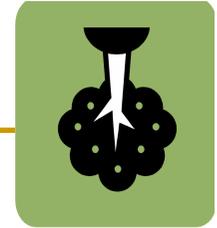# [HR-Join]

## Sum-Max Monotonic Ranked Joins for Evaluating Top-K Twig Queries on Weighted Data Graphs

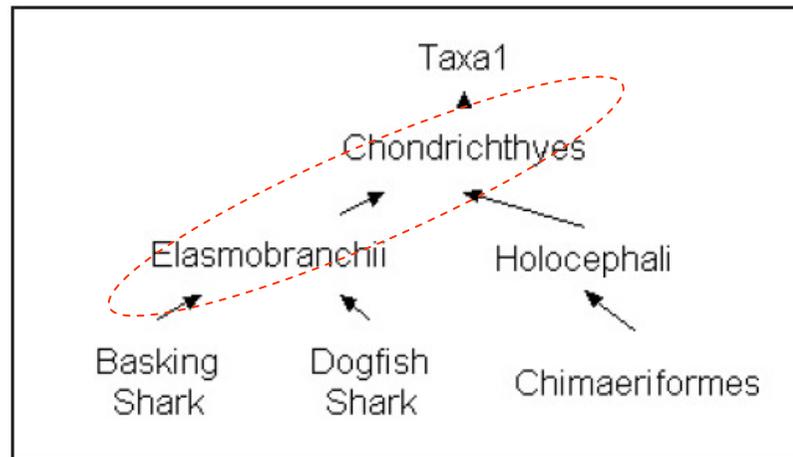Yan Qi                    Arizona State University
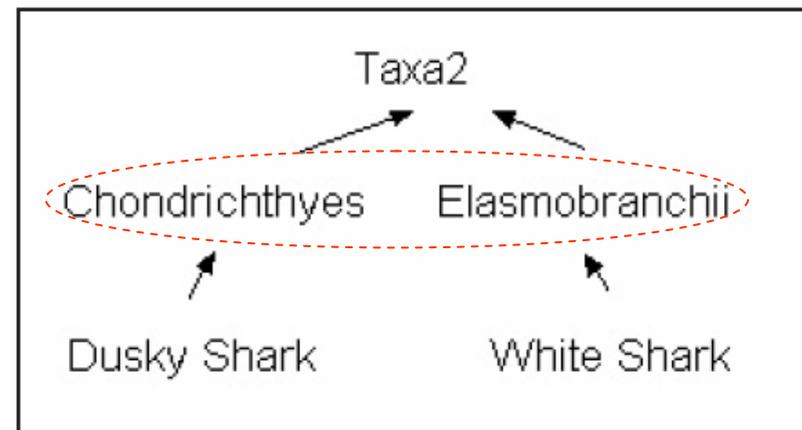
K. Selcuk Candan          Arizona State University

Maria Luisa Sapino        University of Torino

# Motivation: Query Processing on Metadata with Conflicts (FICSR[SIGMOD07])
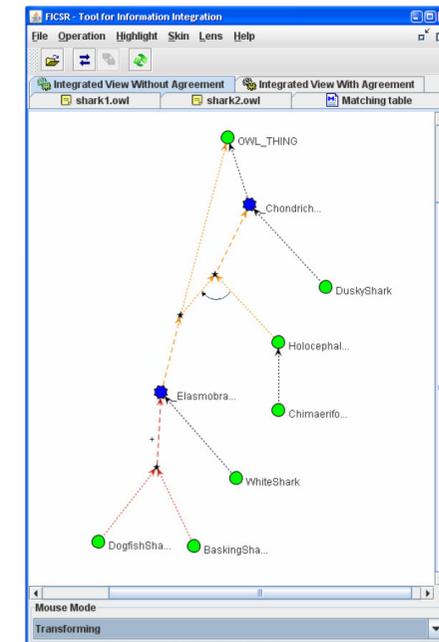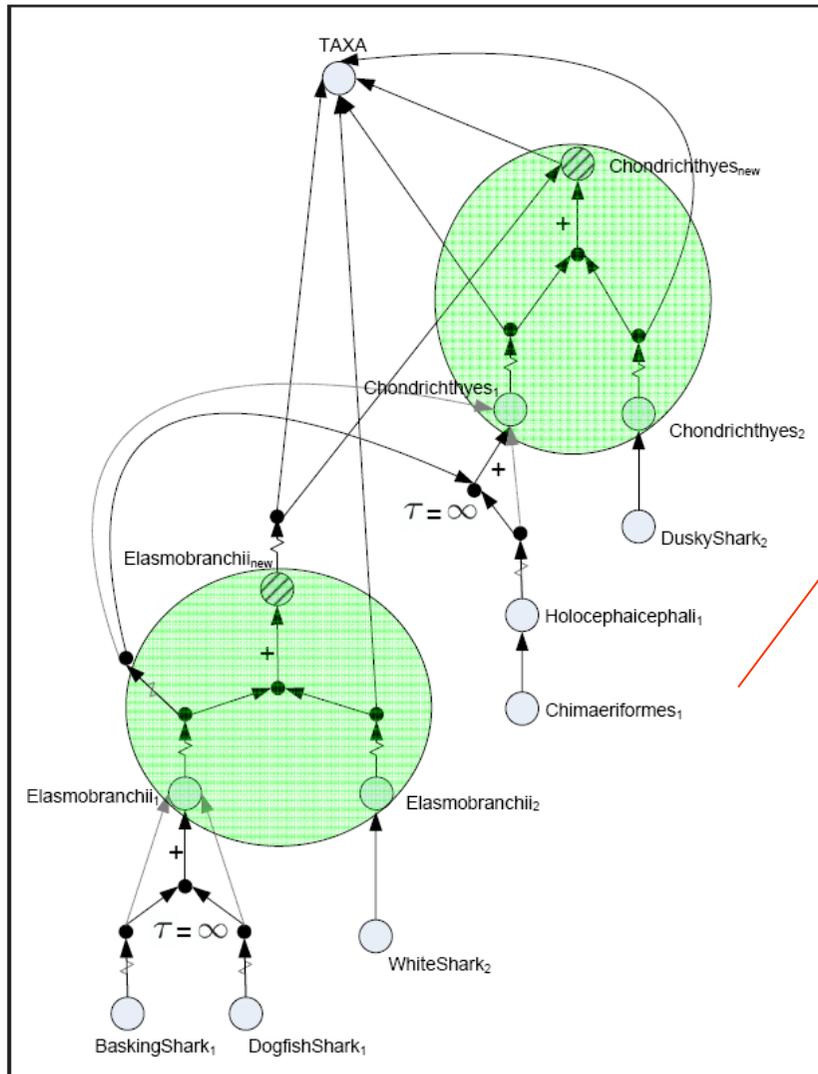


Taxonomy1



Taxonomy2

# FICSR Integrated Representation



Internal FICSR Representation

Simplified Visualization
for the User

# FICSR Agreements [based on source analysis and user feedback]



**Statement of interest:** is it true that "**//Chondrichthyes[//Chimaeriformes]//WhiteShark**"?
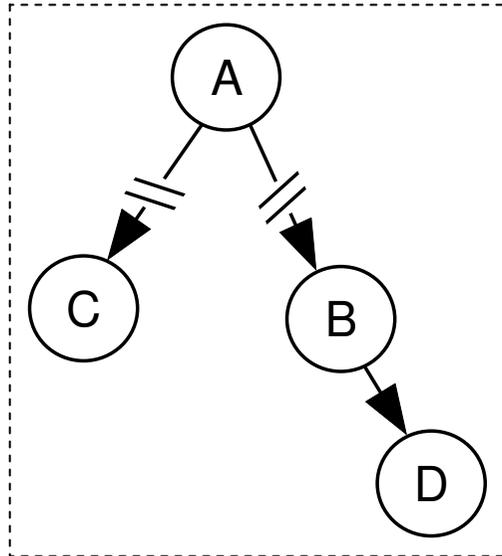
# Related problems

- **Web information units** [RIU ,Banks, DPBF]

- **Keyword search in Relational/OO/XML databases**
  [XRank, ObjectRank, Banks, CP/CV, DPBF]
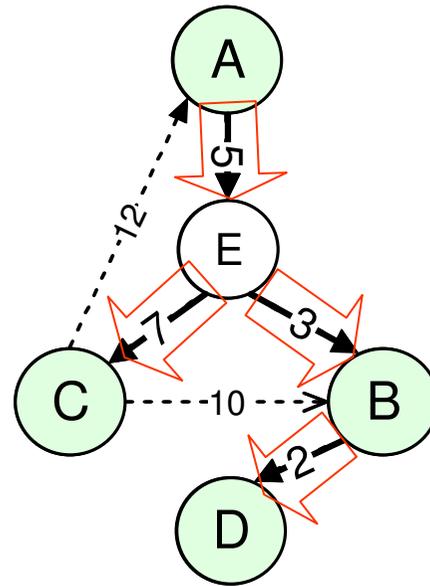
- **Social network analysis** [CDIP]

Common theme:
- Data is a graph…
- …relevant content is distributed across the graph…
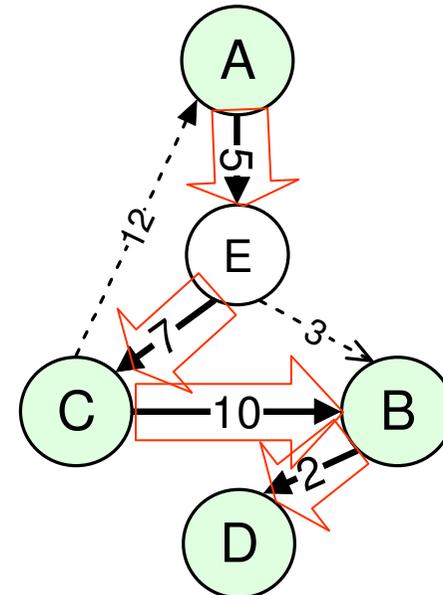- …but, queries (e.g. keyword sets) are not structured.

