

**CS307**

**Database Principles**

Stéphane Faroult  
faroult@sustc.edu.cn

朱悦铭 Zhu Yueming zhuym@sustc.edu.cn

Communication Manager

Query Processor

Transaction & Buffer Manager

Storage Manager

Process Manager

What we have seen last time explains why an enterprise-grade DBMS is a complicated beast, with a lot of subcomponents taking some special tasks in charge, embodied by several background processes that must be coordinated.

Loosely inspired by Hellerstein/Stonebraker/Hamilton  
Architecture of a Database System

**STARTUP**

Something as "simple" as starting a database requires multiple steps, checking that the database was cleanly shut down, and starting a recovery process if it wasn't.


- Read configuration file
- Allocate shared memory
- Start processes
- Identify database files
- Check for recovery
- Accept connections

**SCALING UP**

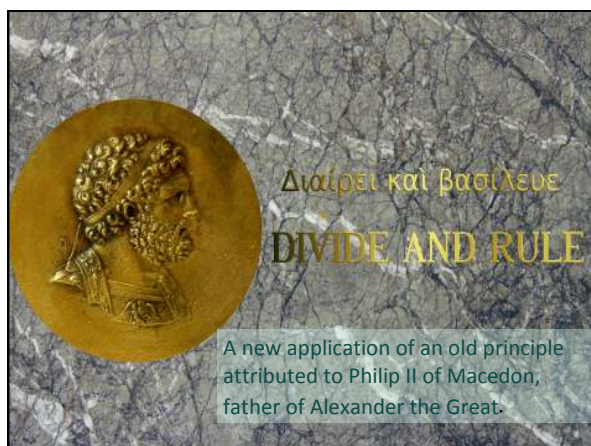
For many years, the answer to a database outgrowing the processing power of its server has been to replace the server by a bigger server.



## SCALING OUT



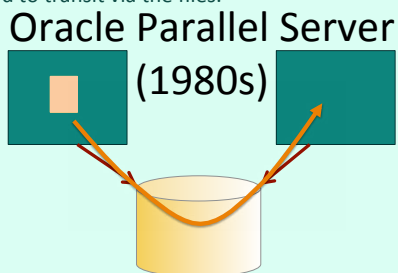
This is why people quickly thought of an alternative, adding more servers and making them share the load.



## Distributed Systems

But here everything becomes more complicated all of a sudden.

In the 1980s, Oracle tried connecting multiple servers to a single database. Complete disaster when the two servers wanted to modify the same data (or simply when one wanted to see what had just been modified by the other), data had to transit via the files.



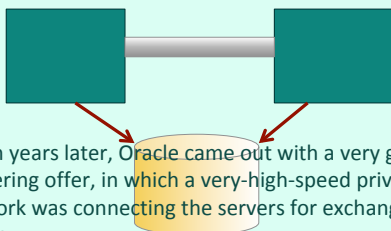
In 1994 Oracle bought from Digital Equipment (which they swallowed whole later), RDB, a respected relational database management system that only worked on Digital systems (SW/HW).

**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)

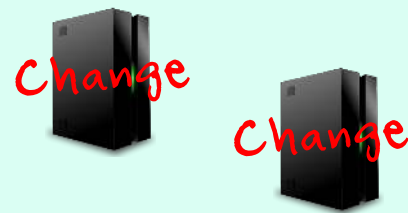
The problem with clustering is that it follows a law of diminishing returns: adding a second server will less than double your capacity, and in practice people have clusters of 2, 3 or 4 machines at most (Neil Gunther is a famous Australian consultant/academic specializing on performance)



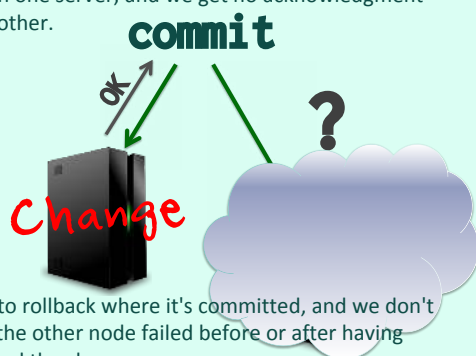
www.perfdynamics.com

One big problem is with transactions that involve several servers. Remember that transactions are meant to be atomic operations.

