

DIPSUM: Distributed Pattern Summaries for Efficient CEP Aggregates

Steven Purtzel¹, Samira Akili², Marc Sebastian Kühne¹, and Matthias Weidlich¹
Humboldt-Universität zu Berlin¹ & BIFOLD, TU Berlin²

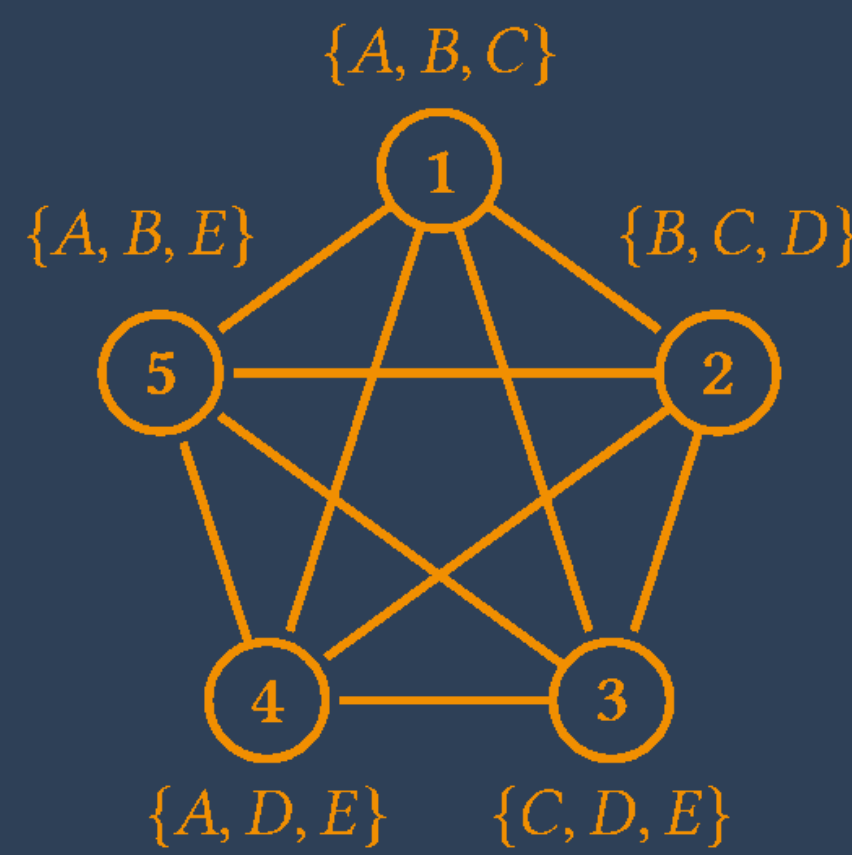
purtzesc@hu-berlin.de

Problem Setting

Evaluate CEP aggregate query...

SEQ(A,B,C⁺,D,E)
WITHIN 1 min
COUNT

...in event-network...



