

Band Joins for Interval Data

Panagiotis Bouros¹ Konstantinos Lampropoulos²

Dimitrios Tsitsigkos^{2,3} Nikos Mamoulis² Manolis Terrovitis³

¹ Johannes Gutenberg University Mainz, Germany

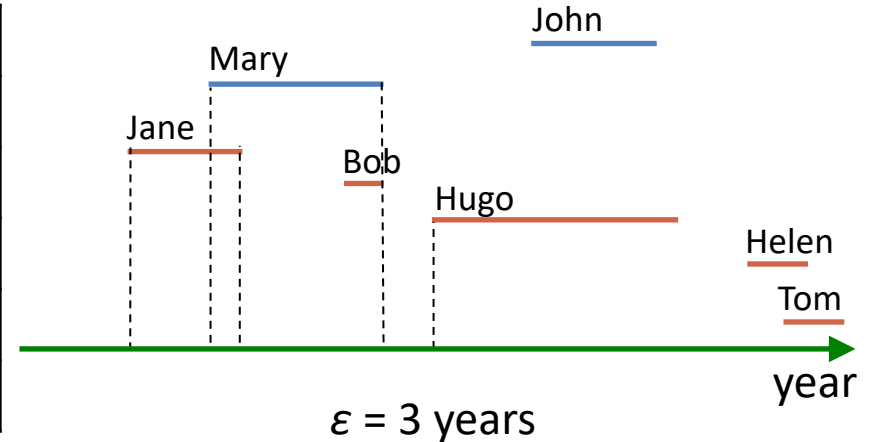
² University of Ioannina, Greece

³ Athena RC, Greece

Joining Interval Data

employee	start	end
John	1999	2002
Mary	1992	1996

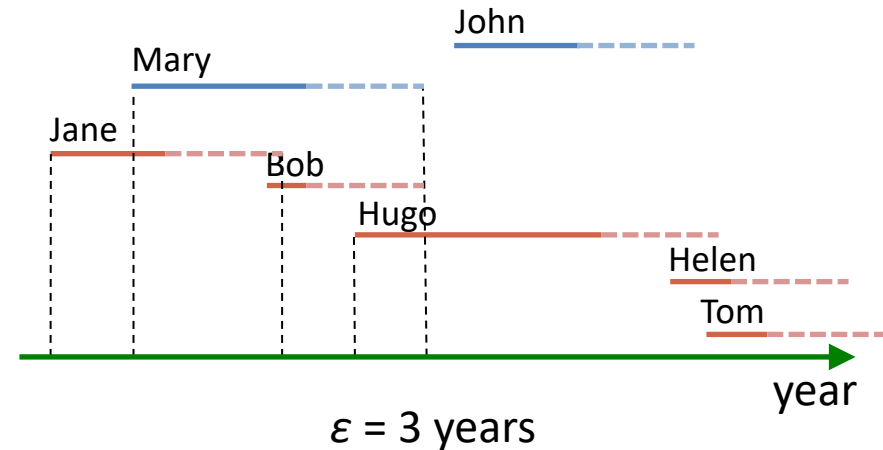
employee	start	end
Jane	1990	1993
Bob	1995	1996
Hugo	1997	2003
Helen	2005	2007
Tom	2006	2008



- **Overlap joins**
 - Find all pairs of employees whose period of work **overlaps**
- **Band joins**
 - Find all pairs of employees whose period of work has **an at most ϵ gap**
 - Overlap joins, special case with ϵ set to 0
- **Applications**
 - **Temporal** data, **uncertain** data, **XML** query processing, **streaming** data

Baseline

- A **straightforward** approach
 - Expand every interval in both inputs by ϵ
 - Compute a **traditional overlap join**
 - Directly apply **domain-based partitioning** from [Bouros and Mamoulis, PVLDB 2017]



Pros

- ✓ Simple and straightforward
- ✓ Capitalize on overlap join **methods** and **optimizations**

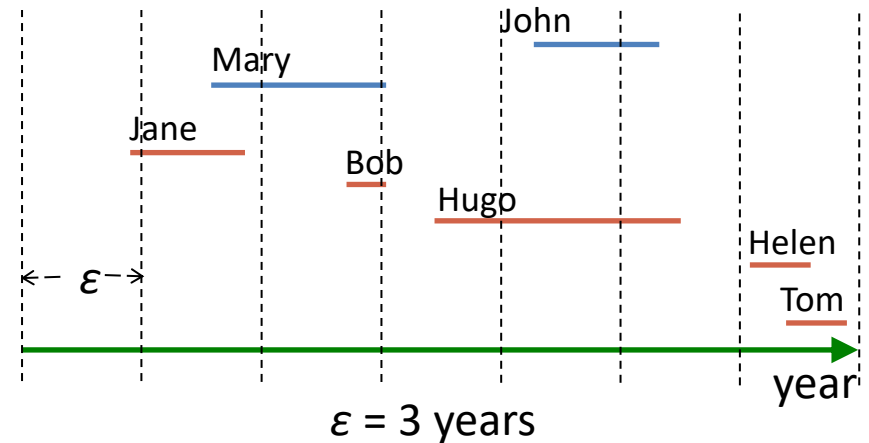
Cons

- ✗ Data replication **increases**, cost for partition-to-partition join **higher**
- ✗ Evaluation **agnostic** to parameter ϵ , **unnecessary comparisons**