## Distributed Transactions



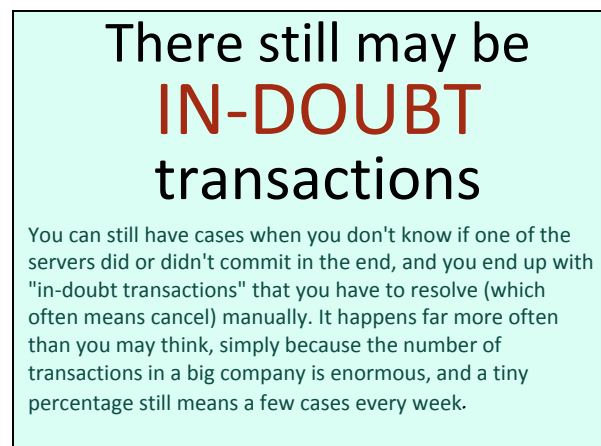
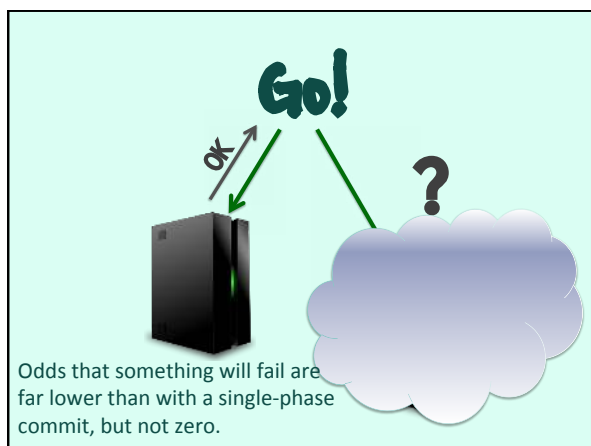
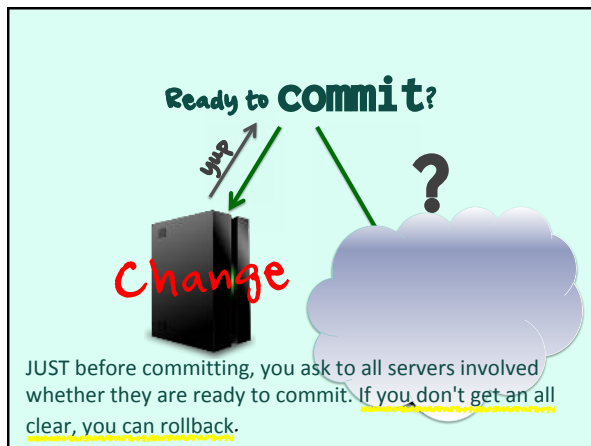
It may happen that when we commit we know for sure it worked on one server, and we get no acknowledgment from the other.



No way to rollback where it's committed, and we don't know if the other node failed before or after having committed the change.

## TWO-PHASE COMMIT

One algorithm was devised (a long time ago), called a "two-phase commit".



## Latency

Additionally, you have latency issues. All machines in a cluster may not be sitting next to each other, they may be a few miles apart in different data centers for security reasons (fire, flood ...)

**1 KM = 0.000005s**

Even if information travels fast, multiplying exchanges (two-phase commit) may become a sensitive issue.

## Synchronous? Asynchronous?

At which point you have to choose between a synchronous mode in which you always wait for an acknowledgment that everything went fine, or an asynchronous mode in which you cross fingers and switch to something else. The type of hardware you have may influence your decision: if all your disks have a big buffer and a reliable battery for instance, you may be more likely to trust them.



**Senior Database Administrator**

What you will do:

**Job posting (US)**

- Manage MySQL in production/QA/dev environments including installation, configuration, upgrades, schema changes, etc.
- Troubleshoot database issues, maintain database systems availability and scalability within production environments
- Perform capacity planning exercises to properly identify required hardware, software, database configuration/architecture necessary to support application needs
- Monitor database performance, identify performance problems and make adjustments to database parameters as needed.
- Monitor key performance indicators and make enhancements to improve/maintain performance/productivity at acceptable levels
- Enforce best practices for improving performance, scalability and operational manageability of production databases
- Part of on-call rotation to respond to and resolve application issues to ensure production applications are online

DBAs come in two flavors (in smaller companies, it's the same guy). They often have completely different backgrounds, and don't necessarily agree on everything.

**PRODUCTION DBA**

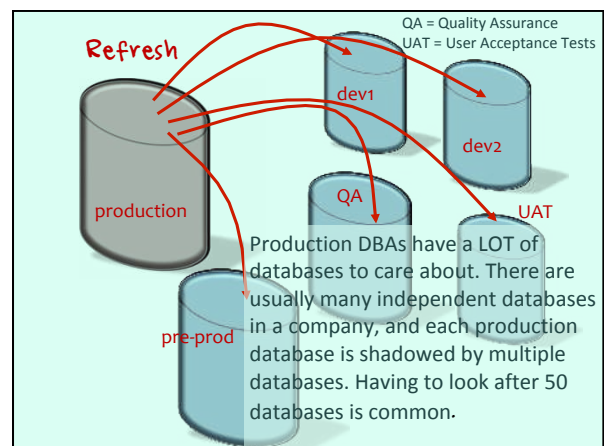
monitoring  
shell scripting  
security

**System Engineer**

planning  
architecture  
SQL scripting

**Developers****DEVELOPMENT DBA**

In many ways, a development DBA is a developer who knows more than the average about databases. We'll focus here a little more on the job of a production DBA, because architects who design information systems aren't always well aware of what their job is, and some "solutions" that look great on the paper are sometimes a hell to maintain in a daily production, with the constraints of keeping systems available and running as much as possible. Keeping everything simple and manageable should be the first concern of every architect.



