

CMU SCS

# 15-721 DB Sys. Design & Impl.

## Distributed DBMSs & R\*

Christos Faloutsos  
[www.cs.cmu.edu/~christos](http://www.cs.cmu.edu/~christos)

CMU SCS

## Roadmap

- 1) Roots: System R and Ingres
- 2) Implementation: buffering, indexing, q-opt
- 3) Transactions: locking, recovery
- 4) Distributed DBMSs
  - intro
  - R\* architecture
  - R\* optimizer
- 5) Parallel DBMSs: Gamma, Alphasort
- 6) OO/OR DBMS
- 7) Data Analysis - data mining

15-721 C. Faloutsos 2

CMU SCS

## Problem – definition

- centralized DB:

15-721 C. Faloutsos 3

CMU SCS

## Problem – definition

- Distr. DB:
- DB stored in many places
- ... connected

15-721 C. Faloutsos 4

CMU SCS

## Problem – definition

now:

connect to LA; exec sql select \* from EMP; ...  
 connect to NY; exec sql select \* from EMPLOYEE; ...

15-721 C. Faloutsos 5


CMU SCS

## Problem – definition

ideally:

connect to distr-LA; exec sql select \* from EMPL;

15-721 C. Faloutsos 6




CMU SCS

## Requirements

- ?

15-721 C. Faloutsos 7




CMU SCS

## Requirements

- location transparency
- performance transparency (-> distr. q-opt)
- copy transparency
- transactions transparency
- fragment transparency
- schema transparency
- local dbms transparency
- (no system has all these features)

15-721 C. Faloutsos 8




CMU SCS

## What's new?

- Q-opt

15-721 C. Faloutsos 9




CMU SCS

## What's new?

- Q-opt
  - communication cost
  - larger search space
  - load balance
  - speed, cost, space, time differences on machines

15-721 C. Faloutsos 10




CMU SCS

## What's new?

- Q-opt
- CC

15-721 C. Faloutsos 11



CMU SCS

## What's new?

- Q-opt
- CC
  - need distributed algorithms (distr. deadlock)

15-721 C. Faloutsos 12

## What's new?

- Q-opt
- CC
- Recovery

## What's new?

- Q-opt
- CC
- Recovery
  - much more complex; more parts that can fail;  
‘2 phase commit’

## What's new?

- Q-opt
- CC
- Recovery
- multiple copies; fragments

## What's new?

- Q-opt
- CC
- Recovery
- multiple copies; fragments
  - (but, at most 2, in practice)

## D-DBMS in practice

Why would one need a D-DBMS?

## D-DBMS in practice

Why would one need a D-DBMS?

- geographic distribution / performance
- off-loading mainframes with local processing
- ‘sins of the past’ - integrating legacy systems

## D-DBMS in practice

- there are products (IBM Data Joiner, Oracle\*)
- BUT: they are not commercially as successful as we would expect! - why?

15-721

C. Faloutsos

19

## D-DBMS - why not?

Speculations:

- data warehouses (copy DBs locally! Sears, Wal-Mart, Kmart)
- D-DBMSs would cut down sales of D/W products
- distr. q-opt is immature

15-721

C. Faloutsos

20

## D-DBMS - other issues?

Integration of data sources: desirable, because of the web - remaining issues:

- semantic consistency (e.g., salaries before/after taxes)
- authentication
- 2-phase-commit on top of legacy databases

15-721

C. Faloutsos

21

## Conclusions

D-DBMS research produced great ideas, useful for

- parallel dbms / “active disks” / sensors
- p2p (peer to peer networks)
- ‘middle-ware’

15-721

C. Faloutsos

22

## Conclusions

Namely:

- 2 phase commit
- distributed q-opt - semi-joins/bloom-joins
- distributed catalogue
- distributed deadlock detection

15-721

C. Faloutsos

23

## Roadmap

- 1) Roots: System R and Ingres
- 2) Implementation: buffering, indexing, q-opt
- 3) Transactions: locking, recovery
- 4) Distributed DBMSs
  - intro
  - ➡ **R\* architecture**
  - R\* optimizer
- 5) Parallel DBMSs: Gamma, Alphasort
- 6) OO/OR DBMS
- 7) Data Analysis - data mining

15-721

C. Faloutsos

24

CMU SCS

## System R\* architecture

Citation  
R. Williams, et al., "R\* : An Overview of the Architecture." IBM Research Report RJ3325

