

15-721 DB Sys. Design & Impl.

GAMMA

Christos Faloutsos
www.cs.cmu.edu/~christos

Roadmap

- 1) Roots: System R and Ingres
- 2) Implementation: buffering, indexing, q-opt
- 3) Transactions: locking, recovery
- 4) Distributed DBMSs
- ➔ 5) Parallel DBMSs: **Gamma**, Alphasort
- 6) OO/OR DBMS
- 7) Data Analysis - data mining
- 8) Benchmarks
- 9) vision statements
 extras (streams/sensors, graphs, multimedia, web, fractals)

Citation

DeWitt, et al. "*The Gamma Database Machine Project.*"
 IEEE TKDE 2(1): 44-63 (March 1990)

Detailed outline

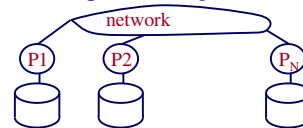
- ➔ Introduction / History
- Hardware
- Software architecture
- Query processing
- Performance evaluation

History

- DIRECT 1977 - 84
 - early database machine project
 - showed parallelism useful for db apps
- Flaws curtailed scalability
 - shared memory
 - central control of execution
- Gamma 1984 - 92

Key Ideas

- Shared-nothing
- Hash-based parallel algorithms
- Horizontal partitioning ('declustering')



CMU SCS

Gamma Hardware (v2.0)

- iPSC/2 Intel hypercube
- 32 x386 processors
- 8MB of memory
- 330MB Maxtor drive / node (45KB cache)
- Routing modules
 - 2.8 Mb/s
 - full duplex, serial, reliable

15-721 C. Faloutsos 7

CMU SCS

Gamma v2.0

- OS: NOSE
- multiple, lightweight processes with shared memory
- Entire DB in one NX/2 process

15-721 C. Faloutsos 8

CMU SCS

Detailed outline

Introduction / History
 Hardware
 ➡ Software architecture
 Query processing
 Performance evaluation

15-721 C. Faloutsos 9

CMU SCS

Storage Organization

- Horizontal partitioning (user selectable)
 - round-robin; hashed; range partitioned
 - all relations, on all units
 - method of partitioning: recorded in catalog
- Clustered index (not necessarily on partitioning attribute)
- (in retrospect: Partitioning relations should have been based on ‘heat’)

15-721 C. Faloutsos 10

CMU SCS

Gamma Process Structure

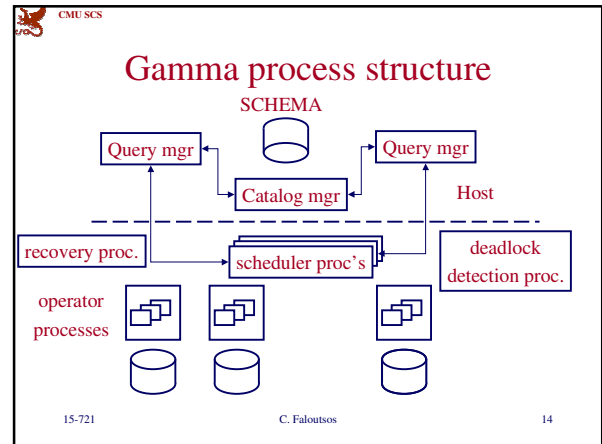
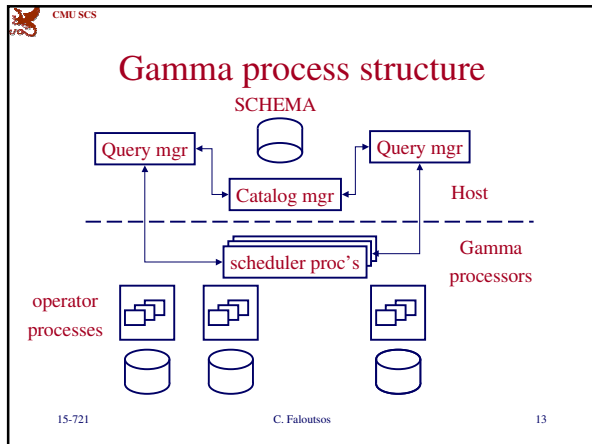
- there is a ‘host’ machine;
- and the Gamma processors

