# PFLOCK Report

Andres Calderon

University of California, Riverside

February 7, 2020
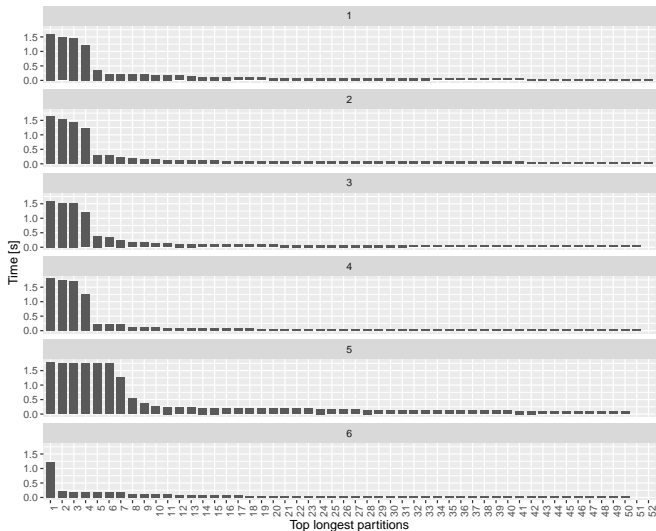
# Experiment setup...

1. Running a distance self-join to find the pairs of points far each other less than $\varepsilon$.

2. Using LA_25K dataset (time instant 19), quadtree partitioner and quadtree local indexing.

3. $\varepsilon = 20$, $partitions \approx 150$ (3x number of available cores).

4. 18 executors, 3 cores each.

# Partitions performace...



Execution time for partitions during a distance self–join

Showing 6 of 10 runs...

# Finding pairs algorithm...

---

**Algorithm 1** FINDPAIRS algorithm

---

**Require:** a dataset of points $\mathcal{P}$, a number of partitions $p$ and a distance threshold $\varepsilon$.

1: **function** FINDPAIRS ( $\mathcal{P}, p, \varepsilon$ )
2:     Partition $\mathcal{P}$ using a Quadtree and $p$ partitions        ▷ Using Algorithm 1 in next page.
3:     Create a circle of radius $\varepsilon$ for each point in $\mathcal{P}$ and store them in $\mathcal{Q}$        ▷ keep same id
4:     Partition $\mathcal{Q}$ using the same partitioner of $\mathcal{P}$
5:     Build local index in $\mathcal{P}$        ▷ Using operations provided by GeoSpark
6:     Execute a distance join query in $\mathcal{P}$ and $\mathcal{Q}$ using $\varepsilon$ as distance        ▷ Using Algorithm 4 in following pages.
7:     Filter those pairs where $p_1.id < p_2.id$
8: **end function**

---

# GeoSpark Partitioning algorithm...

---

**Algorithm 1** SRDD spatial partitioning

---

   **Data**: An original SRDD
   **Result**: A repartitioned SRDD
   /*   **Step** 1**: Build a global grid file at master node**                                  */
1  Take samples from the original SRDD $A$ partitions in parallel;
2  Construct the selected spatial structure on the collected sample at master node;
3  Retrieve the grids from built spatial structures;
   /*   **Step** 2**: Assign grid ID to each object in parallel**                             */
4  **foreach** *spatial object in SRDD A* **do**
5     **foreach** *grid* **do**
6         **if** *the grid intersects the object* **then**
7             Add (grid ID, object) pair into SRDD $B$;
       // Only needed for R-Tree partitioning
8     **if** *no grid intersects the object* **then**
9         Add (overflow grid ID, object) pair into SRDD $B$;
   /*   **Step** 3**: Repartition SRDD across the cluster**                              */
10 Partition SRDD $B$ by ID and get SRDD $C$;
11 Cache the new SRDD $C$ in memory and return it;

---

# GeoSpark GSJoin algorithm...

---

**Algorithm 4** *GSJoin* algorithm for range join and distance join query

---

**Data**: (repartitioned) SRDD A and (repartitioned) SRDD B
**Result**: PairRDD in schema <Left object from A, right object from B>
/* **Step1: Zip partitions** */
1 **foreach** *partition pair from SRDD A and B with the same grid ID i* **do**
2      Merge two partitions to a bigger partition that has two sub-partitions;
3 Return the intermediate SRDD C;
/* **Step2: Run partition-level local join** */
4 **foreach** *partition P in the C* **do**
5      **foreach** *object $O_A$ in the sub-partition from A* **do**
6          **if** *an index exists in the sub-partition from B* **then**
             // Filter phase
7              Query the spatial index of this partition using the $O_A$'s MBR;
             // Refine phase
8              Check the spatial relation using real shapes of $O_A$ and candidate objects $O_B$s;
             /* **Step3: Remove duplicates** */
9              Report $<O_A, O_B>$ pair only if the reference point of this pair is in P;
10          **else**
11              **foreach** *object $O_B$ in the sub-partition from B* **do**
12                  Check spatial relation between $O_A$ and $O_B$;
                 /* **Step3: Remove duplicates** */
13                  Report $<O_A, O_B>$ pair only if the reference point of this pair is in P;
14 Generate the result PairRDD;

---