If you aren't square, you can't survive in this job.  
You must automate as many tasks as possible.

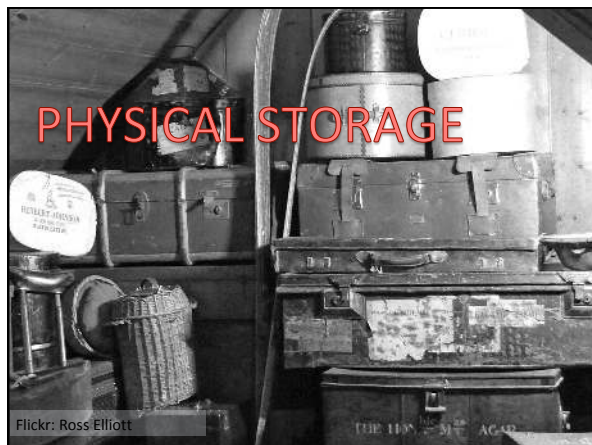
# STANDARDS ORGANIZATION AUTOMATION

Other than performance issues (than we have partly seen) these are the big topics for a DBA. We have talked a bit about privileges in labs, let's talk about storage.

Physical Storage

User Management

Backup / Recovery



## DATABASE?

To start with, the word "database" has a different meaning depending on the DBMS (and sometimes DBMS version) that you are talking about.





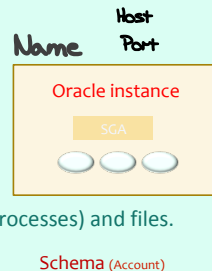
## User's View

ORACLE

&lt; 12c

For instance with Oracle up to 12c a database is basically an instance (shared memory and processes) and files. A schema is a user account that owns tables.

```
select ...
from schema.table_name
synonym_name
```



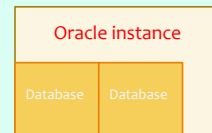
Schema (Account)



## User's View

ORACLE

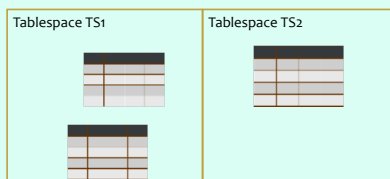
&gt; 12c



From Oracle 12 the vision is altered, and an instance becomes a container that may hold several independent databases. Prior to Oracle 12 it was common to have multiple instances, each with their (big) shared memory and processes on the same server. Oracle 12 pools resources.

## Developer's View

ORACLE

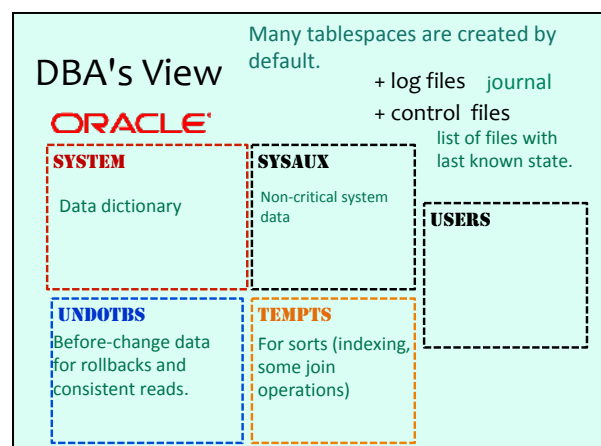
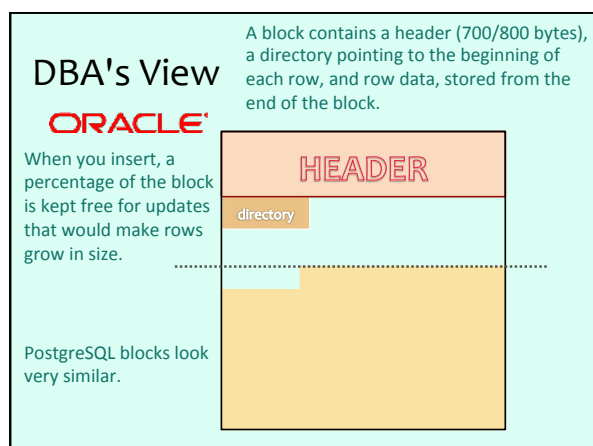
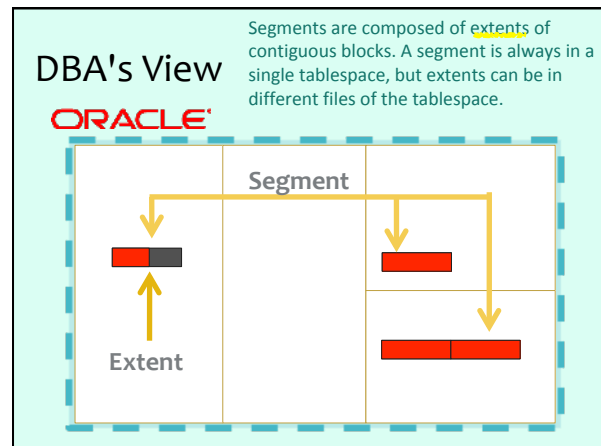
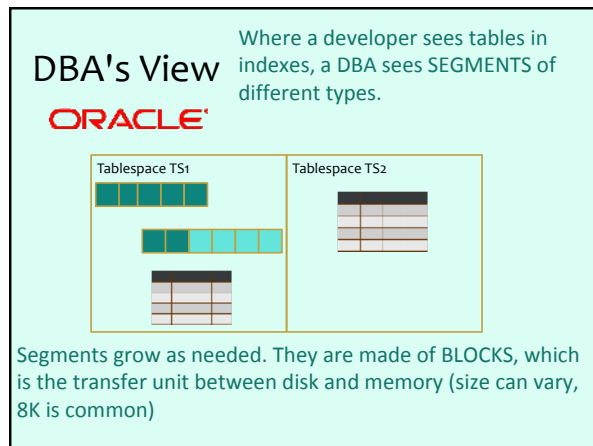


