COS 597D: Principles of Database and Information Systems

Distributed computing on large data sets



Distributing query evaluation

- · Data distribution?
- Computation distribution?
- · Goals
 - Keep all machines busy
 - Be able to replace badly-behaved machines seamlessly!

MapReduce framework

- · programming model
- implementation for large clusters
- Google introduced for index building and PageRank circa 2003
 - "for processing and generating large data sets" The Apache Hadoop project developed open-source
- software





- Map: produce {(term, docID)} for each term appearing in docID
- Input to Reduce: list of all (term, docID) pairs for one term
- · Output of Reduce: (term, sorted list of docIDs containing that term) - postings list!

keys 6





MapReduce Fault Tolerance

- Master fails => restart whole computation
- Worker node fails
 - Master detects failure
 - must redo all Map tasks assigned to worker
 output of completed Map tasks on failed worker's disk
 - for failed Map worker, Master
 - reschedules each Map task
 - notifies reducer workers of change in input location
 - for failed Reduce worker, Master
 - reschedules each Reduce task
 - rescheduling occurs as live workers become available

Communication Cost Rajaraman, Leskovec and Ullman

- Model: algorithm is acyclic network of tasks
 cascaded MapReduce tasks
- Communication cost task = size of input
 bytes versus tuples
- · Commun. cost alg. = sum of cost of all tasks
- Captures
 - Network cost
 - Disk cost
- Why not output cost?
 - input to another algorithmoutput "rarely large"

11

Remarks

- Google built on large collections of inexpensive "commodity PCs"
 - always some not functioning
- Solve fault-tolerance problem in software
 redundancy & flexibility NOT special-purpose hardware
- · Keep machines relative generalists
 - machine becomes free \Rightarrow

assign to any one of set of tasks

13