


15-826: Multimedia Databases and Data Mining


Lecture #23: intro to hadoop
C. Faloutsos



Resources


- Software
 - <http://hadoop.apache.org/>
- Map/reduce paper [Dean & Ghemawat]
 - <http://research.google.com/archive/mapreduce.html>
- Tutorial: see part1, foils #9-20 from
 - videlectures.net/site/normal_dl/tag=75554/kdd2010_papadimitriou_sun_yan_lsdm_01.pdf
 - videlectures.net/kdd2010_papadimitriou_sun_yan_lsdm/

15-826 (c) 2017 C. Faloutsos 2





Motivation: Scalability

- Google: > 450,000 processors in clusters of ~2000 processors each [Barroso, Dean, Hölzle, “*Web Search for a Planet: The Google Cluster Architecture*” IEEE Micro 2003]
- Yahoo: 5Pb of data [Fayyad, KDD’ 07]
- Problem: **machine failures**, on a daily basis
- How to parallelize data mining tasks, then?
- A: map/reduce – hadoop (open-source clone)
<http://hadoop.apache.org/>



15-826 (c) 2017 C. Faloutsos 3

2’ intro to hadoop

- master-slave architecture; n-way replication (default n=3)
- ‘group by’ of SQL (in parallel, fault-tolerant way)
- e.g. find histogram of word frequency
 - compute local histograms
 - then merge into global histogram

```
select course-id, count(*)
from ENROLLMENT
group by course-id
```

15-826 (c) 2017 C. Faloutsos 4

CMU SCS details

2' intro to hadoop

- master-slave architecture; n-way replication (default n=3)
- 'group by' of SQL (in parallel, fault-tolerant way)
- e.g, find histogram of word frequency
 - compute local histograms
 - then merge into global histogram

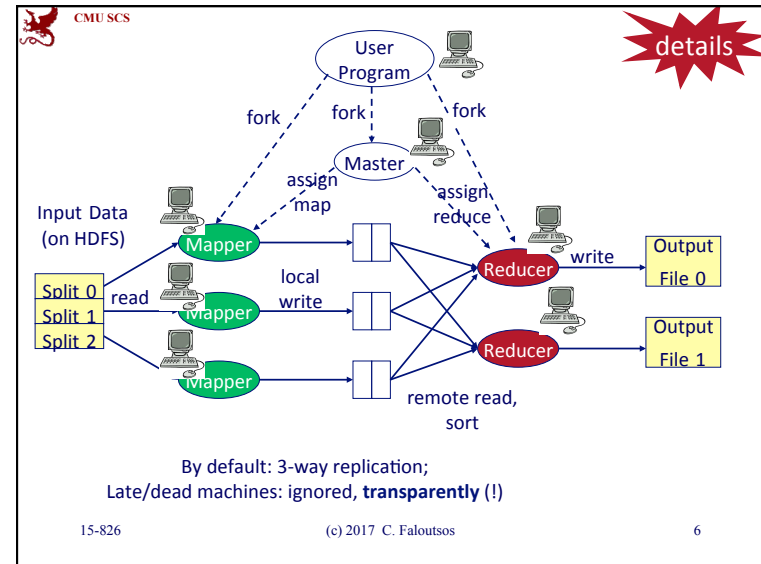
```

select course-id, count(*)
from ENROLLMENT
group by course-id
    
```

reduce

map

15-826 (c) 2017 C. Faloutsos 5



CMU SCS

More details:

- (thanks to U Kang for the animations)



15-826 (c) 2017 C. Faloutsos 7