Rates:
 $r(A) > r(B)$,
 $r(C) \gg \text{OTHERS}$,
 $r(E) \gg r(D)$

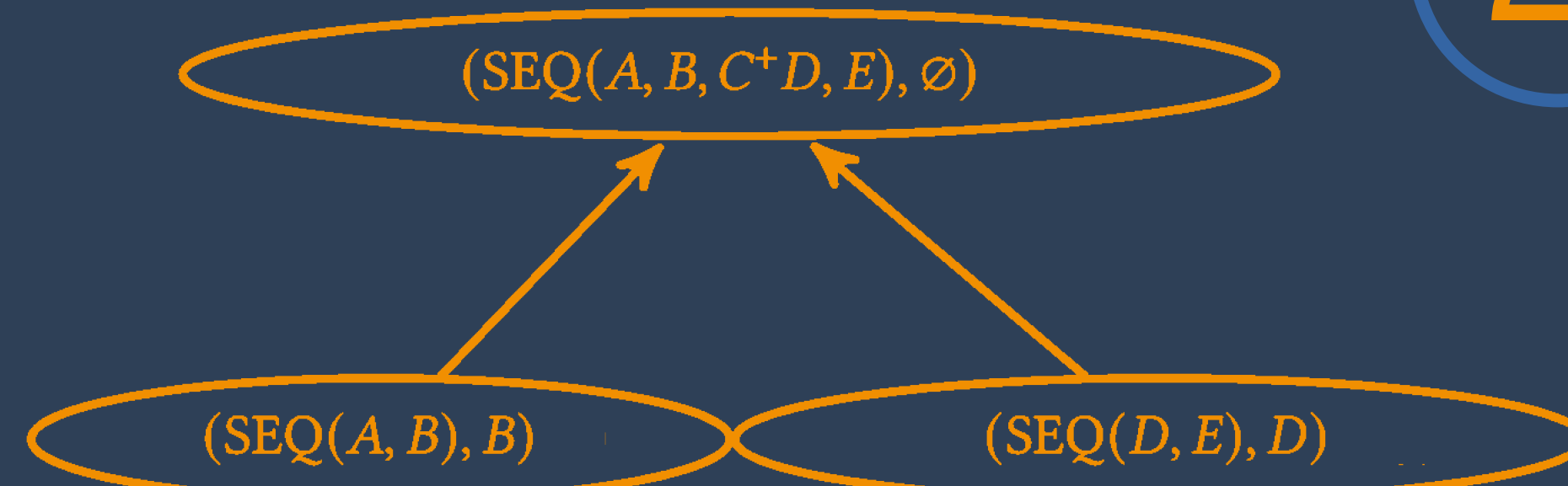
...such that:

- aggregate function is evaluated
- network transmission costs are minimized

DIPSUM Framework

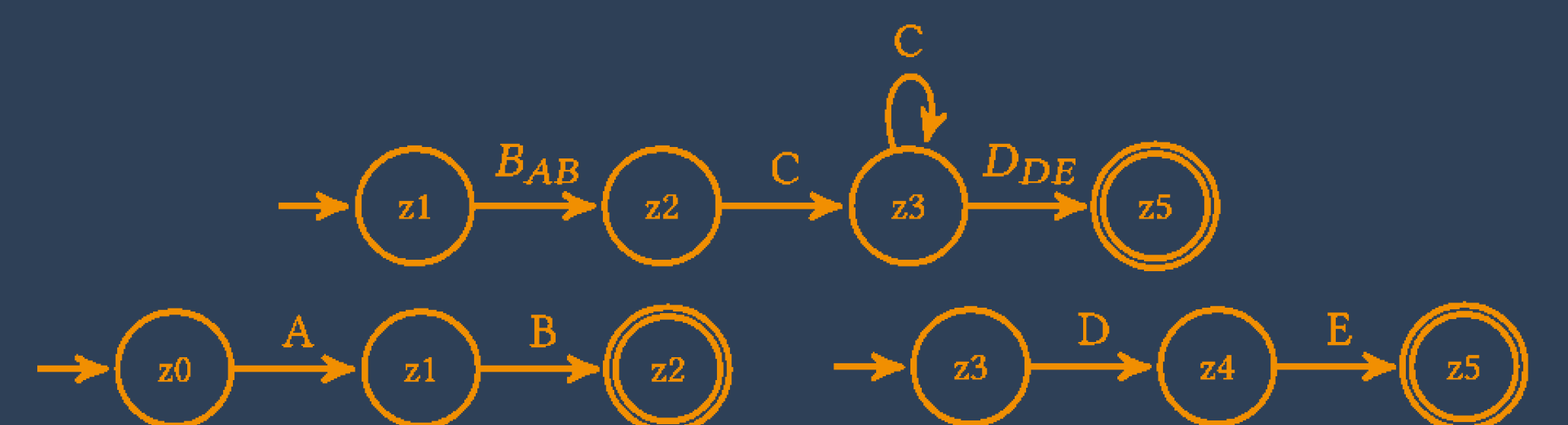
1. Query Decomposition:

- based on event type rates
- establish cut event types (event types connecting sub-queries)
- decompose query into sub-queries



