

Top-k Spatial Joins of Probabilistic Objects

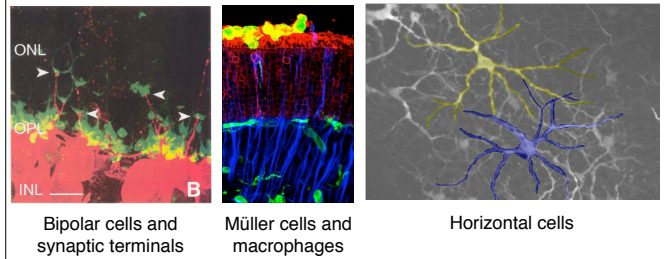
Vebjorn Ljosa
Broad Institute of MIT and Harvard

Ambuj K. Singh
University of California, Santa Barbara

April 9, 2008 , Cancún, México
Geeks Gone Wild: International Conference on Data Engineering (ICDE)

Many biological questions are really spatial joins!

Examples from neuroscience:

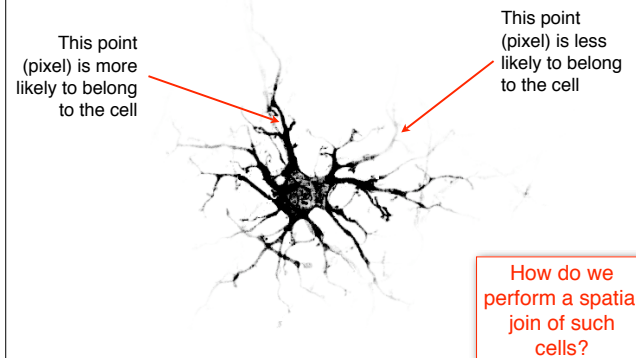


[Geoff Lewis; Mark Verardo] 2

Horizontal cells are hard to segment

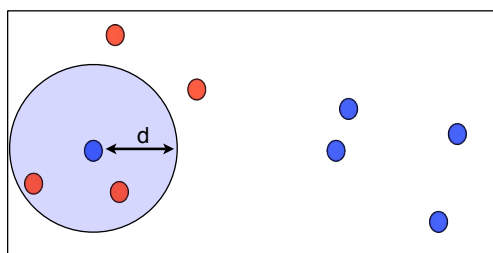
Uncertain extent:
We are not sure
which pixels belong
to the cell

