Factor Windows: Cost-based Query Rewriting for Optimizing Correlated Window Aggregates

Wentao Wu (Microsoft Research)

Philip A. Bernstein (Microsoft Research)

Alex Raizman (Microsoft)

Christina Pavlopoulou (University of California, Riverside)



Scenario: Correlated Window Aggregates

An example "window set" query from **Azure Stream Analytics**: <u>https://azure.microsoft.com/en-us/services/stream-analytics/</u>

SELECT DeviceID, System.Window().Id, Min(T) AS MinTemp, FROM Input TIMESTAMP BY EntryTime GROUP BY DeviceID, Windows(Window('20 min', TumblingWindow(minute, 20), Window('30 min', TumblingWindow(minute, 30),

Window('40 min', TumblingWindow(minute, 40))

Overview of Our Solution



Window Coverage



(b) Window partitioning

Window Coverage Graph (WCG)

- We represent a window as W(r, s), where r and s represent W's "range" and "slide".
 - W is a "tumbling" window if r = s.

W is a "hopping" window if r > s.



(a) Tumbling Window



(b) Hopping Window

- Window coverage graph (WCG)
 - Vertices represent the windows in the window set.
 - There is an edge from W_2 to W_1 if W_1 is covered by W_2 .
 - Example: Consider a window set with four tumbling windows $W_1(10, 10), W_2(20, 20), W_3(30, 30), W_4(40, 40)$.
 - WCG is a directed acyclic graph (DAG).



Cost Modeling

- Assume that the cost of computing an aggregate function *f* is proportional to the number of events processed.
- The "instance/interval cost" of window W(r,s) is $\eta \cdot r$, where η is the input event rate.
- If W_1 is covered by W_2 , the instance cost of W_1 is reduced from $\eta \cdot r_1$ to $M(W_1, W_2) = 1 + \frac{r_1 r_2}{r_2}$.
 - Intuitively, $M(W_1, W_2)$ represents the number of intervals from W_2 that cover an interval from W_1 .
- If W_1 is covered by multiple W_2 's, the instance cost of W_1 is reduced to $\min_{W_2} \{M(W_1, W_2)\}$.

Min-Cost WCG

• <u>Example</u>: Consider a window set with $W_1(10, 10), W_2(20, 20), W_3(30, 30), W_4(40, 40),$.



Factor Windows

- Factor window a window not in the input window set that can help reduce the evaluation cost.
 I.e., it's a window that covers some window(s) in the window set.
- Example: Consider a window set with $W_2(20, 20)$, $W_3(30, 30)$, $W_4(40, 40)$.



(a) Initial WCG (i.e., Min-cost WCG w/o factor window), Cost = 246

(b) Min-cost WCG with factor window $W_1(10, 10)$, Cost = 150

Technical Problem: How to Find the Optimal Factor Windows?

- NP-hard problem (Steiner tree)
- Our "greedy" solution: Find factor windows that are "locally" optimal.
 (1) Candidate generation
 (2) Candidate selection
- Special treatment for "**partitioned by**" semantics: See the paper for details.



Candidate Generation

- Find eligible slides S_f . (1) $s_d = gcd\{s_1, ..., s_K\}$ (2) $S_f = \{s_f: s_d \mod s_f = 0 \text{ and } s_f \mod s_W = 0\}$
- Find eligible ranges R_f for each $s_f \in S_f$. (1) $r_{min} = min\{r_1, ..., r_K\}$ (2) $R_f = \{r_f: r_f \mod s_f = 0 \text{ and } r_f \le r_{min}\}$



• Construct a candidate factor window $W_f = \langle r_f, s_f \rangle$ for each pair $(r_f \in R_f, s_f)$.

Candidate Selection

• Compute the "benefit" δ_f of each candidate factor window W_f .

$$\delta_{f} = \sum_{j=1}^{K} n_{j} \cdot \left(\frac{r_{j} - r_{W}}{s_{W}} - \frac{r_{j} - r_{f}}{s_{f}}\right) - n_{f} \cdot \left(1 + \frac{r_{f} - r_{W}}{s_{W}}\right)$$

$$W_{1} \quad W_{2} \quad W_{K} \quad W_{1} \quad W_{2} \quad W_{K}$$

• Select the factor window with the maximum benefit.

Put Things Together

1 // Construct the set W_f of candidate factor windows. 2 $\mathcal{W}_f \leftarrow \emptyset;$ 3 $s_d \leftarrow \gcd\{s_1, \dots, s_K\};$ Candidate 4 $S_f \leftarrow \{s_f : s_d \mod s_f = 0 \text{ and } s_f \mod s_W = 0\};$ Generation 5 $r_{\min} \leftarrow \min\{r_1, ..., r_K\};$ 6 foreach $s_f \in \mathcal{S}_f$ do $\mathcal{R}_f \leftarrow \{r_f : r_f \mod s_f = 0 \text{ and } r_f \leq r_{\min}\};$ 7 foreach $r_f \in \mathcal{R}_f$ do 8 Construct a candidate factor window $W_f \langle r_f, s_f \rangle$; 9 if $W_f \leq W$ and $W_j \leq W_f$ for $1 \leq j \leq K$ then 10 $\mathcal{W}_f \leftarrow \mathcal{W}_f \cup \{W_f\};$ 11 12 // Find the best factor window from W_f . 13 $\delta_f^{\max} \leftarrow 0, W_f^{\max} \leftarrow \text{null};$ Candidate 14 foreach $W_f \in \mathcal{W}_f$ do Selection Compute the benefit δ_f of W_f using Equation 2; 15 if $\delta_f \geq 0$ and $\delta_f > \delta_f^{\max}$ then 16 $\delta_f^{\max} \leftarrow \delta_f, W_f^{\max} \leftarrow W_f;$ 17 18 return W_f^{max} ;

