



# Scale-out Storage Architectures in the NVM Era

*“Evolution or Revolution?”*

Sandeep Uttamchandani

Chief Architect, Cloud Storage

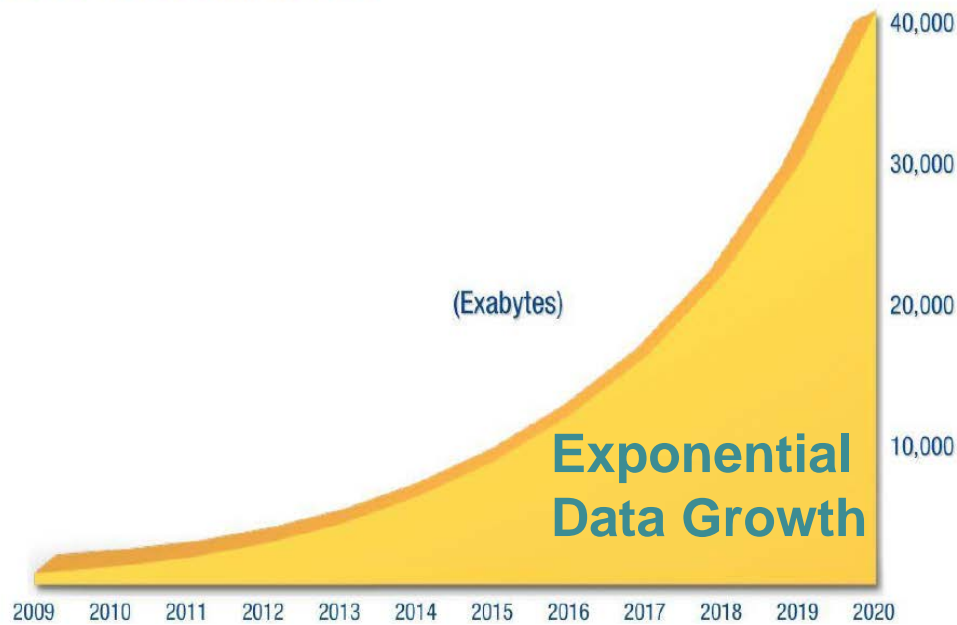
VMWare Inc.

## Agenda: 3 WHYs

- **Why** Enterprise Storage → Scale-out Architectures
- **Why** NVM != *<Yet Another Storage Tier>*
- **Why** New Scale-out Design != Clean slate



# Exponential Data Growth in Enterprises



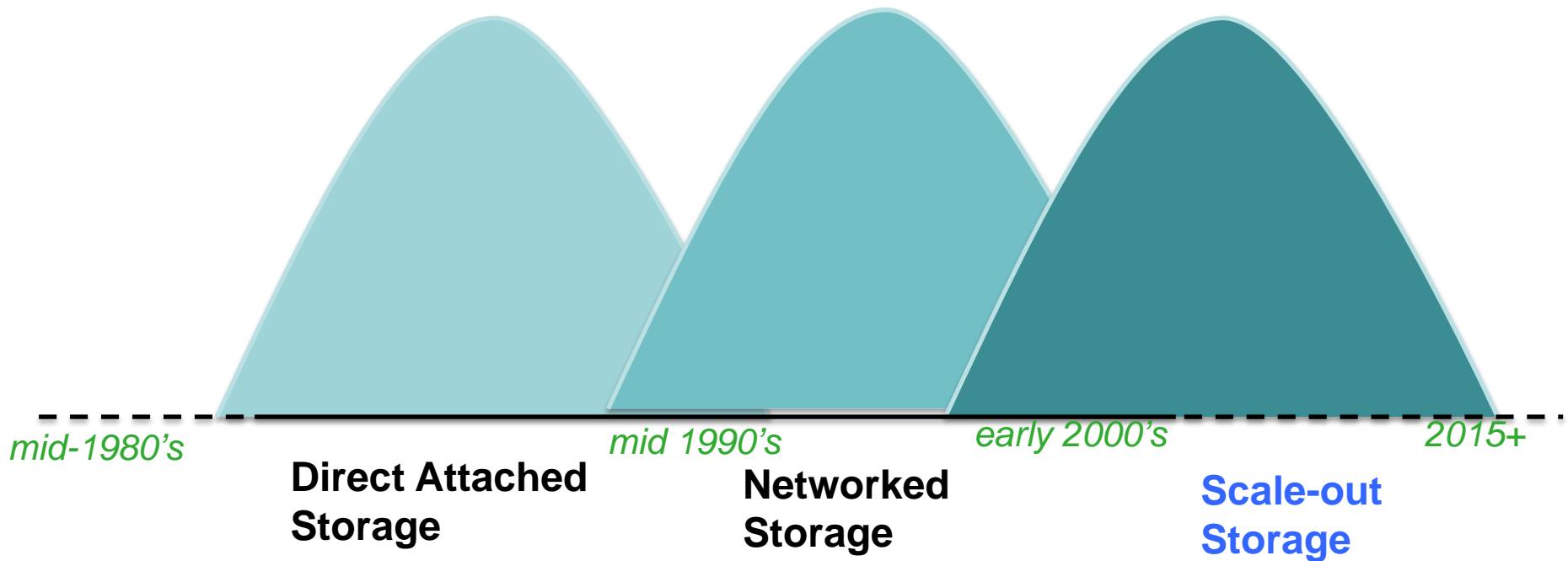
Source: IDC's Digital Universe Study, 2012

***“There were 5 Exabytes of information created between the dawn of civilization through 2003, but that much information is now created every 2 days.”***

