The Case for Learned Spatial Indexes

Varun Pandey TUM pandey@in.tum.de Ibrahim Sabek MIT CSAIL sabek@mit.edu Alexander van Renen TUM renen@in.tum.de Jialin Ding MIT CSAIL jialind@mit.edu Andreas Kipf MIT CSAIL kipf@mit.edu Alfons Kemper TUM kemper@in.tum.de

ABSTRACT

Spatial data is ubiquitous. Massive amounts of data are generated every day from billions of GPS-enabled devices such as cell phones, cars, sensors, and various consumerbased applications such as Uber, Tinder, location-tagged posts in Facebook, Twitter, Instagram, etc. This exponential growth in spatial data has led the research community to focus on building systems and applications that can process spatial data efficiently. In the meantime, recent research has introduced learned index structures. In this work, we use techniques proposed from a state-of-the art learned multi-dimensional index structure (namely, Flood) and apply them to five classical multi-dimensional indexes to be able to answer spatial range queries. By tuning each partitioning technique for optimal performance, we show that (i) machine learned search within a partition is faster by 11.79% to 39.51% than binary search when using filtering on one dimension, (ii) the bottleneck for tree structures is index lookup, which could potentially be improved by linearizing the indexed partitions (iii) filtering on one dimension and refining using machine learned indexes is 1.23x to 1.83x times faster than closest competitor which filters on two dimensions, and (iv) learned indexes can have a significant impact on the performance of low selectivity queries while being less effective under higher selectivities.

1. INTRODUCTION

With the increase in the amount of spatial data available today, the database community has devoted substantial attention to spatial data management. For e.g., NYC Taxi Rides open dataset [31] consists of pickup and drop-off locations of more than 2.7 billion rides taken in the city since 2009. This represents more than 650,000 taxi rides every day in one of the most densely populated cities in the world, but is only a sample of the location data that is captured by many applications today. Uber, a popular ride hailing service available via a mobile application, operates on a global



Figure 1: Machine Learning vs. Binary Search. For low selectivity (0.00001%), the index and refinement phases dominate, while for high selectivity (0.1%), the scan phase dominates (parameters are tuned to favor Binary Search).

scale and completed 10 billion rides in 2018 [47]. The unprecedented rate of generation of location data has led to a considerable amount of research efforts that have been focused on, systems that scale out [1, 2, 8, 9, 15, 40, 41, 42, 49, 51, 52], databases [12, 27, 28, 32, 34], improving spatial query processing [11, 19, 20, 21, 36, 43, 44, 45, 54, 35], or leveraging modern hardware and compiling techniques [6, 7, 38, 37, 39, 53], to handle the increasing demands of applications today.

Recently, Kraska et al. [24] proposed the idea of replacing traditional database indexes with learned models that predict the location of a key in a sorted dataset, and showed that learned models are generally faster than binary search. Kester et al. [17] showed that index scans are preferable over optimized sequential scans in main-memory analytical engines if a query selects a narrow portion of the data.

In this paper, we build on top of these recent research results, and provide a thorough study for the effect of applying ideas from learned index structures (e.g., Flood [29]) to classical multi-dimensional indexes. In particular, we focus on five core spatial partitioning techniques, namely Fixedgrid [4], Adaptive-grid [30], Kd-tree [3], Quadtree [10] and STR [25]. Typically, query processing on top of these partitioning techniques include three phases; index lookup, refinement, and scanning (Details of these phases are in Section 2.1). We propose to replace the typical search techniques used in the refinement phase (e.g., binary search)

This article is published under a Creative Commons Attribution License (http://creativecommons.org/licenses/by/3.0/), which permits distribution and reproduction in any medium as well allowing derivative works, provided that you attribute the original work to the author(s) and AIDB 2020. 2nd International Workshop on Applied AI for Database Systems and Applications (AIDB'20), August 31, 2020, Tokyo, Japan.



with learned models (e.g., RadixSpline [22]).

Interestingly, we found that, by using a learned model as the search technique, we can gain a considerable speedup in the query run-time, especially for low selectivity range queries (Similar to the observation from Kester et al. [17]). Figure 1 shows the average running time of a range query using Adaptive-grid on a Tweets dataset, which consists of 83 million records (Section 3.1), with and without learning. It can be seen that for a low selectivity query (which selects 0.00001% of the data, i.e., 8 records) the index and refinement times dominate the lookup, while for a high selectivity query (which selects 0.1% of the data, i.e., 83 thousand records) the scan time dominates. Another interesting finding from our study is that 1-dimensional grid partitioning techniques (e.g., Fixed-grid) can benefit from the learned models more than 2-dimensional techniques (e.g., Quadtree). Our study will assist researchers and practitioners in understanding the performance of different spatial indexing techniques when combined with learned models.

2. APPROACH

In this section, we first explain how a range query processing has been implemented. Then, we describe the spatial partitioning techniques that we have implemented in our work. We conclude the section by describing the search techniques used within the individual partitions.

