

# Exadata for Oracle DBAs

**Arup Nanda**  
*Longtime Oracle DBA*

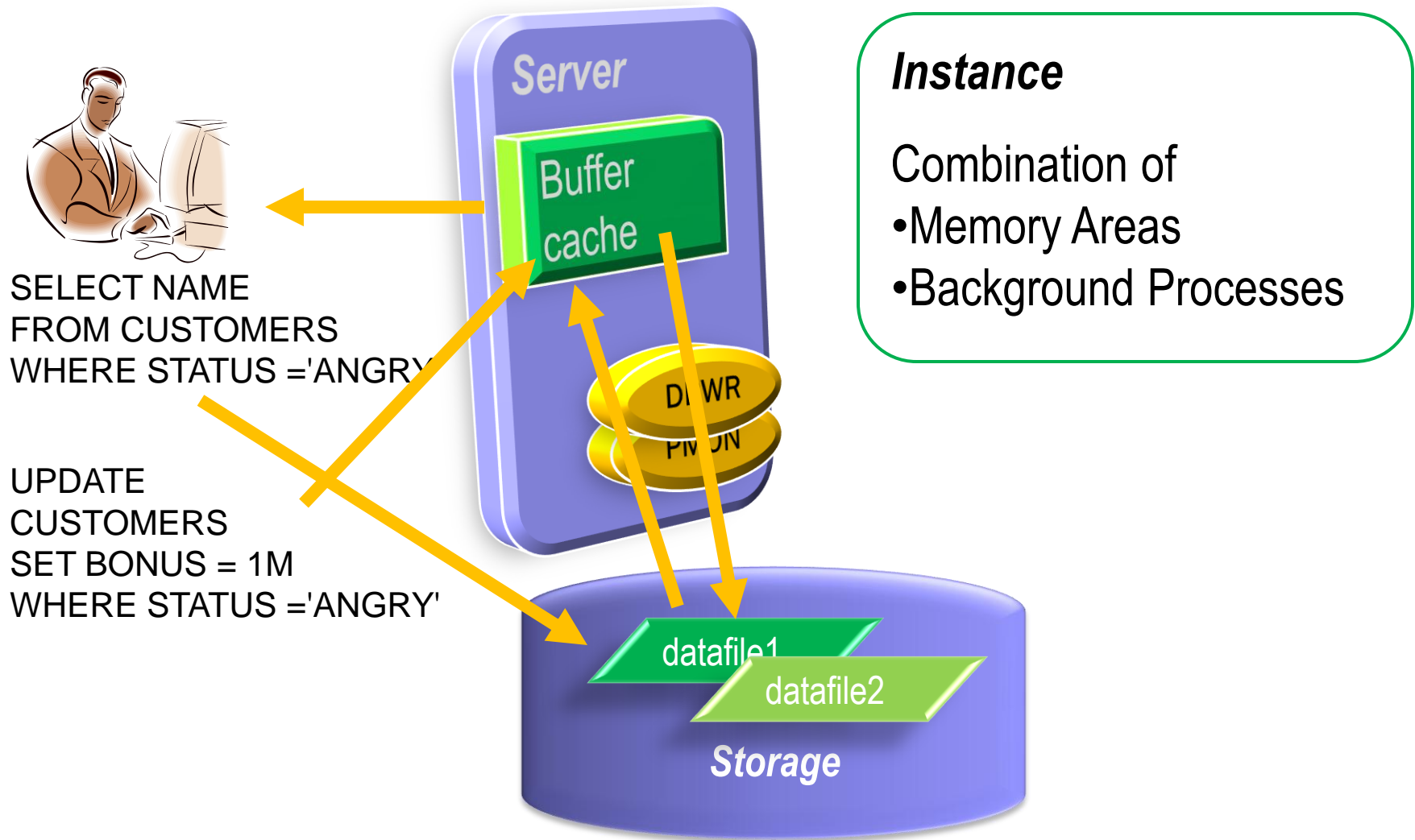
# Why this Session?

- I'm an Oracle DBA
  - Familiar with RAC, 11gR2 and ASM
- About to become a Database Machine Administrator (DMA)
- **How much do I have to learn?**
- How much of my own prior knowledge I can apply?
- What's different in Exadata?
- What makes it special, fast, efficient?
- Do I have to go through a lot of training?