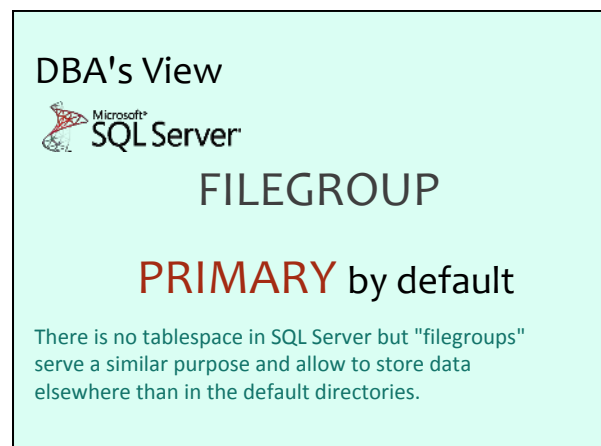
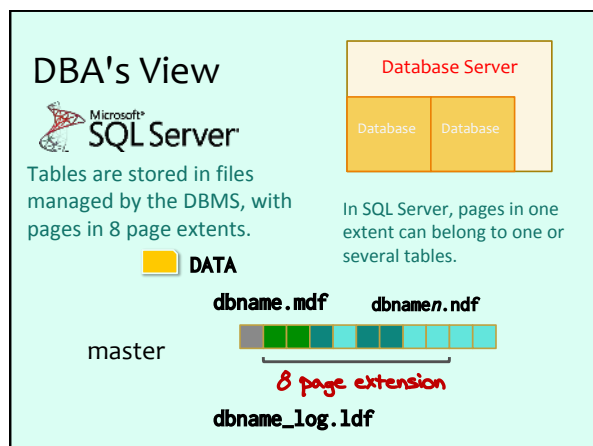
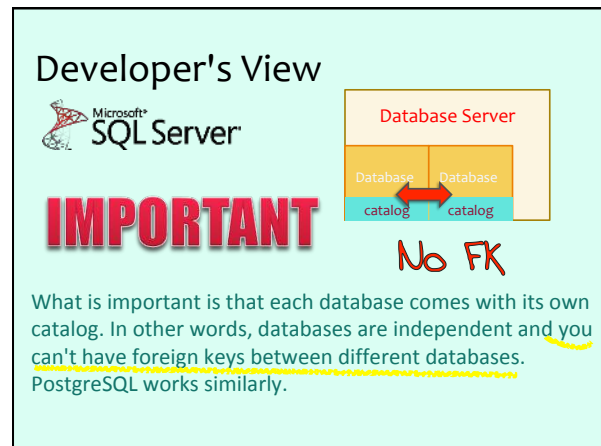
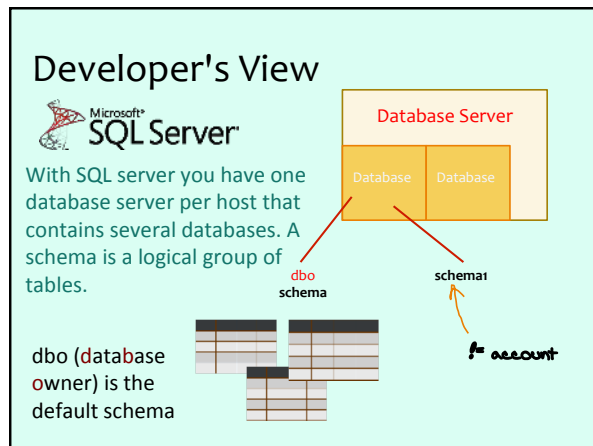
A developer who creates tables or indexes in Oracle create them in a tablespace, which can be specified at the end of the CREATE statement (there is a default tablespace for every account)

