# **Interval Count Semi-Joins**



JOHANNES GUTENBERG UNIVERSITÄT MAINZ

# Panagiotis Bouros<sup>1</sup> and Nikos Mamoulis<sup>2</sup>

<sup>1</sup>Institute of Computer Science, Johannes Gutenberg University Mainz, Germany <sup>2</sup>Department of Computer Science & Engineering, University of Ioannina, Greece bouros@uni-mainz.de, nikos@cs.uoi.gr



Introduction			Baseline Approaches	Smart Counting Algorithm		
Company R	Company S		Naïve Algorithm:	<b>Observation:</b> to compute an <i>r.count</i> we need $\Box$ . The number $ \Delta^{S} $ of active intervals from $S$		
employee start end	employee	start end	Sort and aggregate join results	when <i>r</i> . <i>start</i> is accessed		
John 1993 2003	Jane	1994 1996		The number of intervals form S which		
Mary 1995 2008	Bob	1998 2002	(John, Jane) (John, Jane) (Mary, Jane) (John, Bob)	become active after r start and before r and		
Tom 2000 2015	Hugo	2007 2010	(John, Bob) (Mary, Jane) (John, 2)	became active after <i>i.sturt</i> and before <i>i.end</i>		



## Interval Join

Find all pairs of employees whose working periods at companies **R** and **S** overlap

## Interval Count Semi-Join

For *each* employee *r* from company **R**, count how many employees from company **S** worked at a period of time *overlapping* with *r*'s employment

## Applications

*Temporal* databases
 *Selecting and/or ranking* objects



X High cost: join + aggregation

## Simple Counting Algorithm:

Adapt Interval Join algorithm to count ICSJ results instead of producing IJ results



Approach: use a global counter g to track for how many intervals s ∈ S their start point has been seen so far
When the start point of a s ∈ S is accessed → increase |A<sup>S</sup>|; increase g
When the end point of a s ∈ S is accessed → decrease |A<sup>S</sup>|
When the start point of a r ∈ R is accessed → initialize r.count = |A<sup>S</sup>|-g
When the end point of a r ∈ R is accessed → finalize r.count = r.count + g



## **Related work**

Description Plane-sweep based interval joins [1], [2]

- **Top-k** count semi-joins
  - Relational [3], spatial [4]

## Background

Interval Join based on plane sweep

(1) Preprocessing:

- □ Each interval  $\rightarrow$  a *start* and an *end* domain point
- □ Sort endpoints of each input collection

X Cost similar to interval join
X Cost sensitive to join output

Use a hash table (or array) to track *r.count* 's
 Complexity: O(|R|+|S|), excluding sorting

## Experiments

#### Setup

In-memory processing

Gapless hash map [2]

Lazy optimization [2] for Naïve and Simple Counting

### Datasets

	FLIGHTS	BOOKS	GREEND	WEBKIT
Cardinality	445,827	2,312,602	110, 115, 441	2,347,346
Domain duration (secs)	2,750,280	31, 507, 200	283, 356, 410	461,829,284
Shortest interval (secs)	1,261	1	1	1
Longest interval (secs)	42,301	31,406,400	59,468,008	461,815,512
Avg. interval duration (secs)	8,791	2,201,320	16	33,206,300
Distinct domain points	41,975	5,330	182,028,123	174,471

## Execution time (varying |R|/|S|)

## Breakdown (|R|=|S|)



21st International Conference on Extending Database Technology (EDBT), Vienna, Austria, March 26-29, 2018