


**CS307**

**Database Principles**


Stéphane Faroult  
[faroult@sustc.edu.cn](mailto:faroult@sustc.edu.cn)

Liu Zijian    [liuzijian47@163.com](mailto:liuzijian47@163.com)

**NEXT WEEK !**

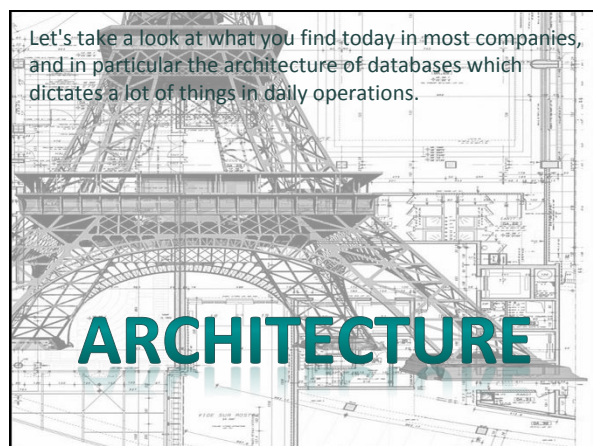


**Guy Harrison**



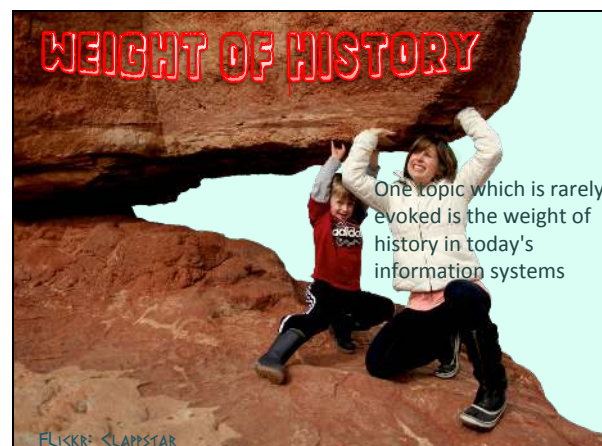
Let's take a look at what you find today in most companies,  
and in particular the architecture of databases which  
dictates a lot of things in daily operations.

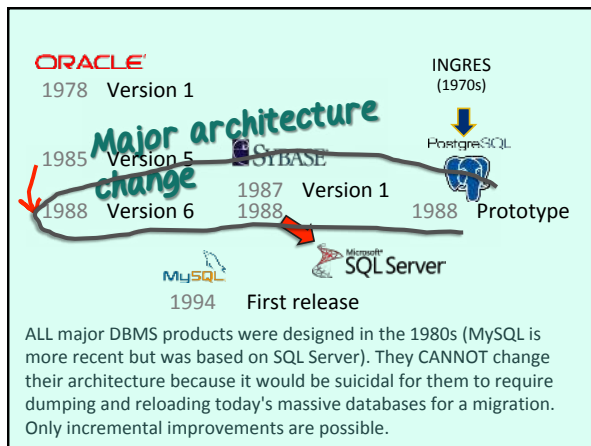
**ARCHITECTURE**



**WEIGHT OF HISTORY**

One topic which is rarely  
evoked is the weight of  
history in today's  
information systems





## Mid 1980s

This is what a \$M 1 machine was looking like in the 1980s.

VAX 8600 (high-end) **digital**

32-bit architecture

Processor, ~10 to 20MHz clock

4 to 256 Mb of memory

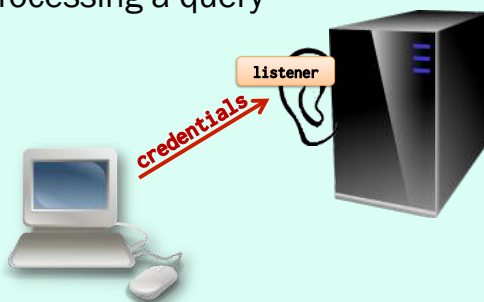
I/Os ~10 to 30 Mb/s

Keep in mind that the big DBMS products were designed for this.

*Dumb hardware* *Beginning of multi-processor computers*

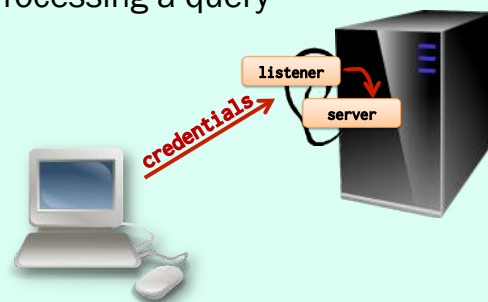
OK, how are we going to process a query? First we must send our credentials to the DBMS server and establish a session.

