# Technical Perspective: Optimized Wandering for Online Aggregation

Jeffrey F. Naughton
University of Wisconsin, Madison
naughton@cs.wisc.edu

There is a rich history in the DBMS research literature involving sampling to estimate the results of queries faster than they can be computed exactly. A particularly interesting example of this is "Online Aggregation" proposed by Hellerstein et al. in 1997 [2]. There the idea is to combine sampling with a creative and intuitive user interface. Briefly, when a query starts to run, Online Aggregation will quickly present an estimate of the result of the query (based on data sampled up to that point) and will also present a confidence interval around the estimate. As query execution continues, the estimate is refined, and the confidence interval shrinks.

Hidden in this attractive idea, however, are some difficult challenges. As an example, for queries that involve joins, the sampling process is in general slow, especially if most of the tuples from one relation participating in the join "match" with only a few tuples in the other relation. For 20 years the state of the art approach to this problem has been the "Ripple Join" [1]. The following paper by Li, Wu, Yi, and Zhao presents a highly effective alternative.

The main idea behind the wander join is to use the presence of indexes to speed the sampling, effectively making a random walk through the data join graph. The details of doing this efficiently (both computationally and statistically) are not obvious. The authors of this paper use a clever combination of sampling strategies from the statistical literature and an on-line optimization process to order the paths chosen for the random walk, in the process achieving much better computational and statistical properties than the previously state of the art algorithm. The authors convincingly prove this through experimentation with an open-source implementation in the Postgres database management system.

## References

[1] P. J. Haas and J. M. Hellerstein. Ripple joins for online aggregation. In *Proceedings of the 1999 ACM SIGMOD International Conference on Management of Data*, SIGMOD '99, pages 287–298, New York, NY, USA, 1999. ACM.

[2] J. M. Hellerstein, P. J. Haas, and H. J. Wang. Online aggregation. In *Proceedings of the 1997 ACM SIGMOD International Conference on Management of Data*, SIGMOD '97, pages 171–182, New York, NY, USA, 1997. ACM.