substream and StateSummary
- StateSummary holds aggregated information on (partial) matches within the substream and for each of its selected elements
- maintain state consistency during element insertions and removals

### summary selector

- queries the StateSummary to decide element insertions/replacements

### 2 StateSummary:



- count matches per state
  - global state counters
    - to quantify the current quality
- count matches per element per state
  - (active) local state counters
    - to quantify an element's current (and future) match participation
- derive count rules; state counter updates

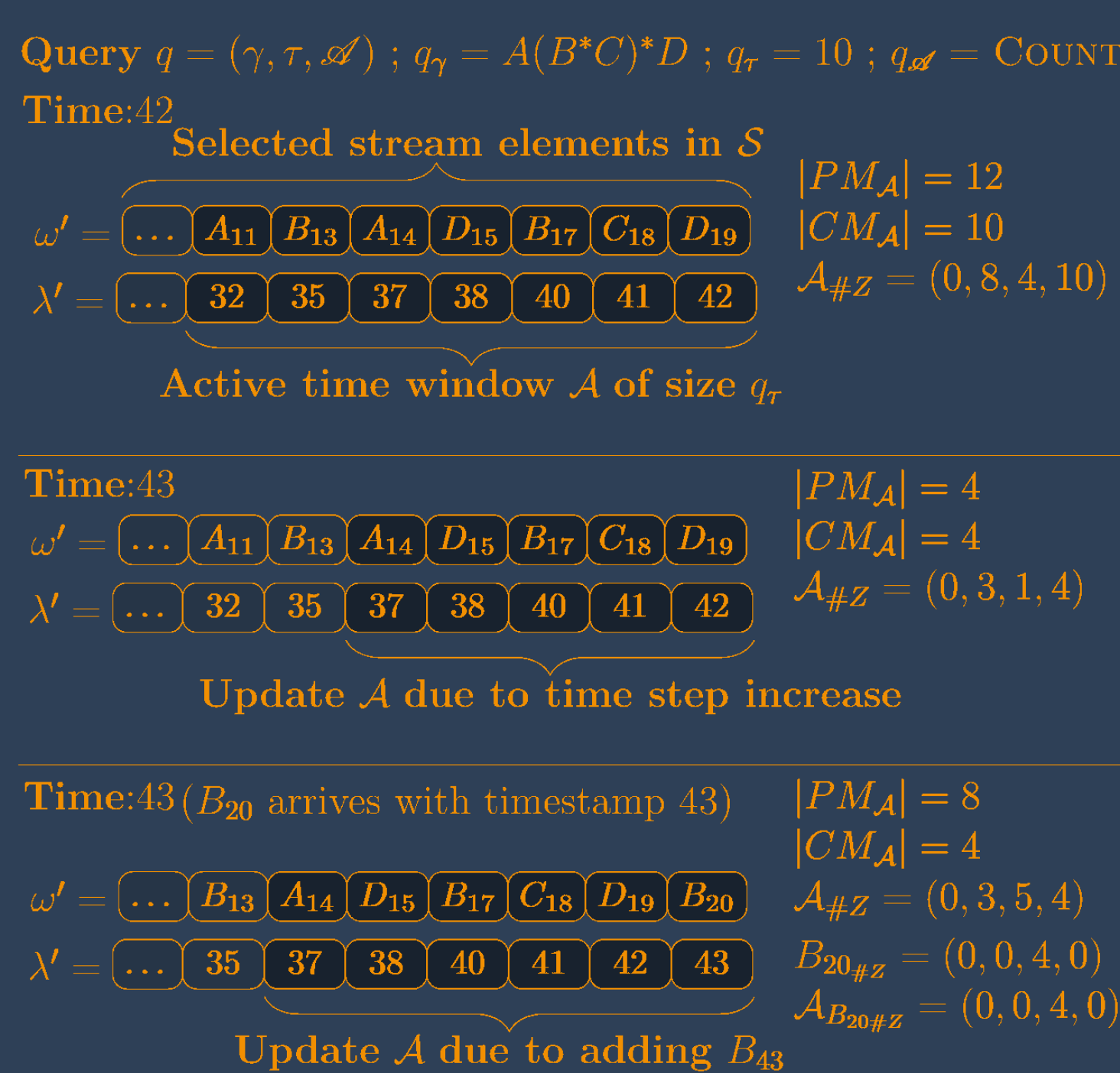
Element Stream	$A_1$	$A_1B_2$	$A_1B_2B_3$	$A_1B_2B_3C_4$
$A_1 : \#z_1 \leftarrow \#z_1 + 1 = 1$	$B_2 : \#z_2 \leftarrow \#z_2 + \#z_2 + \#z_1 = 3$	$B_3 : \#z_3 \leftarrow \#z_3 + \#z_3 + \#z_2 = 3$	$C_4 : \#z_4 \leftarrow \#z_4 + \#z_4 + \#z_3 = 3$	