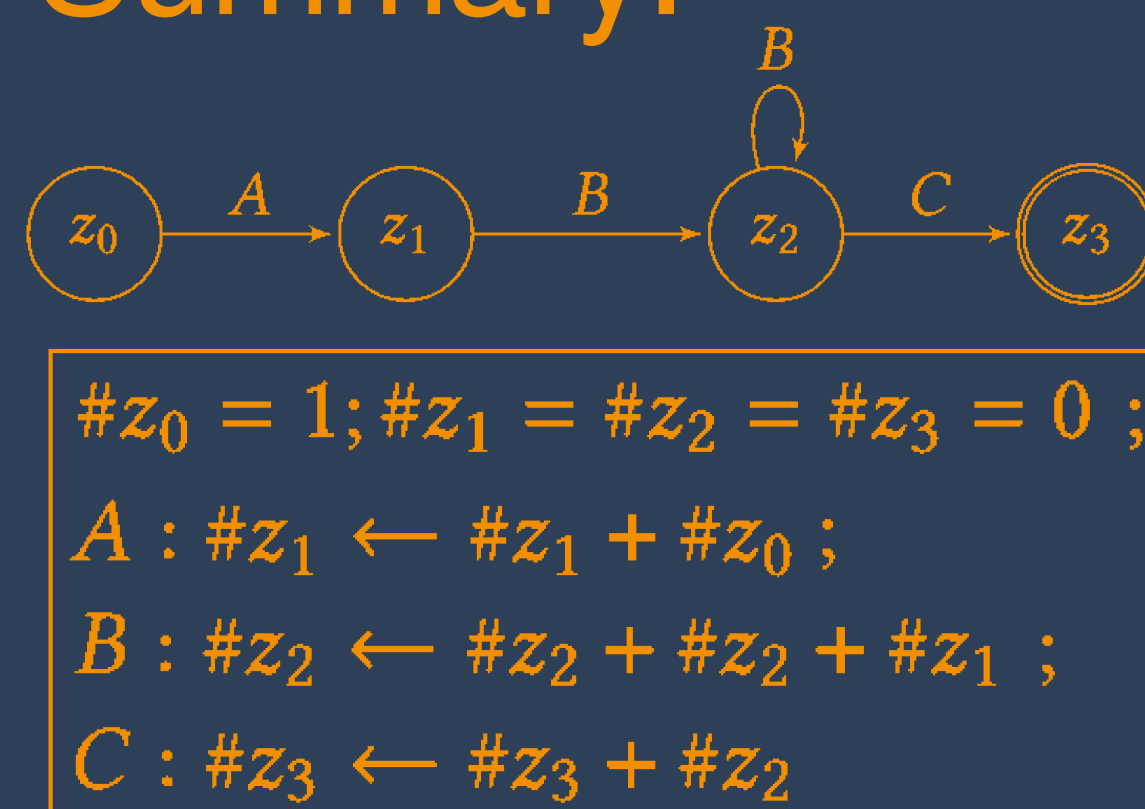
2. (Aggregator) NFA Derivation:

- derive (aggregator) NFA for each sub-query



3. Partial State Summary:

- for each (aggregator) NFA derive a partial State Summary
- example: SEQ(A,B+,C)