2.1 Range Query Processing

A given range query has a lower bound and an upper bound in both dimensions. The task is to materialize all the points that lie within the bounds of the query. Query processing works in three phases:

- Index Lookup: In index lookup, we intersect a given range query using the grid directories (or trees) to find the partitions the query intersects with.
- Refinement: Once the partitions intersected have been determined from the index lookup phase, we use a search technique (Section 2.3) to find the lower bound of the query on the sorted dimension within the partition. There can be various cases on how a query intersects with the partition, and we only consult the search technique when it is actually needed to find the lower bound of the given query on the sorted dimension.. For example, a partition could be fully inside the range query, and in such a case we simply copy all the points in the partition rather than use a search technique.
- Scan: Once the lower bound in the sorted dimension has been determined in refinement, we scan the partition to find the qualifying points on both dimensions.

We stop as soon as we reach the upper bound of the query on the sorted dimension, or we reach the end of the partition.

2.2 Partitioning Techniques

Spatially partitioning a dataset into partitions (or cells), such that the objects within the partitions are also close in space, is known as spatial partitioning. Spatial partitioning techniques can be classified into space partitioning techniques (partitions the embedded space) or data partitioning techniques (partitions the data space). In this paper, we employ Fixed-grid [4], Adaptive-grid [30], and Quadtree [10] as space partitioning techniques; and Sort-Tile-Recursive [25] and K-d tree [3] as data partitioning techniques. Figure 2 illustrates these techniques on a sample of the Tweets dataset used in our experiments (details are in Section 3.1), where sample points and partition boundaries are shown as dots and grid axes respectively.

2.2.1 Fixed and Adaptive Grid

The grid (or cell) methods were primarily designed to optimize retrieval of records from disk and generally they share a similar structure. The grid family imposes a d-dimensional grid on the d-attribute space. Every cell in the grid corresponds to one data page (or bucket) and the data points that fall within a particular cell boundary resides in the data page of that cell. Every cell thus has to store a pointer to the data page it indexes. This mapping of grid cells to data pages is known as the grid directory. The Fixed-grid [4] method requires that the grid subdivision lines to be equidistant. The Grid File [30], or the Adaptive-grid, on the other hand relaxes this restriction. Since the grid subdivision lines are not equidistant in the case of Grid File, it introduces an auxiliary data structure called linear scales, which are a set of d-dimensional arrays and define the partition boundaries of the d-dimensions. Flood [29] is a state-of-the-art learned multi-dimensional index for d-dimensional data, which partitions the data using a grid over d-1 dimensions and uses the last dimension as the sort dimension. In our implementation, the grid partitioning techniques use a similar approach where the space is divided in one dimension and the other dimension is used as the sort dimension.

2.2.2 Quadtree

Quadtree [10] along with its many variants is a tree data structure that also partitions the space like the k-d tree. The term quadtree is generally referred to the two-dimensional variant, but the basic idea can easily be generalized to d dimensions. Like the k-d tree, the quadtree decomposes the space using rectilinear hyperplanes. The important distinction is that quadtree is not a binary tree, and the interior nodes in the tree have 2^d children for d-dimensions. For d = 2, each interior node has four children, each corresponding to a rectangle. The search space is recursively decomposed into four quadrants until the number of objects in each quadrant is less than a predefined threshold (usually the page size). Quadtrees are generally not balanced as the tree goes deeper for the areas with higher densities.

2.2.3 K-d tree

K-d tree [3] is a binary search tree that recursively subdivides the space into equal subspaces by means of rectilinear (or iso-oriented) hyperplanes. The subdivision alternates between the k dimensions to be indexed. The splitting hyperplanes at every level are known as the discriminators. For k = 2, for example, the splitting hyperplanes are alternately perpendicular to the x-axis and the y-axis, and are called the x-discriminator and the y-discriminator respectively. The original K-d tree partitioned the space into equal partitions, for example if the input space consists of GPS co-ordinate system (-90.0, -180 to 90, 180) the space would be divided into equal halves (-45, -90 to 45, 90). Kd trees are thus not balanced if the data is skewed (most of which might only lie in one partition). K-d tree can be made data-aware by selecting a median point from the data and dividing the data into two halves. This ensures that both partitions in the binary tree are balanced. We have implemented the data-aware k-d tree in our work.

2.2.4 Sort-Tile-Recursive (STR) packed R-tree

Sort-Tile-Recursive [25] is a packing algorithm to fill Rtree [13] and aims to maximize space utilization. The main idea behind STR packing is to tile the data space into $S \times S$ grid. For example, consider the number of points in a dataset to be P and N be the capacity of a node. The data space can then be divided into $S \times S$ grid where $S = \sqrt{P/N}$. The points are first sorted on the x-dimension (in case of rectangles, the x-dimension of the centroid) and then divided into S vertical *slices*. Within each vertical slice, the points are sorted on the y-dimension, and packed into nodes by grouping them into runs of length N thus forming S horizontal slices. The process then continues recursively. Packing the R-tree in this way packs all the nodes completely, except the last node which may have fewer than N elements.