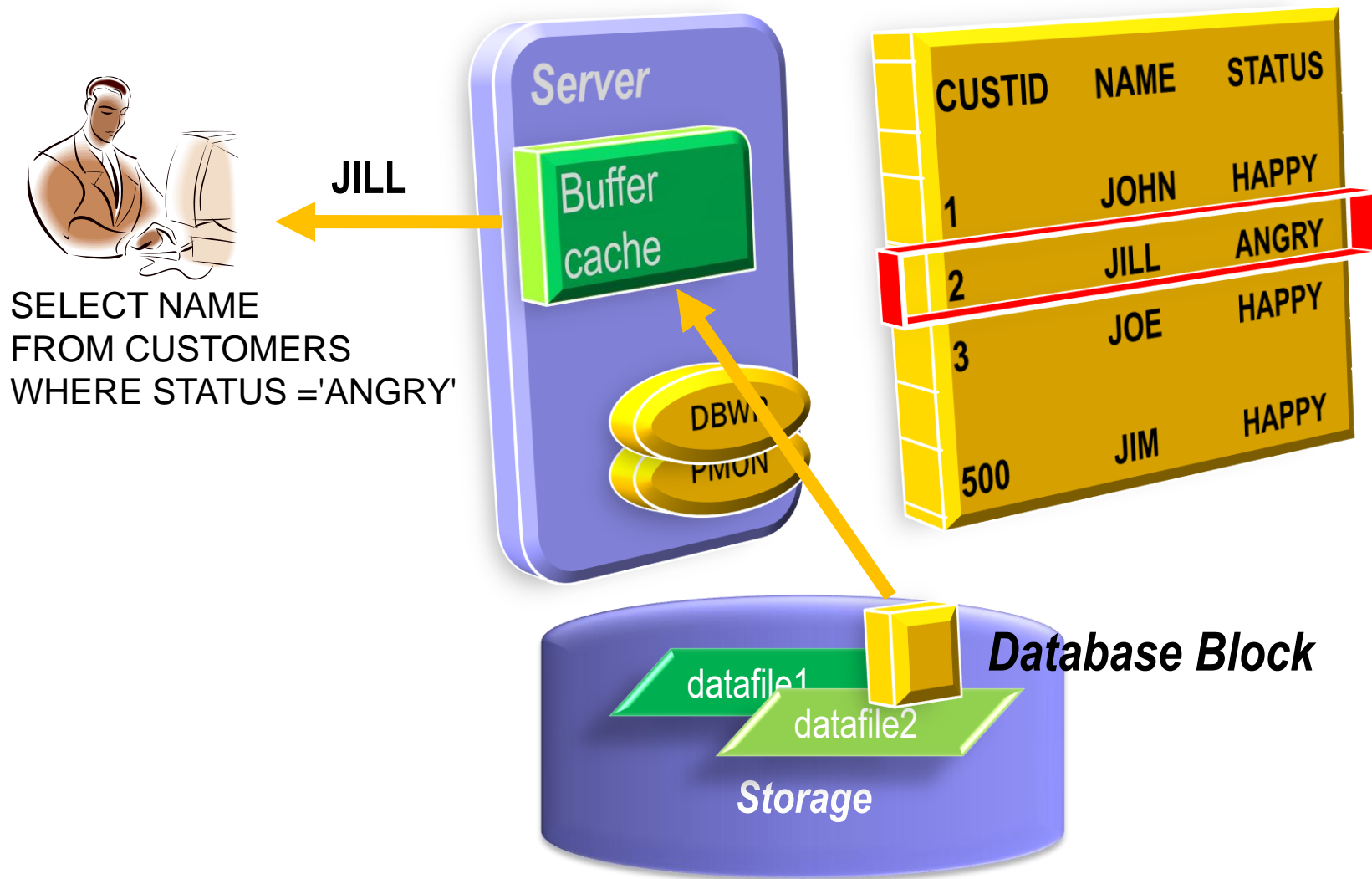
# What is Exadata

- Is an *appliance* containing
  - Storage
  - Flash Disks
  - Database Servers
  - Infiniband Switches
  - Ethernet Switches
  - KVM (some models)
- But is *not* an appliance. Why?
  - It contains additional software to make it a better database machine
- That's why Oracle calls it a **Database Machine** (DBM)
- And **DMA** – Database Machine Administrator

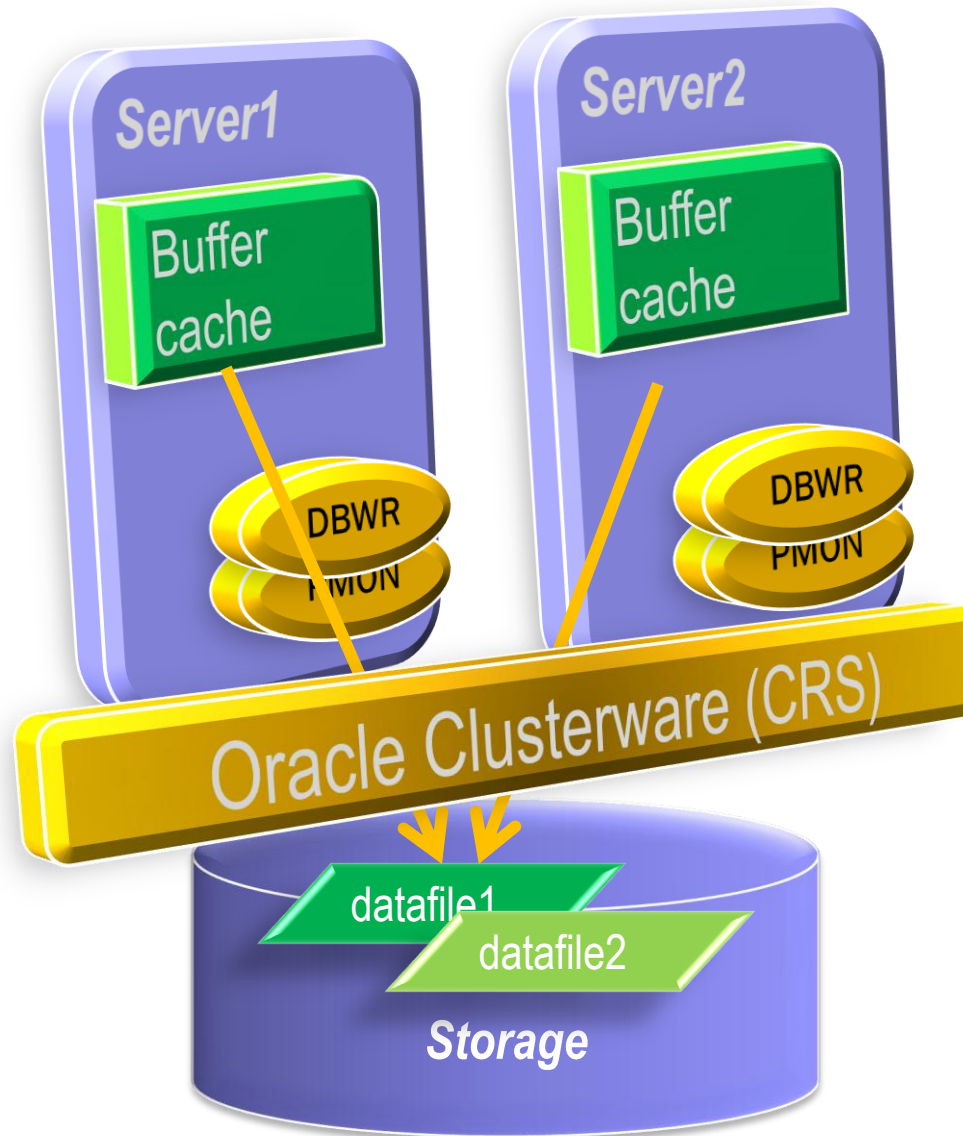
# Anatomy of an Oracle Database



# Query Processing



# RAC Database



# Components for Performance

CPU

Memory

Network

I/O Controller

Disk

Storage has been and will continue to be the bane of all databases. Simply put, less I/O → better performance

# What about SAN Caches?

- Success of SAN caches is built upon predictive analytics
- They work well, if a small percentage of *disk* is accessed most often
  - The emphasis is on *disk*; not *data*
- Most database systems
  - are way bigger than caches
  - need to get the data to the memory to process
    - ➔ I/O at the disk level is still high
- Caches are excellent for filesystems
  - ➔ or very small databases