*Eric Schmidt, Google 2010 Convention*



# Enterprise Storage Evolution



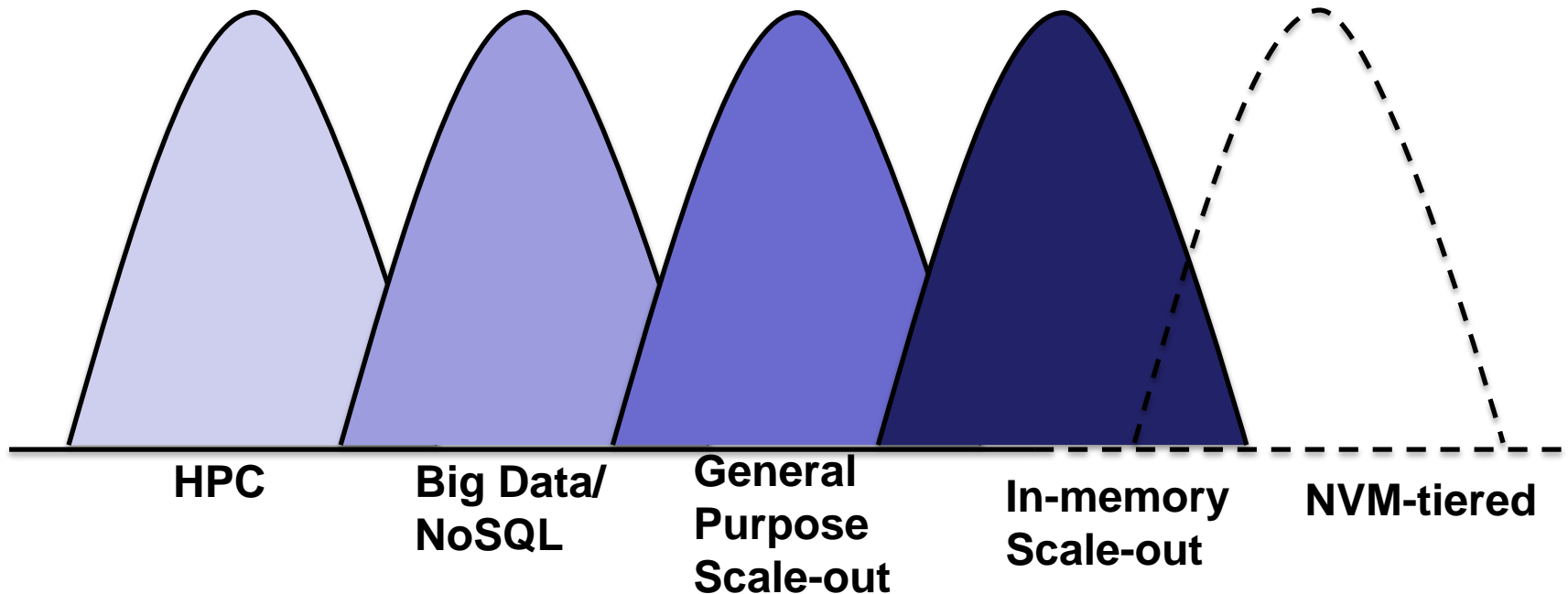
# Scale-out Storage Evolution

- Lustre
- IBM GPFS
- VMWare VMFS
- Veritas CFS
- ...

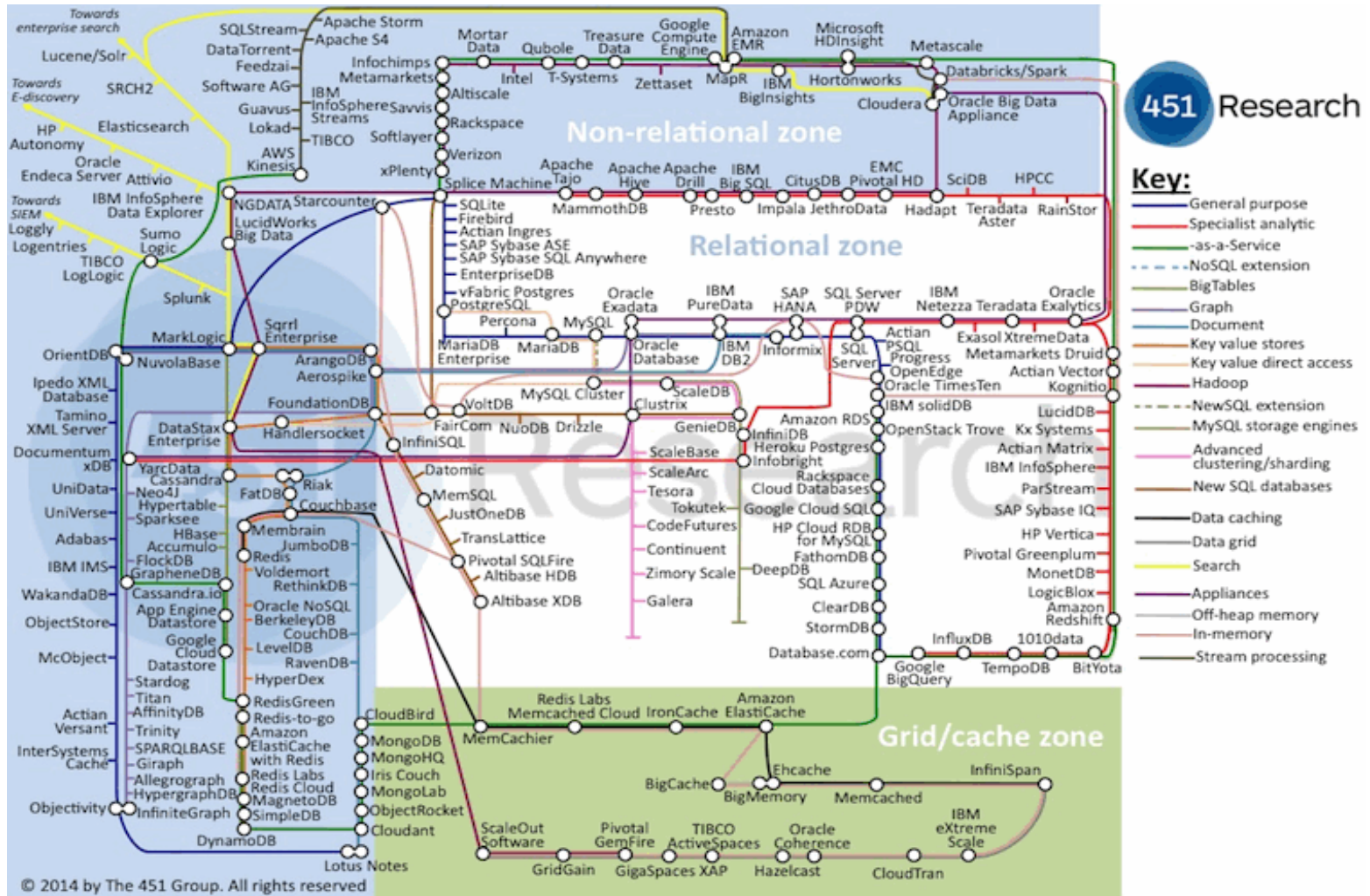
- Cassandra
- HDFS
- MongoDB
- HBase
- ...

