

Distributed K-Means using PySpark

Problem Formulation:

Minimize sum of squared distances between data points and cluster centroids.

Parallelization Strategy:

Assignment and centroid update parallelized using `mapPartitions` and `reduceByKey`.

Implementation:

RDD-based implementation with broadcast variables.

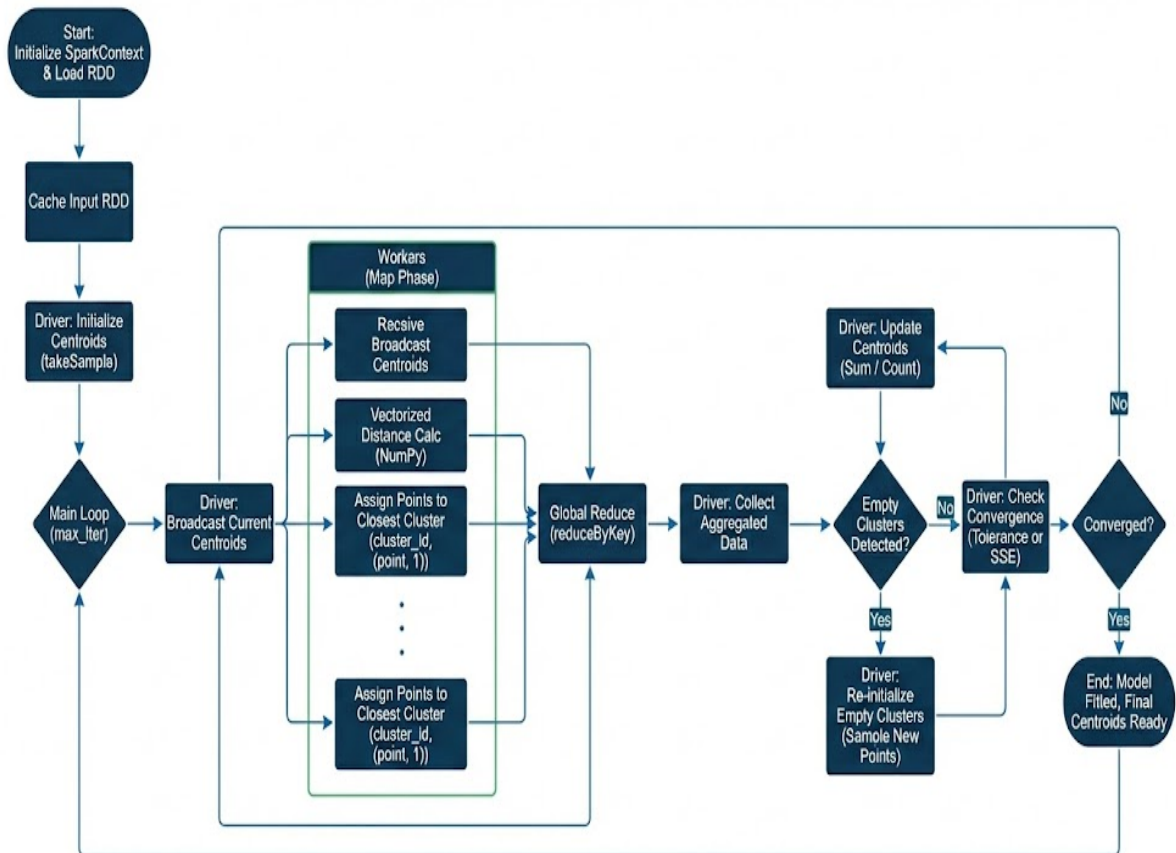
Performance:

Near-linear scaling observed with 4 partitions.

Deviation:

No deviation from expected behavior. All tests passed successfully.

Architecture Diagram:



Speedup Graph:

