

Distributed – What to expect

The next trend (hype?) seems to be : Clustered, serverless, distributed, sharded, replicated, databases. That is a lot of buzzwords together, and I have omitted “multi-cloud”, “raft” and “vector”.

We’ll try to find out how various vendors interpret those (buzz)words and how/when these are useful. We’ll explore what users, administrators and application-programmers can expect when data is stored in one of those “new and shiny” databases.

And then, when your boss tells you to “Use This”, maybe you know what to watch out for. In the end your aim is probably to keep your data Safe and Available.

Distributed Databases – What to Expect

Piet de Visser
The Simple (Oracle) DBA



What are they and how do they behave ?
Should you use those databases ?



Piet de Visser - PDVBV



For some, I am like a Dinosaur
And I spend a lot of time on that motorcycle. – need picture with Frecciarossa...

Agenda

(approx 45 minutes)



Agenda. I'll try to cover a few topics and talk about 2 of the products.
List isn't complete, but I looked at those two (bcse I had time)

Distributed – the Buzzwords + "why".

1/2

Cluster / Clustered : Using multiple machines, nodes – e.g. 42...
(shared Nothing!)

Distributed: multiple machines, nodes.
=> "spreading out the data"

Serverless: Not using any specific servers.
e.g. docker, k8s, lambdas, services

Replicated: Keeping copies of data.
Mostly RF = 3 or RF = 5. Overhead !

Quorum: Majority of votes (e.g. 2 out of 3)



Distributed – the Buzzwords + "why".

2/2

Sharded:	Divide data over Shards (range, hash, list...)
Sharding-Key: (=! PK)	Data, column(s) to use for Sharding. Sharding-Key determines "Where data Goes"
Leader / leaseholder : Followers, replicas:	The process in charge of an obj or data-item. The (keepers of the) copies...
Raft :	Algorithm, to elect and re-elect a leader. Notably to be Fast and avoid race/lock conditions frequent heartbeats...



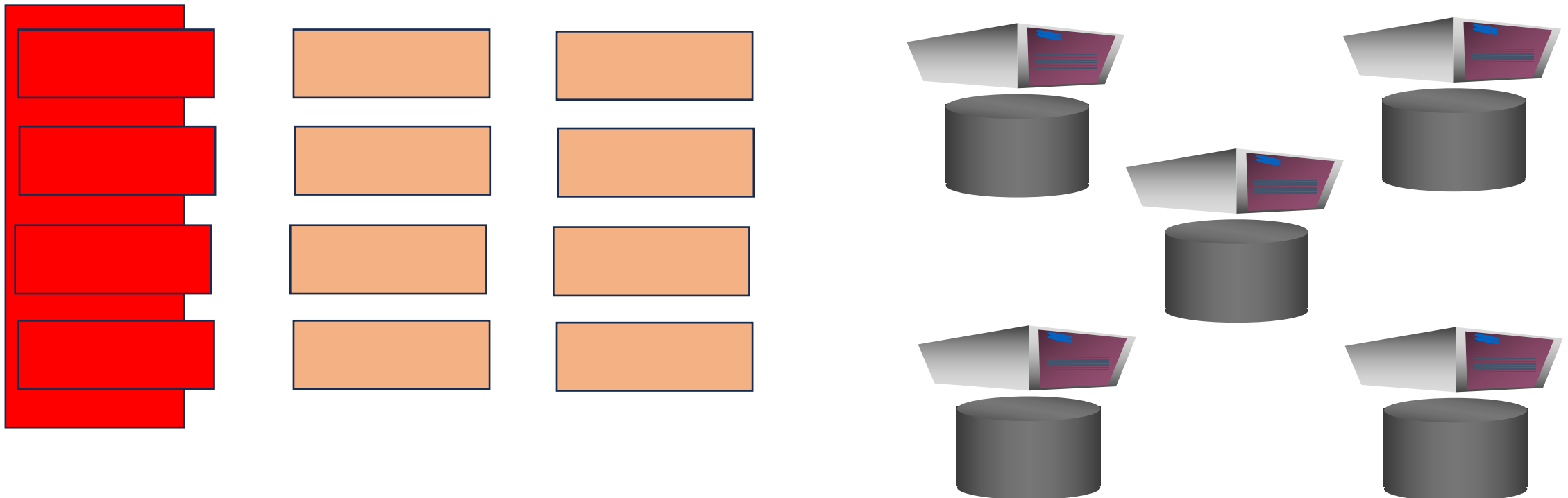
[Remove Notes...]