- VMWare VSAN
- Ceph
- GlusterFS
- Swift
- ...

- Redis Cluster
- Stanford RamCloud
- Spark/Tachyon
- VoltDB/H-Store
- ...



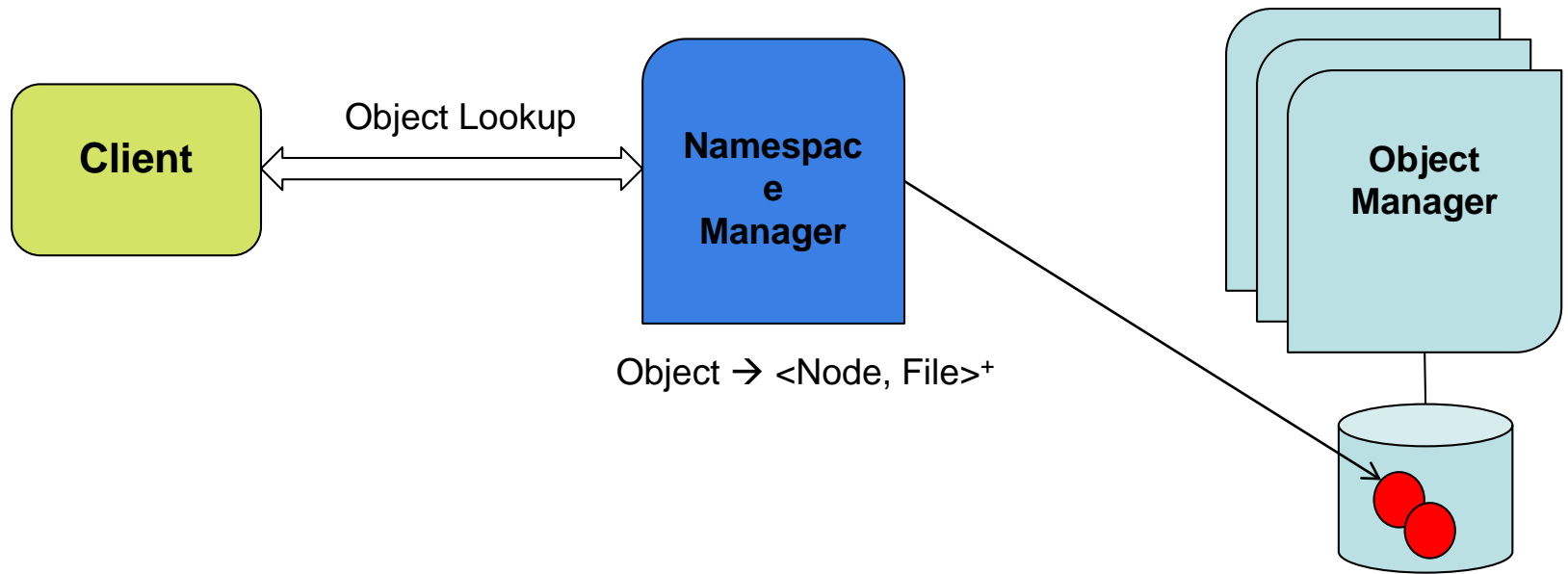
# Scale-out Architectures: End of *one-size-fits-all*



## Agenda: 3 WHYs

- Why Enterprise Storage → Scale-out Architectures
- Why NVM != <Yet Another Storage Tier>
- Why New Scale-out Design != Clean slate