# Twig queries (i.e., structure of interest) and top-k results
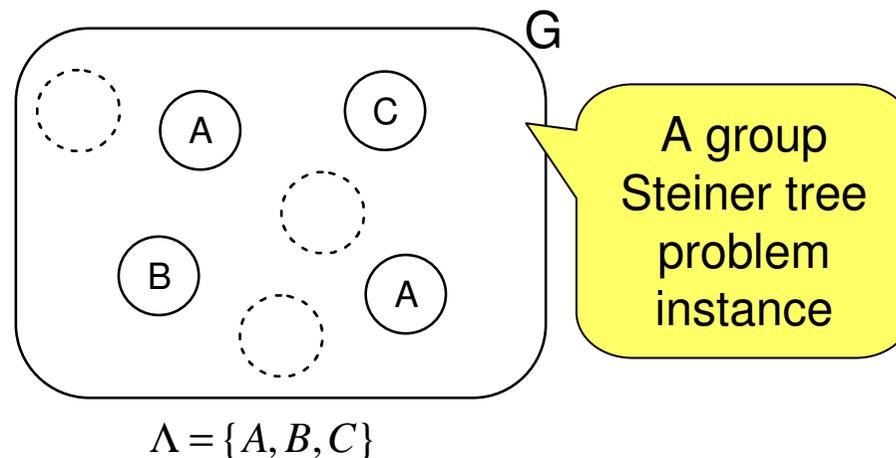


A[//C]//B/D

Cost = 17

Cost = 24

More desirable….can we find it before the other??
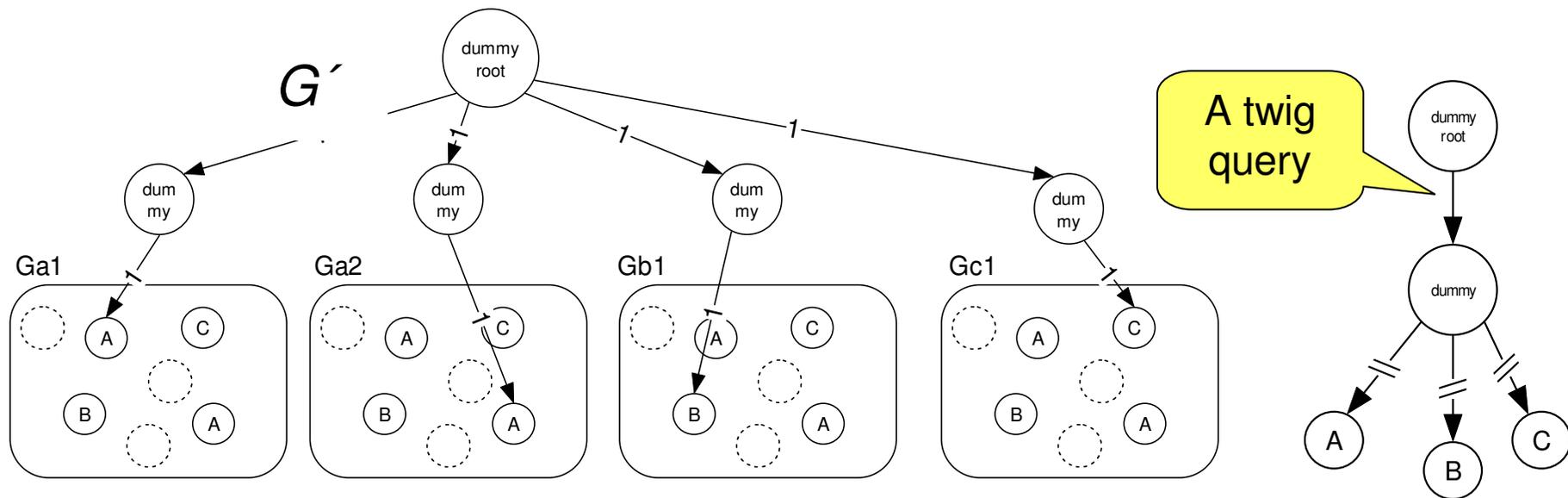
# Answering twig queries on weighted graphs

- So, how hard is the "min-cost twig query problem"?
- NP-complete (by reduction from the "*group Steiner tree* problem*")

# Answering twig queries on weighted graphs

- So, how hard is the "min-cost twig query problem"?
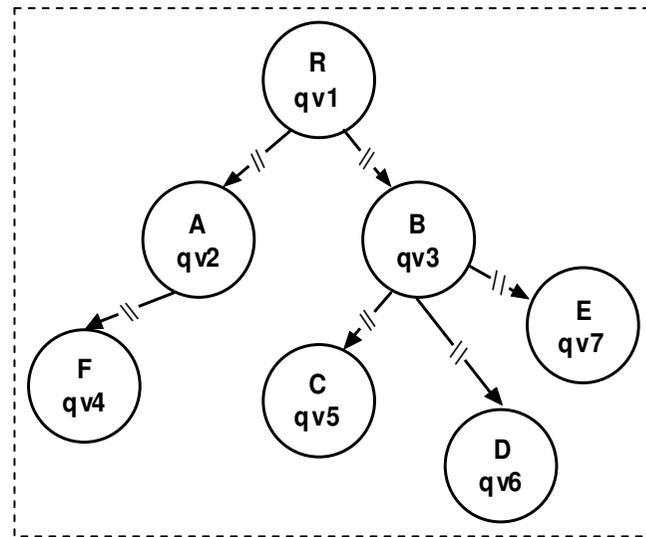- NP-complete (by reduction from the "*group Steiner tree* problem*")



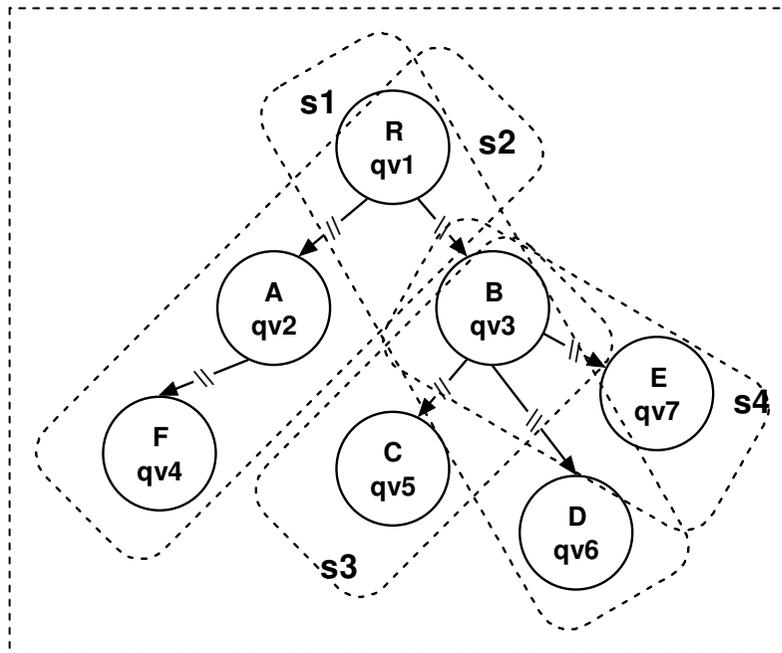Also see DBTwig results (Kimelfeld and Sagiv, 06)

# So what can we do?

- Keyword search on graph data [RIU, BANKS] ?
  - No…we need to enforce query structure..

# So what can we do?

- **Keyword search on graph data** [RIU, BANKS] ?
  - No…we need to enforce query structure..
- **Ranked-join algorithms** (FA,TA, NRA) for top-k queries

$$Q = s_1 \bowtie s_2 \bowtie s_3 \bowtie s_4$$

# So what can we do?

- Keyword search on graph data [RIU, BANKS] ?
  - No…we need to enforce query structure..
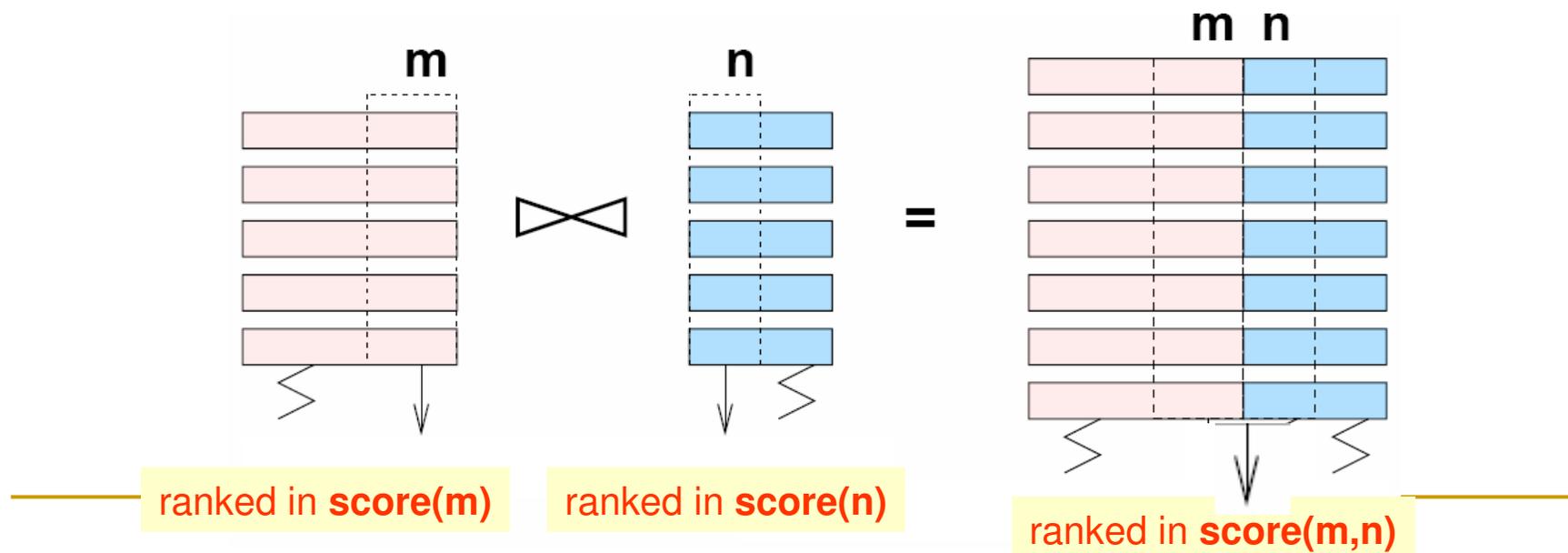- Ranked-join algorithms (FA,TA, NRA) for top-k queries
  - ….score combination function must be monotonic.