2.3 Search Within Partition

The learned index structures require the underlying data to be sorted. In multi-dimensions, there is no inherent sort order over all dimensions. Thus, after partitioning the data, a sort ordering on some dimension is required for the learned indexes to work. To achieve that, within each partition we sort the data using one dimension. Since spatial data consists of two dimensions (in two-dimensional space), either of the two dimensions can be selected as the sort dimension. Once the data within the partition has been sorted, either a learned index or binary search (hereby search technique) can be used on the sorted dimension. In all our experiments, we have sorted on the longitude value of the location.

We use a RadixSpline [22, 23] over the sorted dimension which consists of two components: 1) a set of spline points, and 2) a radix table to quickly determine the spline points to examine for a lookup key (in our case the dimension over which the data is sorted). At lookup time, first the radix table is consulted to determine the range of spline points to examine. In the next step, these spline points are searched over to determine the spline points surrounding the lookup key. In the last step, linear interpolation is used to predict the position of the lookup key in the sorted array. Unlike the RMI [24], the RadixSpline only requires one pass over the data to build the index, while retaining competitive lookup times. The RadixSpline and the RMI, at the time of writing, only work on integer values, and we adapted the open-source implementation of RadixSpline to work with floating-point values (spatial datasets generally contain floating point values). In our implementation, we have set the spline error to 32 in all experiments.

It is important to make a distinction between how we use RadixSpline and binary search for refinement. In case of binary search, we do a lookup for the lower bound of the query on the sorted dimension. As learned indexes come with an error, usually a local search is done to find the lookup point (in our case the query lower bound). For range scans, as we do, there can be two cases. The first case is that the estimated value from the spline is lower than the actual lower bound on the sorted dimension. In this case, we scan up until we reach the lower bound on the sorted dimension. In the second case, the estimated value is higher than the actual lower bound, hence, we first scan down to the lower bound, materialize all the points in our way until we reach this bound, and after that we scan up until the query upper bound (or the partition end). In case the estimated value is lower than the upper bound of the query (i.e. the estimated value is within both query bounds), the second case incurs zero cost for local search as we can scan in both directions until we reach the query bounds within the partition.

3. EVALUATION

All experiments were run single threaded on an Ubuntu 18.04 machine with an Intel Xeon E5-2660 v2 CPU (2.20 GHz, 10 cores, 3.00 GHz turbo)¹ and 256 GB DDR3 RAM. We use the *numactl* command to bind the thread and memory to one node to avoid NUMA effects. CPU scaling was also disabled during benchmarking using the *cpupower* command.

3.1 Datasets

For evaluation, we used three datasets, the New York City Taxi Rides dataset [31] (NYC Taxi Rides), geo-tagged tweets in the New York City area (NYC Tweets), and Open Streets Maps (OSM). NYC Taxi Rides contains 305 million taxi rides from the years 2014 and 2015. NYC Tweets data was collected using Twitter's Developer API [46] and contains 83 million tweets. The OSM dataset has been taken from [33] and contains 200M records from the All Nodes (Points) dataset. Figure 3 shows the spatial distribution of the three datasets. We further generated two types of query workloads for each of the three datasets: skewed queries (which follows the distribution of the underlying data) and uniform queries. For each type of query workload, we generated six different workloads ranging from 0.00001% to 1.0%selectivity. For example, in the case of Taxi Rides dataset (305M records), these queries would materialize 30 records to 3 million records. These query workloads consist of one million queries each. To generate skewed queries, we select a record from the data, and expand its boundaries (using

¹CPU: https://ark.intel.com/content/www/us/en/ark/ products/75272/intel-xeon-processor-e5-2660-v2-25m-cache-2-20-ghz.html



Figure 3: Datasets: (a) Tweets are spread across New York, (b) NYC Taxi trips are clustered in central New York, and (c) All Nodes dataset from OSM.

a random ratio in both dimensions) until the selectivity requirement of the query is met. For uniform queries, we generated points uniformly in the embedding space of the dataset and expand the boundaries similarly until the selectivity requirement of the query is met. The query selectivity and the type of query are mostly application dependent. For example, consider the application Google Maps, and a user issues a query to find the popular pizzeria near the user. The expected output for this query should be a handful of records, i.e. a low selectivity query (a list of 20-30 restaurants near the user). On the other hand a query on an analytical system, would materialize many more records (e.g. find average cost of all taxi rides originating in Manhattan).

3.2 Tuning Partitioning Techniques

Recent work in learned multi-dimensional and spatial indexes have focused on learning from the data and the query workload. The essential idea behind learning from both data and query workload is that a particular usecase can be instance-optimized. To study this effect, we conducted multiple experiments on the three datasets by varying the sizes of the partitions, tuning them on two workloads with different selectivities (to cover a broad spectrum we tune the indexes on queries with low and high selectivity) for both skewed and uniform queries.