# What about In-Memory DBs

- Memory is still more expensive
- How much memory is enough?
- You have a 100 MB database and 100 MB buffer cache
- The whole database will fit in the memory, right?
- NO!
- Oracle database fills up to 7x DB size buffer cache  
<http://arup.blogspot.com/2011/04/can-i-fit-80mb-database-completely-in.html>

# The Solution

- A typical query may:
  - Select 10% of the entire storage
  - Use only 1% of the data it gets
- To gain performance, the DB needs to shed weight
- It has to get less from the storage
  - Filtering at the storage level
  - The storage must be cognizant of the data

```
SELECT NAME  
FROM CUSTOMERS  
WHERE STATUS = 'ANGRY'
```

**Filtering  
should be  
Applied Here**

CPU

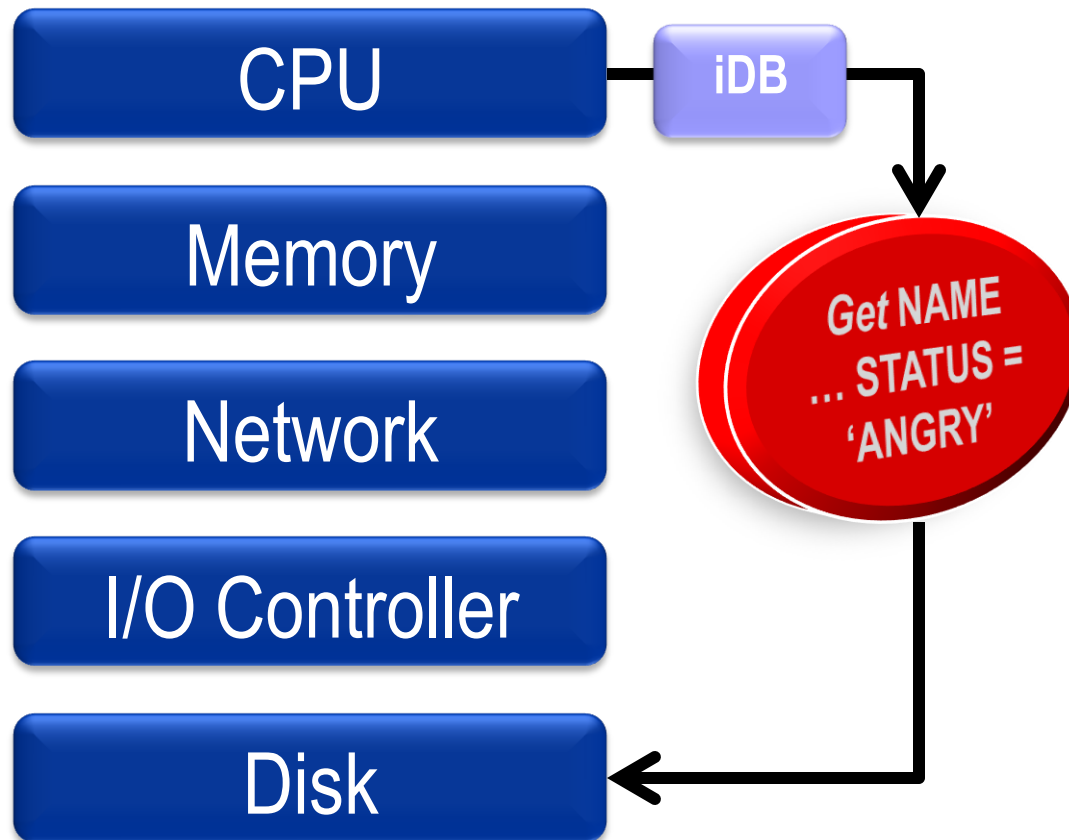
Memory

Network

I/O Controller

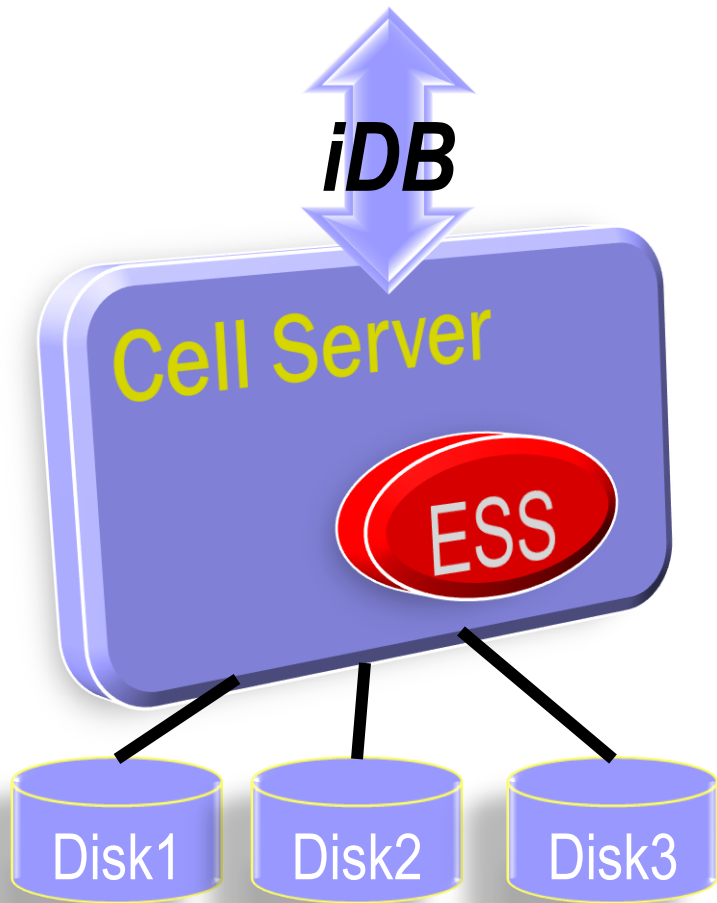
Disk

# The Magic #1



The communication between CPU and Disk carries the information on the query – columns and predicates. This occurs as a result of a special protocol called iDB.