Evaluation on ϵ -wide Partitioning

- Principles

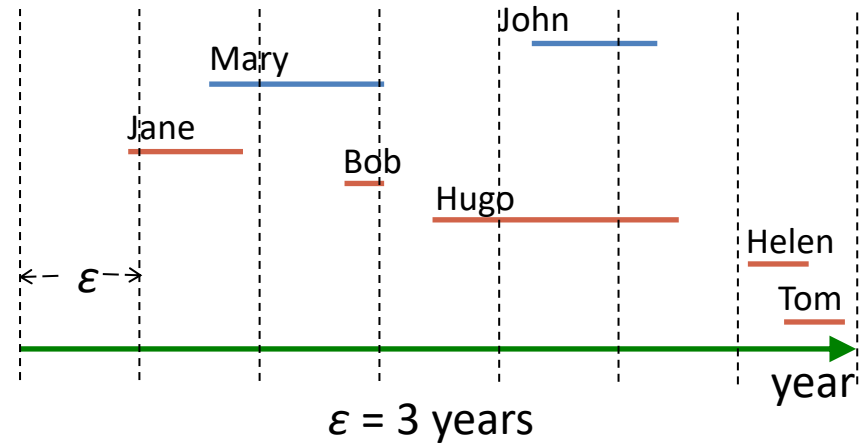
- Divide domain in ϵ -wide stripes
- Replicate intervals spanning multiple stripes
- Every partition is joined with exactly two partitions, e.g., for R_i from R
 - With S_j in a cross product
 - With S_{i+1}



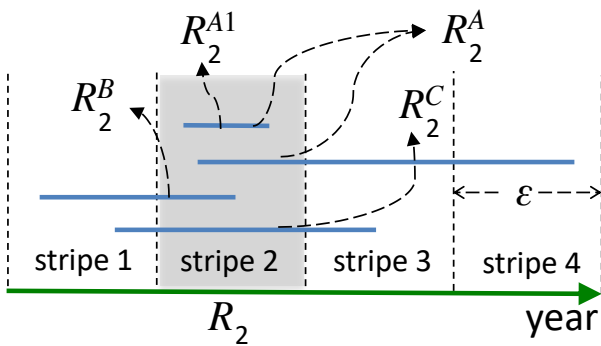
Evaluation on ϵ -wide Partitioning

- Principles

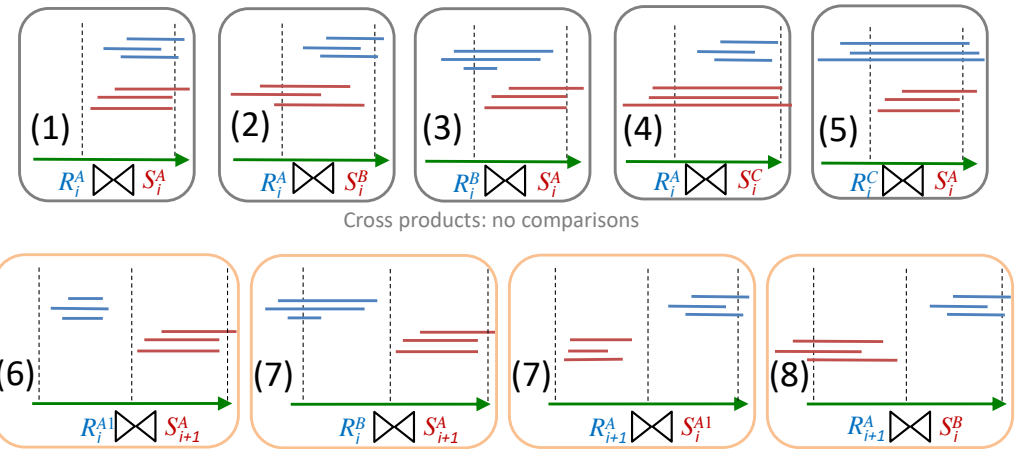
- Divide domain in ϵ -wide stripes
- Replicate intervals spanning multiple stripes
- Every partition is joined with exactly two partitions, e.g., for R_i from R
 - With S_j in a cross product
 - With S_{i+1}
- Divide partition contents to define mini-join tasks



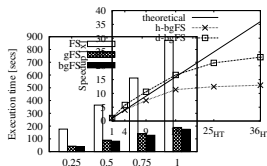
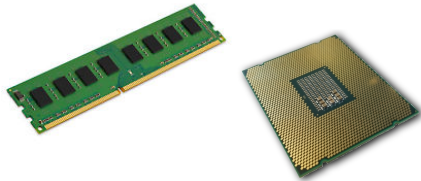
Mini-partitions



Mini-joins



Experimental Analysis



- Setup
 - In-memory, parallel processing
 - C++ with OpenMP for multi-threading
- Methods
 - BSL
 - ϵ -BSL, i.e., BSL with ϵ -wide stripes
 - ϵ -WIDE
- Experiments
 - 4 real-world datasets
 - Vary threshold ϵ , number of threads, input size ratio $|R|/|S|$
- Key finding
 - ϵ -WIDE most efficient method, unless ϵ too small compared to the domain

Thank you

Q & A