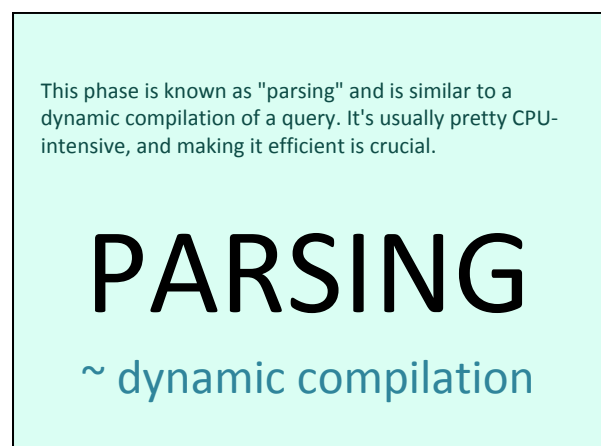
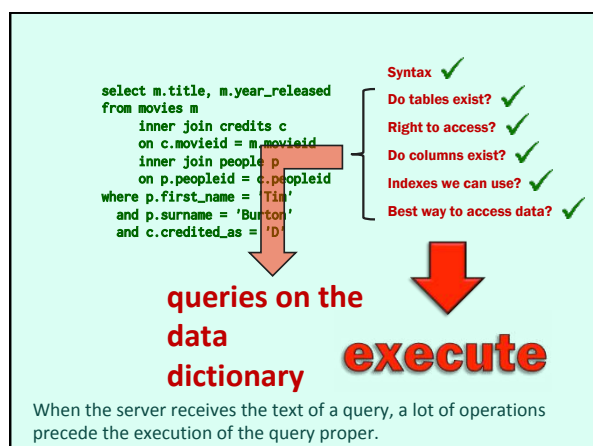
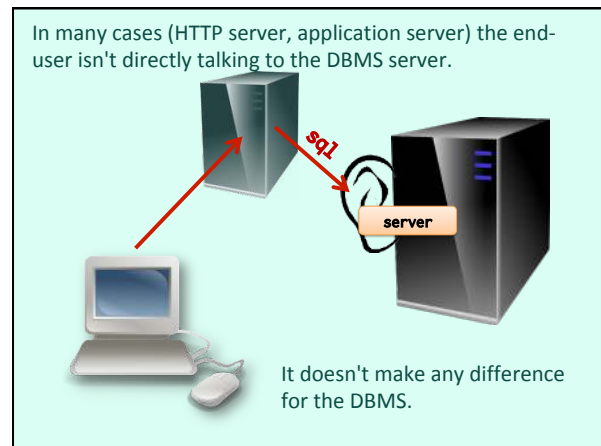
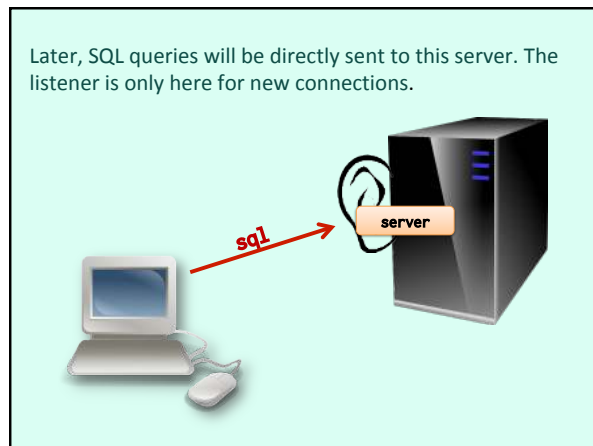
## Processing a query



The listener may fork a process that will become our personal server and execute queries for us.

## Processing a query





```

select m.title, m.year_released
from movies m
  inner join credits c
    on c.movieid = m.movieid
  inner join people p
    on p.peopleid = c.peopleid
where p.first_name = 'Tim'
and p.surname = 'Burton'
and c.credited_as = 'D'

```

- Syntax ✓
- Do tables exist? ✓
- Right to access? ✓
- Do columns exist? ✓
- Indexes we can use? ✓
- Best way to access data? ✓

One way to improve efficiency is to keep data dictionary information (meta-data) in a shared cache to avoid additional queries.

Kept in <sup>shared</sup> memory → **meta data**

```

select m.title, m.year_released
from movies m
  inner join credits c
    on c.movieid = m.movieid
  inner join people p
    on p.peopleid = c.peopleid
where p.first_name = 'Tim'
and p.surname = 'Burton'
and c.credited_as = 'D'


```

- Syntax ✓
- Do tables exist? ✓
- Right to access? ✓
- Do columns exist? ✓
- Indexes we can use? ✓
- Best way to access data? ✓

Another crucial phase is the one when the optimizer tries to determine the most efficient way to access data.