- Insert drawing of multiple servers:
- Connected, but not shared.
- Data is sharded over nodes
- Data is replicated for resilience..
- (ELES ... need Quorum... $RF=3 \Rightarrow N-1$)
- Nature of shared nothing clusters...
- Big Q: How much can we loose... (hint: not much)

Distributed = Sharded + Replicated...

1/2

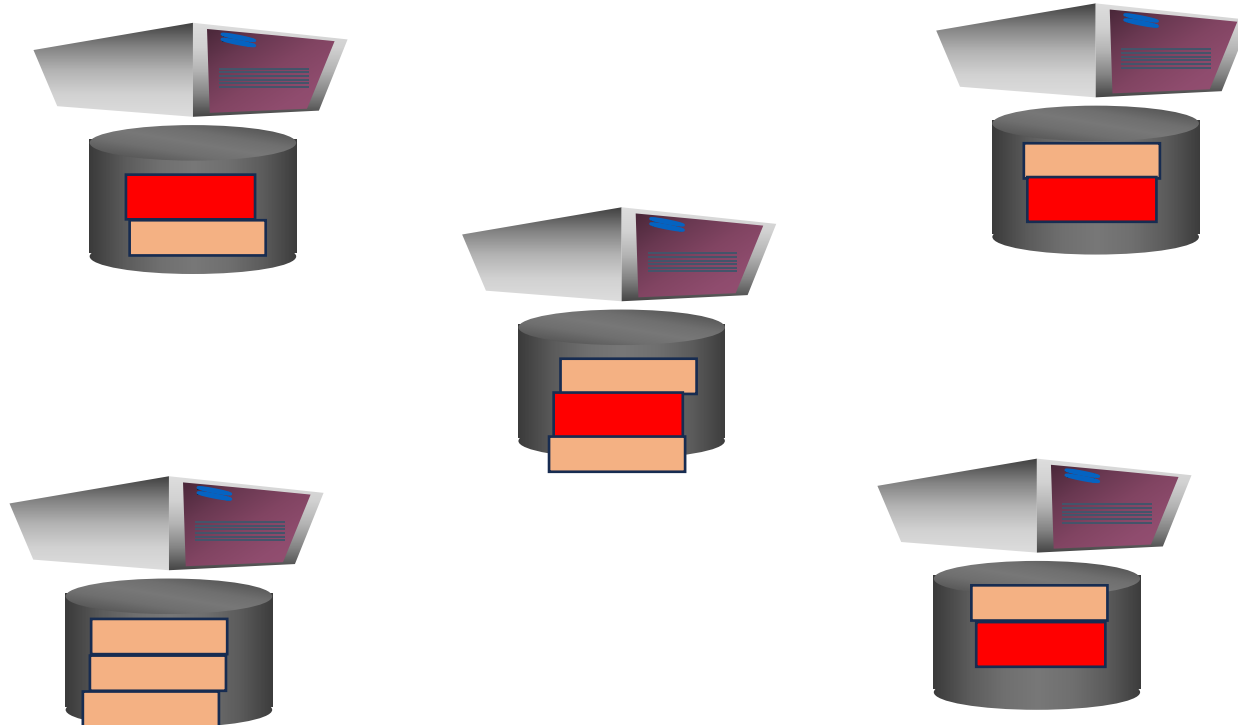
- Sharding : Split table into... Shards, then spread over “servers”
- Replicate: make copies (of shards), RF = 3..., then spread...
- Example: 4 Shards x 3 Replicas = 12 items (tablets, partitions, files)



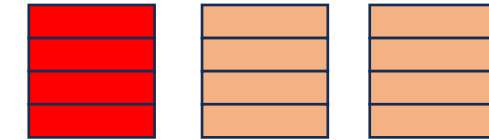
Sharding = partitioning, Replication = copies..
And then imagine spreading that over a park of servers or containers

Distributed – Sharding (spread load)

- Multiple Server (VMs, Containers)
- Sharded Data over ... many.
- And RF = 3 ..