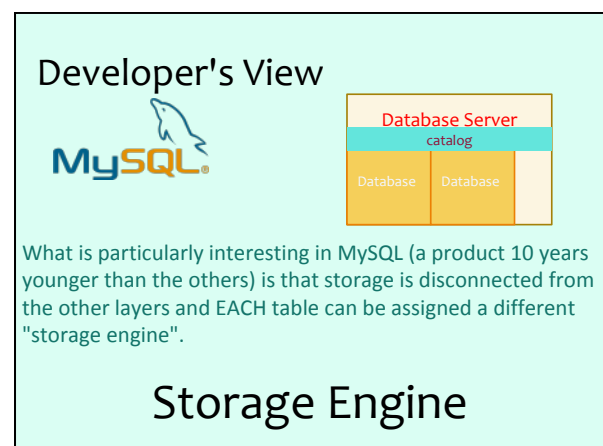
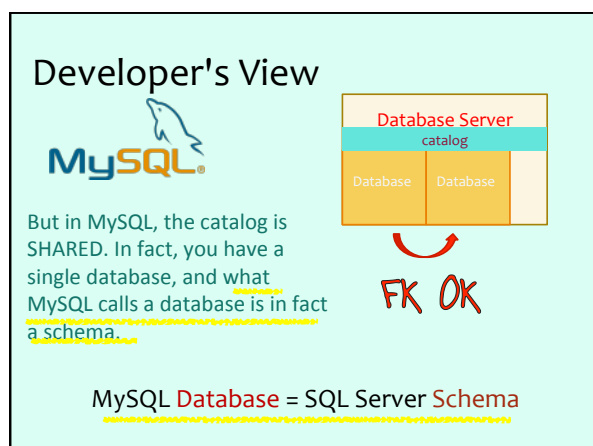
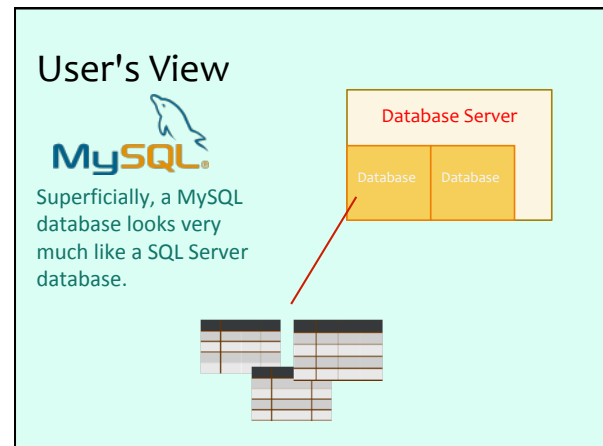
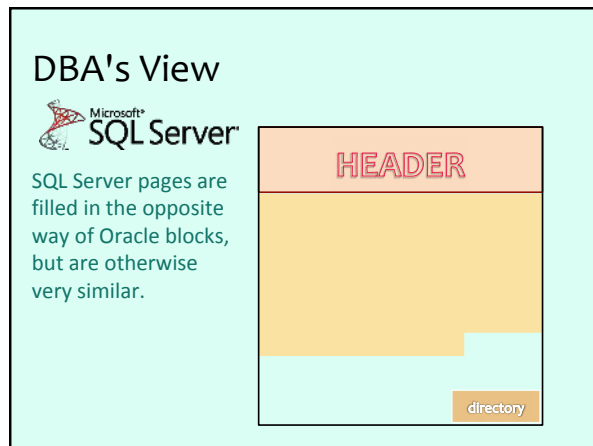
## DBA's View