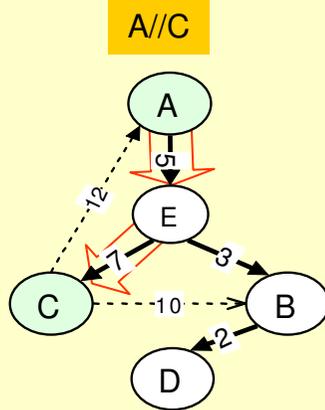ranked in **score(m)**          ranked in **score(n)**          ranked in **score(m,n)**

# Sum-Max Monotonicity

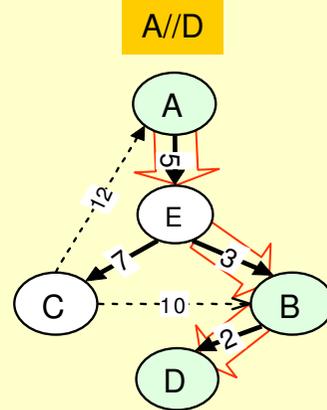- Ranked joins is a good idea…
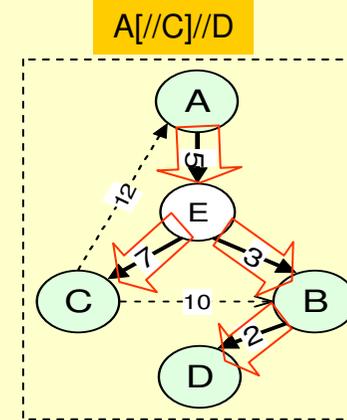  - ..but, monotonicity does not hold.
- Good news: Sum-Max monotonicity

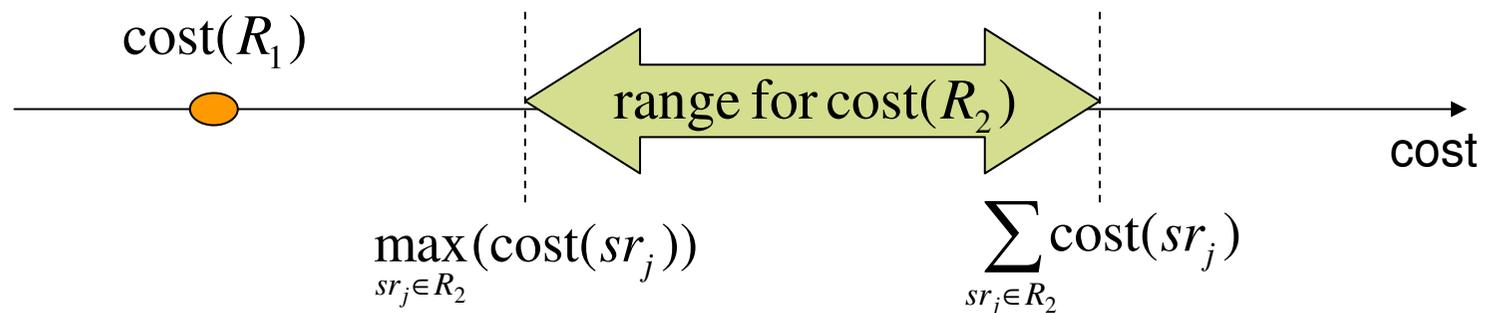$$\max_{sr_i \in R}(\text{cost}(sr_i)) \leq \text{cost}(R) \leq \sum_{sr_i \in R} \text{cost}(sr_i)$$



Cost = 12    Cost = 10    max(10,12)=12  <  Cost =17  <  12+10 = 22

# Sum-Max Monotonicity

- **Ranked joins is a good idea…**
  - ..but, monotonicity does not hold.
- **In fact, we can also see that**
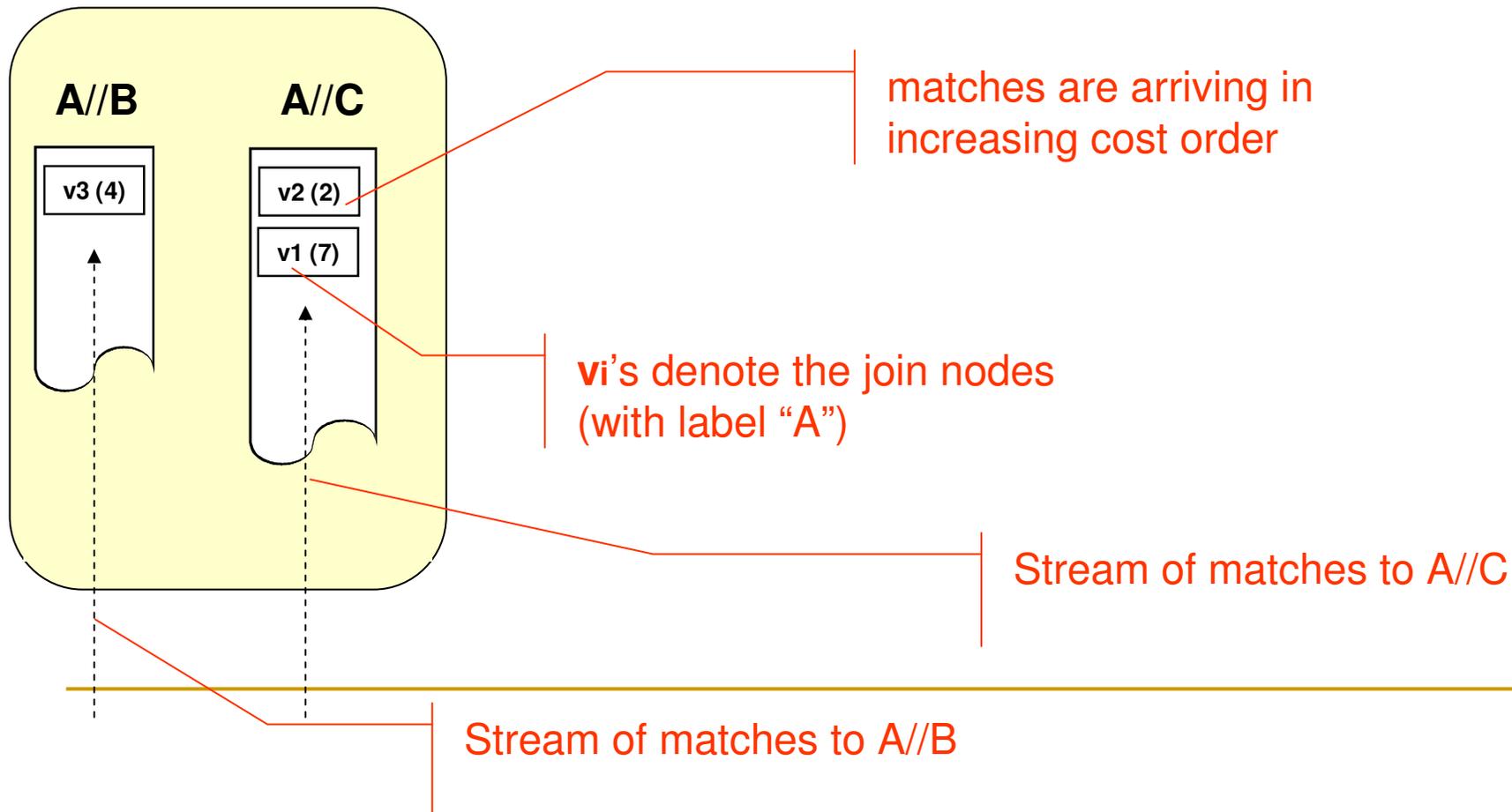
$$\left( \mathrm{cost}(R_1) \leq \max_{sr_j \in R_2}(\mathrm{cost}(sr_j)) \right) \rightarrow \mathrm{cost}(R_1) \leq \mathrm{cost}(R_2)$$

# Progressive enumeration based on Sum-Max monotonicity

- Query: A[//B]//C

# Progressive enumeration based on Sum-Max monotonicity

- Horizon -> Stopping criterion



(a) Horizon = ∞

(b) Horizon = 14

(c) Horizon = 11

(d) Horizon = 14

No

Match with cost 11

Match with cost 11 returned

**HORIZON TIGHTENING**

**HORIZON RELAXATION**

# Pruning and data overhead



$$dataoverhead = \frac{\#all\_submatches - \#necessary\_submatches}{\#all\_submatches}$$

# HR-Join: Horizon based Ranked Join

**HR-Join**

O

Ordered output stream

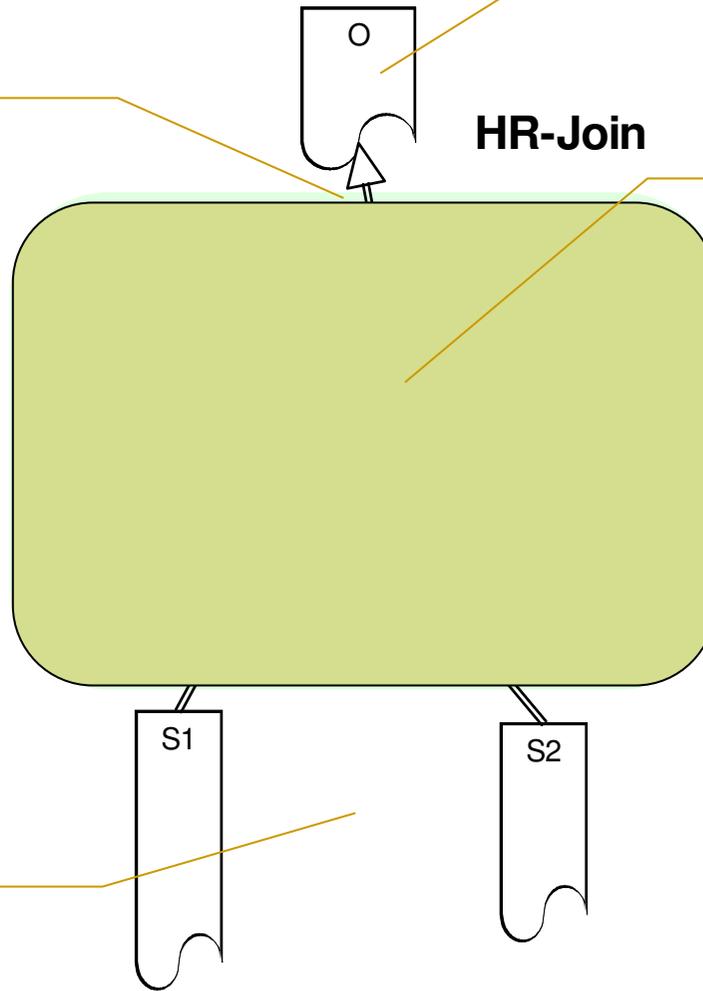Result sieve controls when a candidate match be declared "result"

Implements sum-max pruning

**Relies on the "punctuations" received from the horizon valves**

S1    S2

Horizon valve regulates the availability of the incoming data to the hash join.