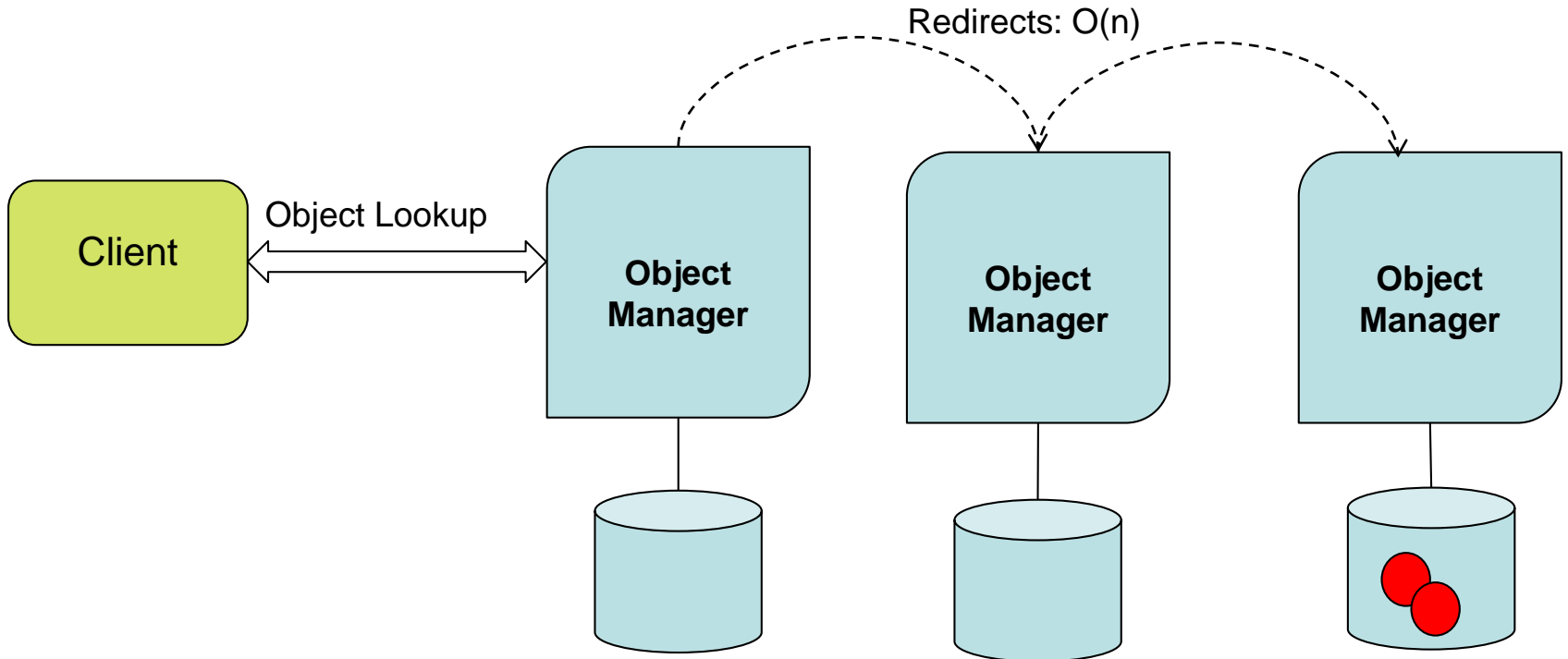
# Scale-out Design: Object Lookup



**Directory Lookup Design Pattern**

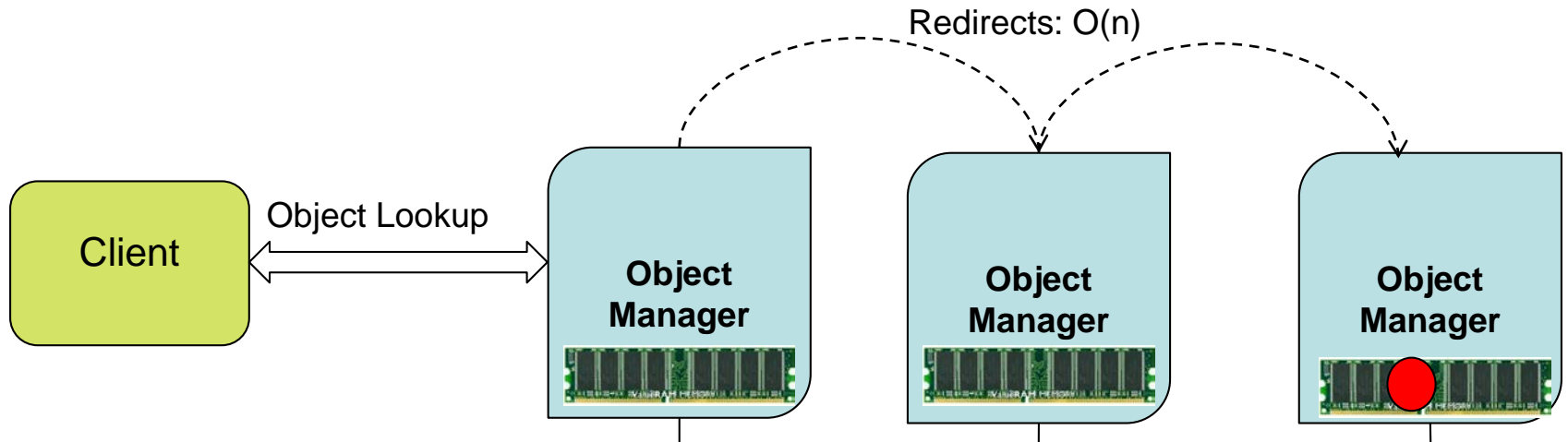


# Scale-out Design: Object Lookup



**Range-based Lookup Design Pattern**

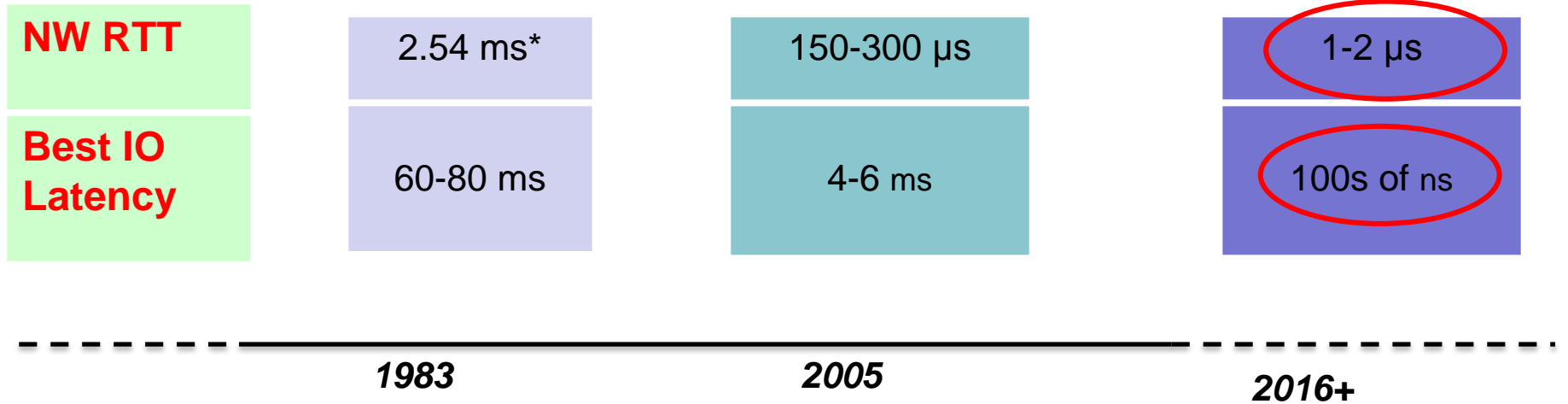
# Scale-out Design: Object Lookup



**Does this deliver the true value of NVM?**

**Range-based Lookup Pattern**

# Network is the new bottleneck



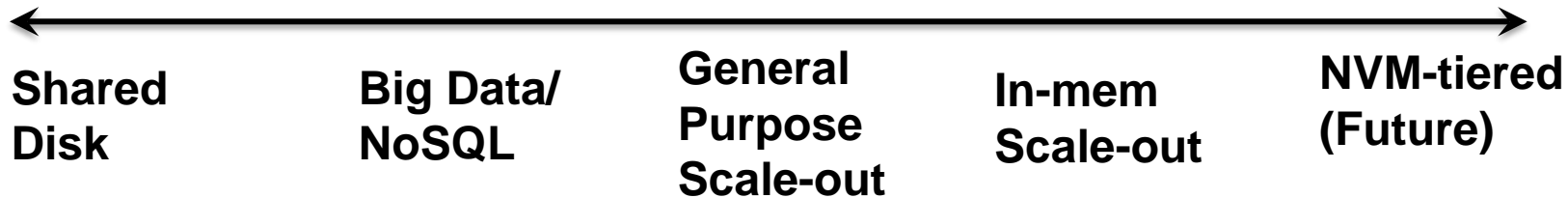
\*S. Rumble, et. al. Its time for low Latency. HotOS 2011