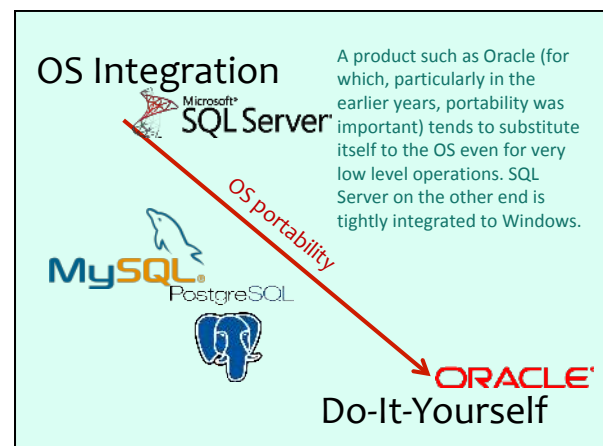
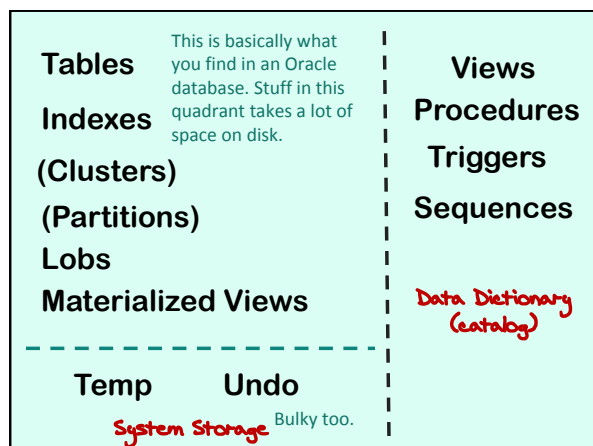
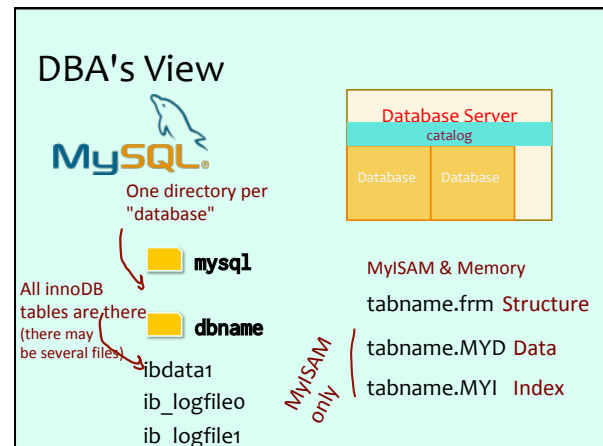
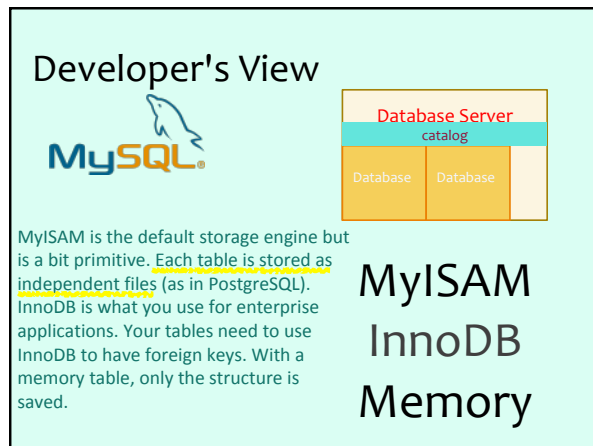
For a DBA, a tablespace is a set of files. It can also be a partition without a filesystem on a disk (raw device)









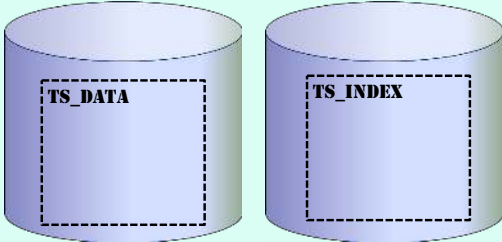


Historically, DBAs used to care about where files were going.

**Physical file placement  
USED to be important**

**ORACLE'**

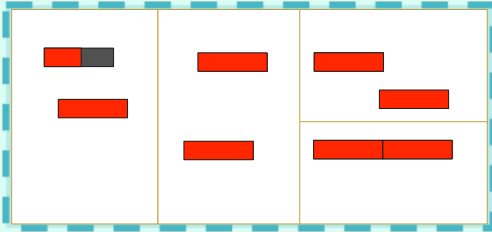
Most Oracle applications still separate data and indexes in different tablespaces.



Defragmentation used to be a concern. SQL Server DBAs periodically rebuild their indexes.

**Physical file placement  
USED to be important**

**ORACLE'**



# TOO **BIG** TO REORG

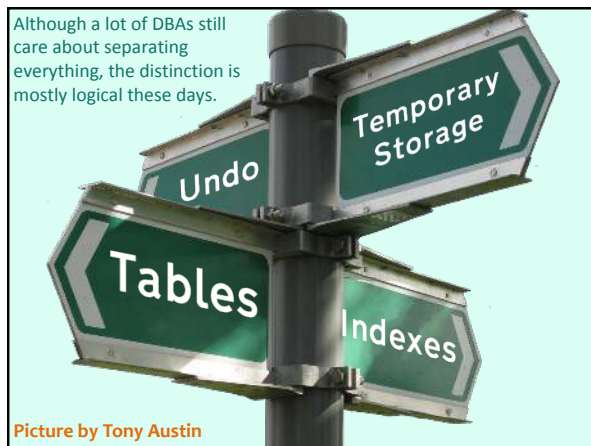
Came one day when bases became so big and availability requirements so strong that defragmenting a database simply was out of the question. Today focus is mostly on finding a way to store data that doesn't degrade too much over time. DBMS vendor also introduced tools for reorganizing an active database with not too much impact.