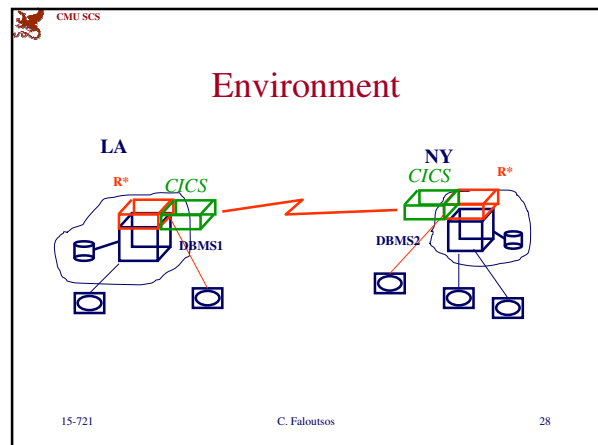
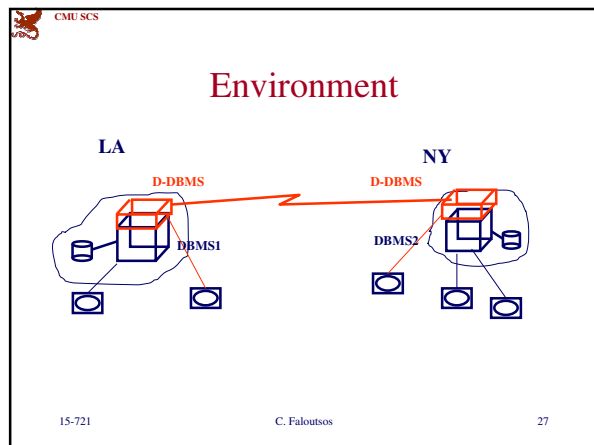
15-721 C. Faloutsos 25

CMU SCS

## Detailed outline

- Environment
- Object naming
- Distributed Catalogues
- Xact management - commit protocols
- Q-opt
- CC- recovery
- SQL changes

15-721 C. Faloutsos 26



CMU SCS

## Environment

Relations:

- dispersed;
- replicated
- fragmented (hor./vert.);
- snapshots

15-721 C. Faloutsos 29

CMU SCS

## Object naming

No global naming system (why?)  
Instead: System Wide Names (SWN)

- by attaching 'site' on user-names
- by attaching 'birth-site' on tables

e.g.:

bruce.EMPLOYEE ->  
bruce@san-jose . EMPLOYEE@yorktown

15-721 C. Faloutsos 30

## Distributed catalogues

- Q: where and how should we store the schema?
- A1: fully replicated (but:....)
- A2: single copy (but:....)
- A3?

## Distributed catalogues

- Q: where and how should we store the schema?
- A1: fully replicated (but:....)
- A2: single copy (but:....)
- A3: only birth site keeps moving info - thus

## Distributed catalogues

- A3: only birth site keeps moving info - thus each site has
  - local schema +
  - moving info (for items 'born' here) and
  - birth sites of global objects
- thus:  $\leq 2$  messages are enough to locate non-local object

## Detailed outline


- Environment
- Object naming
- Distributed Catalogues
- ➡ Xact management - commit protocols
- Q-opt
- Query execution
- SQL changes

## Xact management

- Q: how to give xact-ids?

## Xact management


- Q: how to give xact-ids?
- A: site-id & sequence#
  - ordered (to break deadlocks)


CMU SCS

## Transactions – recovery

- Problem: eg., a transaction moves \$100 from NY -> \$50 to LA, \$50 to Chicago
- 3 sub-transactions, on 3 systems, with 3 W.A.L.s
- how to guarantee atomicity (all-or-none)?
- Observation: additional types of failures (links, servers, delays, time-outs ....)


15-721
C. Faloutsos
37


CMU SCS

## Transactions – recovery

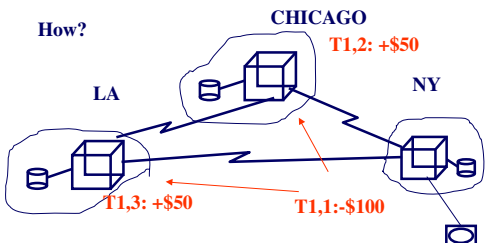
- Problem: eg., a transaction moves \$100 from NY -> \$50 to LA, \$50 to Chicago

15-721
C. Faloutsos
38



CMU SCS

## Distributed recovery

How?