Kept in <sup>shared</sup> memory → **meta data**

## QUERY OPTIMIZER



The query optimizer is a vital element in a DBMS product. You have seen that we can obtain the same result by writing the query in different ways, but that "entry points" and indexes used can be completely different.

flickr: Brian Hilligas

```

SELECT Customers.CustomerName, Orders.OrderID
FROM Customers
INNER JOIN Orders
ON Customers.CustomerID=Orders.CustomerID
ORDER BY Customers.CustomerName

```

## TEXTBOOK

The problem is that if SQL queries are simple in textbooks, they are far more complicated in real life and, like in a chess game, you cannot explore all possibilities. Query optimizing is part of the overall response time.

REAL LIFE

Complex view

## QUERY OPTIMIZER



Ha  
Th  
can  
lot

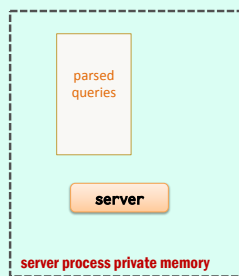
hlickr:Brian Hillegas

As a result

# PARSING takes time

5

## Keep parsed queries in memory

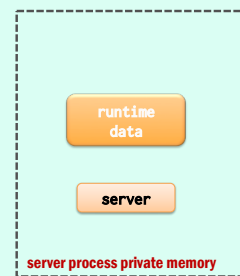


As most applications run exactly the same SQL statements again and again, a DBMS will cache a parsed query for reuse. For MySQL, it will be cached for a session.

MySQL

For Oracle and SQL Server, cached parsed queries will be shared between sessions.

## Keep parsed queries in memory



SQL Server  
ORACLE

## Query cache management

# LRU

## Least Recently Used

Of course we cannot hold in cache zillions of parsed queries. We need to manage the cache, and replace queries that haven't been executed in a while with new ones.

```
select m.title, m.year_released
from movies m
  inner join credits c
    on c.movieid = m.movieid
  inner join people p
    on p.peopleid = c.peopleid
where p.first_name = 'Tim'
   and p.surname = 'Burton'
   and c.credited_as = 'D'
```

Checksum

We primarily recognize identical queries by computing a text checksum.

+ check tables are same  
and context identical

```
select m.title, m.year_released
from movies m
  inner join credits c
    on c.movieid = m.movieid
  inner join people p
    on p.peopleid = c.peopleid
where p.first_name = 'Howard'
and p.surname = 'Hawks'
and c.credited_as = 'D'
```

And here we have a problem because this looks like a different query from the preceding one, yet there is no reason why the execution plan should be any different.

Basically the same query

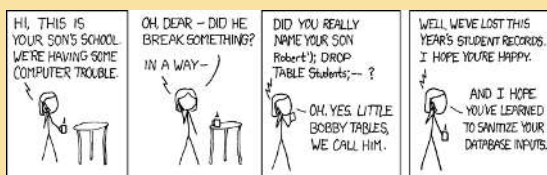
DIFFERENT CHECKSUM!

```
select m.title, m.year_released
from movies m
  inner join credits c
    on c.movieid = m.movieid
  inner join people p
    on p.peopleid = c.peopleid
where p.first_name = :first_name
and p.surname = :surname
and c.credited_as = 'D'
```

This is the main reason why we should use parameters and "bind" them at run time. The text of the query is unchanged whichever name we attach to it. Note that I keep 'D' because, as we have more actors, the plan might be different for them.

"variable binding"

## (Parenthesis)



<http://xkcd.com/327/>


Additionally, of course, variable binding protects against SQL injection.

## Some issues:

~~select ...  
from :my\_table  
select ...  
from ...  
order by :user\_choice  
select ...  
from ...  
where col in (:user\_list)~~

And unfortunately you cannot bind a table or column name, nor (which would be quite useful) a whole list of values.

**Some issues:  
Selectivity?**



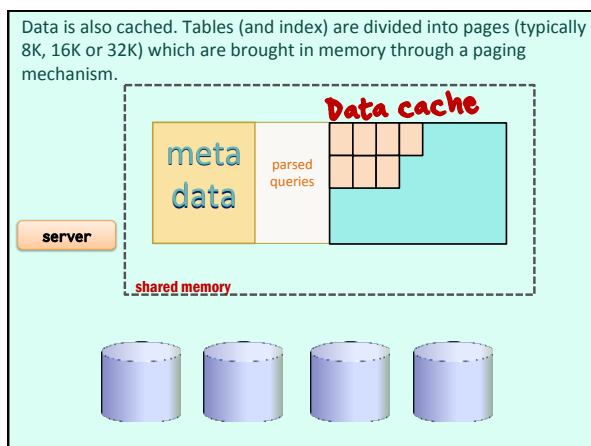
There is another problem with parameters: the choice of using or not using an index depends mostly on data selectivity. The optimizer will base its decision on the values passed when the query is first parsed. Nothing guarantees that the decision will still be adequate with other parameter values.

flickr: Brian Hillegas

## More radical Cache Query Result

Read Only/Slowly Changing

Some products can go further than caching parsed queries and can even cache results associated with a query and its parameters. It's OK if the database is a read-only database or changes only very slowly (for instance, for a query such as retrieving the last articles in a blog).