Figure 4 shows the effect of tuning when the indexes are tuned for the lowest selectivity workload for the two query types. It can be seen in the figure that it is essential to tune the grid partitioning techniques for a particular workload. Firstly, they are susceptible to the size of the partition. As the size of the partition increases, we notice an improvement in the performance until a particular partition size is reached which corresponds to the optimal performance. After this point, increasing the size of the partitions only degrades performance. It can be seen that, usually, for grid (singledimension) partitioning techniques the partition sizes are much larger compared to partitioning techniques which filter on both dimensions (only Quadtree is shown in the figure but the same holds for the other partitioning techniques we have covered in this work, we do not show the other trees because the curve is similar for them). Due to the large partition sizes in grid partitioning techniques, we notice a large increase in performance while using a learned index compared to binary search. This is especially evident for skewed queries (which follow the underlying data distribution). We encountered a speedup from 11.79% up to 39.51% compared to binary search. Even when we tuned a learned index to a partition size which corresponds to the optimal performance for binary search, we found that in multiple cases learned index frequently outperformed binary search. Learned indexes do not help much for partitioning techniques which filter on both dimensions, instead the performance of Quadtree (and STRtree) dropped in many cases, see Table 1. The reason is that the optimal partition sizes for these techniques is very low (less than 1,000 points per partition for most configurations). The refinement cost for learned indexes is an overhead in such cases. K-d tree on the other hand, contains more points per partition (from 1200 to 7400) for the optimal configuration for Taxi Trips and OSM datasets and thus learned indexes perform faster by 2.43% to 9.17% than binary search. For Twitter dataset, the optimal configuration contains less than 1200 points per partition, and we observed a similar drop in performance using learned indexes.

Figure 5 shows the effect of number of cells and number of points that are scanned in each partition on query runtime for Fixed-grid on Taxi Trips dataset for lowest selectivity. As the number of points per partitions increases (i.e. fewer number of partitions), the number of cells decreases. At the same time, the number of points that need to be scanned for the query increases. The point where these curves meet is the optimal configuration for the workload which corresponds to the lowest query runtime. For tree structures, the effect is different. Figure 6 shows that the structures that filter on both dimensions do most of the pruning in the index lookup. The dominating cost in these structures is the number of points scanned within the partition and the query runtime is directly proportional to this number. To minimize the number of points scanned, they do most of the pruning during index lookup which require more partitions (i.e. less number of points per partition), but then they pay more for index lookup.

3.3 Range Query

Figure 7 shows the query runtime for all learned index structures. It can be seen that Fixed-grid along with Adaptivegrid performs (1D schemes) perform the best for all the cases except uniform queries on Taxi and OSM datasets. For skewed queries, Fixed-grid is $1.23 \times$ to $1.83 \times$ faster than the closest competitor, Quadtree (2D), across all datasets and selectivity. The slight difference in performance between Fixed-grid and Adaptive-grid comes from the index lookup. For Adaptive-grid, we use binary search on the linear scales to find the first partition the query intersects with. For Fixed-grid, the index lookup is almost negligible as only an offset computation is needed to find first intersecting partition. It can also be seen in the figure that the Quadtree is significantly better for uniform queries in case of Taxi Rides dataset $(1.37 \times)$ and OSM dataset $(2.68 \times)$ than the closest competitor Fixed-grid. There are two reasons for this, firstly the Quadtree intersects with fewer number of partitions than the other index structures, see Table 2. Secondly, for uniform queries, the Quadtree is more likely to traverse the sparse and low-depth region of the index. This is in conformance with an earlier research [18], where the authors report similar findings while comparing the Quadtree to the R^{*}-tree and the Pyramid-Technique.

3.4 Indexing Costs

Figure 8 shows that Fixed-grid and Adaptive-grid are faster to build than the tree based indexes. Fixed-grid is $2.11\times$, $2.05\times$, and $1.90\times$ faster to build than closest competitor STR tree. Quadtree is the slowest to build because it



Figure 4: Configuration Experiments - ML vs. BS for low selectivity (0.00001%).

	Taxi Trips (Skewed Queries)						Taxi Trips (Uniform Queries)					
	Fixed		Adaptive		Quadtree		Fixed		Adaptive		Quadtree	
Selectivity (%)	ML	BS	ML	BS	ML	BS	ML	BS	ML	BS	ML	BS
0.00001	1.78	2.35	1.86	2.40	2.77	2.51	2.02	2.58	81.4	10.54	1.48	1.31
0.0001	4.54	5.82	4.67	6.12	6.12	5.82	5.85	6.91	228.1	27.69	3.69	3.42
0.001	14.97	18.83	15.32	19.49	20.84	19.47	22.87	24.34	708.8	87.49	13.59	12.98
0.01	90.13	97.04	89.48	95.96	117.01	104.37	141.24	151.47	2634.4	309.62	98.85	112.77
0.1	678.12	698.39	675.14	696.49	922.67	793.96	988.35	922.96	9609.9	1174.79	891.24	1101.95
1.0	8333.94	8408.15	8301.56	8399.69	10678.04	9512.29	8843.71	8753.68	8574.84	8836.28	10647.97	12377.14

Table 1: Total query runtime (in microseconds) for both RadixSpline (ML) and binary search (BS) for Taxi Rides dataset on skewed and uniform query workloads (parameters are tuned for selectivity 0.00001%).