# Namespace Sharding

Object Lookup: ★ Client-aware    ☆ Client-opaque

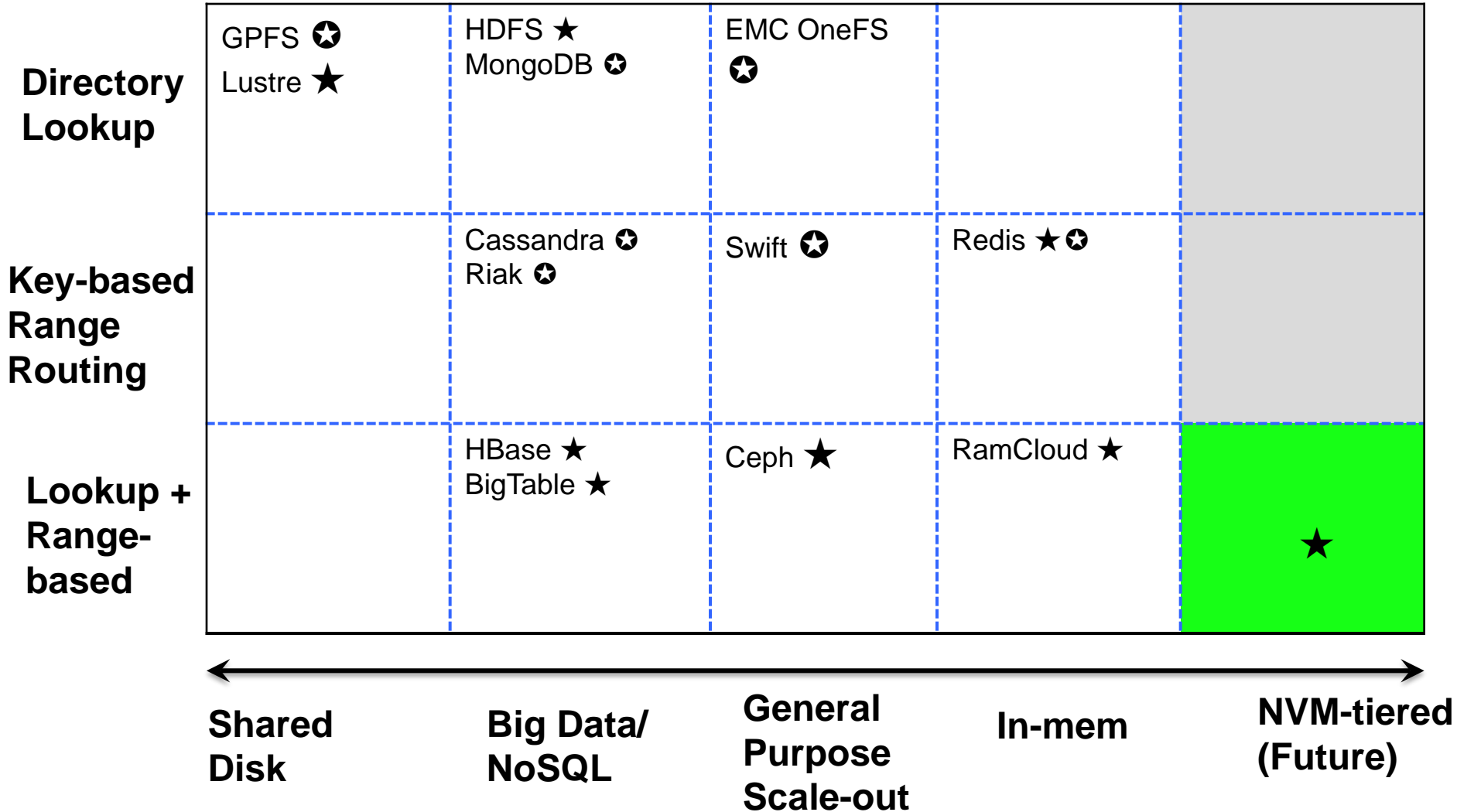
<b>Directory Lookup</b>	GPFS Lustre	HDFS MongoDB	EMC OneFS	
<b>Key-based Range Routing</b>		Cassandra Riak	Swift	Redis
<b>Lookup + Key-based Routing</b>		HBase BigTable	Ceph	RamCloud