Probabilistic mask of one horizontal cell



[Ljosa and Singh, ICDM 2006] 4

Spatial joins



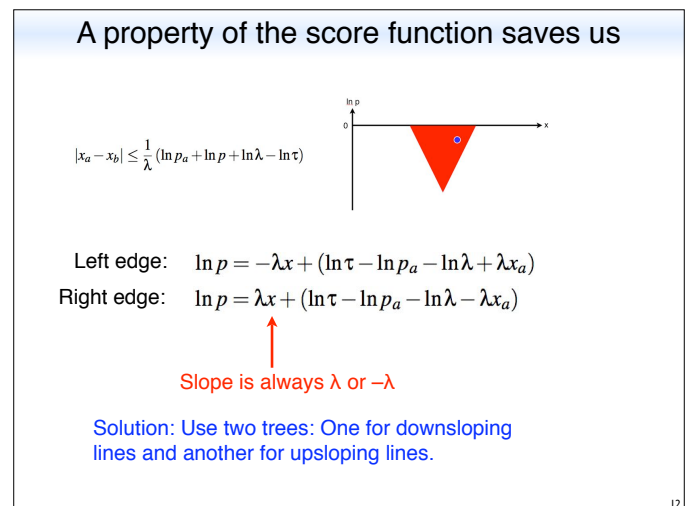
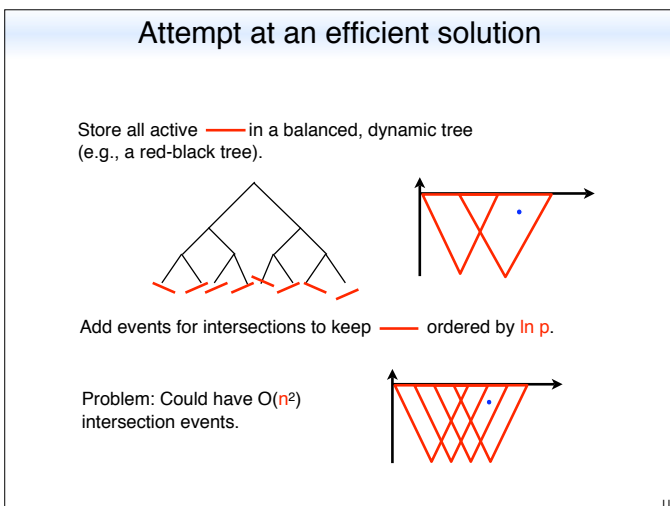
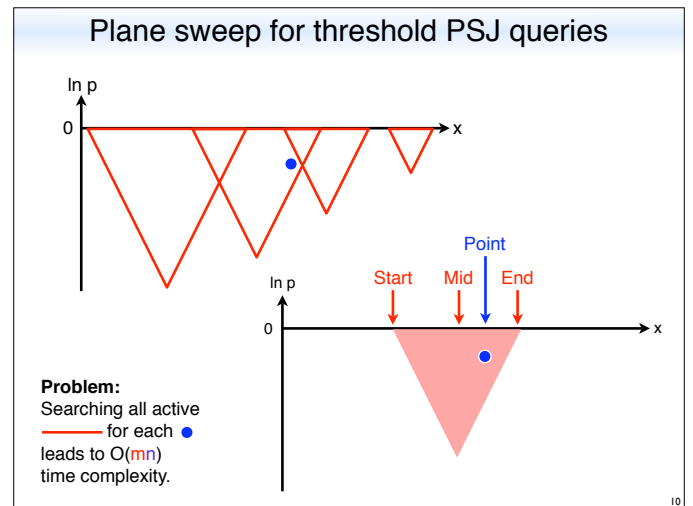
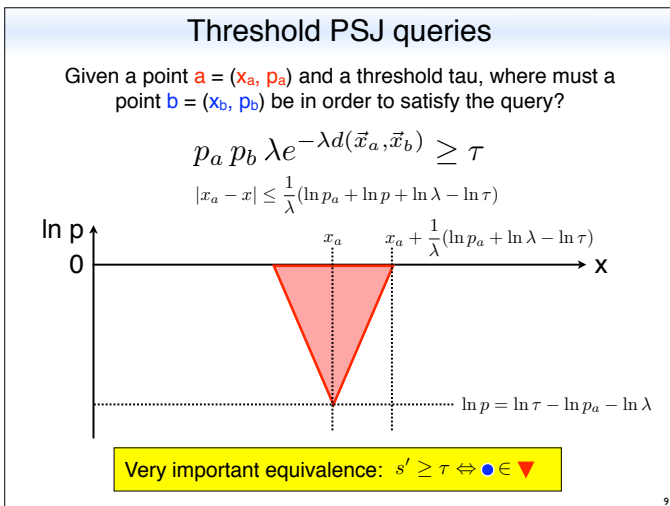
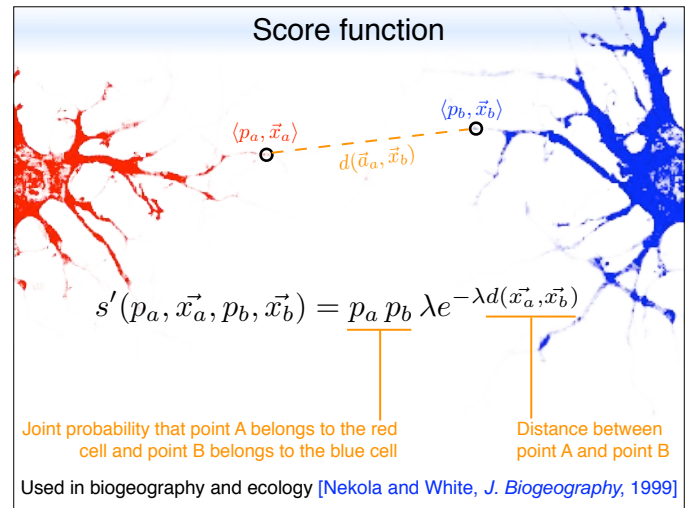
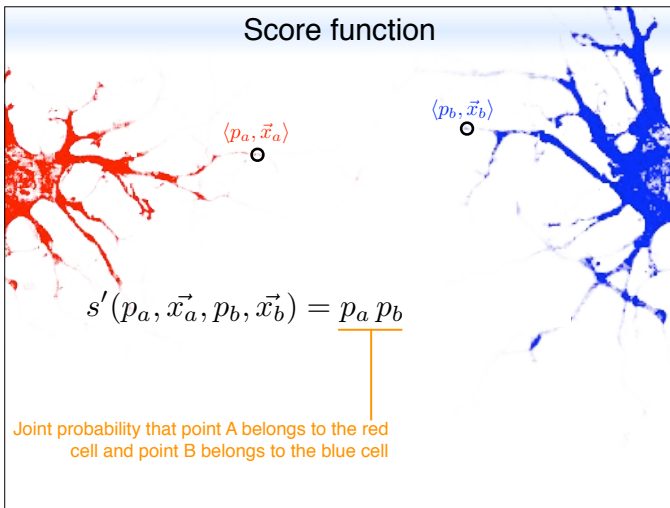
Spatial join of objects with certain extent:
Find all pairs of red and blue points less than d apart