**Global State Counters**  
 $(0,0,0,0) \ A_1 \rightarrow (0,1,0,0) \ B_2 \rightarrow (0,1,1,0) \ B_3 \rightarrow (0,1,3,0) \ C_4 \rightarrow (0,1,3,3)$

**Local State Counters  $B_2$**   
 $(0,0,1,0) \ B_3 \rightarrow (0,0,2,0) \ C_4 \rightarrow (0,0,2,2)$

### 3 Active Time Window:

- substream segment within the current (sliding) time window
- contains (active) partial matches with match potential
- used to initialize (active) local state counters
- crucial for state consistency and estimating future matches



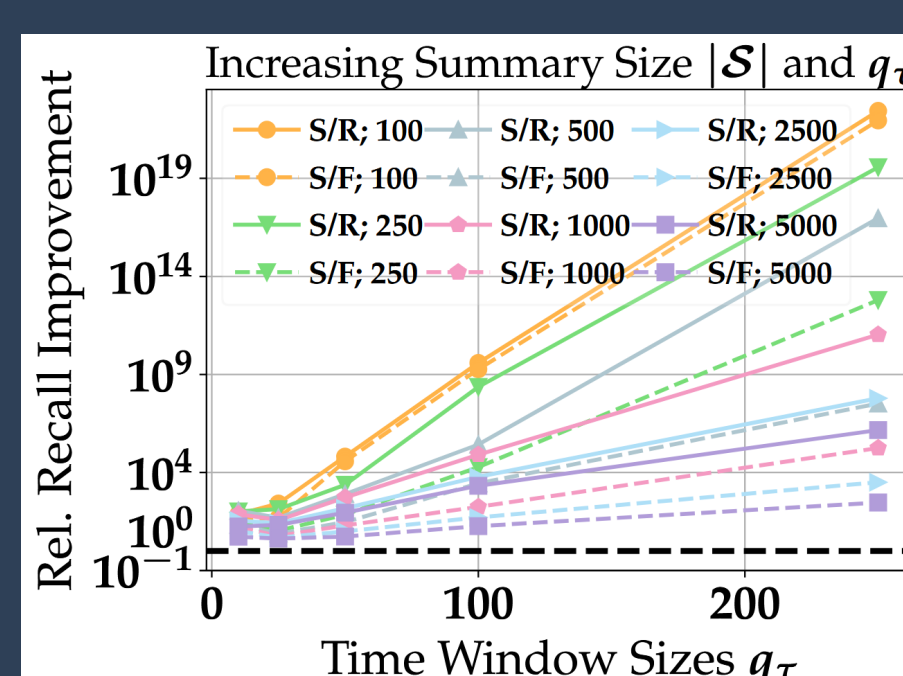
### 4 Selection Strategy

- goal: keep elements in substream that maximize RegEx matches
- benefit(e) = present benefit(e) + expected benefit(e)
  - present: #complete matches e participates in
  - expected: #complete matches e \*will\* participate in
    - partial matches including e with future match potential
- idea: simulate how e's (active) partial matches evolve with future arrivals
  - for the remaining time span where e stays within the window, e.g.,  $\Delta t = 2$ 
    - t=0: (0, 5, 10, **15**); t=1: (0, 7.75, 13, **15.75**); t=2: (0, 10, 16.5, **17**)
    - accumulate expected match completions via counters on final states