Two ordered streams

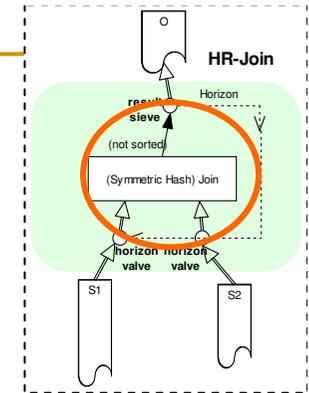**Relies on the "horizon" value from result sieve**

# Operation of the horizon valve

•Punctuation: blocking due to the horizon limit

17

19

20

21

27

29

# Punctuations are propagated by the symmetric-hash join
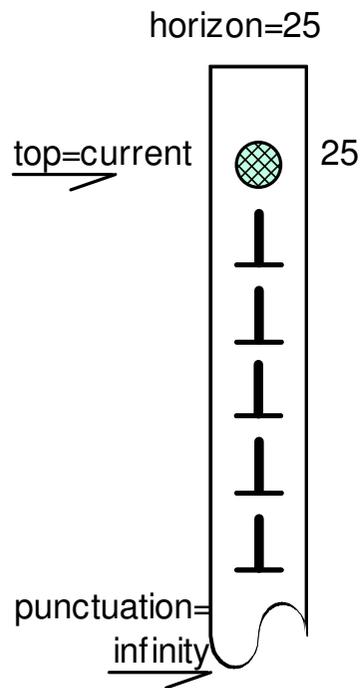
**HR-Join**

Horizon

result
sieve
(not sorted)

(Symmetric Hash) Join

horizon horizon
valve valve

S1  S2

O

punctuation

Symettric Hash Join

Punctuations
in the input
streams are
not met yet

used

used

punctuation
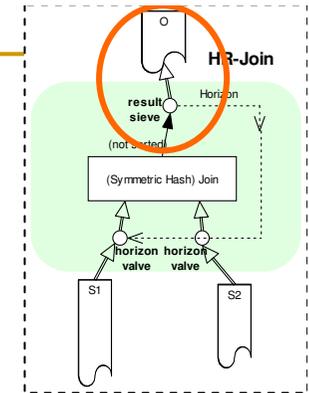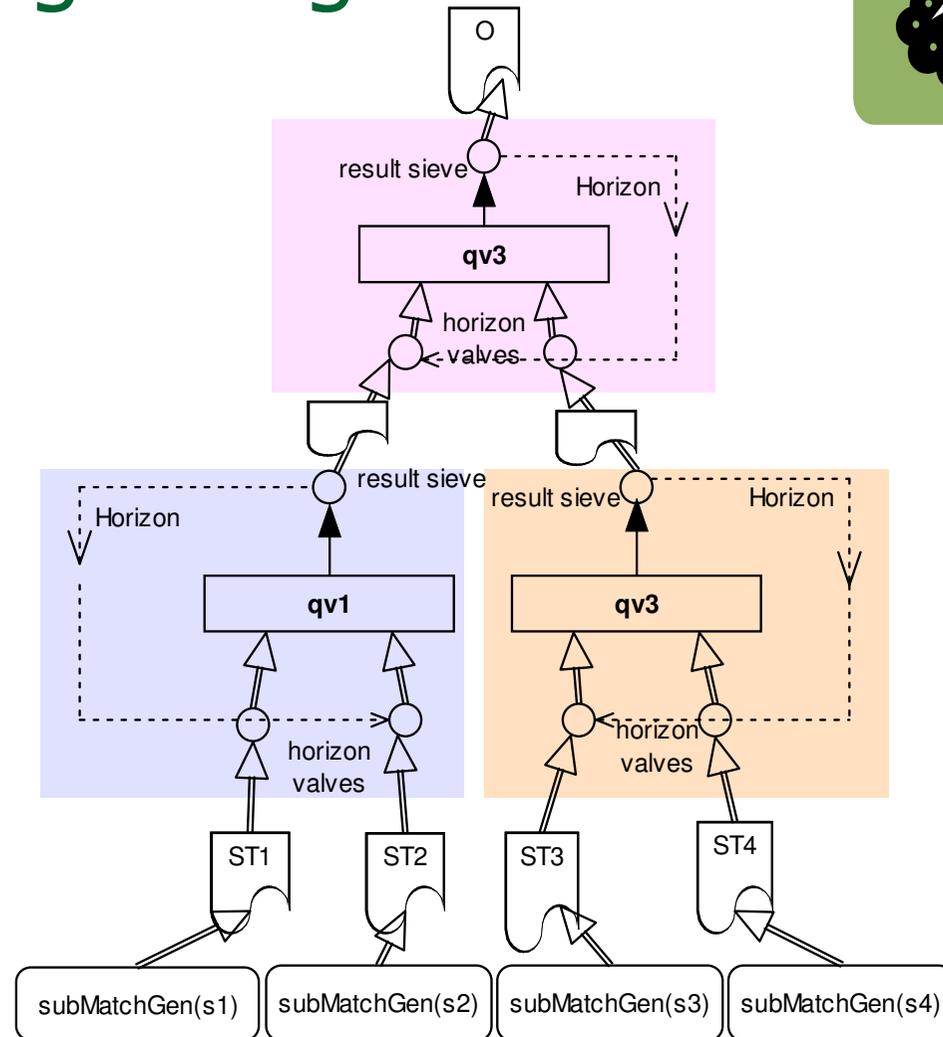
punctuation

(a)

# Operation of the result sieve

- Top: The current best
- Punctuation: indicates that the input streams are punctuated

horizon=25

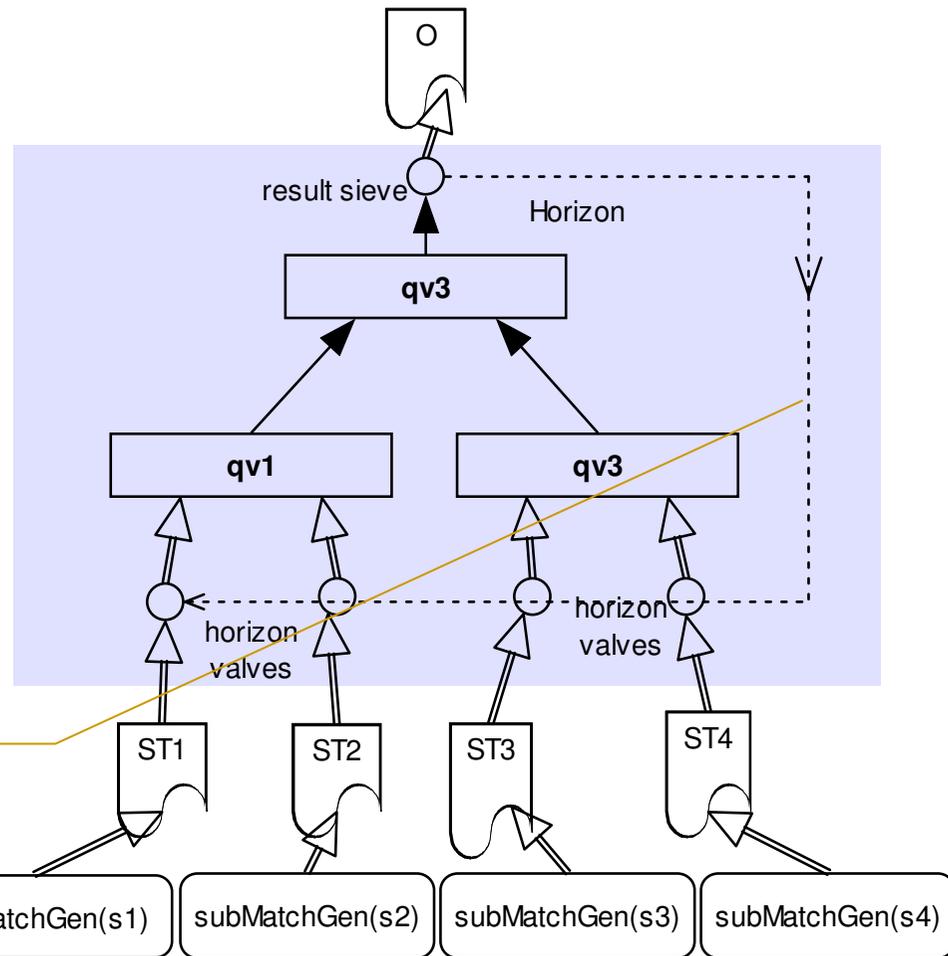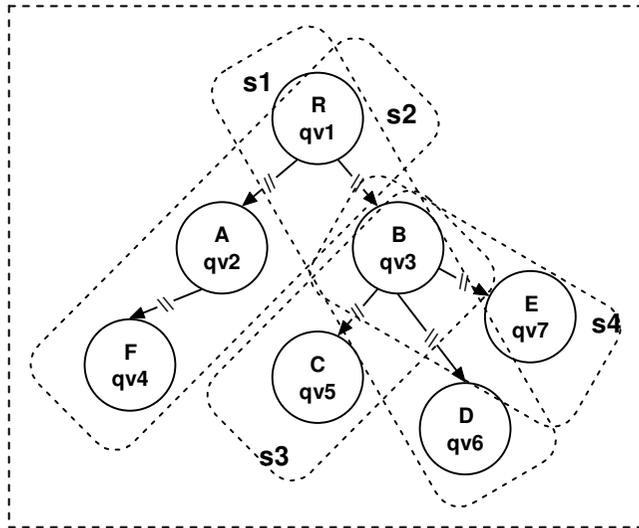top=current          25

punctuation=
infinity

Indirectly regulates its own input stream by updating the horizon value

# Query processing using HR-Join



not ordered →     cost-ordered ⇒

# M-way HR-Joins



Only one result sieve regulating all valves of the four input streams.