**Physical file placement  
USED to be important**

**ORACLE'**

Today  
**Automatic Storage Management**

Oracle has for instance introduced ASM, which is a volume manager and file system specially designed for Oracle files, which automatically takes care of a lot of thorny issues. Oracle DBAs no longer care much about physical files.



## Storage Area Network (SAN)

Seen as a disk

EMC<sup>2</sup>

Hitachi Data Systems

## Network-Attached Storage

Seen as a file-server



## Exadata ("Database Appliance")

ORACLE

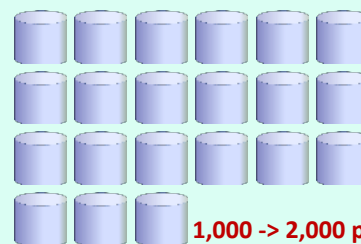
Especially as a lot of intelligent (and expensive) "storage devices" are usually plugged behind computers (and shared by them).

# Redundant Array of ~~Inexpensive~~ Disks

All started with RAID and the idea of putting a bunch of disks together, making them look like a single unit and reading/writing them in parallel to massively increase throughput. 'I' in RAID used to mean inexpensive (because each disk IS inexpensive), but as you have a LOT of inexpensive disks the total isn't cheap ... It was conveying the wrong idea and 'I' now stands for "independent".

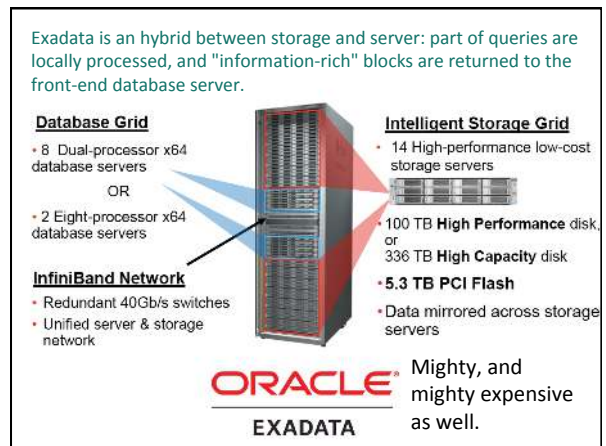
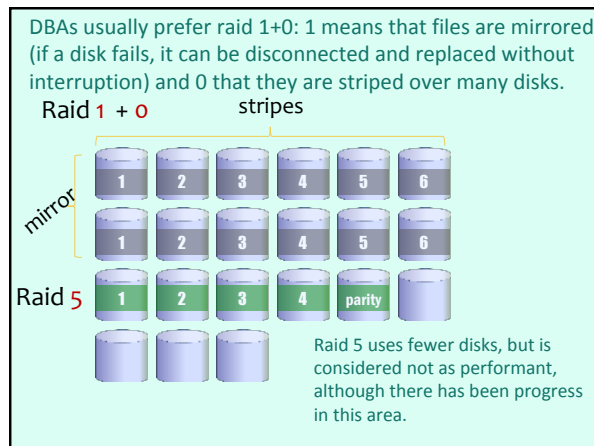
Cache

512GB-2TB



That's what a SAN looks like.

1,000 -> 2,000 physical disks



## FILE PLACEMENT

⇒ **Storage Specialist**

With a SAN at the back of a computer, a DBA these days no longer really knows where everything is. A storage specialist usually knows, and a DBA may have to work with this specialist for files that are heavily accessed, including log files that are written sequentially and should ideally not be on disks that undergo other heavy I/O activity.

Things that **STILL** matter  
To DBAs

Directory structures (scripts!)

Table organization

The logical structure though still matters (having database files everywhere is a bad idea, especially in /tmp) and care must be taken of inside table organization.





## Do we need this?

**A**

- aborting transactions, 130
- ABS function, 48
- ACOS function, 50
- ADD\_MONTHS function, 45
- aggregate functions, 56
- aliases
  - column, 110-114
  - UNION queries and, 133
  - qualifying join columns, 82
- table, 35
  - flashback queries and, 40
  - in FROM clauses, 116
- ALL keyword, 56, 114
  - comparison operators and, 97
  - vs. FIRST keyword, 76
- ANSI/ISO CAST function, 9
- MySQL, 29, 34
- SQL Server, 23
- ANSI/ISO EXTRACT function, 10, 29
- ANSI/ISO WITH clause
  - correlated subqueries, factoring out, 123-125
  - noncorrelated subqueries, factoring out, 122
  - recursive, 66
- ANSI\_NULLS setting, 95
- ANY keyword and comparison operators, 97
- APPEND hint and direct-path inserts, 74
- AS OF keyword, 40
- ASCENDING keyword, 120
- ASIN function, 50
- asterisk ( \* ) regular-expression metacharacter, 104, 106
- COUNT function and, 57
  - qualifying with table names, 108
  - returning all columns from a table, 107
- ATAN function, 50
- ATAN2 function, 50
- ATANH function, 50
- ATN2 function (SQL Server), 50
- autocommit mode, 35, 125
- AVG function, 57

You are going to turn to a book index only when you are looking for very specific information.