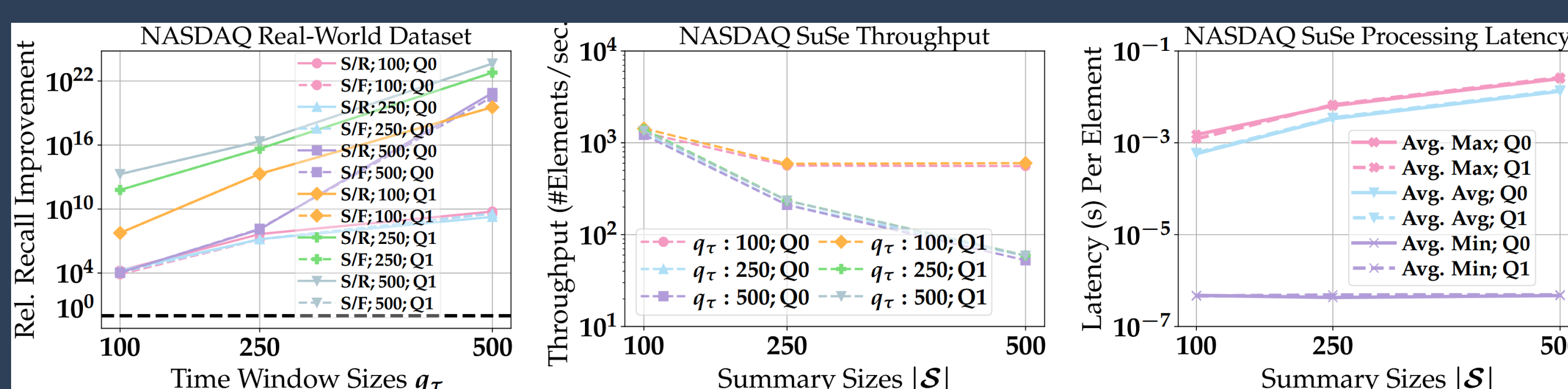
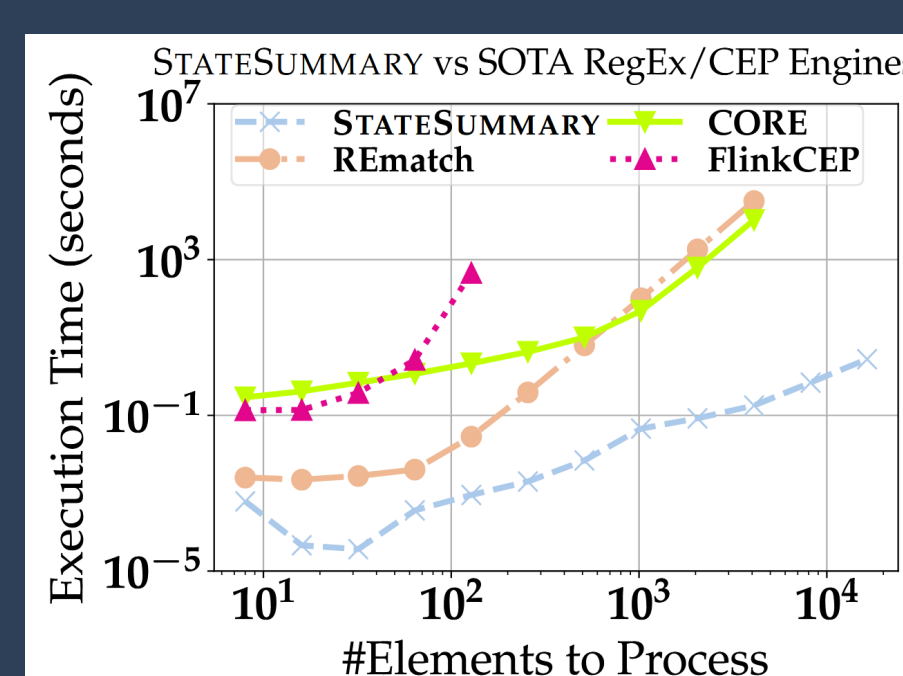
## (Some) Evaluation Results

- Random & FIFO baselines
  - $RRI = \frac{1}{|\mathcal{E}|} \sum_{\varepsilon \in \mathcal{E}} \frac{|CM_{SuSe}(\varepsilon)|}{|CM_{base}(\varepsilon)|}$
- FlinkCEP, CORE, REmatch
- NASDAQ ~ 450k elements

Effectiveness



Efficiency



## Take Away

- SuSe: architecture for efficient RegEx subsequence summarization over streams
- leverages StateSummary to maintain query-specific, aggregated match summaries
- enables substream selection that minimizes aggregation loss by maximizing subsequence matches
- SuSe is orders of magnitude faster than leading RegEx and CEP engines, while producing richer substream aggregates than baselines