1/2



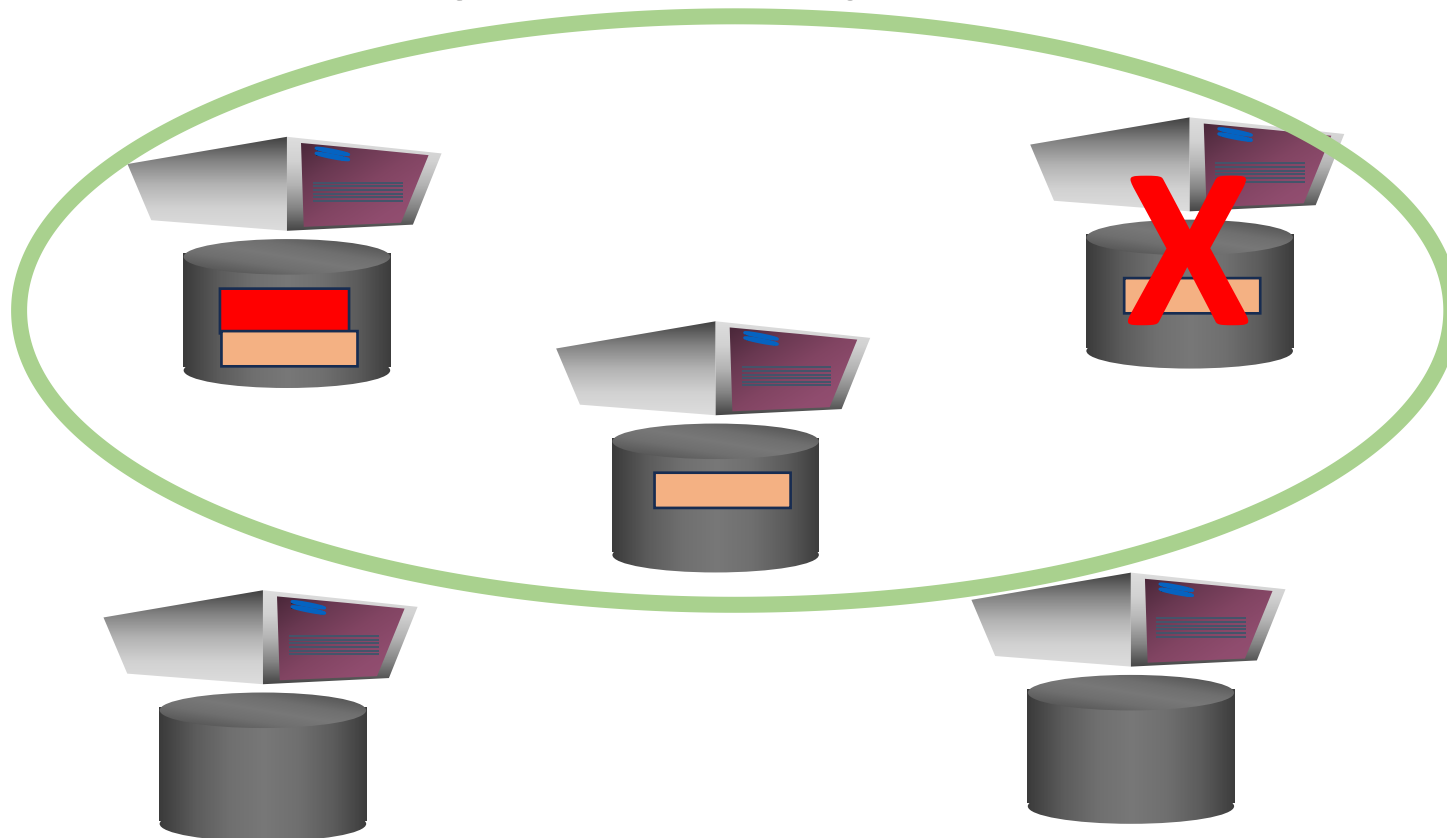
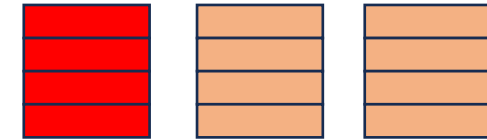
- 4 shards...
- Plus 2 copies of each
- Leaders + Copies everywhere..
- Can afford to lose 1 (just 1)
- Any loss = Work (copies!)
- What if you loose 2 ?
- Should have had ..
- “Many” servers, still @ Risk.