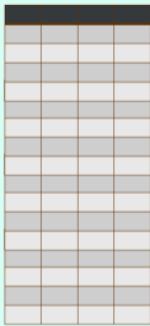
## Or do we need this?

<b>Functions</b>	<b>120</b>
Date Functions	120
Numeric and Math Functions	122
Trigonometric Functions	124
String Functions	125
Miscellaneous Functions (Oracle)	130
<b>Grouping and Summarizing</b>	<b>132</b>
Aggregate Functions	132
GROUP BY	133
Useful GROUP BY Techniques	135
HAVING	137
GROUP BY Extensions (Oracle)	139
GROUP BY Extensions (SQL Server)	141
<b>Hierarchical Queries</b>	<b>143</b>
ANSI/ISO Recursive WITH (DB2)	143
CONNECT BY Syntax (Oracle)	146

When there are many entries in the index for what you are looking for, you rather turn to the table of contents.

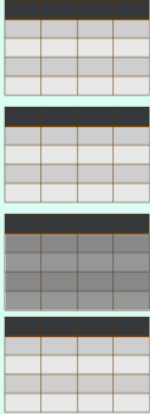
## Partitioning

The idea of partitioning is very similar. Suppose that a table holds data for say over one year, and that we are interested in data for one month (say 10% of the table to simplify). If the table contains 500,000,000 rows, 10% are still 50,000,000 rows and fetching them one by one with an index will take ages. Scanning the full table may be more efficient, but perhaps that instead of discarding 90% of what we read, we could regroup data by month and only read what we want.



That's the idea behind partitioning. You see one table but physically data is regrouped in semi-independent tables called partitions. The DBMS know the organization and direct you to the right table and ignore the others.

## Partition key



When you create the table, you also create partitions (they can also be added and removed on the fly) and specify a set of columns (partition key) that determines in which partition a row should be inserted. When the partition key is a criterion in a query, only the right partitions are considered.

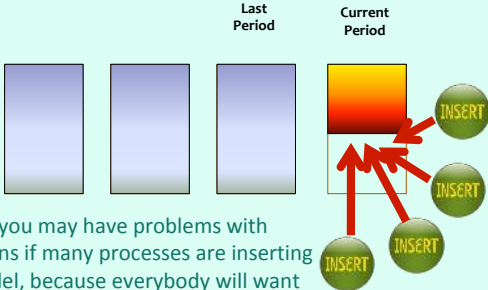
There are different ways to partition. When your concern is to control data grouping, you usually partition by range (so as to have rows grouped by month, week, or whatever interval), or by list (If the row contains this value in the partition key, then it goes into this partition). When your concern is spreading inserts over a table, for instance, you may opt for hash partitioning, and the system will compute partition placement for you.

**By range** (*dates usually*)

**By list**

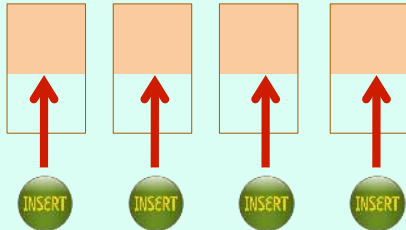
**By hashing**

Partitioning may want to address very different problems. If you group your rows by period (good for querying) ....

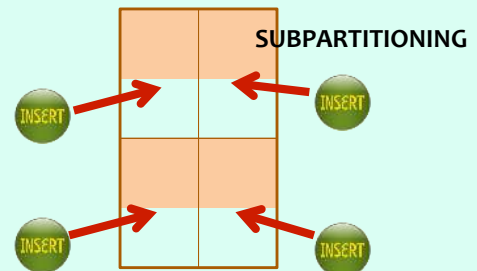


... then you may have problems with insertions if many processes are inserting in parallel, because everybody will want to insert at the same place and you'll need some serialization (two processes cannot write bytes in the same block at the same time).

For parallel insertions, you'd rather see each process inserting into separate partitions. But then if you want to retrieve all these rows at once in a query, you must



query n partitions instead of just one, which isn't ideal.



You can sometimes have your cake and eat it by using subpartitioning, for instance subpartitioning by hash partitions defined by range.