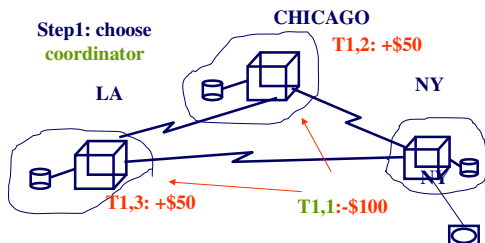


15-721
C. Faloutsos
39



CMU SCS

## Distributed recovery

Step1: choose coordinator




15-721
C. Faloutsos
40


CMU SCS

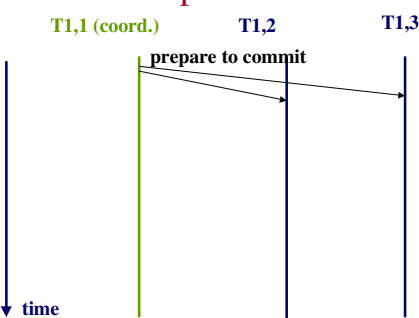
## Distributed recovery

- Step 2: execute a protocol, eg., “2 phase commit”

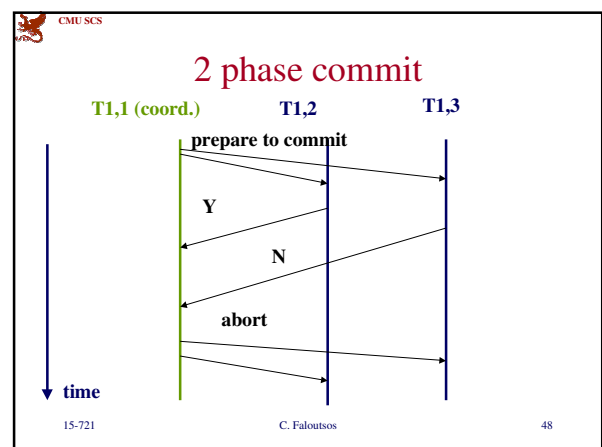
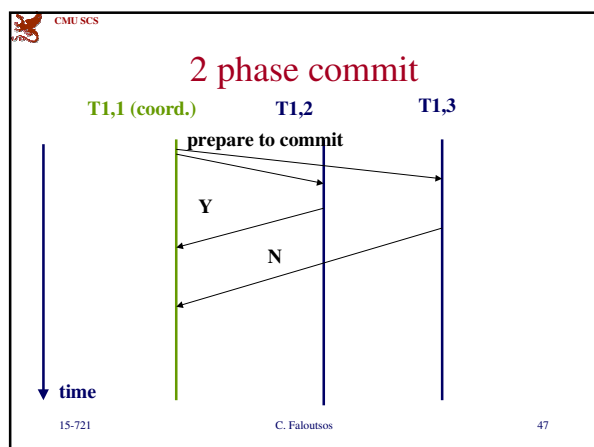
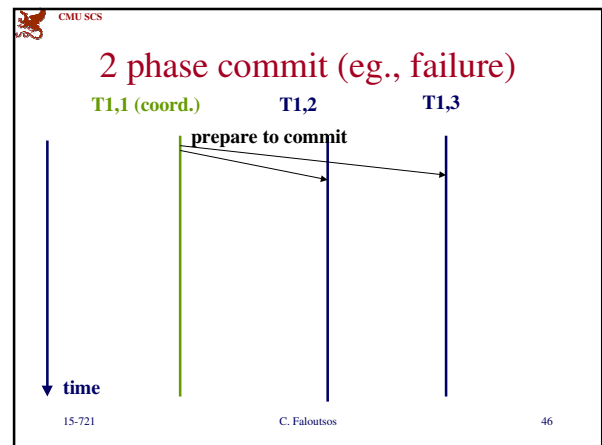
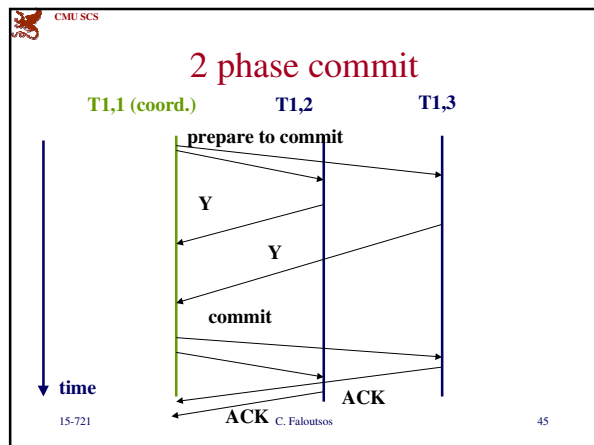
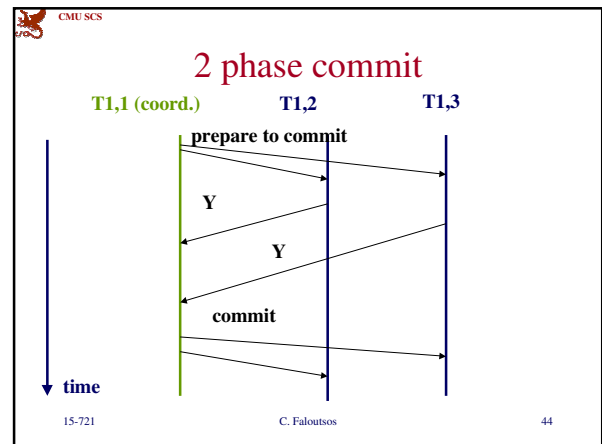
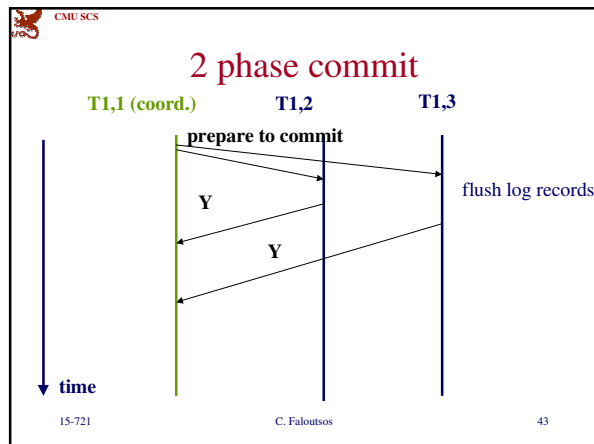
15-721
C. Faloutsos
41