Each ● is either a match or not to ●.

Two types of probabilistic spatial join (PSJ) queries

Threshold PSJ: Given two sets A and B of probabilistic objects, and a score threshold τ , find all pairs (a, b) in $A \times B$ such that $s(a, b) \geq \tau$

Top-k PSJ: Given two sets A and B and a natural number k , find a set $R \subseteq A \times B$ of size k such that other pairs in $A \times B$ score no higher than the lowest-scoring pair in R .

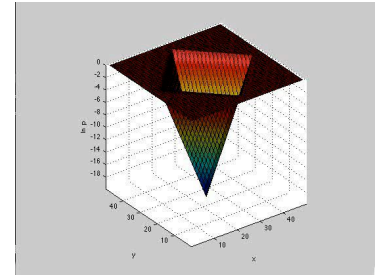


Time complexity

- Sort events: $O(n \log n) + O(m \log m)$
- There are $O(n + m)$ events
 - Processing a start/mid/end event: $O(\log n)$
 - Processing a point event: $O(\log n + k')$
 - k' is the number of results for this point
- Time complexity: $O(m \log m + (n + m) \log n + k)$
 - k is total number of results
 - If we assume that $m = n$: $O(n \log n + k)$

13

Algorithm generalizes to multiple dimensions

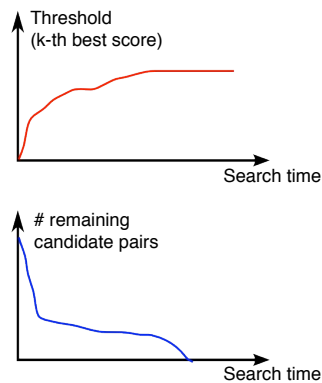


In 2D: Pyramid instead of triangle

14

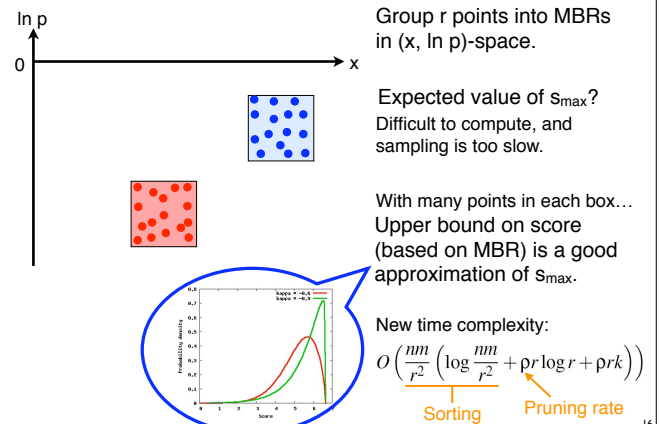
Plane sweep for top-k PSJ

- Query: Find the k top-scoring pairs.
- Plane sweep algorithms adapt easily
 - Move start and end events as threshold increases.
- Key to efficiency is to find some good pairs early
 - Brings the threshold up
 - Prunes most of the dataset