# What is missing?

- ## How to enumerate (subresult) paths in cost order?
  - K-shortest simple paths problem [Qi et al. SIGMOD 07]
    $$O(k \,|V\,|\,(|\,E\,|+|V\,|\log|V\,|))$$
    - ..details are in the paper

- ## How to deal with "*" wildcards in twigs??
  - can be expensive (too many matches and joins)
  - query rewriting…......details are in the paper

# Can we do better?

- Horizon values are set based on

$$\left(\operatorname{cost}(R_1) \le \max_{sr_j \in R_2}(\operatorname{cost}(sr_j))\right) \to \operatorname{cost}(R_1) \le \operatorname{cost}(R_2)$$

which assumes the worst case:

i.e., subresults may overlap fully.

- Horizon tightening factor (tf) can be used when overlaps are known to be bounded

$$\left(\operatorname{cost}(R_1) \le tf \max_{sr_j \in R_2}(\operatorname{cost}(sr_j))\right) \to \operatorname{cost}(R_1) \le \operatorname{cost}(R_2)$$
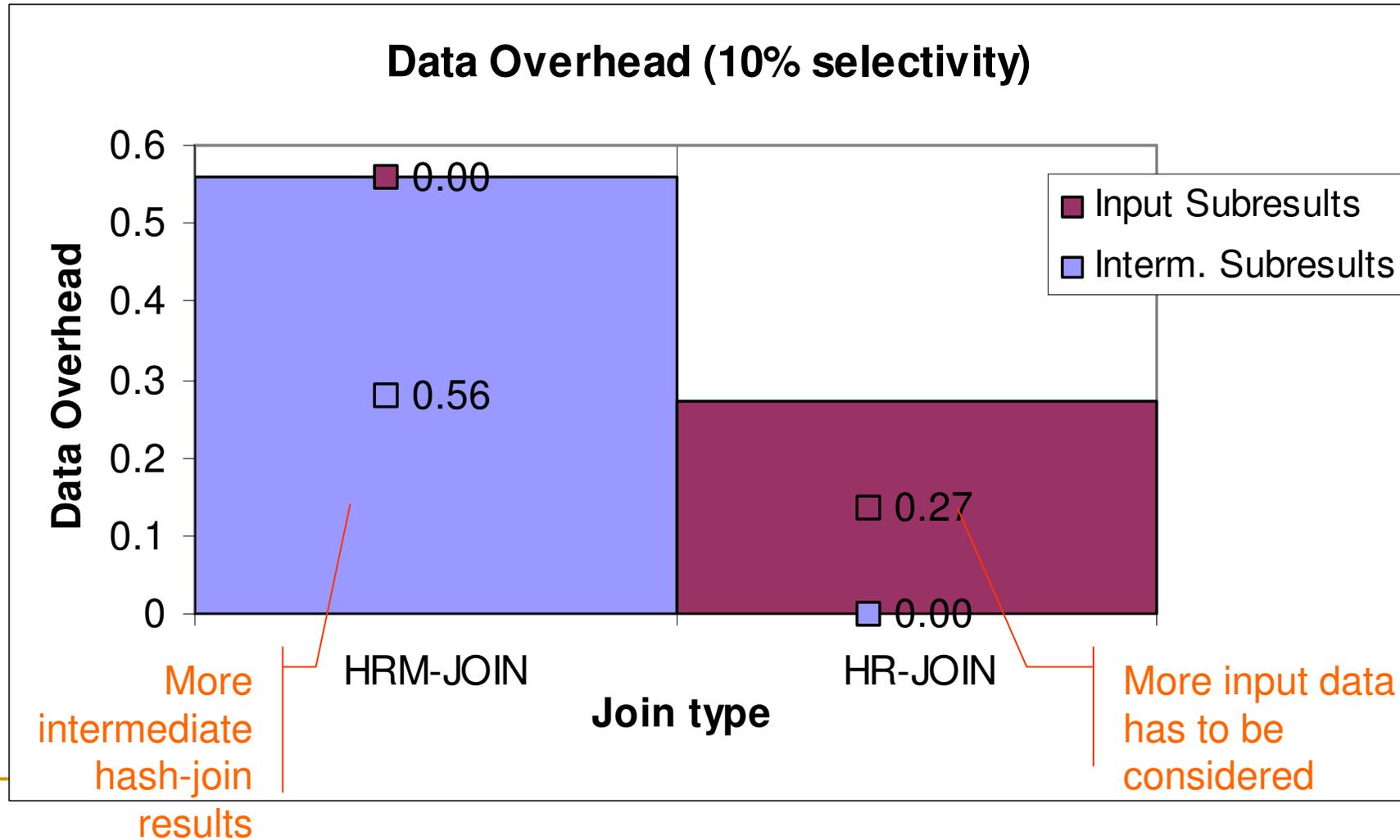
# Experiments

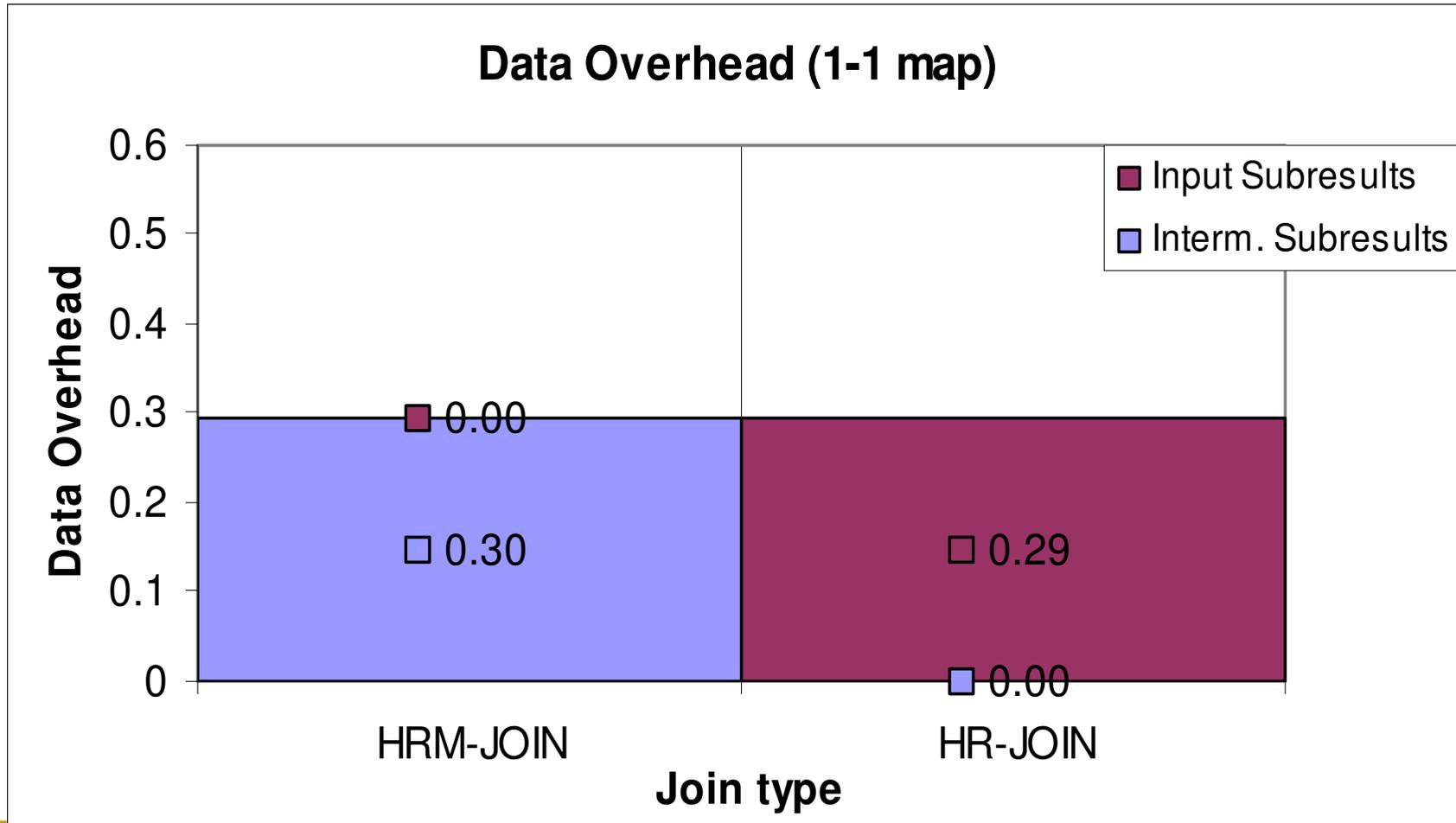- ## 2GHz Pentium with 1GB main memory.

- ## Query plans:
  - HR-Join and M-Way HR-Join (MHR-Join)
  - 2 significantly different join-selectivity distributions: ~10% and 1-to-1.

- ## Data
  - FICSR weighted graph data
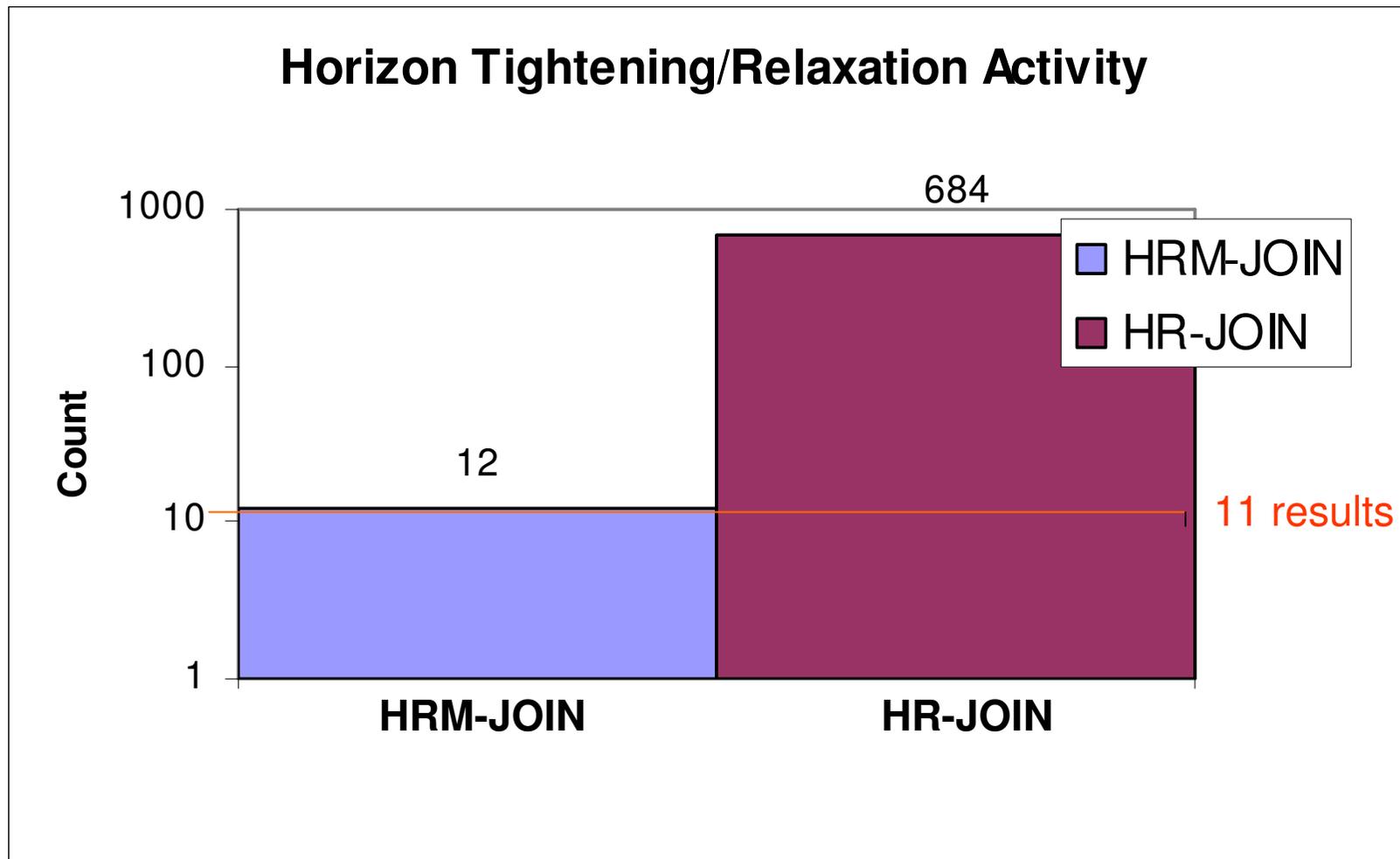    - Zipfian-like distribution of edge weights
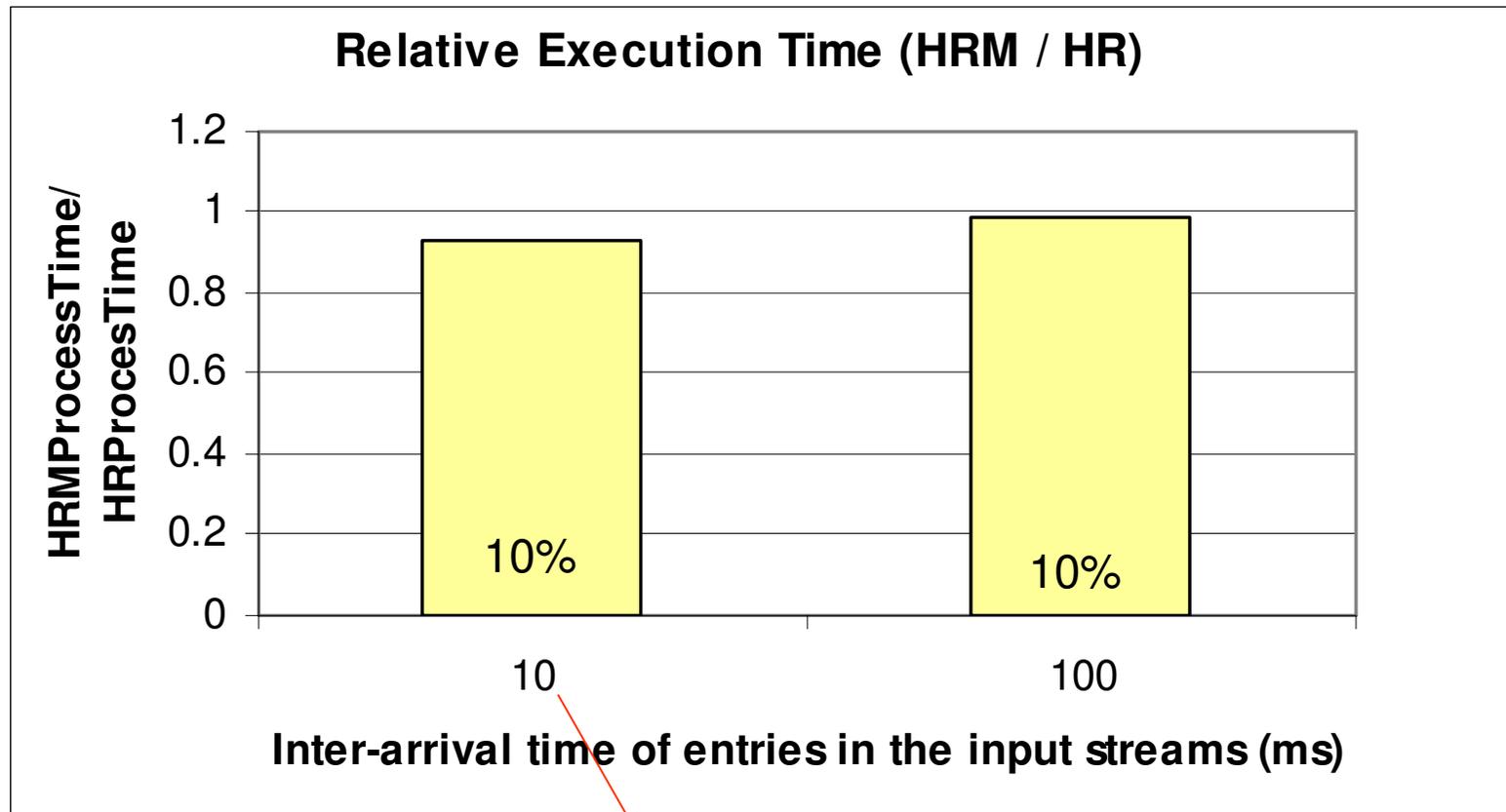
# Data overhead of HRM-Join versus HR-Join