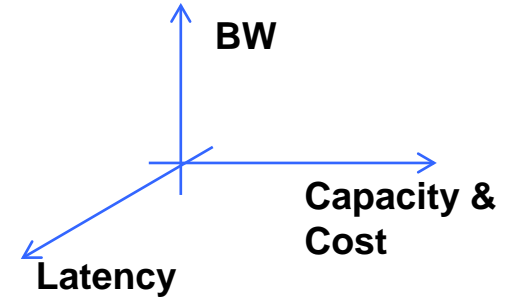


# Namespace Sharding

Object Lookup: ★ Client-aware    ☆ Client-opaque



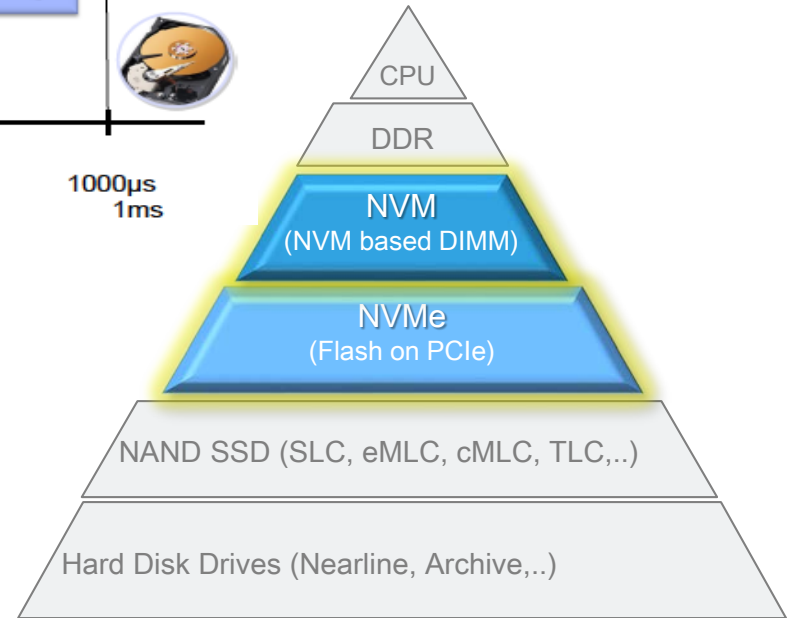
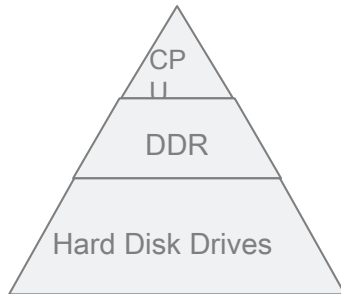
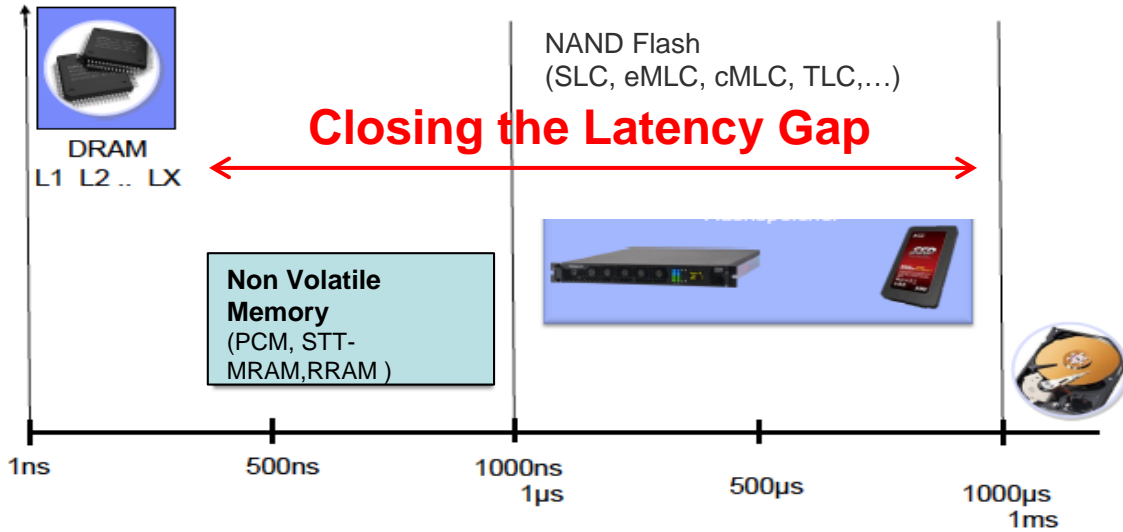
# Holistic Shifts in Design Constraints



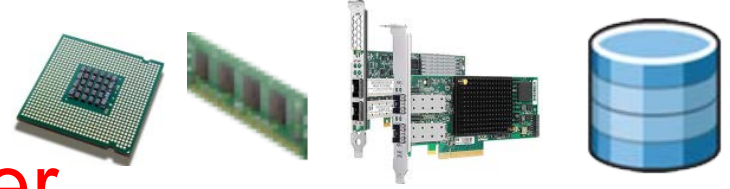
<b>Application</b>	POSIX/ACID Batch-oriented	Non-POSIX/BASE Real-time	<ul style="list-style-type: none"> <li>Beyond traditional block/file POSIX applications</li> </ul>
<b>CPU</b>	1x10MHz	16x3GHz	<ul style="list-style-type: none"> <li>Slow-down of Moore's law</li> </ul>
<b>Memory</b>	< 2 MB	>16GB	<ul style="list-style-type: none"> <li>Bigger &amp; Cheaper</li> </ul>
<b>Network</b>	3Mbps	10Gbps	<ul style="list-style-type: none"> <li>Network is becoming the new latency bottleneck</li> </ul>
<b>Storage</b>	<30MB	>4TB	<ul style="list-style-type: none"> <li>Emergence of distinct Capacity and Performance Storage Tiers</li> </ul>