Figure 5: Effect of number of cells and number of points scanned for Fixed-grid on Taxi Trip dataset for skewed queries (0.00001% selectivity).

	Taxi	Rides	OSM			
Partitioning	Skewed	Uniform	Skewed	Uniform		
Fixed Adaptive k-d tree Quadtree STR	$1.97 \\ 1.74 \\ 1.70 \\ 1.79 \\ 2.60$	$7.98 \\ 31.57 \\ 21.62 \\ 2.12 \\ 47.03$	$1.72 \\ 1.51 \\ 1.56 \\ 1.37 \\ 1.90$	$23.73 \\ 24.80 \\ 30.95 \\ 7.96 \\ 11.05$		

Table 2: Average number of partitions intersected for each partitioning scheme for selectivity 0.00001% on Taxi Rides and OSM datasets.



Figure 6: Effect of number of cells and number of points scanned for Quadtree on Taxi Trip dataset for skewed queries (0.00001% selectivity).

generates a large number of cells for optimal configuration. Not all partitions in Quadtree contain an equal number of points as it divides space rather than data, thus leading to an imbalanced number of points per partition. Fixed-grid and Adaptive grid do not generate large number of partitions, as the partitions are quite large for optimal configuration. They are lower in size for similar reasons. The index size in Figure 8 also includes the size of data being indexed.

4. RELATED WORK

Recent work by Kraska et al. [24] proposed the idea of replacing traditional database indexes with learned models. Since then, there has been a corpus of work on extending the



Figure 7: Total query runtime with parameters tuned on selectivity 0.00001%.



Figure 8: Index build times and sizes for the three datasets.

ideas of the learned index to spatial and multi-dimensional data. Flood [29] is an in-memory read-optimized multidimensional index that organizes the physical layout of ddimensional data by dividing each dimension into some number of partitions, which forms a grid over *d*-dimensional space and adapts to the data and query workload. Learning has also been applied to the challenge of reducing I/O cost for disk-based multi-dimensional indexes. Qd-tree [50] uses reinforcement learning to construct a partitioning strategy that minimizes the number of disk-based blocks accessed by a query. LISA [26] is a disk-based learned spatial index that achieves low storage consumption and I/O cost while supporting range queries, nearest neighbor queries, and insertions and deletions. The ZM-index [48] combines the standard Z-order space-filling curve with the RMI from [24] by mapping multi-dimensional values into a single-dimensional space, which is then learnable using models. The ML-index [5] combines the ideas of iDistance [16] and the RMI to support range and KNN queries. [14] augment existing indexes with

light-weight models to accelerate range and point queries.

5. CONCLUSIONS AND FUTURE WORK

In this work, we implemented techniques proposed in a state-of-the-art multi-dimensional index, namely, Flood [29], which indexes points using a variant of the Grid-file and applied them to five classical spatial indexes. We have shown that replacing binary search with learned indexes within each partition can improve overall query runtime by 11.79% to 39.51%. As expected, the effect of using a more efficient search within a partition is more pronounced for queries with low selectivity. With increasing selectivity, the effect of a fast search diminishes. Likewise, the effect of using a learned index is larger for (1D) grid partitioning techniques (e.g., Fixed-grid) than for (2D) tree structures (e.g., Quadtree). The reason is that the partitions (cells) are less representative of the points they contain in the 1D case than in the 2D case. Hence, 1D partitioning requires more refinement within each cell.

In contrary, finding the qualifying partitions is more efficient with 1D than with 2D partitioning, thus contributing to lower overall query runtime (1.23x to 1.83x times faster). Currently, we are using textbook implementations for Quadtree and K-d tree. Future work could study replacing these tree structures with learned counterparts. For example, we could linearize Quadtree cells (e.g., using a Hilbert or Z-order curve) and store the resulting cell identifiers in a learned index.

So far we have only studied the case where indexes and data fit into RAM. For on-disk use cases, performance will likely be dominated by I/O and the search within partitions will be of less importance. We expect partition sizes to be performance-optimal when aligned with the physical page size. To reduce I/O, it will be crucial for partitions to not contain any unnecessary points. Hence, we expect 2D partitioning to be the method of choice in this case. We refer to LISA [26] for further discussions on this topic.