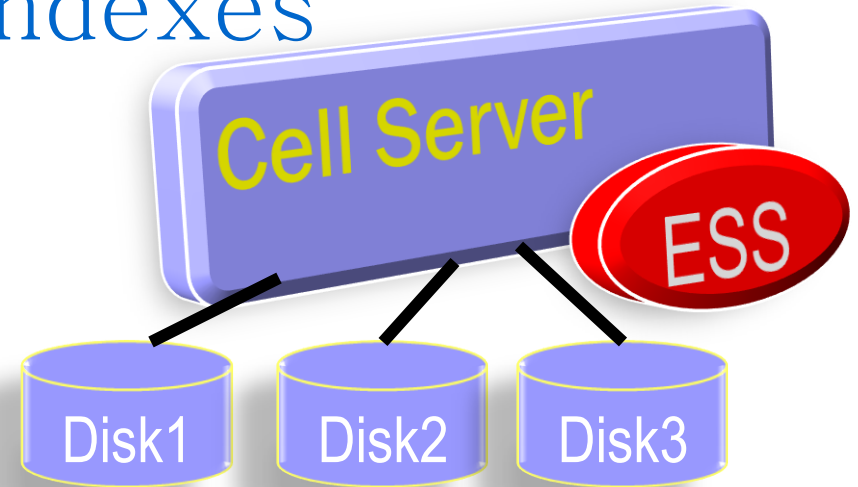
# Magic #2 Storage Cell Server



- Cells are Sun Blades
- Run Oracle Enterprise Linux
- Software called Exadata Storage Server (ESS) which understands iDB

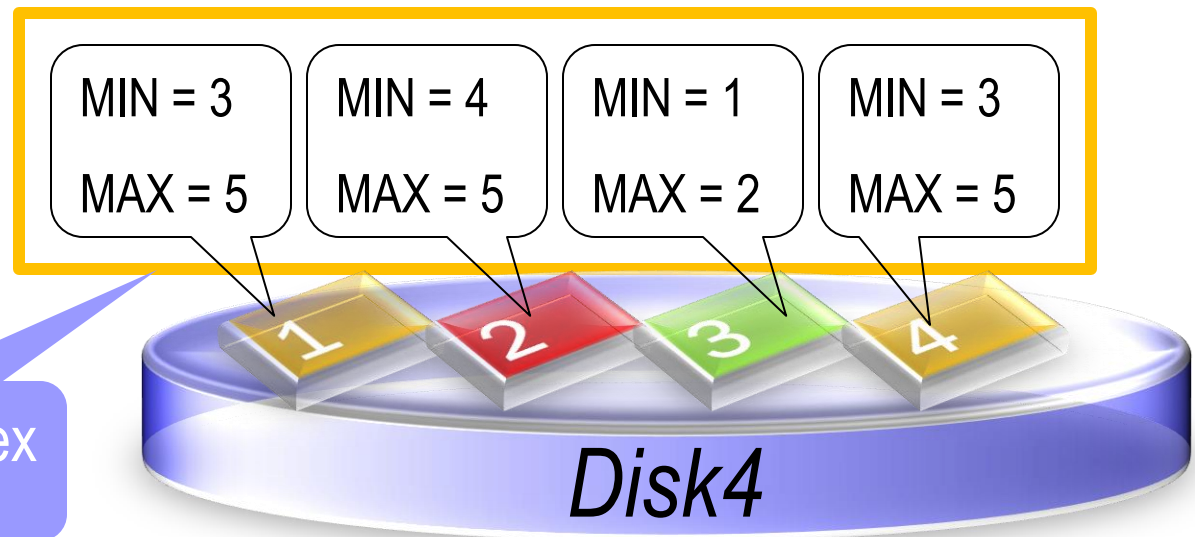
# Magic #3 Storage Indexes

Storage Indexes store in memory of the Cell Server the areas on the disk and the MIN/MAX value of the column and whether NULL exists. They eliminate disk I/O.

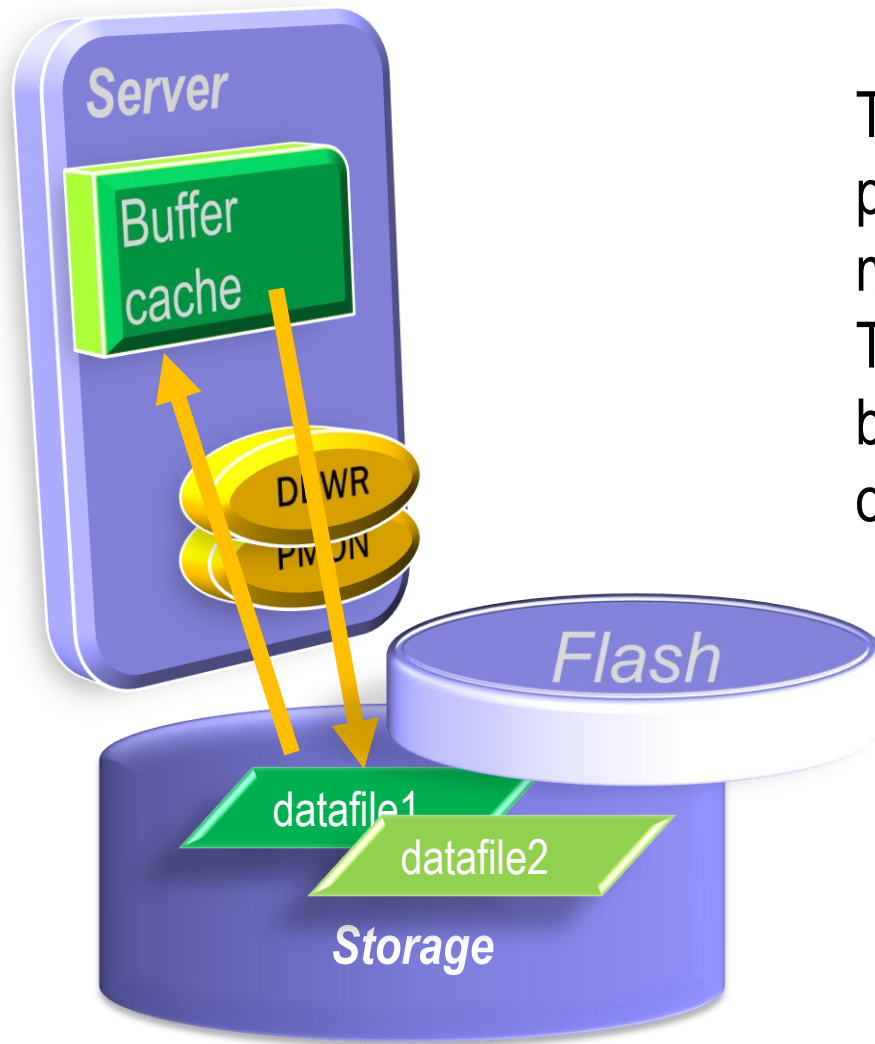


**SELECT ...  
FROM TABLE  
WHERE COL1 = 1**

Storage Index



# Magic #4 Flash Cache



These are flash cards presented as disks; not memory to the Storage Cells. They are similar to SAN cache; but Oracle controls what goes on there and how long it stays.