With an RF=3, even with many, 42, servers, any loss of “Two” means some component loses Quorum
And any “loss of server incurs the “copy-jobs” – an un-stable cluster will incur WORK – (what about GEO...)

Distributed – Replication + Raft

2/2

- Multiple Servers, Many.. (VMs, Containers)
- Replicate Data over ...
- RF = 3, every item has 3 copies on diff servers...



- Lose 1 server
 - Detect loss
 - [RAFT: Elect Leader if needed]
 - **Still Quorum (2/3)? OK!**
 - Copy to re-store RF=3
-
- What if you loose 2 ?
 - Should have had .. RF = 5 !
 - What if .. GEO-constraint ?

Distributed:	Every Node holds <u>Some</u> data...
Shared-Nothing:	You “need all nodes”
Storage is Not Shared:	Hence, every node “leads” <u>some</u> data.
a weakness...?	add/remove nodes = Re-Distribution.

IMHO: Data is not just a “Stateless app”.

You can scale an “app” by adding/removing Computing power.

You can not “quickly” scale(-down) a database by removing >1 nodes...

Replication – to always be safe + available

1/3

RF=3 (or 5):

**Three (or 5) copies of every “item”
every shard, every partition, 3 (or 5) copies...**

Raft:

**On “lost leaders”: need to Re-Elect (millisec)
A “Shard” needs a “Master” somewhere.**

Lost one node?

**No problem, 2=quorum, All still works
timeout: objects are re-replicated to RF=3
Cluster is thus “self healing”**

Lost 2 nodes ?

**Should have had RF= ...
5 or more, to keep quorum.**



RF=3/5:

**Hardware: You need 3 (or 5)++ “Racks”.
There is a lot of Communication between...**

Storage:

**At Least 3x (or 5x)
Most systems will do some compress...**

Processing:

**Ditto. Notably Writing (commits).
More CPU ticks, more time.
(even parallel-work needs “some coord.”)**

GEO-Location:

**Need “Quorum” in Every Location => servers++!
+ Latency will hit you... (more comm...)**

Does it work?

Hell Yes!

Just dont loose 2 or more nodes (15-20min...)

Just keep distance/latency Low.

Bonus Question:

Losing 2 nodes at once, out of 42 ?

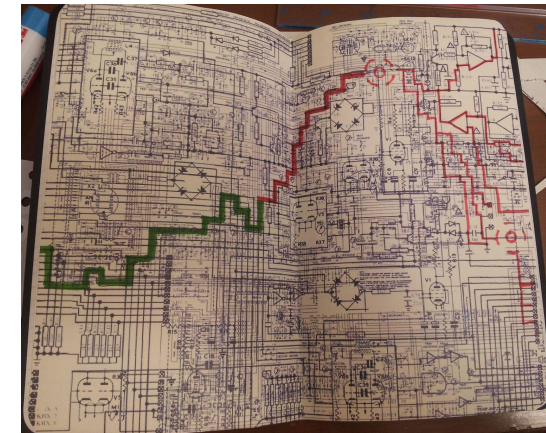
**Only “some shards” are affected.
Only those that lost “quorum”
(but generally Not Acceptable!)**

Each vendor has (more) tricks...