6. **REFERENCES**

- A. Aji, F. Wang, H. Vo, R. Lee, Q. Liu, X. Zhang, and J. H. Saltz. Hadoop-gis: A high performance spatial data warehousing system over mapreduce. *PVLDB*, 6(11):1009–1020, 2013.
- [2] K. Amemiya and A. Nakao. Layer-integrated edge distributed data store for real-time and stateful services. In NOMS 2020 - IEEE/IFIP Network Operations and Management Symposium, Budapest, Hungary, April 20-24, 2020, pages 1–9. IEEE, 2020.
- [3] J. L. Bentley. Multidimensional binary search trees used for associative searching. *Commun. ACM*, 18(9):509–517, 1975.
- [4] J. L. Bentley and J. H. Friedman. Data structures for range searching. ACM Comput. Surv., 11(4):397–409, 1979.
- [5] A. Davitkova, E. Milchevski, and S. Michel. The ML-Index: A Multidimensional, Learned Index for Point, Range, and Nearest-Neighbor Queries. In 2020 Conference on Extending Database Technology (EDBT), 2020.
- [6] H. Doraiswamy and J. Freire. A gpu-friendly geometric data model and algebra for spatial queries. In D. Maier, R. Pottinger, A. Doan, W. Tan, A. Alawini, and H. Q. Ngo, editors, *Proceedings of the* 2020 International Conference on Management of Data, SIGMOD Conference 2020, online conference [Portland, OR, USA], June 14-19, 2020, pages 1875–1885. ACM, 2020.
- [7] H. Doraiswamy and J. Freire. A gpu-friendly geometric data model and algebra for spatial queries: Extended version. *CoRR*, abs/2004.03630, 2020. https://arxiv.org/abs/2004.03630.
- [8] A. Eldawy and M. F. Mokbel. Spatialhadoop: A mapreduce framework for spatial data. In 31st IEEE International Conference on Data Engineering, ICDE 2015, Seoul, South Korea, April 13-17, 2015, pages 1352–1363, 2015.
- [9] A. Eldawy, I. Sabek, M. Elganainy, A. Bakeer, A. Abdelmotaleb, and M. F. Mokbel. Sphinx: Empowering impala for efficient execution of SQL queries on big spatial data. In Advances in Spatial and Temporal Databases - 15th International Symposium, SSTD 2017, Arlington, VA, USA, August 21-23, 2017, Proceedings, pages 65–83, 2017.
- [10] R. A. Finkel and J. L. Bentley. Quad trees: A data structure for retrieval on composite keys. Acta Inf., 4:1–9, 1974.
- [11] F. García-García, A. Corral, L. Iribarne, and M. Vassilakopoulos. Improving distance-join query processing with voronoi-diagram based partitioning in spatialhadoop. *Future Gener. Comput. Syst.*, 111:723–740, 2020.
- [12] D. Gomes. MemSQL Live: Nikita Shamgunov on the Data Engineering Podcast, 2019. https: //www.memsql.com/blog/memsql-live-nikitashamgunov-on-the-data-engineering-podcast/.
- [13] A. Guttman. R-trees: A dynamic index structure for spatial searching. In B. Yormark, editor, SIGMOD'84, Proceedings of Annual Meeting, Boston, Massachusetts, USA, June 18-21, 1984, pages 47–57. ACM Press, 1984.

- [14] A. Hadian, A. Kumar, and T. Heinis. Hands-off model integration in spatial index structures. CoRR, abs/2006.16411, 2020. https://arxiv.org/abs/2006.16411.
- [15] S. Hagedorn, P. Götze, and K. Sattler. The STARK framework for spatio-temporal data analytics on spark. In Datenbanksysteme für Business, Technologie und Web (BTW 2017), 17. Fachtagung des GI-Fachbereichs ,,Datenbanken und Informationssysteme" (DBIS), 6.-10. März 2017, Stuttgart, Germany, Proceedings, pages 123–142, 2017.
- [16] H. V. Jagadish, B. C. Ooi, K.-L. Tan, C. Yu, and R. Zhang. Idistance: An adaptive b+-tree based indexing method for nearest neighbor search. ACM Trans. Database Syst., 30(2):364–397, June 2005.
- [17] M. S. Kester, M. Athanassoulis, and S. Idreos. Access path selection in main-memory optimized data systems: Should I scan or should I probe? In Proceedings of the 2017 ACM International Conference on Management of Data, SIGMOD Conference 2017, Chicago, IL, USA, May 14-19, 2017, pages 715–730. ACM, 2017.
- [18] Y. J. Kim and J. M. Patel. Rethinking choices for multi-dimensional point indexing: Making the case for the often ignored quadtree. In CIDR 2007, Third Biennial Conference on Innovative Data Systems Research, Asilomar, CA, USA, January 7-10, 2007, Online Proceedings, pages 281–291. www.cidrdb.org, 2007.
- [19] A. Kipf, H. Lang, V. Pandey, R. A. Persa, C. Anneser, E. T. Zacharatou, H. Doraiswamy, P. A. Boncz, T. Neumann, and A. Kemper. Adaptive main-memory indexing for high-performance point-polygon joins. In Proceedings of the 23nd International Conference on Extending Database Technology, EDBT 2020, Copenhagen, Denmark, March 30 - April 02, 2020, pages 347–358. OpenProceedings.org, 2020.
- [20] A. Kipf, H. Lang, V. Pandey, R. A. Persa, P. A. Boncz, T. Neumann, and A. Kemper. Adaptive geospatial joins for modern hardware. *CoRR*, abs/1802.09488, 2018. http://arxiv.org/abs/1802.09488.
- [21] A. Kipf, H. Lang, V. Pandey, R. A. Persa, P. A. Boncz, T. Neumann, and A. Kemper. Approximate geospatial joins with precision guarantees. In 34th IEEE International Conference on Data Engineering, ICDE 2018, Paris, France, April 16-19, 2018, pages 1360–1363, 2018.
- [22] A. Kipf, R. Marcus, A. van Renen, M. Stoian, A. Kemper, T. Kraska, and T. Neumann. Radixspline: A single-pass learned index. *CoRR*, abs/2004.14541, 2020. https://arxiv.org/abs/2004.14541.
- [23] A. Kipf, R. Marcus, A. van Renen, M. Stoian, A. Kemper, T. Kraska, and T. Neumann. Radixspline: a single-pass learned index. In Proceedings of the Third International Workshop on Exploiting Artificial Intelligence Techniques for Data Management, aiDM@SIGMOD 2020, Portland, Oregon, USA, June 19, 2020, pages 5:1–5:5. ACM, 2020.
- [24] T. Kraska, A. Beutel, E. H. Chi, J. Dean, and N. Polyzotis. The case for learned index structures. In Proceedings of the 2018 International Conference on Management of Data, SIGMOD Conference 2018,