Now what happens when data is changed?

## SEVERAL problems

Atomicity	Transactions (rollback)
Durability	Persistence
	Consistency
Isolation	Concurrency


32



The ACID acronym predates relational database and symbolizes everything that is desirable to have with a transactional databases.

**ACID**

The dogma that every database management system should be "ACID" has been questioned by the NOSQL movement (more about this later)



33

**A**tomicity

Atomicity is related to a transaction: seen as a single unit.

**C**onsistency

Data should always be seen as consistent (constraints)

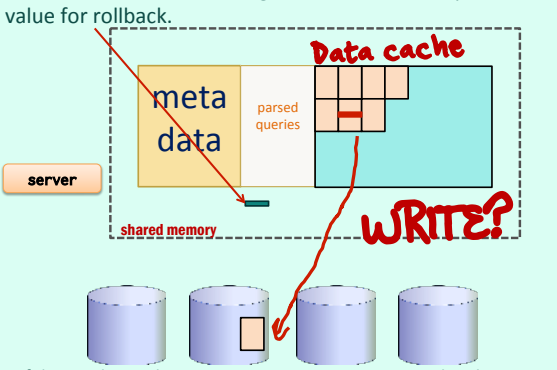
**I**solation

When one user changes data, it should not make it inconsistent to others.

**D**urability

When something is saved to a database, it should become persistent, whatever happens.

For instance, before a change we must save the previous value for rollback.



server

meta data

parsed queries

data cache

shared memory

WRITE?

If data is changed in memory, are we going to immediately write it to file for durability?

35

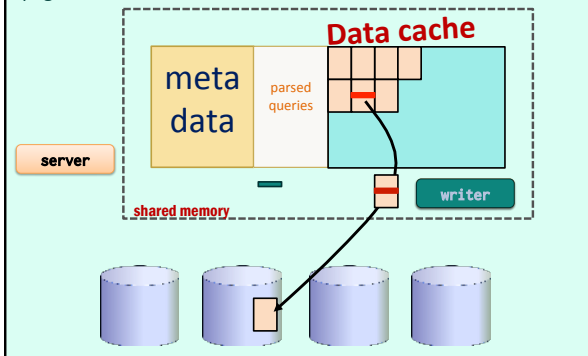
In fact, this would be a massive performance concern (especially in the 1980s, when there was no cache in disk units). With a massive update (massive load, for instance) we would just keep on waiting for data to be written to disk, and it would become a massive bottleneck.

**MASSIVE UPDATE**

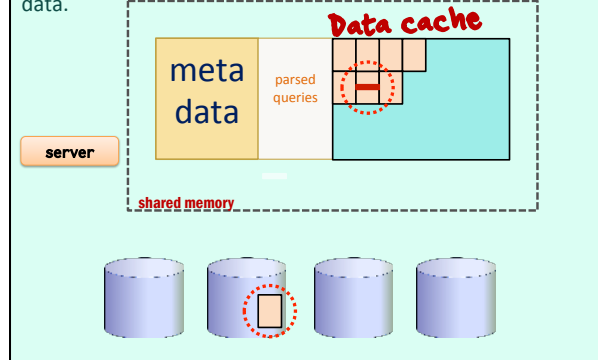
**I/O STORM**

36

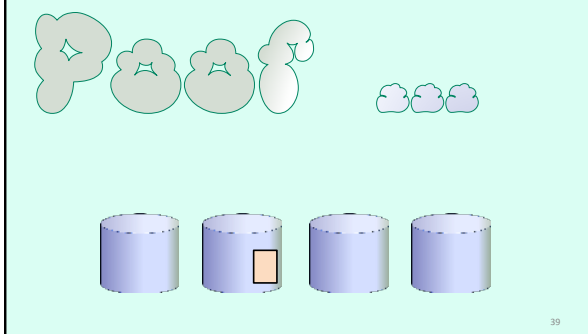
What will happen with a product such as Oracle is that a dedicated writer process will, **ASYNCHRONOUSLY**, write "dirty" pages to files from time to time.



So, at any given time, memory (representing the true database state) and files will not exactly contain the same data.



What happens if the system crashes and everything that was in memory simply vanishes?



## Not committed?

If the change wasn't committed, then losing it is just what we would expect: a crash should rollback any on-going transaction.

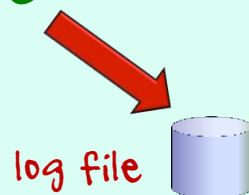
## Committed?

If the change was committed it's a different story. The D in ACID means that we expect any committed change to survive a system crash.