Experimental Settings

- Experimental Setup
 - (1) Prototype: C# implementation that produces optimized query plans (with and without factor windows) represented by Trill expressions.
 - (2) Performance metric: We measure the throughput of the original plan and the optimized query plans.
 - (3) Runtime configuration: Single-core execution with 2.2 GHz Intel CPUs and 128 GB main memory.
- Datasets
 - (1) Synthetic Datasets: Synthetic-1M (1 million events), Synthetic-10M (10 million events).
 - (2) Real Dataset: Real-32M (from DEBS 2012 Grand Challenge, with 32 million events).
- Window Sets
 - (1) RandomGen: Windows with randomly generated slides and ranges.
 - (2) SequentialGen: Windows whose ranges follow a "sequential" pattern (e.g., 10s, 20s, 30s, 40s, ...).

Throughput with and without Factor Windows

<u>Synthetic10M</u>, Window set size = 5



Summary of Throughput Results

Synthetic10M

Window set size = {5, 10}

<u>Real32M</u>
Window set size

= {5, 10}

Setup	w/o FW	w/oFW	w/ FW	w/ FW
	(Mean)	(Max)	(Mean)	(Max)
R-5-tumbling	1.21×	1.92×	1.85×	$2.54 \times$
R-10-tumbling	$1.34 \times$	1.77×	$1.88 \times$	3.38×
R-5-hopping	$1.18 \times$	$1.82 \times$	3.26×	$4.29 \times$
R-10-hopping	1.34×	1.71×	3.20×	6.15×
S-5-tumbling	1.63×	1.67×	4.28×	4.81×
S-10-tumbling	$1.98 \times$	$2.05 \times$	7.91×	9.38×
S-5-hopping	$1.34 \times$	$1.48\times$	2.17×	$2.81\times$
S-10-hopping	$1.58 \times$	1.73×	$2.92\times$	3.79×

Setup	w/o FW	w/o FW	w/ FW	w/ FW
_	(Mean)	(Max)	(Mean)	(Max)
R-5-tumbling	1.19×	1.78 imes	$1.43 \times$	1.91×
R-10-tumbling	$1.30 \times$	$1.71 \times$	$1.53 \times$	$2.86 \times$
R-5-hopping	$1.09 \times$	$1.39 \times$	$1.54 \times$	$2.63 \times$
R-10-hopping	$1.18 \times$	1.39×	1.46×	3.53×
S-5-tumbling	1.63×	$1.67 \times$	$4.12 \times$	$4.85 \times$
S-10-tumbling	$1.90 \times$	$1.97 \times$	$7.53 \times$	9.14 ×
S-5-hopping	$1.12 \times$	$1.30 \times$	$1.22 \times$	$1.77 \times$
S-10-hopping	$1.22 \times$	$1.51 \times$	$1.45 \times$	$2.31 \times$

Scalability w.r.t. Window Set Size

	Setup	w/o FW (Mean)	w/oFW (Max)	w/ FW (Mean)	w/ FW (Max)
etic10M v set size)	R-15-tumbling R-20-tumbling R-15-hopping R-20-hopping	$1.55 \times 1.49 \times 1.55 \times 1.68 \times$	$1.96 \times 2.29 \times 1.95 \times 2.20 \times$	$\begin{array}{c} 2.97 \times \\ 2.10 \times \\ 4.67 \times \\ 4.23 \times \end{array}$	$4.34 \times 4.83 \times 6.59 \times 7.65 \times$
	S-15-tumbling S-20-tumbling S-15-hopping S-20-hopping	$2.43 \times$ $2.42 \times$ $1.85 \times$ $1.91 \times$	$2.49 \times 2.53 \times 2.09 \times 2.15 \times$	$11.29 \times 14.28 \times 3.51 \times 4.02 \times$	13.83× 16.82 × 4.68× 5.32×

Synthe Window = {15, 20

Query Optimization Overhead



Vary window set size from 5 to 20.

Comparison with Window Slicing

- Window slicing
 - Chop the window into smaller chunks (i.e., <u>slices</u>).
 - Compute the aggregate over the window by aggregating sub-aggregates over the slices.
- Scotty (TODS'21): "General stream slicing"
 - We compare our factor-window based optimization against Scotty on top of Apache Flink.
 - We use the same data generator developed by Scotty for benchmarking its own performance: <u>https://github.com/TU-Berlin-DIMA/scotty-window-processor</u>.

Comparison Results



Window set size = 10

Conclusion: Summary of Contributions

- We proposed the "<u>window coverage graph</u>" (WCG) abstraction that captures the overlapping relationships between correlated windows.
- We proposed a <u>cost-based optimization framework</u> on top of the WCG abstraction, to minimize the computation cost of multi-window aggregate queries.
- We extended the cost-based framework by considering "<u>factor windows</u>" and developed technologies to find beneficial factor windows.
- We evaluated our techniques on both synthetic and real datasets and demonstrated that the throughput of optimized plans can outperform the original plans by <u>up to 16.8x</u>.

Thank you for viewing our presentation! We'd be happy to hear your comments and questions: wentao.wu@microsoft.com.