Storage: LSM and SST

LSM:	Log Structured Merge (memory)
SST:	Sorted String Table (file)
What/Why:	Efficient use of Memory and Disk Write-Once - Append Only. Compressable (compaction) = Nice!
Implications... (?)	Disk usage growing ... and ...Shrinking! (for now: assume >3x space, RF + ...) Need to investigate, and YMMV!



Thank You !

(stay tuned)



Serverless

Not strikltly (you can deploy on a server)

Dev:

Developers prefer containers... (duh)

Arch:

Some shops “Require” it (k8s only...)

Scale-Up:

Re-deploy (restart) on bigger box. Works fine.

Scale-Out:

Deploy to more containers...

Not immediate: need re-distribution of data.

Scale-Down ?

**Less nodes ?
(not ideal.....)**



Scale-up/down

**Not Stateless ... (a DB is Not a Micro-service)
Re-Distribution = WORK (not ideal)**

Scale-Back

Better: re-deploy (restart) on smaller “nodes”

Note:

We once thought VMs unsuitable for RDBMS...

IMHO:

Abstraction, Layers ? Complexity ?

**Data should have least possible layers of processing,
Try to deploy “On the Server”, or “on Iron” (says the dinosaur).
Many shops/clouds: Containerization is a Must... ?**



Geo-Sharded – Legal or Latency...

Geo-Located: Each data item (record) has a “home”

Legal ? Data “must reside” at location (?)
But still see it as One Single Database.
Example: Medical records, Tax-records.

Latency ! Data close to user
But still see it as One Single Database.
Example: Superbowl-tickets, YB showcase.

Quorum ? RF=3: nr of servers++ in Each Location.

IMHO: “Legal requirement” will never be “Waterproof”



Errors made... (so you don't have to...)

1/2

Sharding: Too Many shards... (20 regioni, 107 province)
Check the defaults (42 servers... ???)

Sharding Too Few shards => re-sharding starts
Big-Data? - Prevent re-sharding on bulk-load.

Sharding: Hash-sharding on Date...
Hash-sharding on name-lookup

**IMHO: Sharding is a Good Mechanism,
but You Need to Know “What happens under the hood”**



Errors made... (so you don't have to...)

2/2

Sequence:

Cache !

Nodes on/off

Outages, No Problem...

Short (<15min): node can/will come back.

Long (>15min): Data will redistrib !

Storage:

No clear Errors... Disk Full !

CPU/Mem:

No Metrics ! (go see Franck)

**IMHO: There is a lot to Learn still,
notably in Metrics “under the hood”**



Main Messages – What we Learned...

This will happen (dino: who needs it?... not relevant..)

Observe, Learn, Think, Measure.

Sharded: range or hash.. (date ranges...)

Replicated: how + where (Sequences!)

Impact of LSM + SST: need to learn...

Go out and Experiment !





Long term: DB-aaS.. And “serverless” ...
Beware of the learning-curve ahead... (dinos are still here, and picture of stork, for luck)

Don't Take my word for it...

Try for Yourself... and tell me next conference.

Simplicity

- **In case of doubt: Simplify!**
- **Less components**
- **Less complexity**
- **Less tricks...**

Goethe _____ (simplicity)



Quick Q & A (3 min) 3 .. 2 .. 1 .. Zero

- Questions ?
- Reactions ?
- Experiences from the audience ?
- @pdevisser (twitter..)



This is POUG... Let's Demo !

Demo_dd.sql : **overhead of (many)shards...**

Mk_longt.sql : **insert 200M, and check tablesize...**
 - use lots of space
 - compression can help, automatically

Do_fill.sh : **run process and kill nodes...**

Chk sizes again to see compression.

Blank

End - This slide intentionally left Blank...

Distributed – Where does your data go?

1/2

- Multiple Server (VMs, Containers)
- Sharding : split table into... Shards.
- Distribute Shards over ... servers/containers

