# Inside MapReduce

For the cloud:

1.  Parallelize Map: easy! each map task is independent of the other!
    - All Map output records with same key assigned to same Reduce
2.  Transfer data from Map to Reduce:
    - Called Shuffle data
    - All Map output records with same key assigned to same Reduce task
    - use partitioning function, e.g., hash(key)%number of reducers
3.  Parallelize Reduce: easy! each reduce task is independent of the other!
4.  Implement Storage for Map input, Map output, Reduce input, and Reduce output
    - Map input: from distributed file system
    - Map output: to local disk (at Map node); uses local file system
    - Reduce input: from (multiple) remote disks; uses local file systems
    - Reduce output: to distributed file system

    local file system = Linux FS, etc.

    distributed file system = GFS (Google File System), HDFS (Hadoop Distributed File System)