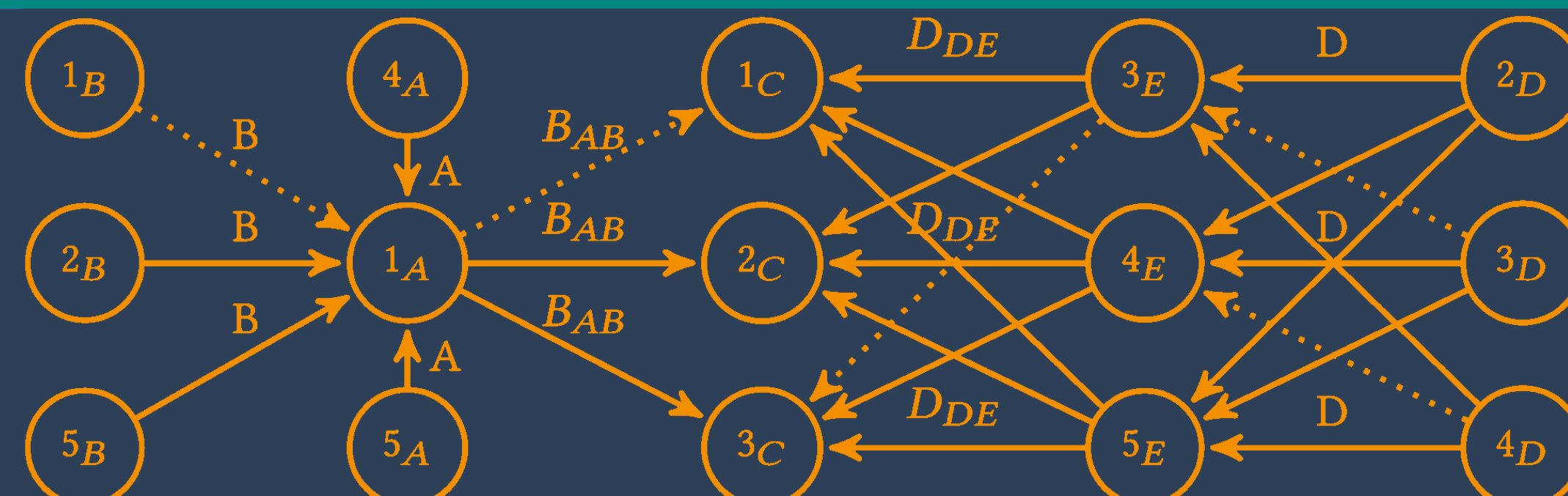
#z₀ = 1; #z₁ = #z₂ = #z₃ = 0;
A : #z₁ ← #z₁ + #z₀;
B : #z₂ ← #z₂ + #z₂ + #z₁;
C : #z₃ ← #z₃ + #z₂

- count matches per state
- global state counters
- initialize: (1,0,0,0)
- count matches per event per state
- local state counters
- represent partial aggregates
- e.g., #matches B occurred in SEQ(A,B)
- derive count rules; state counter updates

Event Stream	A ₁	A ₁ B ₂	A ₁ B ₂ B ₃	A ₁ B ₂ C ₄
	A ₁	B ₂	B ₃	C ₄
A ₁ : #z ₁ ← #z ₁ + #z ₀ = 1				
B ₂ : #z ₂ ← #z ₂ + #z ₂ + #z ₁ = 1				
C ₄ : #z ₃ ← #z ₃ + #z ₂ = 3				
Global State Counters	(1,0,0,0)	A ₁ → (1,1,0,0)	B ₂ → (1,1,1,0)	B ₃ → (1,1,3,0)
Local State Counters B ₂	(0,0,1,0)	B ₃ → (0,0,2,0)	C ₄ → (0,0,2,2)	

4. Partial Summary Placement & Partial Aggregate Exchange:

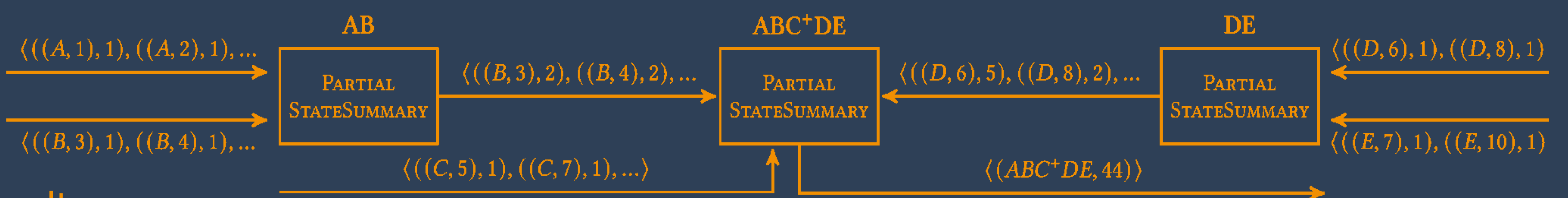
- based on rates & #nodes producing an event type