15-721 C. Faloutsos 11

CMU SCS

Gamma process structure

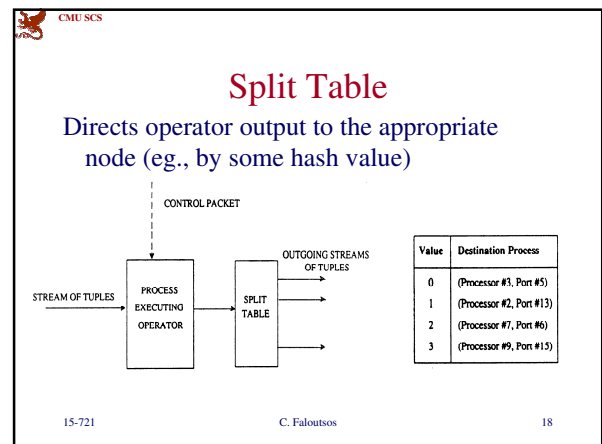
15-721 C. Faloutsos 12

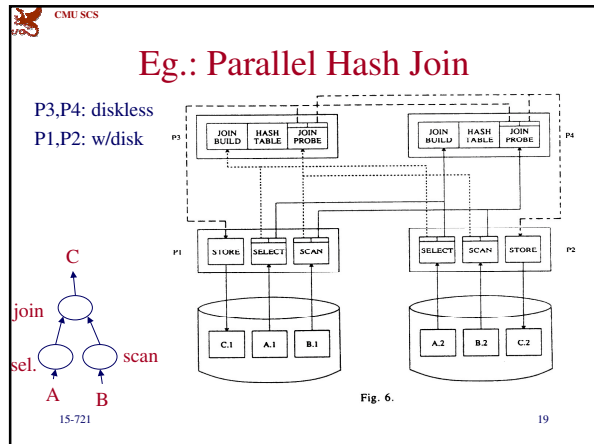


- CMU SCS
- ## Gamma Process Structure
- Catalog manager
 - Query manager
 - one associated with each user (on host)
 - Scheduler processes
 - coordinates multi-site queries (spread out - why?)
 - Operator processes
 - each executes a single relational operator
- 15-721 C. Faloutsos 15

- CMU SCS
- ## Query Processing
- ad-hoc and embedded query interfaces
 - standard parsing, optimization, code gen.
 - left deep trees only
 - hash joins only
 - at most two join operators active simultaneously
- 15-721 C. Faloutsos 16

- CMU SCS
- ## Operator and process structure
- data-flow: each operator process
 - reads input stream
 - output tuples in one or more output streams:
 - split tables (partitioning, joining)
- 15-721 C. Faloutsos 17





- CMU SCS
- ### Detailed outline
- Introduction / History
 - Hardware
 - Software architecture
 - ➔ Query processing
 - Performance evaluation
- 15-721 C. Faloutsos 20

- CMU SCS
- ### Selections
- How?
- 15-721 C. Faloutsos 21

- CMU SCS
- ### Selections
- Start selection operator on each node
 - Exclusion of nodes for hash and range partitioning
 - Throughput considerations
 - one page read-ahead
- 15-721 C. Faloutsos 22

- CMU SCS
- ### Joins
- Partition into buckets, join buckets
 - Implemented: sort-merge, Grace, Simple, Hybrid
 - Parallel Hybrid
- 15-721 C. Faloutsos 23

- CMU SCS
- ### More operations
- Aggregation - How?
- 15-721 C. Faloutsos 24

More operations

- Aggregation
 - Compute partial results for each partition
 - Hash on “group-by” attribute
- Updates
 - Standard techniques
 - How to update partitioning attribute?

Concurrency Control

- How?

Concurrency Control

- 2PL
- Granularity:?

Concurrency Control

- 2PL
- Granularity: file and page
- Modes: S, X, IS, IX, SIX
- Local lock manager and deadlock detector
- (how?)

Concurrency Control

- 2PL
- Granularity: file and page
- Modes: S, X, IS, IX, SIX
- Local lock manager and deadlock detector
- wait-for graph
- Centralized multi-site lock detector

Concurrency Control

- Centralized multi-site lock detector
 - Q: How often to check for deadlocks?

Concurrency Control

- Centralized multi-site lock detector
 - Q: How often to check for deadlocks?
 - A: Period halves/doubles on deadlock / no-deadlock

Recovery

- How?