Houston, TX, USA, June 10-15, 2018, pages 489–504. ACM, 2018.

- [25] S. T. Leutenegger, J. M. Edgington, and M. A. López. STR: A simple and efficient algorithm for r-tree packing. In *Proceedings of the Thirteenth International Conference on Data Engineering, April 7-11, 1997, Birmingham, UK*, pages 497–506. IEEE Computer Society, 1997.
- [26] P. Li, H. Lu, Q. Zheng, L. Yang, and G. Pan. LISA: A Learned Index Structure for Spatial Data. In Proceedings of the 2020 International Conference on Management of Data, SIGMOD '20, New York, NY, USA, 2020. Association for Computing Machinery.
- [27] A. Makris, K. Tserpes, G. Spiliopoulos, and D. Anagnostopoulos. Performance evaluation of mongodb and postgresql for spatio-temporal data. In P. Papotti, editor, *Proceedings of the Workshops of the EDBT/ICDT 2019 Joint Conference, EDBT/ICDT* 2019, Lisbon, Portugal, March 26, 2019, volume 2322 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2019.
- [28] MongoDB Releases New Geo Features in MongoDB 2.4, 2013. https://www.mongodb.com/blog/post/ new-geo-features-in-mongodb-24/.
- [29] V. Nathan, J. Ding, M. Alizadeh, and T. Kraska. Learning multi-dimensional indexes. In *Proceedings of the 2020 International Conference on Management of Data*, SIGMOD '20, New York, NY, USA, 2020. Association for Computing Machinery.
- [30] J. Nievergelt, H. Hinterberger, and K. C. Sevcik. The grid file: An adaptable, symmetric multikey file structure. ACM Trans. Database Syst., 9(1):38–71, 1984.
- [31] NYC Taxi and Limousine Commission (TLC) TLC Trip Record Data, 2019. https://www1.nyc.gov/ site/tlc/about/tlc-trip-record-data.page.
- [32] Oracle Spatial and Graph Spatial Features, 2019. https://www.oracle.com/technetwork/database/ options/spatialandgraph/overview/ spatialfeatures-1902020.html/.
- [33] V. Pandey, A. Kipf, T. Neumann, and A. Kemper. How good are modern spatial analytics systems? *Proc. VLDB Endow.*, 11(11):1661–1673, 2018.
- [34] V. Pandey, A. Kipf, D. Vorona, T. Mühlbauer, T. Neumann, and A. Kemper. High-performance geospatial analytics in hyperspace. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD Conference 2016, San Francisco, CA, USA, June 26 - July 01, 2016*, pages 2145–2148, 2016.
- [35] K. Richly. Optimized spatio-temporal data structures for hybrid transactional and analytical workloads on columnar in-memory databases. In I. Bartolini and F. Li, editors, Proceedings of the VLDB 2019 PhD Workshop, co-located with the 45th International Conference on Very Large Databases (VLDB 2019), Los Angeles, California, USA, August 26-30, 2019, volume 2399 of CEUR Workshop Proceedings. CEUR-WS.org, 2019.
- [36] D. Sidlauskas, S. Chester, E. T. Zacharatou, and A. Ailamaki. Improving spatial data processing by clipping minimum bounding boxes. In 34th IEEE International Conference on Data Engineering, ICDE 2018, Paris, France, April 16-19, 2018, pages

425–436. IEEE Computer Society, 2018.