# Components

CPU

Memory

Network

I/O Controller

Disk

***Database Node***

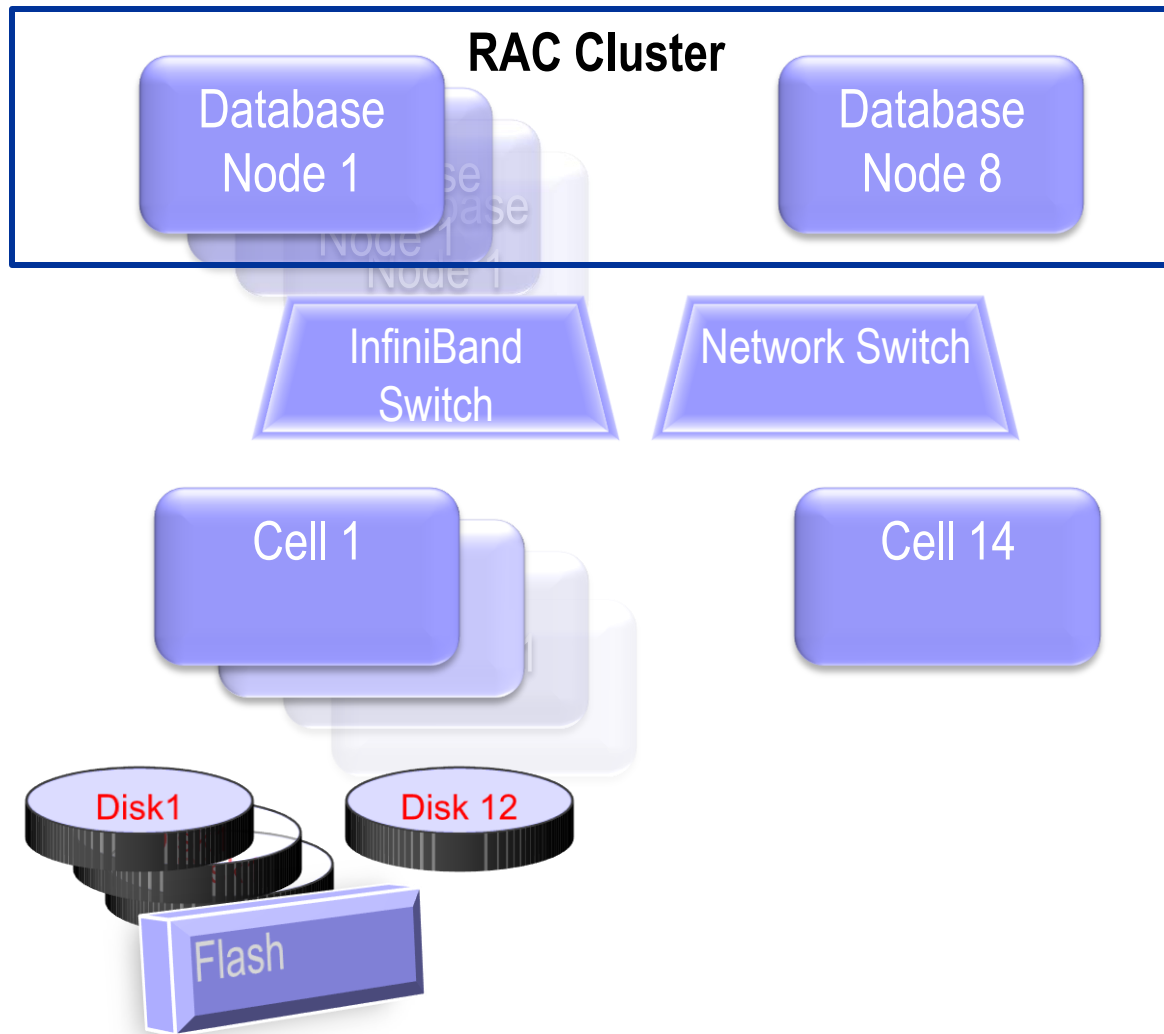
(Sun Blade. OEL)  
Oracle 11gR2 RAC

***InfiniBand Switch***

***Storage Cell***

Exadata Storage Server  
Disks, Flash

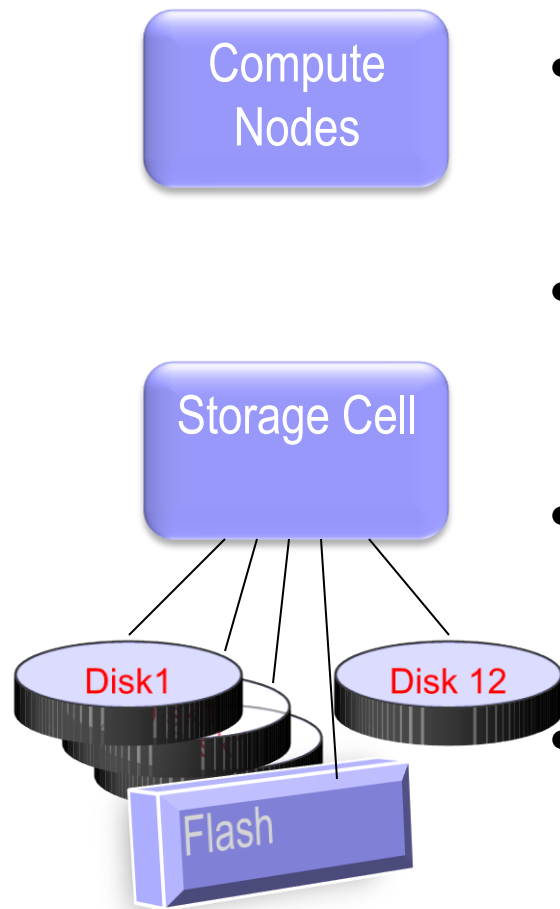
# Put Together: One Full Rack



**Clients  
connect to  
the database  
nodes.**

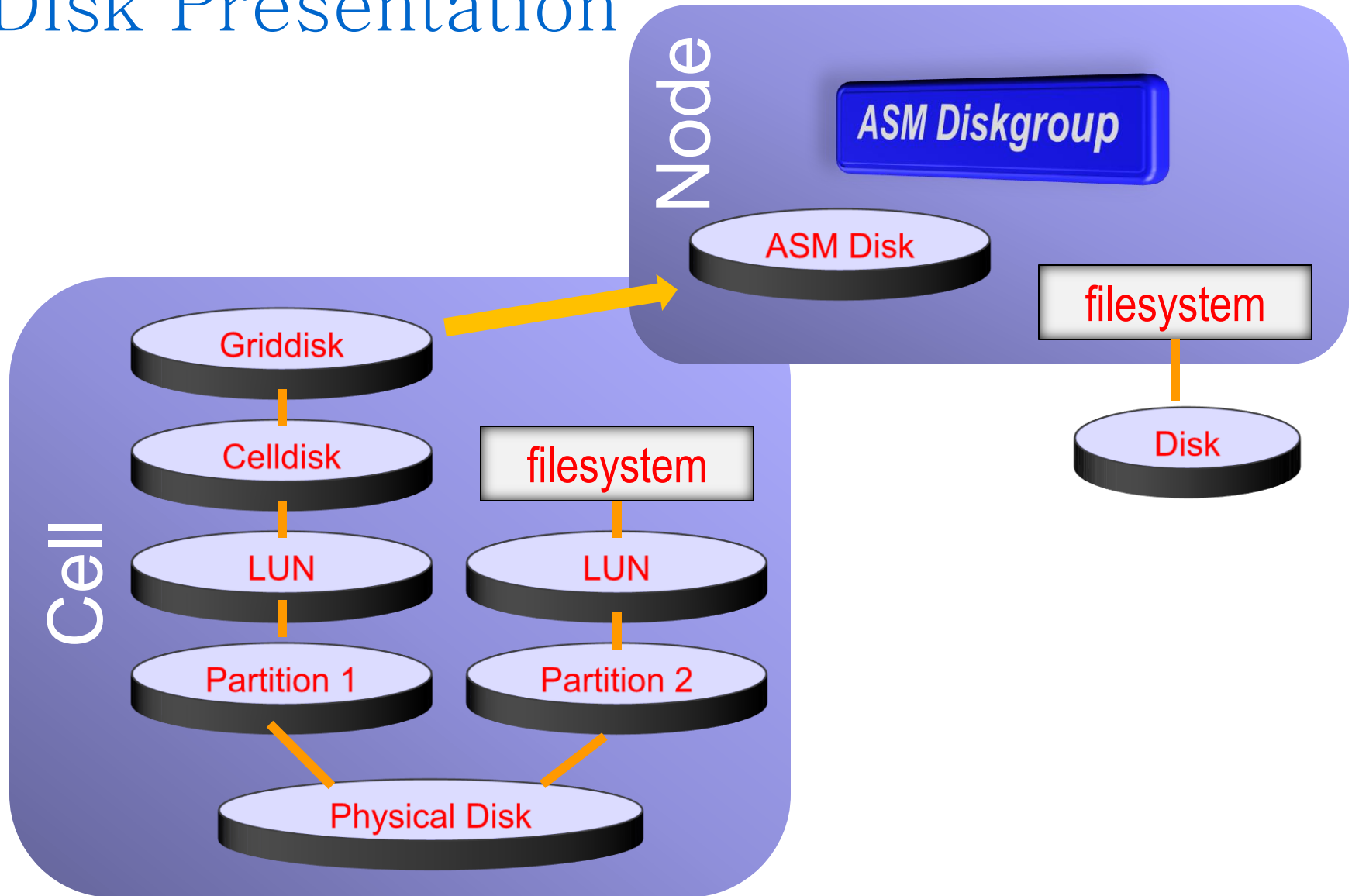


# Disk Layout



- Disks (hard and flash) are connected to the cells.
- The disks are partitioned at the cell
- Some partitions are presented as filesystems
- The rest are used for ASM diskgroups
- All these disks/partitions are presented to the compute nodes

# Disk Presentation



# Command Components

**Linux Commands – vmstat, mpstat, fdisk, etc.**  
**ASM Commands – SQL\*Plus, ASMCMD, ASMCA**  
**Database Commands – startup, alter database, etc.**  
**Clusterware Commands – CRSCTL, SRVCTL, etc.**

