# BIRCH (Balanced Iterative Reducing and Clustering Using Hierarchies)

❑ A multiphase clustering algorithm (Zhang, Ramakrishnan & Livny, SIGMOD'96)

❑ Incrementally construct a CF (Clustering Feature) tree, a hierarchical data structure for multiphase clustering

    ❑ Phase 1: Scan DB to build an initial in-memory CF tree (a multi-level compression of the data that tries to preserve the inherent clustering structure of the data)

    ❑ Phase 2: Use an arbitrary clustering algorithm to cluster the leaf nodes of the CF-tree

❑ Key idea: Multi-level clustering

    ❑ Low-level micro-clustering: Reduce complexity and increase scalability

    ❑ High-level macro-clustering: Leave enough flexibility for high-level clustering

❑ *Scales linearly*:  Find a good clustering with a single scan and improve the quality with a few additional scans