**Data Overhead (10% selectivity)**

Data Overhead vs Join type

- Input Subresults
- Interm. Subresults

HRM-JOIN: ■ 0.00, □ 0.56

HR-JOIN: □ 0.27, ■ 0.00

More intermediate hash-join results

More input data has to be considered

# Data overhead of HRM-Join versus HR-Join

# HR-Join has higher horizon-management cost (10%)



**Horizon Tightening/Relaxation Activity**
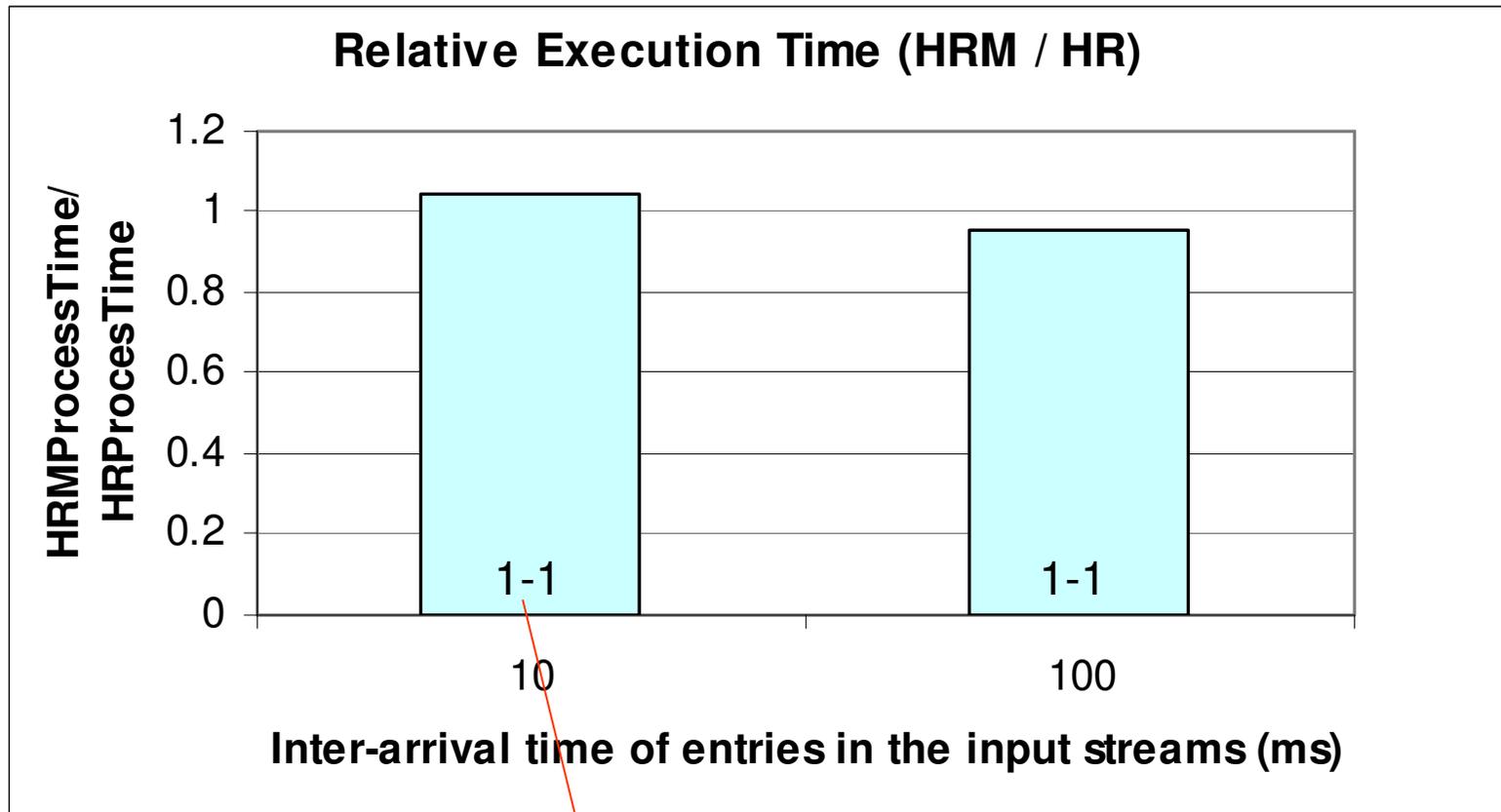
# Inter-arrival time of stream inputs



**Relative Execution Time (HRM / HR)**

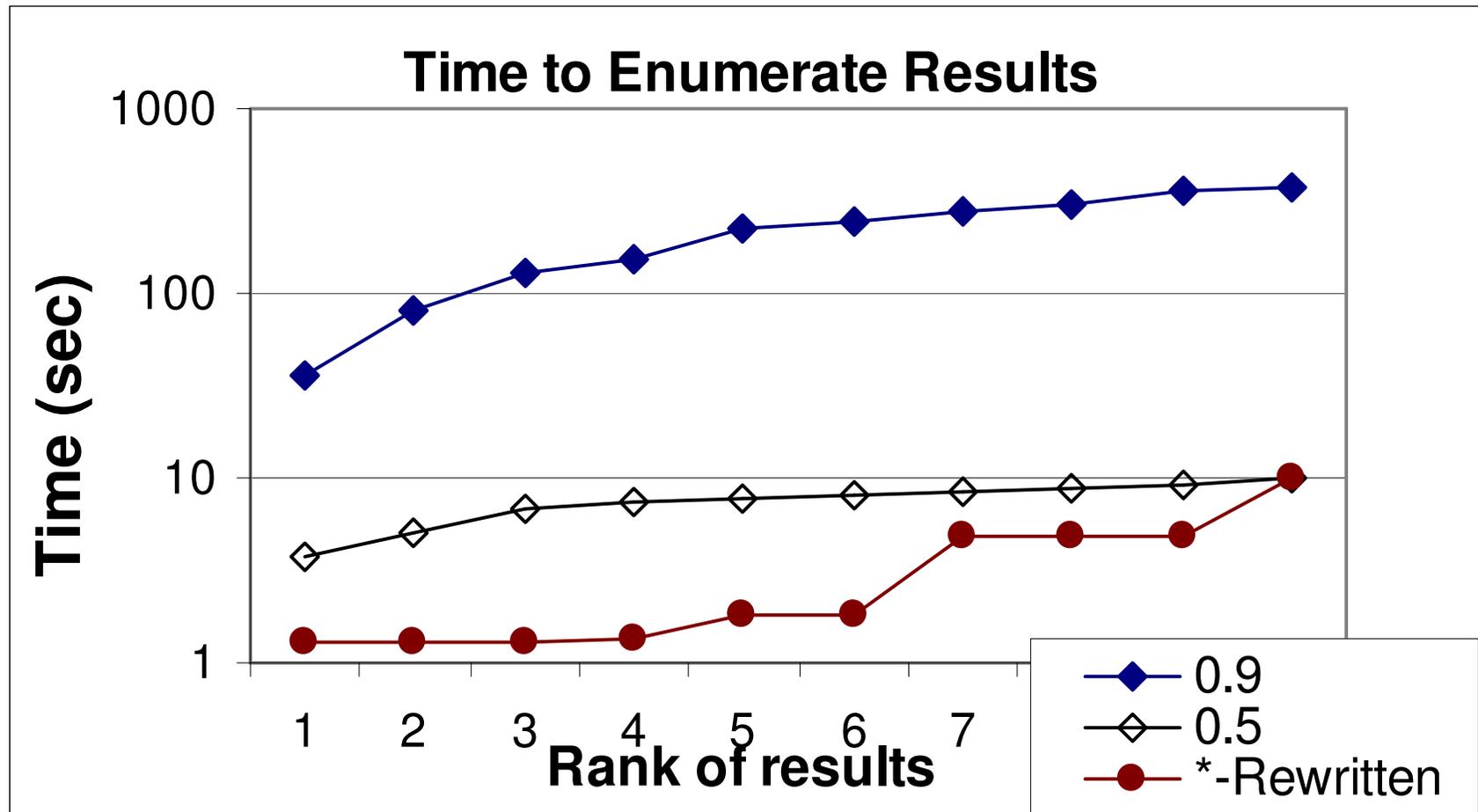When inter-arrival is tight, HRM benefits from reduced horizon management

# Inter-arrival time of stream inputs

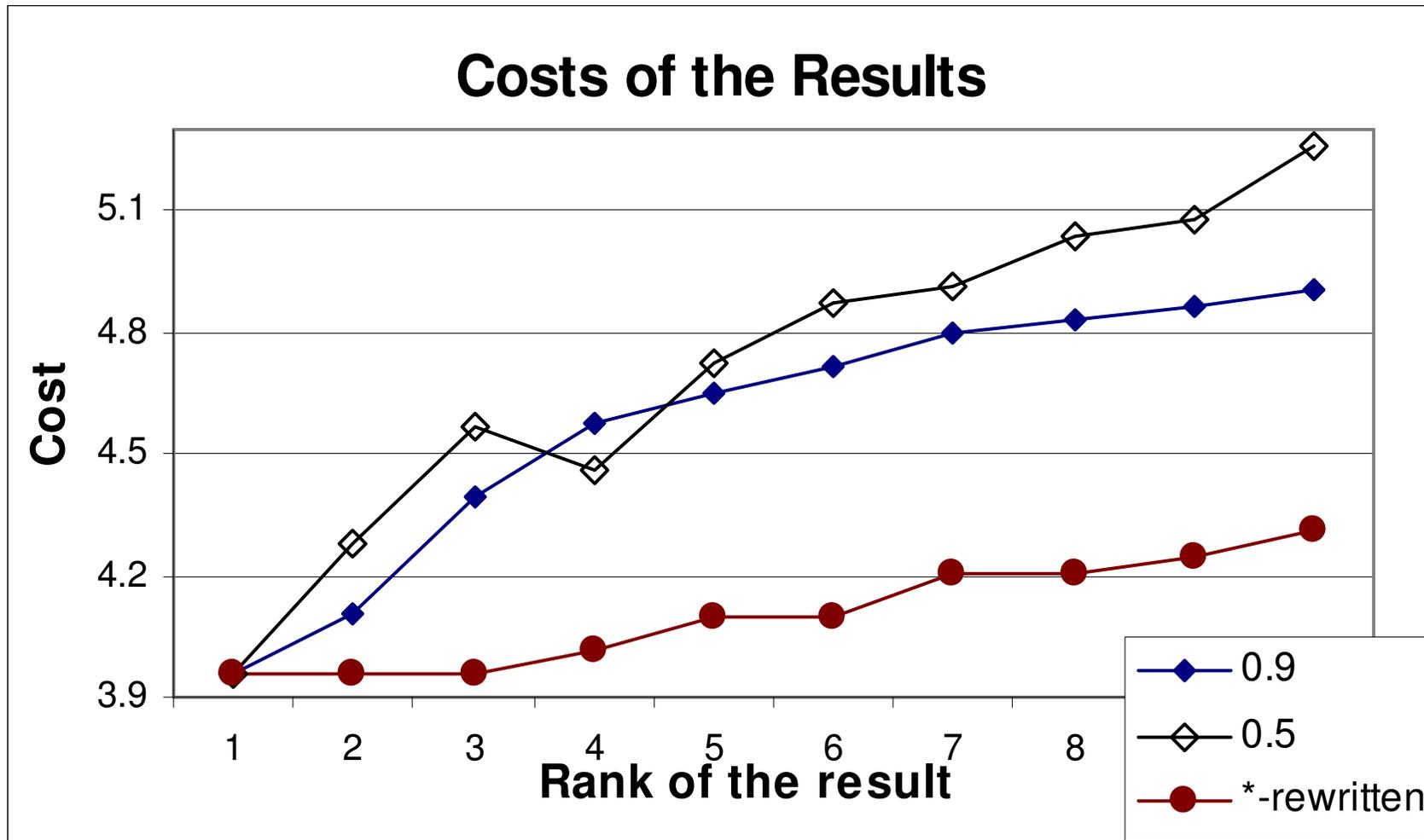**Relative Execution Time (HRM / HR)**



When selectivity is tight, there is less horizon tightening, thus HR-Join performs OK