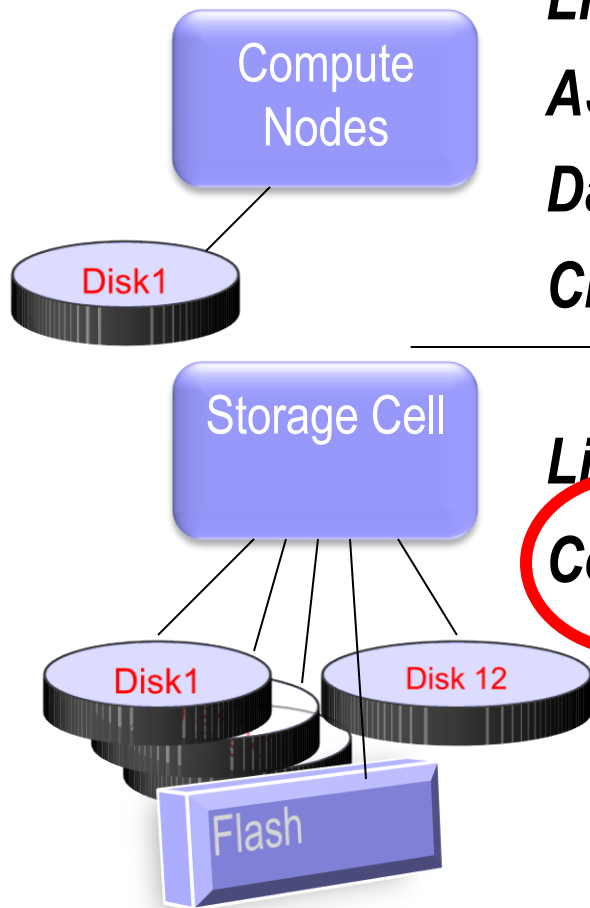
---

**Linux Commands – vmstat, mpstat, fdisk, etc.**  
**CellCLI – command line tool to manage the Cell**

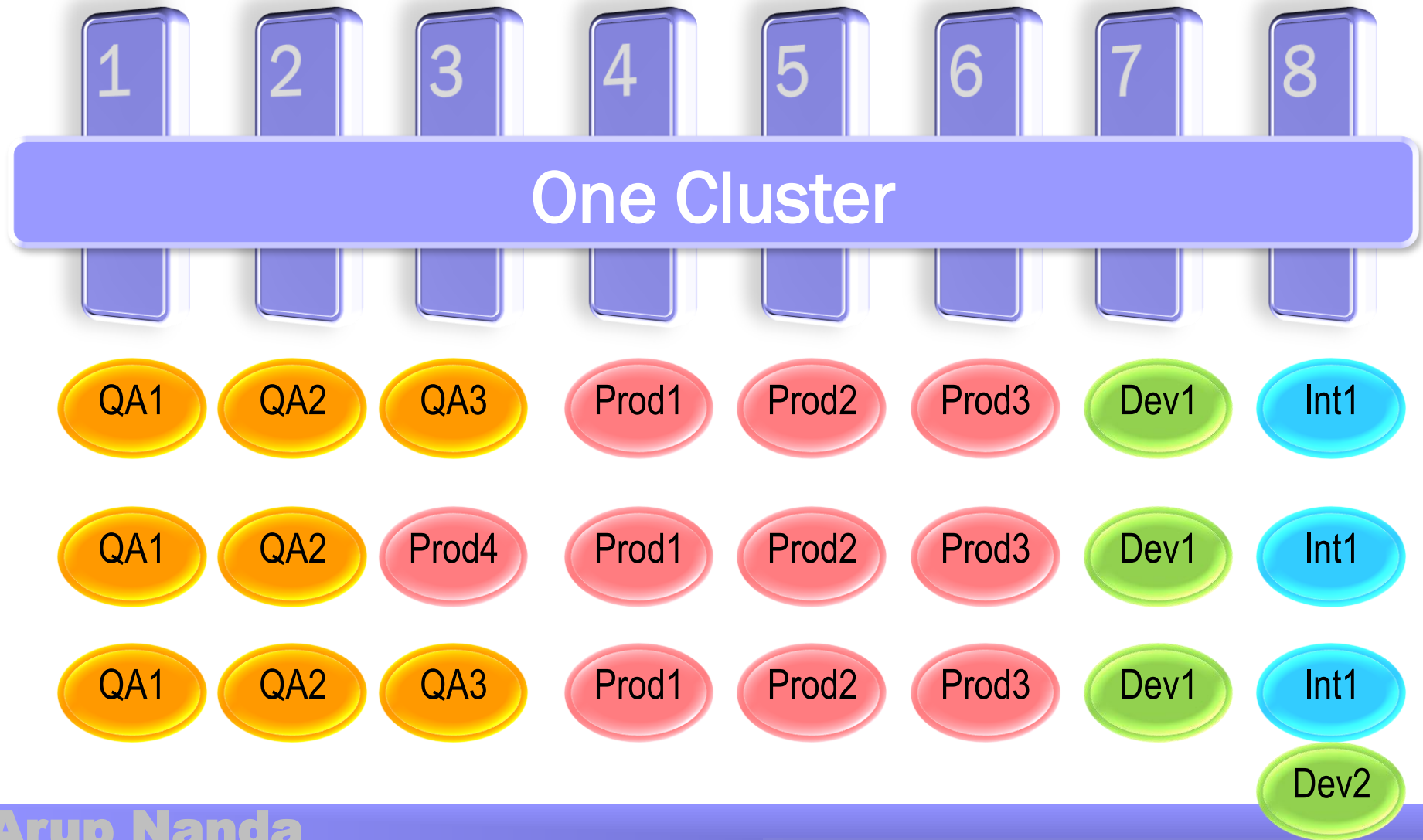
5-part Linux Commands article series

<http://bit.ly/k4mKQS>

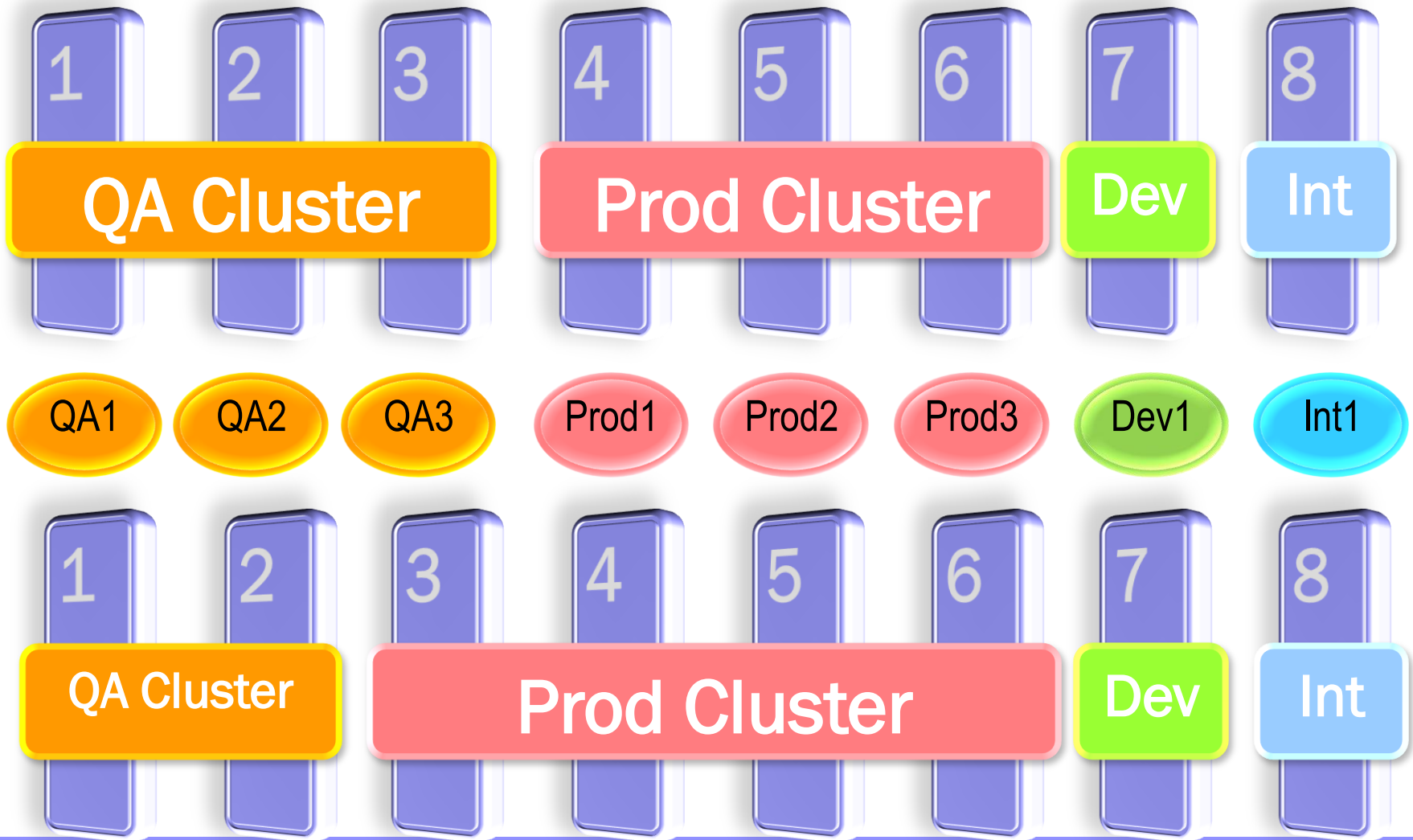
4-part Exadata Command Reference article series <http://bit.ly/IljFIO>