# DURABILITY

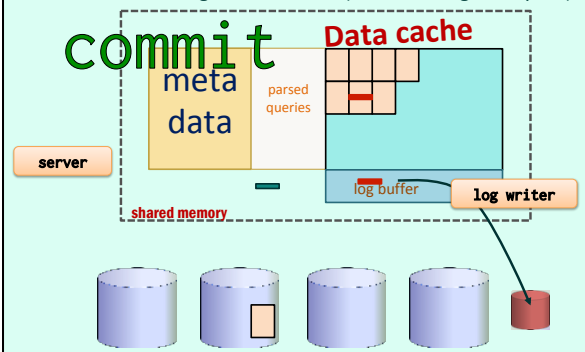
Practically, it means that every COMMIT should result into a SYNCHRONOUS file write: the commit call only returns when changes are safely written to disk. We aren't going to update all the blocks affected by the transaction, which would take too much time, but only save the new values to an additional log (journal) file.

# commit



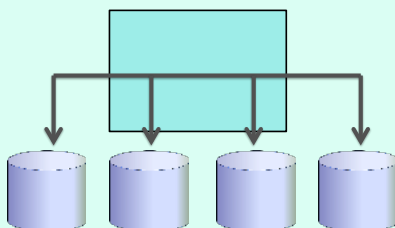
41

In Oracle, every pending change (new values, as opposed to old values kept for rollback) will also be saved to a "log buffer" that will be flushed to a log file on commit (or when dangerously full)



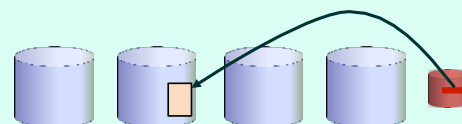
When you shutdown the database cleanly, all the buffers are flushed to files, which represent the last state in memory.

Clean shutdown :  
Flush buffers to files



When restarting after a crash, the log is scanned for all committed changes, and they are applied if need be to the file blocks, until the files contain the exact image of the database at the last COMMIT before the crash.

RESTART AFTER CRASH

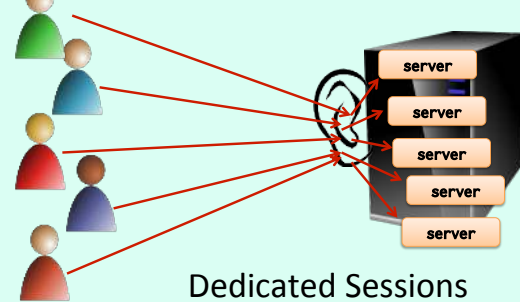


44

### What about **several** sessions?

So far we have, kind of selfishly, considered only the interaction of our session with the database. One of the reasons why databases were created was to share data, which of course implies many concurrent sessions and many interesting problems.

First of all, we may assign to all connected users their own personal server process running SQL for them.



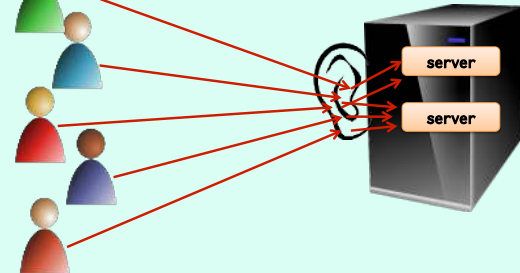
### Dedicated Sessions

It works well, it has been the Oracle standard for years.

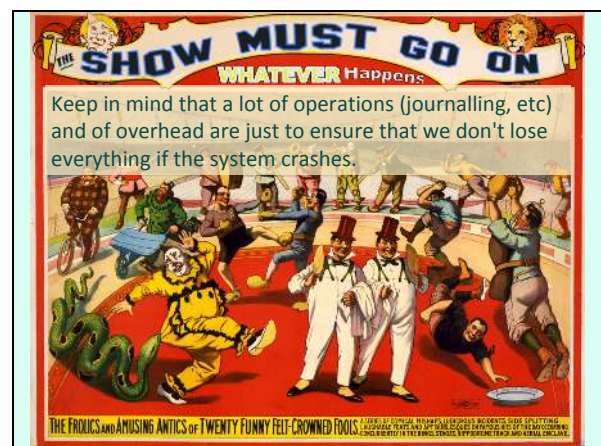
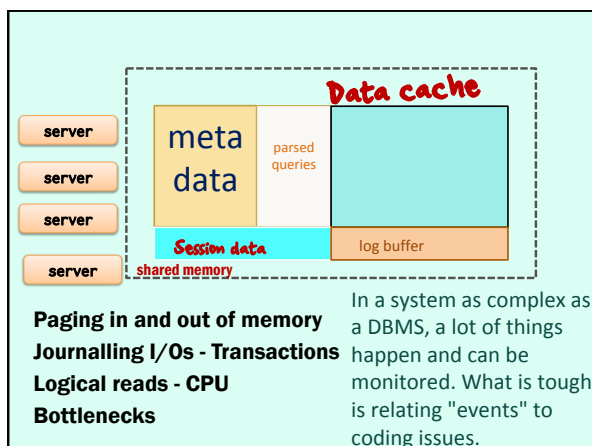
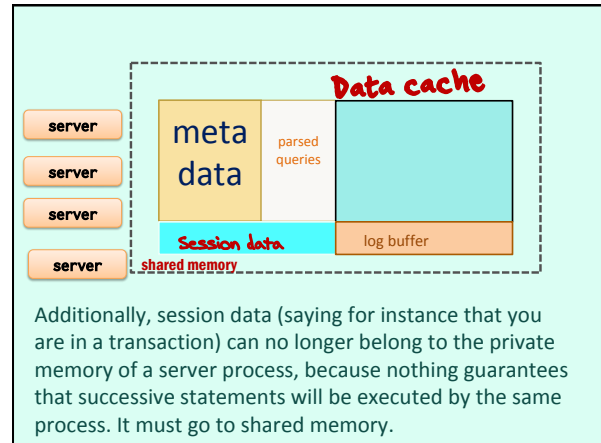
The only problem is that when you have hundreds or thousands of currently connected users (it happens), it's kind of heavy for the operating system, even if all processes aren't active at the same time.