-----  
**1984**
Avg. Node Configuration
**2012**
-----

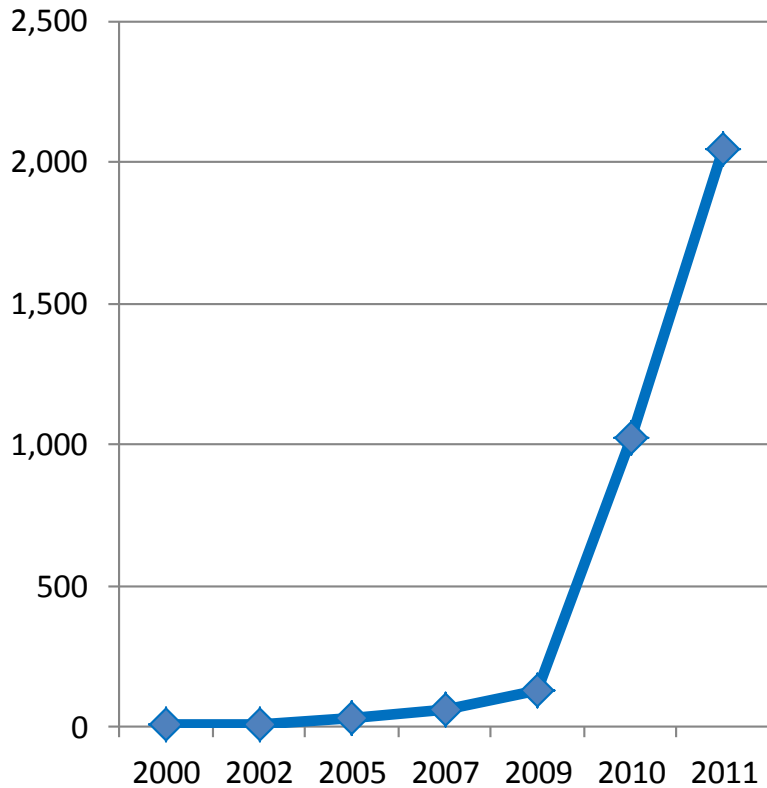
# The New Storage Hierarchy



# Bigger and Cheaper

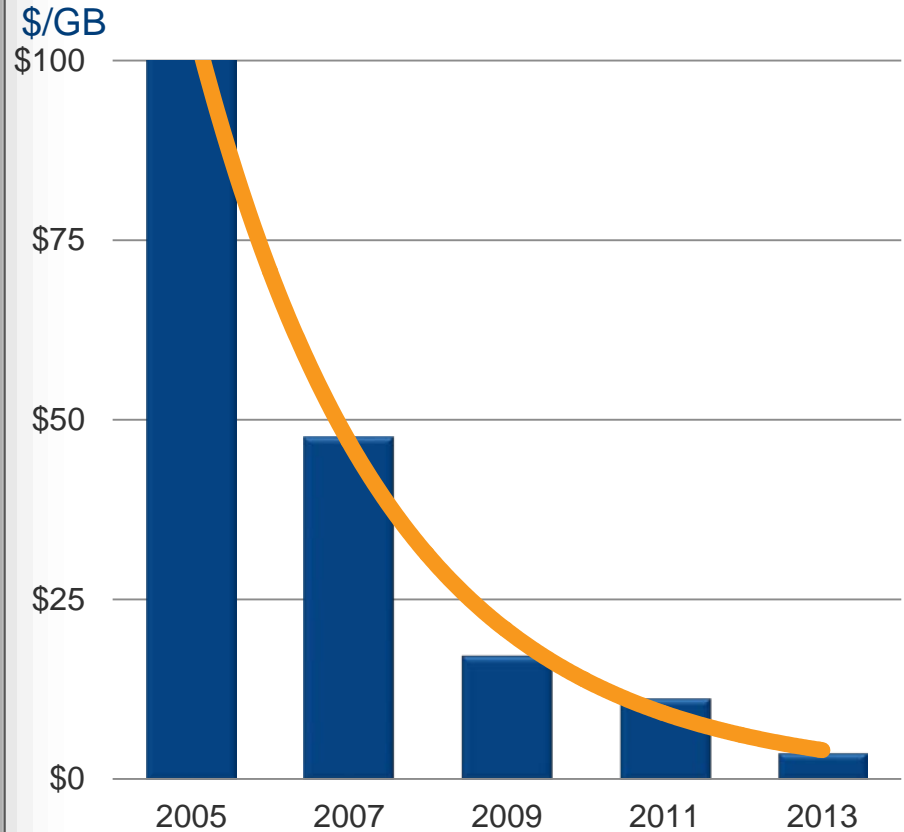


**Maximum DRAM in GB**



Source: Forrester Research, *The x86 Server Grows Up And Out* (October 8, 2010)

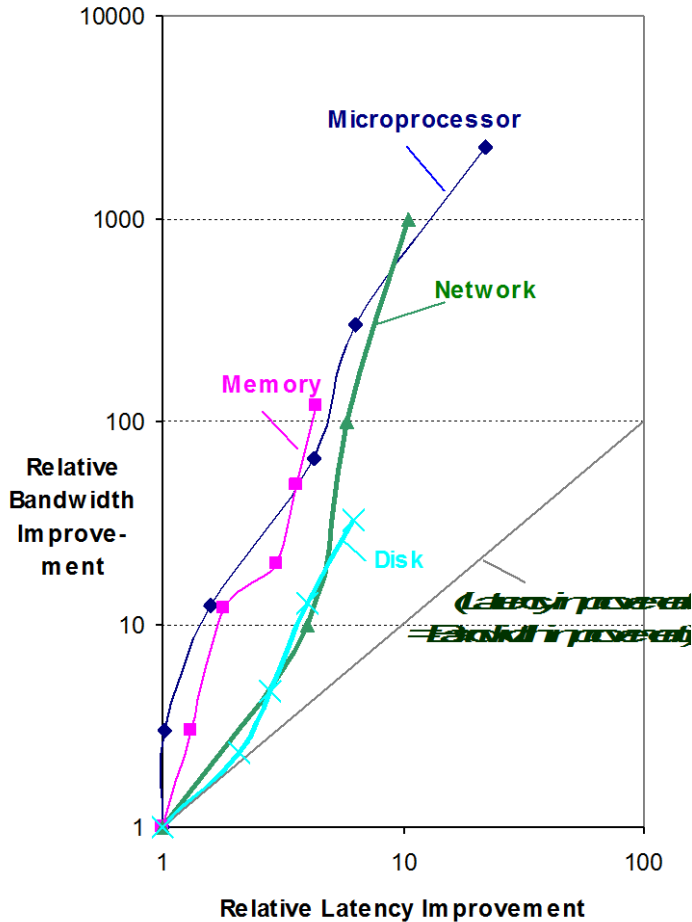
**DRAM**



Source: Gartner Dataquest, *Forecast: DRAM Market Statistics (1Q11)*



# Latency lags Bandwidth\*



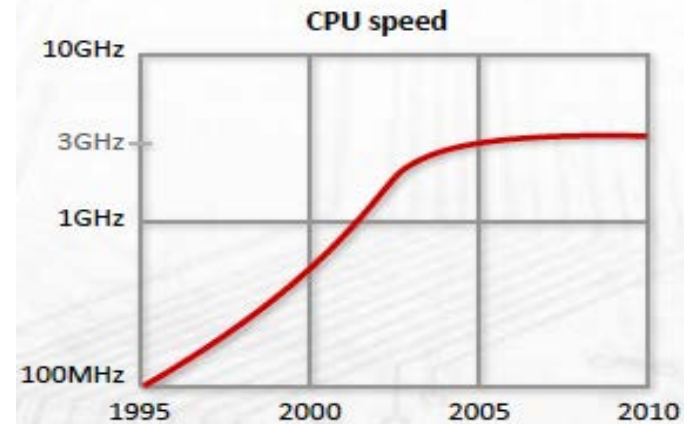
	Mid-1980s	2009	Change
Disk capacity	30 MB	500 GB	16667x
Max. transfer rate	2 MB/s	100 MB/s	50x
Latency (seek & rotate)	20 ms	10 ms	2x