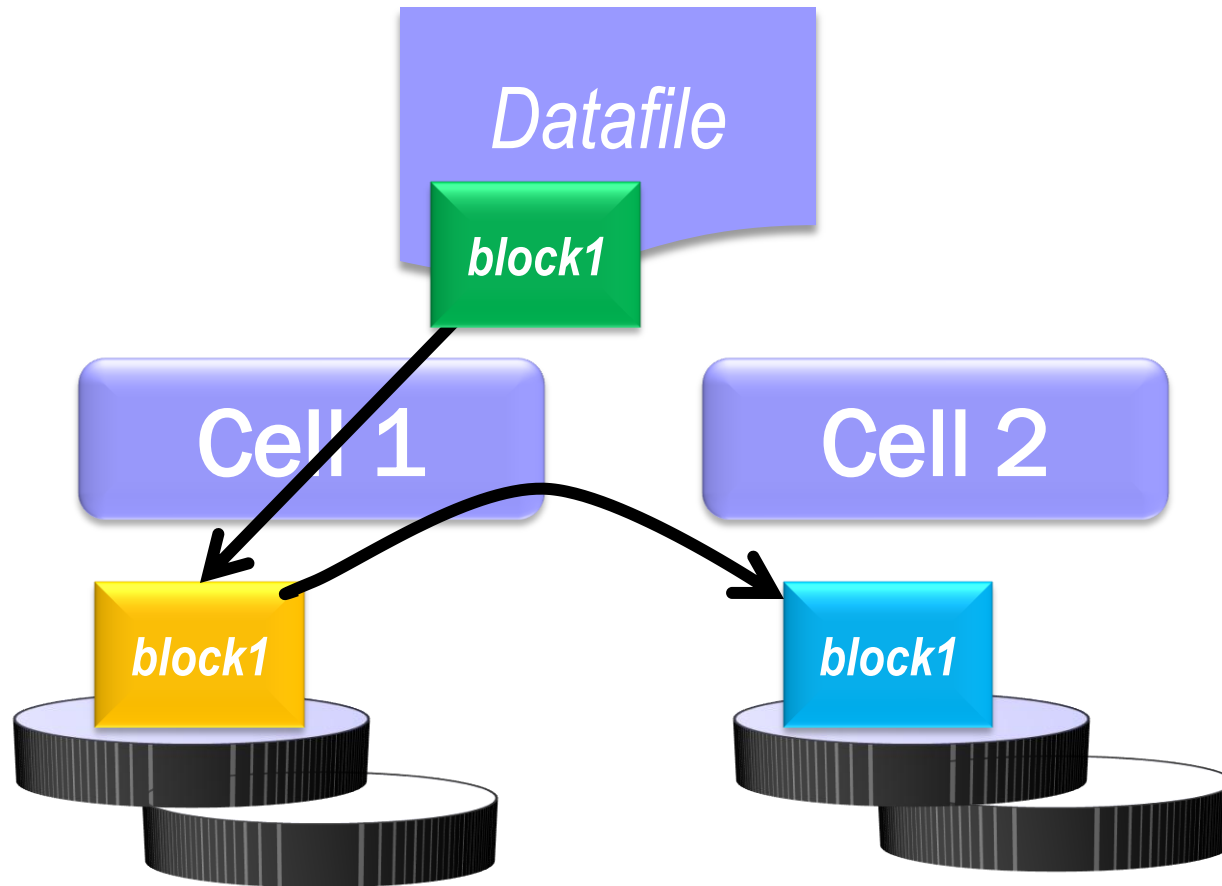
# One Cluster?



# Many Clusters?



# Disk Failures



# Other Questions

Q: Do clients have to connect using Infiniband?

A: No; Ethernet is also available

Q: How do you back it up?

A: Normal RMAN Backup, just like an Oracle Database

Q: How do you create DR?

A: Data Guard is the only solution

Q: Can I install any other software?

A: Nothing on Cells. On nodes – yes

Q: How do I monitor it?

A: Enterprise Manager, CellCLI, SQL Commands

# Summary

- Exadata is an Oracle Database running 11.2
- The storage cells have added intelligence about data placement
- The compute nodes run Oracle DB and Grid Infra
- Nodes communicate with Cells using iDB which can send more information on the query
- Smart Scan, when possible, reduces I/O at cells even for full table scans
- Cell is controlled by CellCLI commands
- DMA skills = 60% RAC DBA + 15% Linux + 20% CellCLI + 5% miscellaneous



# *Thank You!*

Arup Nanda

Blog: [Arup.Blogspot.com](http://Arup.Blogspot.com)

Twitter: [arupnanda](https://twitter.com/arupnanda)

Email: [arup@prolignence.com](mailto:arup@prolignence.com)