An alternative is to pre-start (which makes connection slightly faster) a number of servers, and to assign requests to them with a bit of load-balancing.



### Session Pooling



The real fun, though, begins when concurrent sessions start modifying data required by several of them.

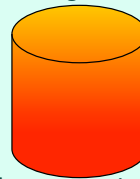
What about a session **querying** data **being changed** by another session?

## ISOLATION

53

## Data Change

The main problem is that a transaction can last "a certain time". You may need to run a lot of code before deciding on commit or rollback.



**Begin Transaction**

**Commit**

During the transaction, the database will be in a kind of transient state of which you cannot say whether it will become permanent.

## CONSISTENT STATE

Between two consistent states, you have a kind of black hole



**CONSISTENT STATE**

## ATOMICITY

## 4 levels for ISOLATION

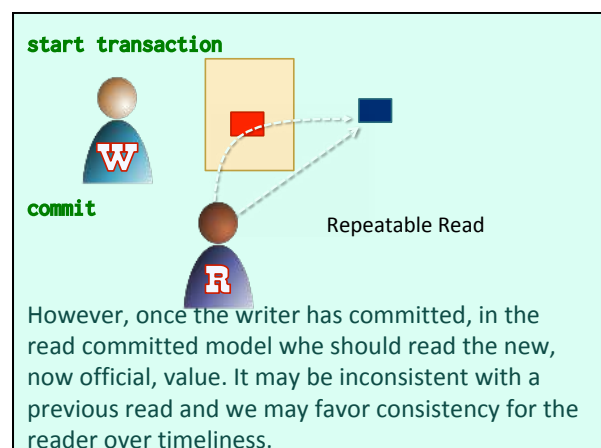
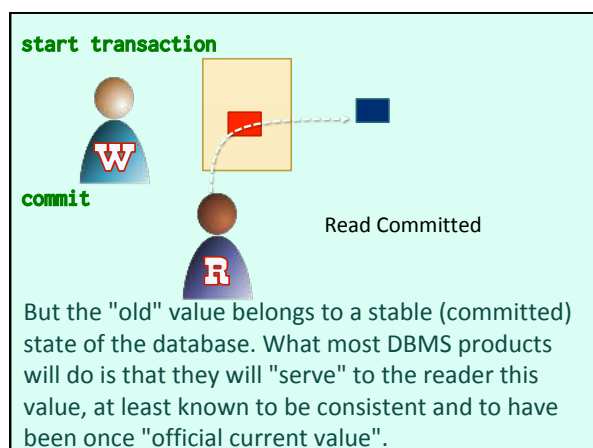
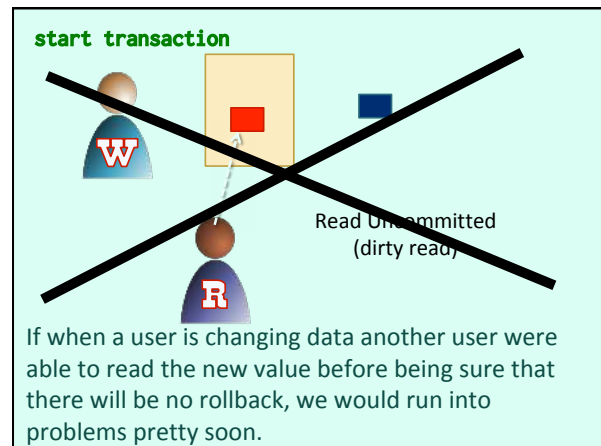
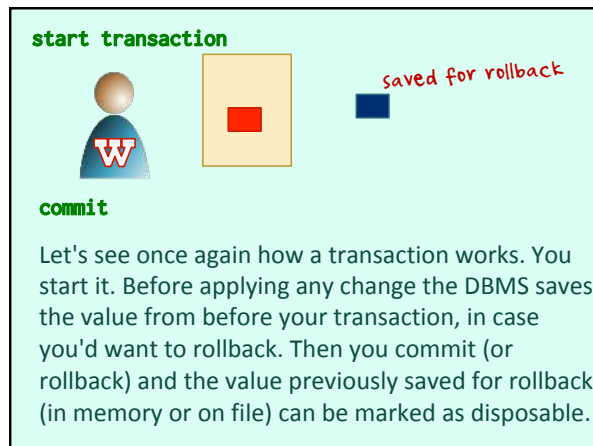
~~dirty reads~~