Recovery

- Standard WAL protocol
 - local log manager generates log records
- one or more log processors
 - collect these log records and write them to the disk

Node Failure

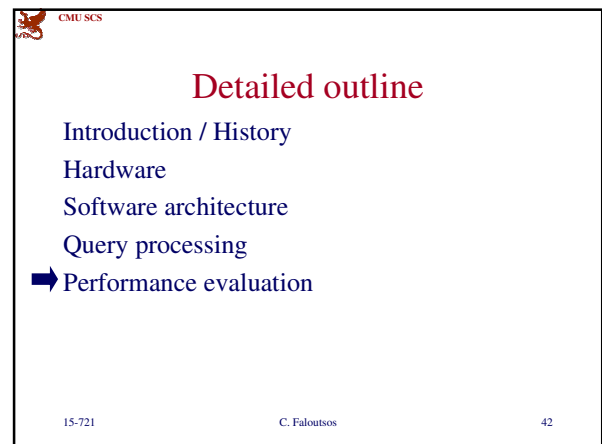
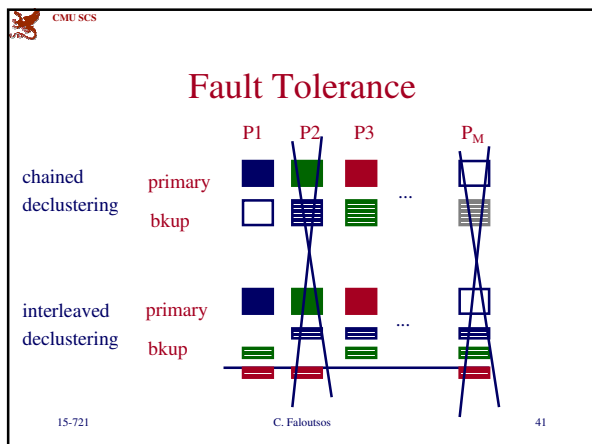
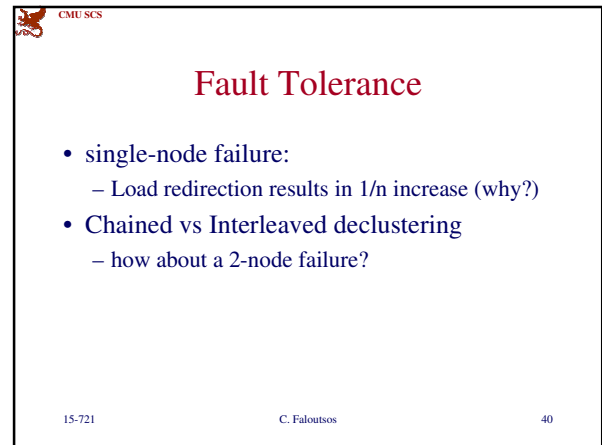
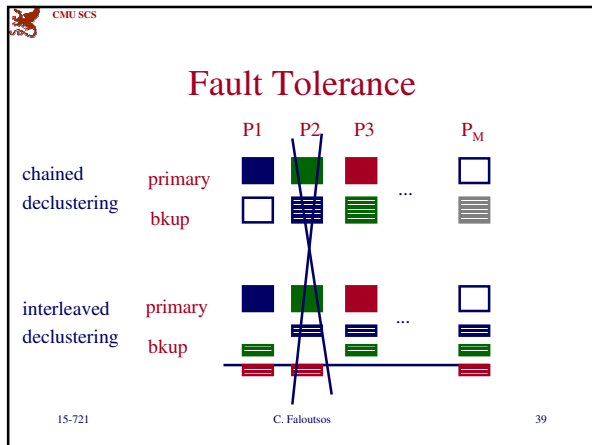
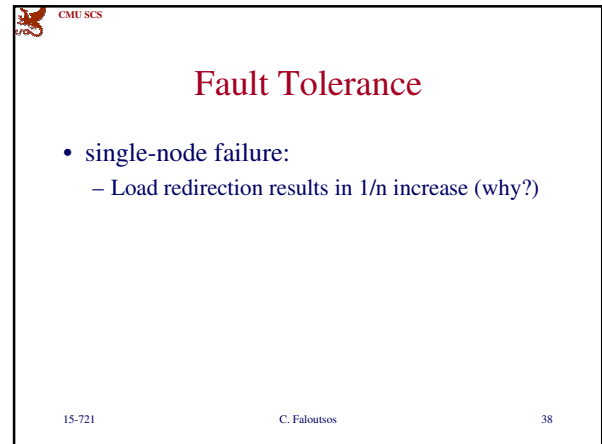
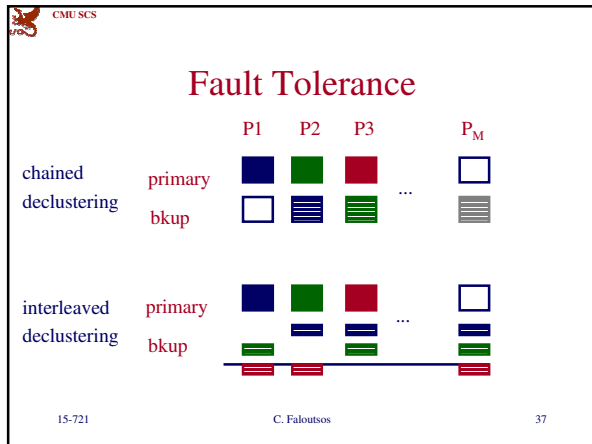
- Goal: Availability in spite of a processor or disk failure
- Goal#2: in that case, spread load as uniformly as possible
- How to handle a (single) node failure?