15

Global scheduling of top-k PSJs



16

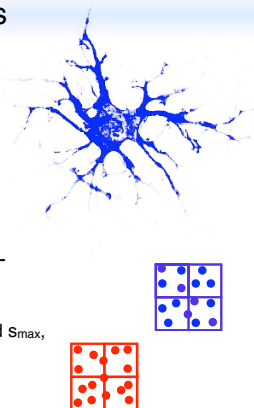
Experiments

Datasets (43k and 52k points) based on horizontal cell images

Increased size synthetically (copy & shift) up to 300 times (from 10^9 pairs to 10^{14} pairs)

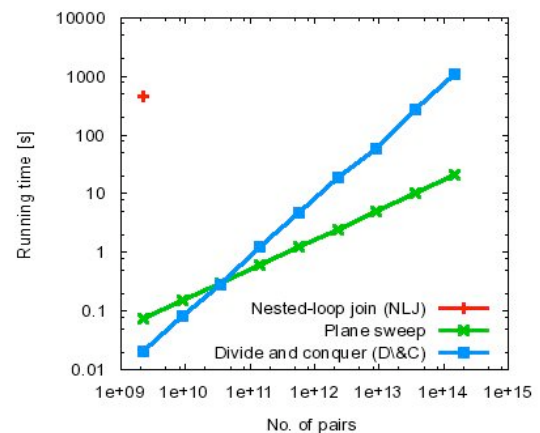
Compared to NLJ and a simple divide-and-conquer technique:

- Split boxes recursively until they
 - can be pruned based on the threshold and s_{\max} ,
 or
 - are small enough to be joined with NLJ



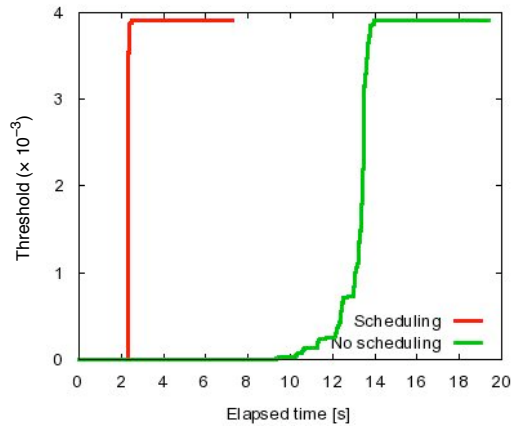
17

Experiments: Scalability of threshold PSJ queries



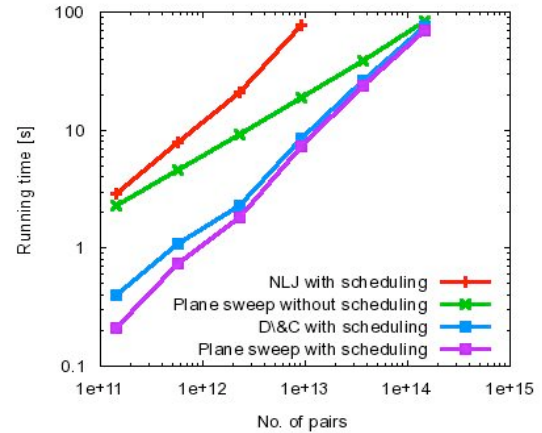
18

Threshold increases faster with scheduling



19

Scalability of top-k PSJ queries



20

Conclusion

Probabilistic spatial joins

- Geographical information systems
- Biomedical image analysis

Technically challenging

- Score depends on not only distance, but on both probabilities
- Finding top-ranking results: spatial join and top-k query at once

Efficient algorithms

- Threshold PSJs and top-k PSJs
- Plane sweeps in $O(n \log n + k)$ time
- Global scheduling: faster top-k by finding high-scoring pairs early

Future work

- Efficient algorithms for more than 2 dimensions
- Compare experimentally to Kriegel et al. [DSFAA 2006]

Acknowledgements: This work was supported in part by grant no. ITR-0331697 from the National Science Foundation. Horizontal cell micrographs were provided by Geoffrey P. Lewis from the laboratory of Steven K. Fisher at UCSB.

21