read committed

repeatable read

serialization

The SQL standard defines four isolation levels, from no isolation at all to paranoid. Some products let you set it, others impose it.



In practice, the problem is more a problem of data consistency of foreign keys when we are scanning big related tables. A single SELECT will usually be consistent (a product such as Oracle ignores any change having happened since the start of the SELECT, even if it was committed). But if we SELECT twice (two different queries) from two tables with an FK relationship, changes that may have happened (and have been committed) between the time when we started reading the first table and the time when we were reading the second table may lead to problems such as orphaned rows.

Repeatable read

Point-in-time read



If I haven't read an uncommitted order when reading ORDER, then I should ignore rows referring to this order in ORDER\_DETAIL even if they have been committed by the time I reach them.



start transaction



commit



Serializable

LOCK

The last level is to block readers when data is being modified. Guaranteed to be consistent but bad for performances.

When you have two writers you cannot do otherwise than locking.  
Table-level  
Row-level

## BACKUP ISSUES

All of this, the fact that files and memory may not been completely in synch, plus the fact that people may be modifying the data, leads to most interesting issues when you want to backup a database without stopping it.

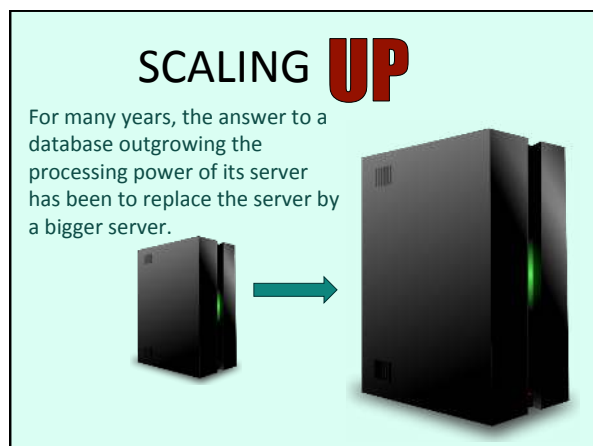
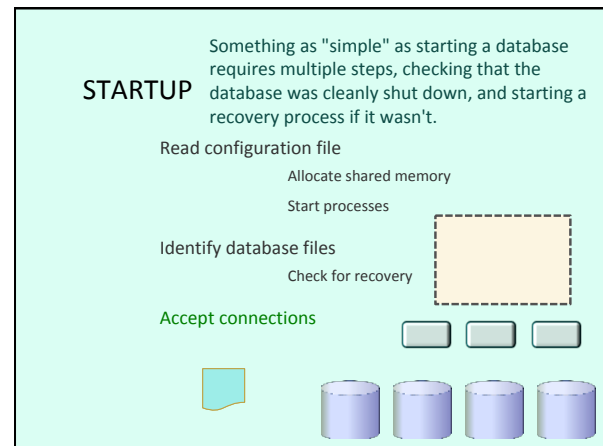
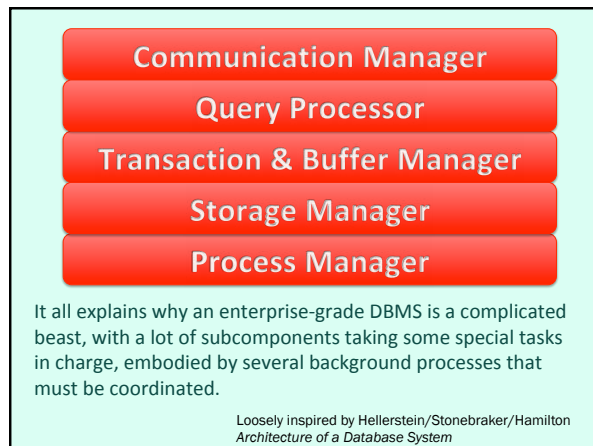
**Backing up something consistent?**