CMU SCS

## 2 phase commit



15-721
C. Faloutsos
42





## Distributed recovery

- Many, many additional details (what if the coordinator fails? what if a link fails? etc)
- and many other solutions (eg., 3-phase commit)

15-721

C. Faloutsos

49

## Detailed outline

- Environment
- Object naming
- Distributed Catalogues
- Xact management - commit protocols
- ➡ • Q-opt
- Query execution
- SQL changes

15-721

C. Faloutsos

50

## Distributed Q-opt

- Steps:
  - parse
  - resolve names
  - authorization
  - compilation + plan generation
    - binding? (eg., an access path may be dropped mid-flight!)

15-721

C. Faloutsos

51

## Distributed Q-opt

- Q: how to do binding?
  - A1: at a chosen site (-> ~ centralized)
  - A2: at the originating site (but: needs much info, which may be out-dated)
  - A3: distributed binding

15-721

C. Faloutsos

52

## Distributed binding

- master site decides inter-site issues + high level binding
  - local sites do low-level decisions
- Local optimality: OK  
 global optimality: NOK  
 Solution: Master sites sends global plan; local sites complain, if things changed

15-721

C. Faloutsos

53

## Distributed q-opt

- cost to minimize?

15-721

C. Faloutsos

54

CMU SCS

## Distributed q-opt

- cost to minimize?
- cost = CPU + I/O + communication
  - comm. cost =
    - msg-cost \* #messages +
    - byte-cost \* #bytes
- (could minimize elapsed time, instead...)

15-721 C. Faloutsos 55

CMU SCS

## Distr. Q-opt –joins

SUPPLIER

s#	...
s1	
s2	
s5	
s11	

SHIPMENT

s#	p#
s1	p1
s2	p1
s3	p5
s2	p9

SUPPLIER Join SHIPMENT = ?

15-721 C. Faloutsos 56

CMU SCS

## Distr. q-opt - join plans

Joins: join order + join method + LOCATION

15-721 C. Faloutsos 57

CMU SCS

## Distr. q-opt - join plans

SEVERAL choices - R\* chooses one of 5:

- (a) ship inner to S1; join there
- (b) ship outer to S2, tuple-at-a-time
- (c) ('semi-join'): reduce inner; ship that to S1
- (d) ship both tables to a third site
- (e) ship outer to a third site; do (c)

15-721 C. Faloutsos 58

CMU SCS

## Semijoins

- Idea: reduce the tables before shipping

SUPPLIER

s#	...
s1	
s2	
s5	
s11	

SHIPMENT

s#	p#
s1	p1
s2	p1
s3	p5
s2	p9

SUPPLIER Join SHIPMENT = ?