# Horizon tightening and *-rewriting help with "wildcard" queries



Query: A//*[//B]//C

# The costs of the distinct results of the rewritten query are significantly better
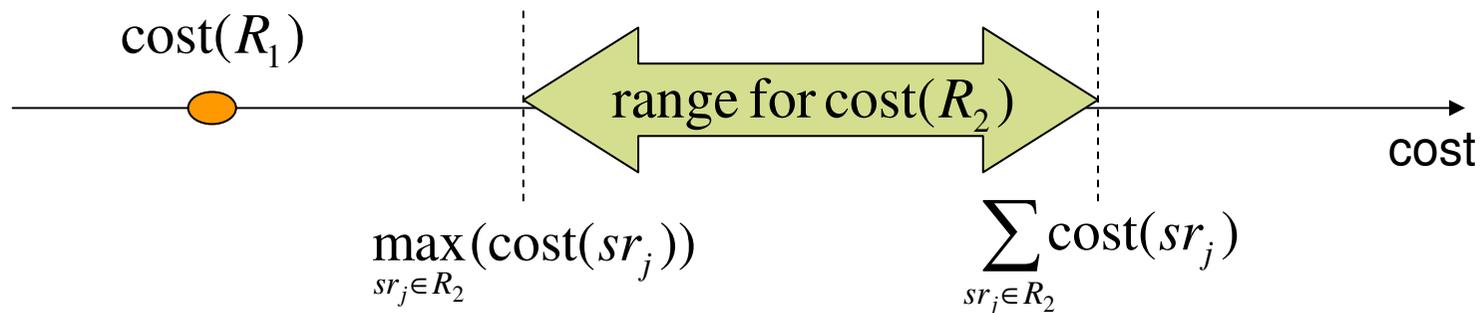


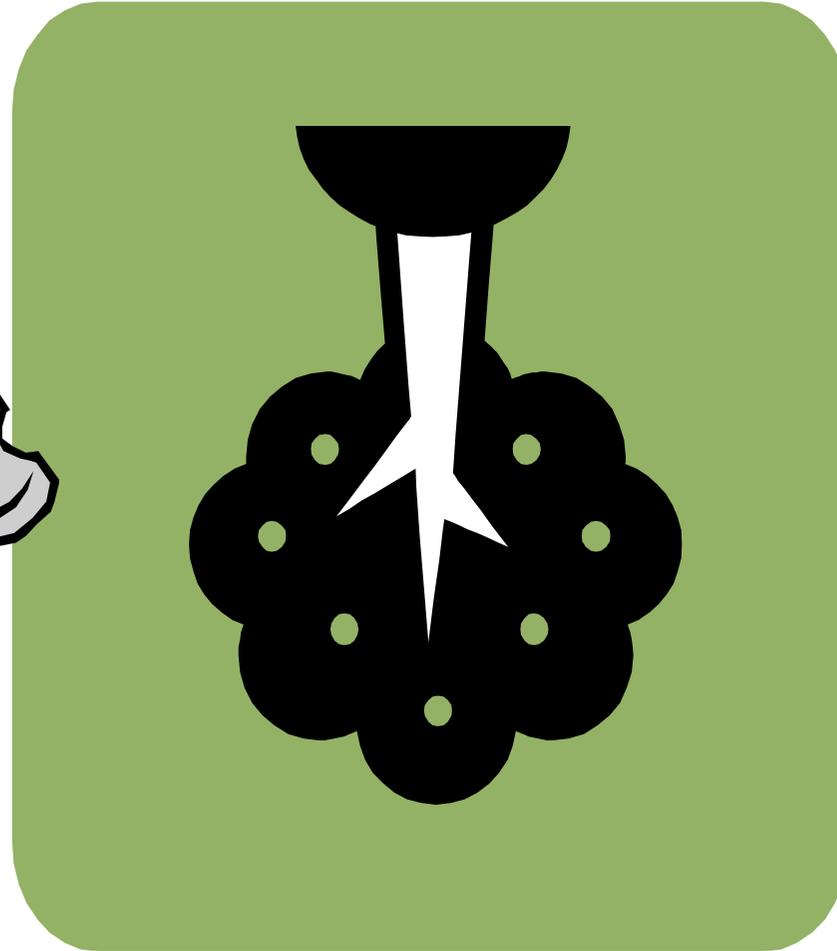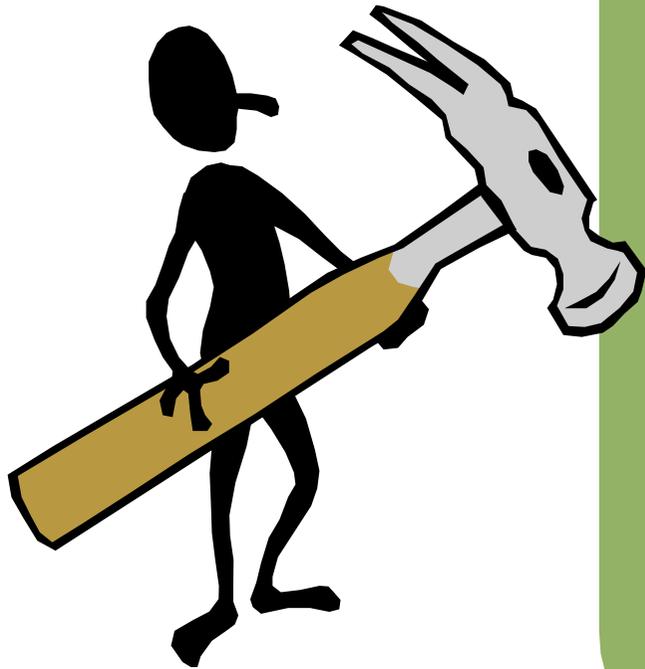**Costs of the Results**

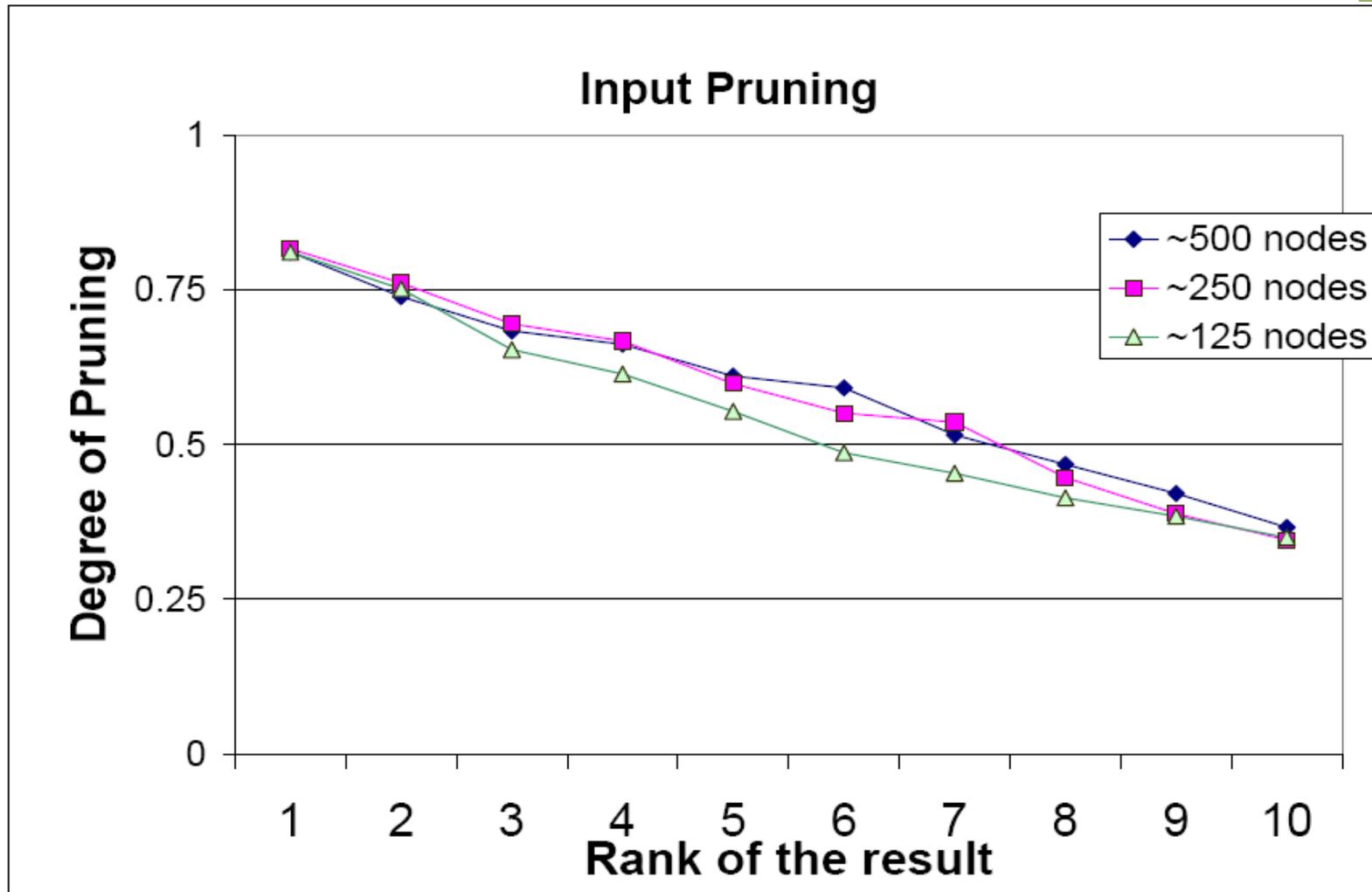Query: A//*[//B]//C

# Conclusion

- Sum-max monotonicity…



- …..a self-punctuating, horizon-based ranked join operator (binary, m-way)…

- …optimizations…

- Twig query processing over weighted data graphs

# Questions?

# The degree of pruning is directly correlated with the size of k



**Input Pruning**

# The degree of input pruning is more important in bigger graphs



Time Between Consecutive HR-Join Results