Node Failure

- Chained declustering (Gamma)
 - backup copy on 'next' node
- Mirrored disk (Tandem)
- Interleaved declustering (Teradata)
 - backup copy: spread over rest of nodes

Fault Tolerance

- specifically:
- Chained Declustering
 - Primary: $i \bmod M$; Backup: $(i+1) \bmod M$
- Interleaved declustering
 - Divide fragments into $N-1$ parts
 - Store parts in all disks but the one containing primary



CMU SCS

Experimental setup

- Wisconsin benchmark (100K, 1M, 10M tuples);
- tables: hash partitioned
- selections (1%, 10%) x (non-indexed, clustered index)
- joins
- wallclock time; speedup; scale-up

15-721 C. Faloutsos 43

CMU SCS

Selections

- non-indexed, 1%, 10%

response time

#processors

15-721 C. Faloutsos 44

CMU SCS

Selections

- non-indexed, 1%, 10%

response time

#processors

15-721 C. Faloutsos 45

CMU SCS

Selections

- non-indexed, 1%, 10%

speed-up

#processors

15-721 C. Faloutsos 46

CMU SCS

Selections

- non-indexed, 1%, 10%
- 10%: superlinear! (why?)

speed-up

#processors

15-721 C. Faloutsos 47

CMU SCS

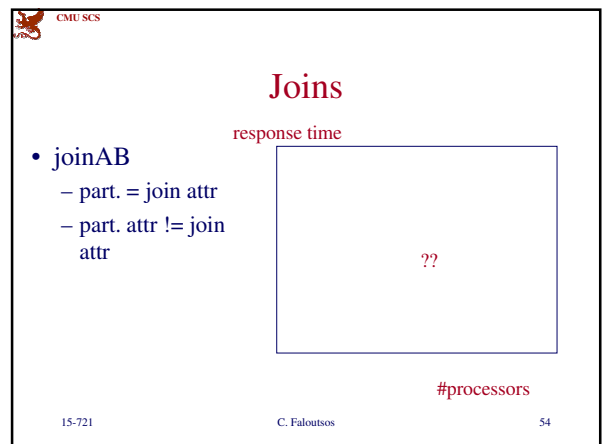
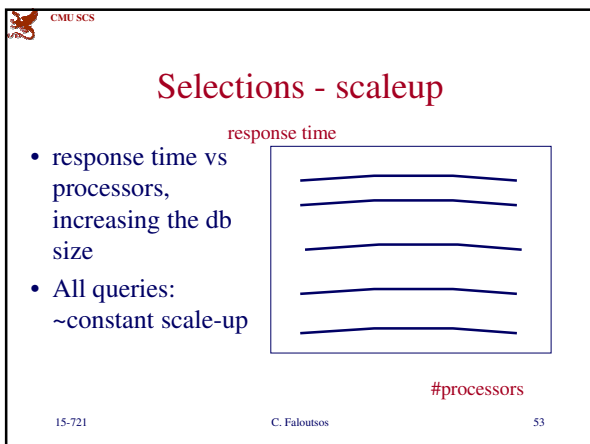
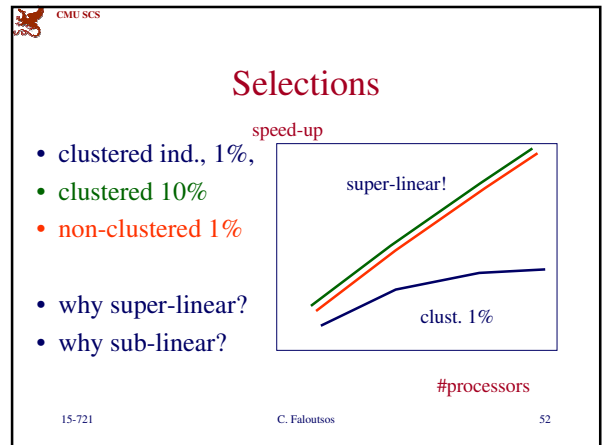
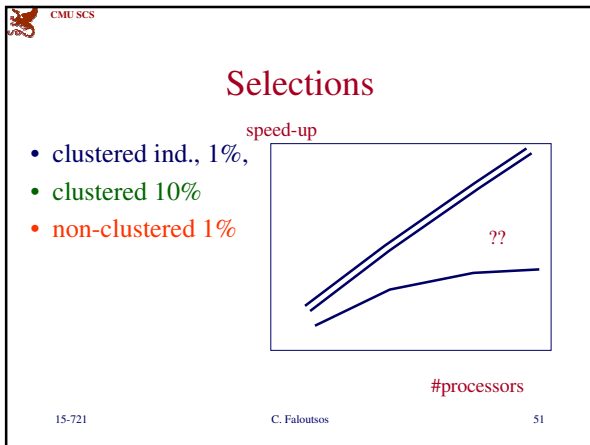
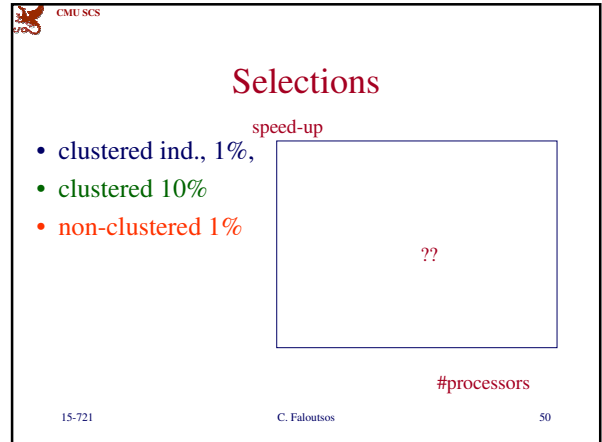
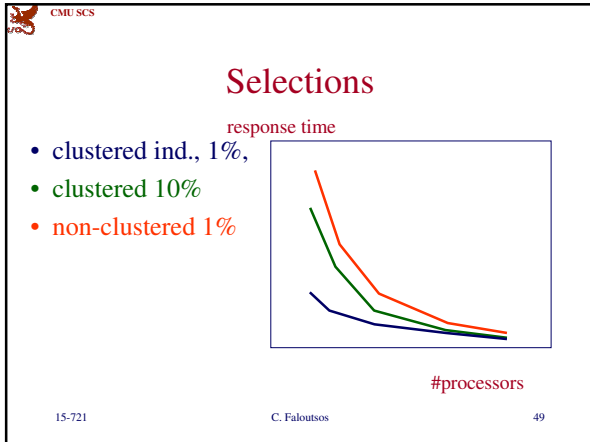
Selections

- clustered ind., 1%, 10%
- clustered 10%
- non-clustered 1%

response time

#processors

15-721 C. Faloutsos 48



CMU SCS

Joins

- joinAB
 - part. = join attr
 - part. attr != join attr

response time

#processors

15-721 C. Faloutsos 55

CMU SCS

Joins

- joinAB
 - part. = join attr
 - part. attr != join attr

speed-up

#processors

15-721 C. Faloutsos 56

CMU SCS

Joins

- joinAB
 - part. = join attr
 - part. attr != join attr
- (scale-up is good, too)

speed-up

#processors

15-721 C. Faloutsos 57

CMU SCS

Conclusions

- Shared-nothing architecture
- horizontal partitioning
- parallel hashing
- data-flow scheduling
- good speed-up and scale-up

15-721 C. Faloutsos 58