- leveraging inequalities for single- and multi-sink placements

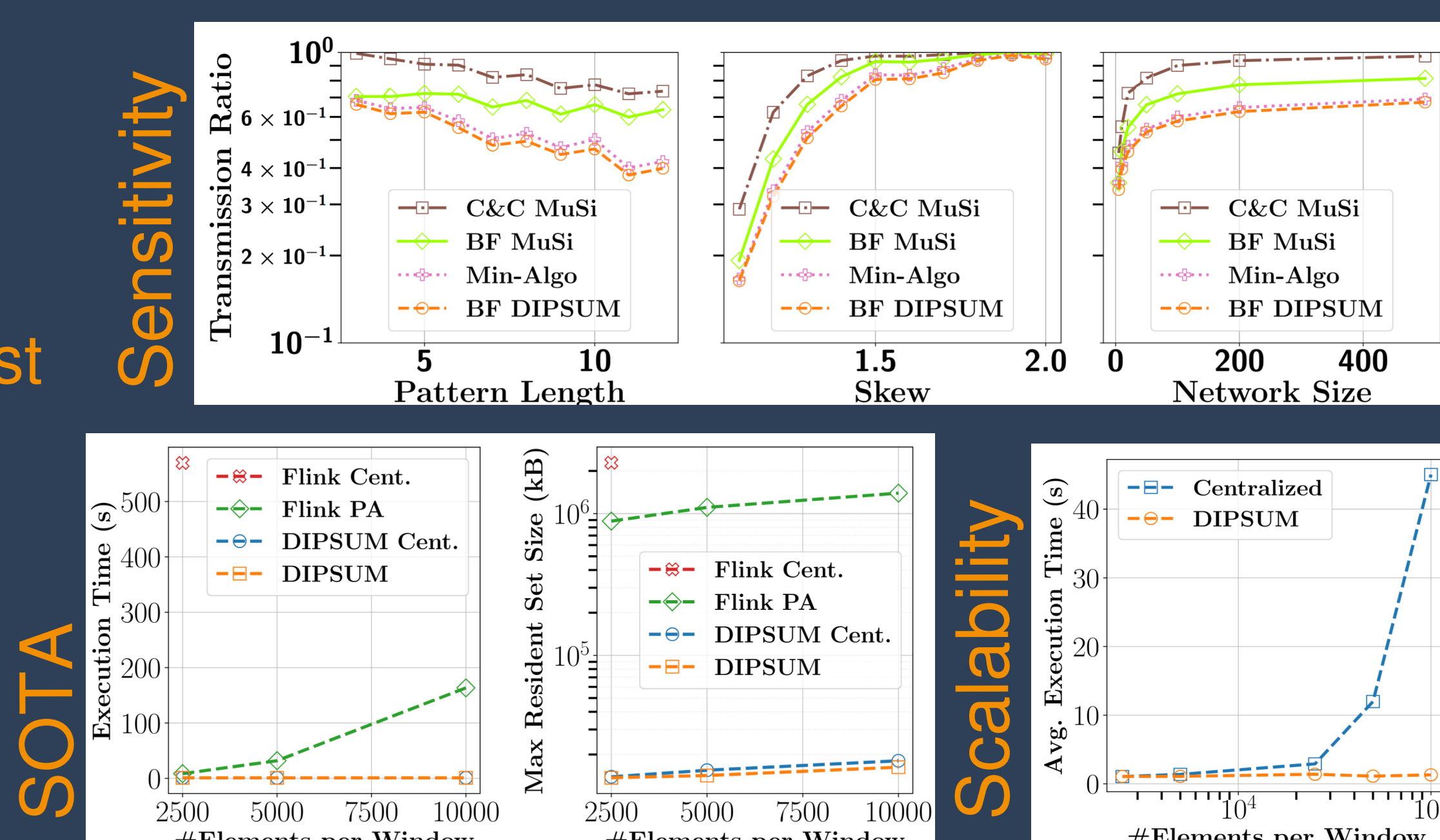
5. Distributed State Summary:

- stitching partial aggregates; obtain result



(Some) Evaluation Results

- synthetic data
- varying parameters
- transmission ratio
- prod. cost/cent. cost
- FlinkCEP baselines
- simul. environment
- execution time & memory usage



Take Away

- DIPSUM plans **decompose** queries into sub-queries, **derive** and **place** pattern summaries, and **coordinate** partial aggregate exchange
- DIPSUM improves **transmission costs**, **execution times**, **throughput**, and **memory usage** by multiple orders of magnitude