15-721 C. Faloutsos 59

CMU SCS

## Semijoins

- Formally:
- SHIPMENT' = SHIPMENT  $\bowtie$  SUPPLIER

15-721 C. Faloutsos 60

CMU SCS

## Detailed outline

- Environment
- Object naming
- Distributed Catalogues
- Xact management - commit protocols
- Q-opt
- • Query execution
- SQL changes

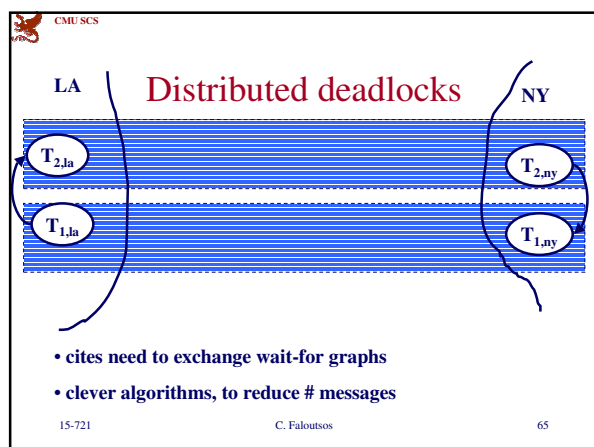
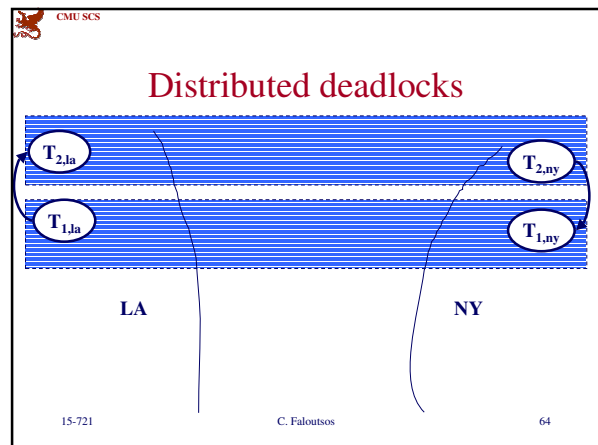
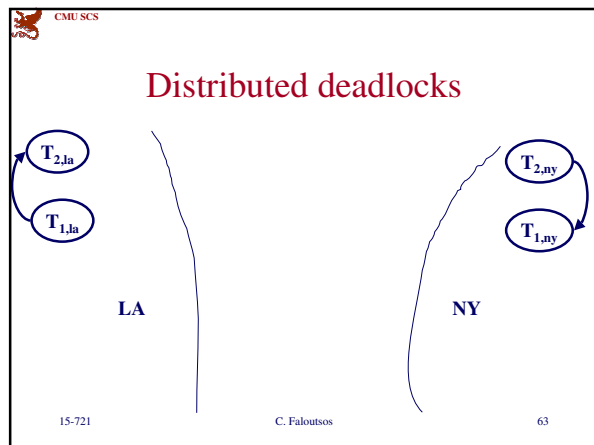
15-721 C. Faloutsos 61

CMU SCS

## Detailed outline

- Environment
- Object naming
- Distributed Catalogues
- Xact management - commit protocols
- Q-opt
- • Query execution
- SQL changes

15-721 C. Faloutsos 62



CMU SCS

## Distributed deadlocks

- cites need to exchange wait-for graphs
- clever algorithms, to reduce # messages
  - naively: each site ships its wait-for strings, until all have all
  - anything better?

15-721 C. Faloutsos 66

## Distributed deadlocks

- anything better?
- A: each site ships ONLY the strings where 'first-xact-id' < 'last-xact-id'
  - (any other ordering, is fine!)
- Eg: LA: T1-> T2; NY T2->T1
  - only NY will send

15-721

C. Faloutsos

67

## Detailed outline

- Environment
- Object naming
- Distributed Catalogues
- Xact management - commit protocols
- Q-opt
- Query execution
- ➡ • SQL changes

15-721

C. Faloutsos

68

## SQL extensions

- DEFINE SYNONYM <rel-name> AS <SWN>
- DISTRIBUTE TABLE <t-name>  
HORIZONTALLY | VERTICALLY |  
REPLICATED ...
- DEFINE SNAPSHOT ...
- REFRESH SNAPSHOT
- MIGRATE TABLE ...

15-721

C. Faloutsos

69

## Conclusions

- 2 phase commit
- distributed q-opt; distr. deadlock detection
- distributed catalogue

15-721

C. Faloutsos

70