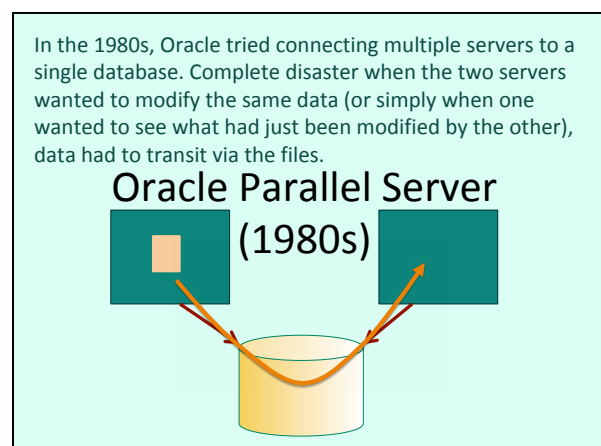
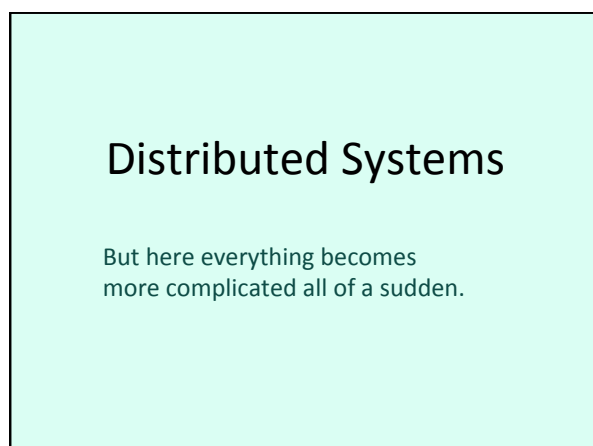
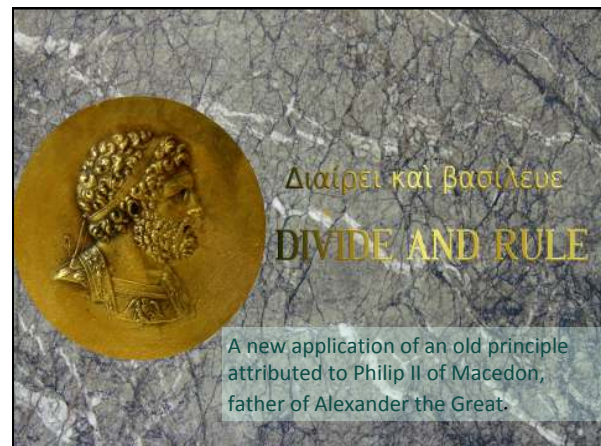
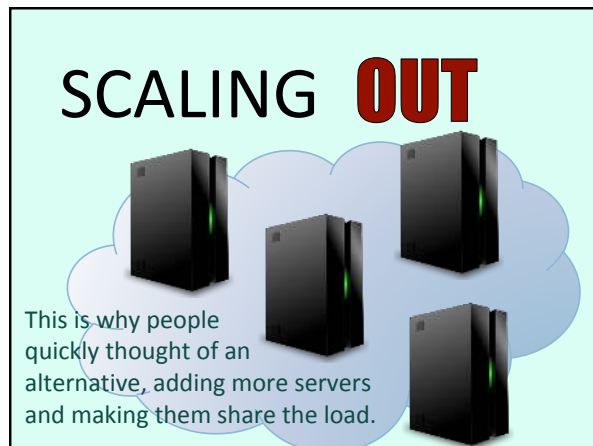
**Files image of memory?**

**No file being written while copied?**

**Need to prevent changes?**







In 1994 Oracle bought from Digital Equipment (which they swallowed whole later), RDB, a respected relational database management system that only worked on Digital systems (SW/HW).

**1994** RDB acquired by Oracle

One of the strengths of RDB was working with clusters of machine; it was the only product doing it well. All of a sudden, Oracle Development gained strong competencies in this area.

## Oracle RAC (2001)

### Real Application Clusters



Seven years later, Oracle came out with a very good clustering offer, in which a very-high-speed private network was connecting the servers for exchanging data blocks.

**Coordinator node**  
**Node "owns" blocks**  
**Owner can change**

RAC implements a complicated "cache-fusion" algorithm in which a coordinator nodes always knows who is working on what, and blocks can be exchanged quickly. Of course it works better when the workload is fairly different for all servers, but when different servers can work in parallel on disjoint data, it can be extremely efficient and it handles conflicts well.

Neil Gunther's

### Universal Law of Computational Scalability

Relative capacity  $C(N)$  of a computational platform:

$$C(N) = \frac{N}{1 + \alpha (N-1) + \beta N(N-1)}$$

$$0 \leq \alpha, \beta < 1$$

$\alpha$  Level of contention

$\beta$  Coherency delay (latency for data to become consistent)



[www.perfdynamics.com](http://www.perfdynamics.com)

The problem with clustering is that it follows a law of diminishing returns: adding a second server will less than double your capacity, and in practice people have clusters of 2, 3 or 4 machines at most (Neil Gunther is a famous Australian consultant/academic specializing on performance)