- [37] R. Y. Tahboub, G. M. Essertel, and T. Rompf. How to architect a query compiler, revisited. In G. Das, C. M. Jermaine, and P. A. Bernstein, editors, *Proceedings of the 2018 International Conference on Management of Data, SIGMOD Conference 2018, Houston, TX, USA, June 10-15, 2018*, pages 307–322. ACM, 2018.
- [38] R. Y. Tahboub and T. Rompf. On supporting compilation in spatial query engines: (vision paper). In S. Ravada, M. E. Ali, S. D. Newsam, M. Renz, and G. Trajcevski, editors, Proceedings of the 24th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems, GIS 2016, Burlingame, California, USA, October 31 - November 3, 2016, pages 9:1–9:4. ACM, 2016.
- [39] R. Y. Tahboub and T. Rompf. Architecting a query compiler for spatial workloads. In Proceedings of the 2020 International Conference on Management of Data, SIGMOD Conference 2020, online conference [Portland, OR, USA], June 14-19, 2020, pages 2103–2118. ACM, 2020.
- [40] M. Tang, Y. Yu, Q. M. Malluhi, M. Ouzzani, and W. G. Aref. Locationspark: A distributed in-memory data management system for big spatial data. *PVLDB*, 9(13):1565–1568, 2016.
- [41] K. Theocharidis, J. Liagouris, N. Mamoulis, P. Bouros, and M. Terrovitis. SRX: efficient management of spatial RDF data. VLDB J., 28(5):703–733, 2019.
- [42] T. Toliopoulos, N. Nikolaidis, A. Michailidou,
 A. Seitaridis, A. Gounaris, N. Bassiliades,
 A. Georgiadis, and F. Liotopoulos. Developing a real-time traffic reporting and forecasting back-end system. In Research Challenges in Information Science 14th International Conference, RCIS 2020, Limassol, Cyprus, September 23-25, 2020, Proceedings, volume 385 of Lecture Notes in Business Information Processing, pages 58–75. Springer, 2020.
- [43] D. Tsitsigkos, P. Bouros, N. Mamoulis, and M. Terrovitis. Parallel in-memory evaluation of spatial joins. *CoRR*, abs/1908.11740, 2019. http://arxiv.org/abs/1908.11740.
- [44] D. Tsitsigkos, P. Bouros, N. Mamoulis, and M. Terrovitis. Parallel in-memory evaluation of spatial joins. In Proceedings of the 27th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems, SIGSPATIAL 2019, Chicago, IL, USA, November 5-8, 2019, pages 516–519. ACM, 2019.
- [45] D. Tsitsigkos, K. Lampropoulos, P. Bouros, N. Mamoulis, and M. Terrovitis. A two-level spatial in-memory index. *CoRR*, abs/2005.08600, 2020. https://arxiv.org/abs/2005.08600.
- [46] Tutorials: Filtering Tweets by location, 2020. https://developer.twitter.com/en/docs/ tutorials/filtering-tweets-by-location.
- [47] Uber. Uber newsroom: 10 billion, 2018. https://www.uber.com/newsroom/10-billion/.
- [48] H. Wang, X. Fu, J. Xu, and H. Lu. Learned index for spatial queries. In 2019 20th IEEE International Conference on Mobile Data Management (MDM), pages 569–574, 2019.
- [49] D. Xie, F. Li, B. Yao, G. Li, L. Zhou, and M. Guo.

Simba: Efficient in-memory spatial analytics. In Proceedings of the 2016 International Conference on Management of Data, SIGMOD Conference 2016, San Francisco, CA, USA, June 26 - July 01, 2016, pages 1071–1085, 2016.

- [50] Z. Yang, B. Chandramouli, C. Wang, J. Gehrke, Y. Li, U. F. Minhas, P.-A. Larson, D. Kossmann, and R. Acharya. Qd-tree: Learning Data Layouts for Big Data Analytics. In *Proceedings of the 2020 International Conference on Management of Data*, SIGMOD '20, New York, NY, USA, 2020. Association for Computing Machinery.
- [51] S. You, J. Zhang, and L. Gruenwald. Large-scale spatial join query processing in cloud. In 31st IEEE International Conference on Data Engineering Workshops, ICDE Workshops 2015, Seoul, South Korea, April 13-17, 2015, pages 34–41, 2015.
- [52] J. Yu, J. Wu, and M. Sarwat. Geospark: a cluster

computing framework for processing large-scale spatial data. In Proceedings of the 23rd SIGSPATIAL International Conference on Advances in Geographic Information Systems, Bellevue, WA, USA, November 3-6, 2015, pages 70:1–70:4, 2015.

- [53] E. T. Zacharatou, H. Doraiswamy, A. Ailamaki, C. T. Silva, and J. Freire. GPU rasterization for real-time spatial aggregation over arbitrary polygons. *PVLDB*, 11(3):352–365, 2017.
- [54] E. T. Zacharatou, D. Sidlauskas, F. Tauheed, T. Heinis, and A. Ailamaki. Efficient bundled spatial range queries. In F. B. Kashani, G. Trajcevski, R. H. Güting, L. Kulik, and S. D. Newsam, editors, Proceedings of the 27th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems, SIGSPATIAL 2019, Chicago, IL, USA, November 5-8, 2019, pages 139–148. ACM, 2019.