Source: Stanford RamCloud Talk, Feb 2010

\*David Patterson, Latency Lags Bandwidth, CACM, 2004

# CPU Scaling: Slow-down of Moore's Law

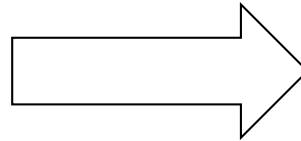
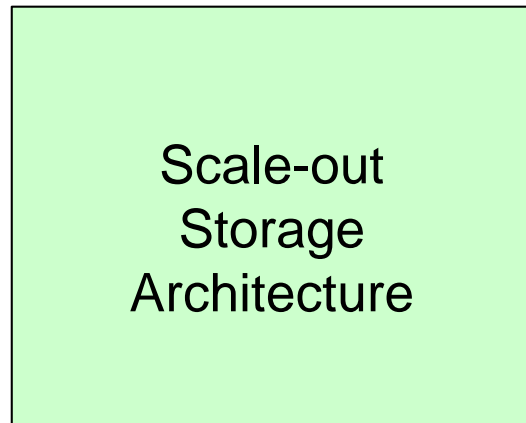
- CPU Scaling is not exponential anymore
  - 3000% increase from 1995-2004
  - 30-40% increase since 2004
- Multi-core scaling is linear
  - NUMA, locking, sharing latencies, programming models...
- Compute per unit of data is decreasing



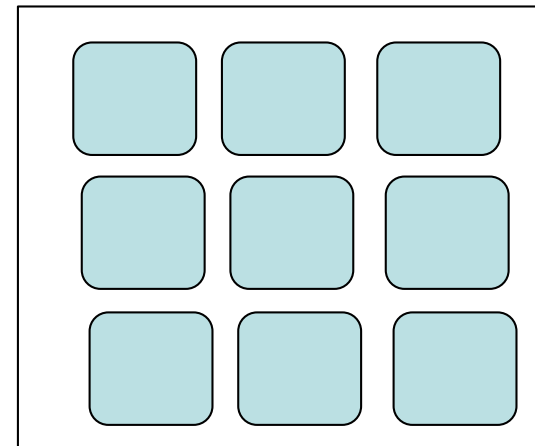
## Agenda: 3 WHYs

- Why Enterprise Storage → Scale-out Architectures
- Why NVM != <Yet Another Storage Tier>
- Why New Scale-out Design != Clean slate Re-design

# What does re-design mean?

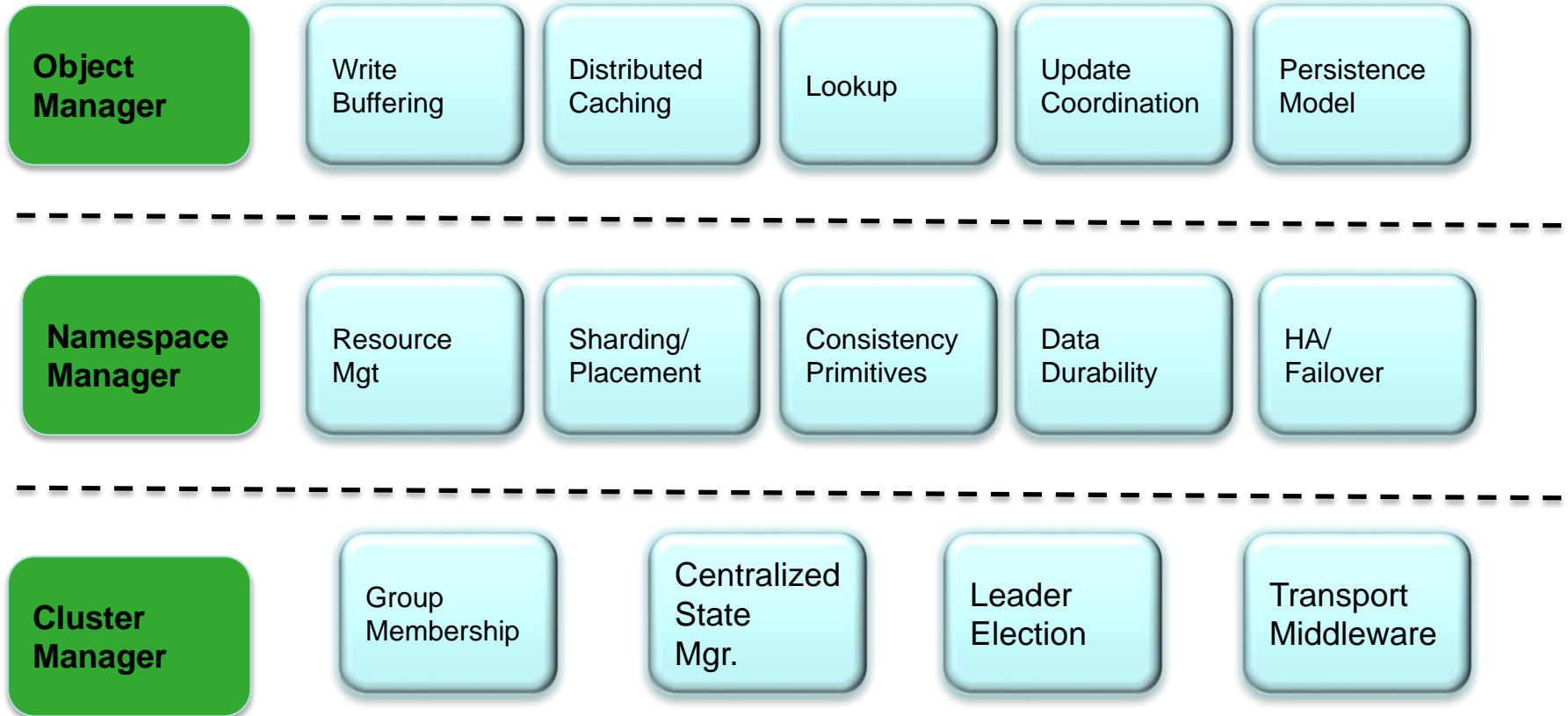


Design of a collection of  
Micro-Services





# Blueprint of a Scale-out Storage Architecture

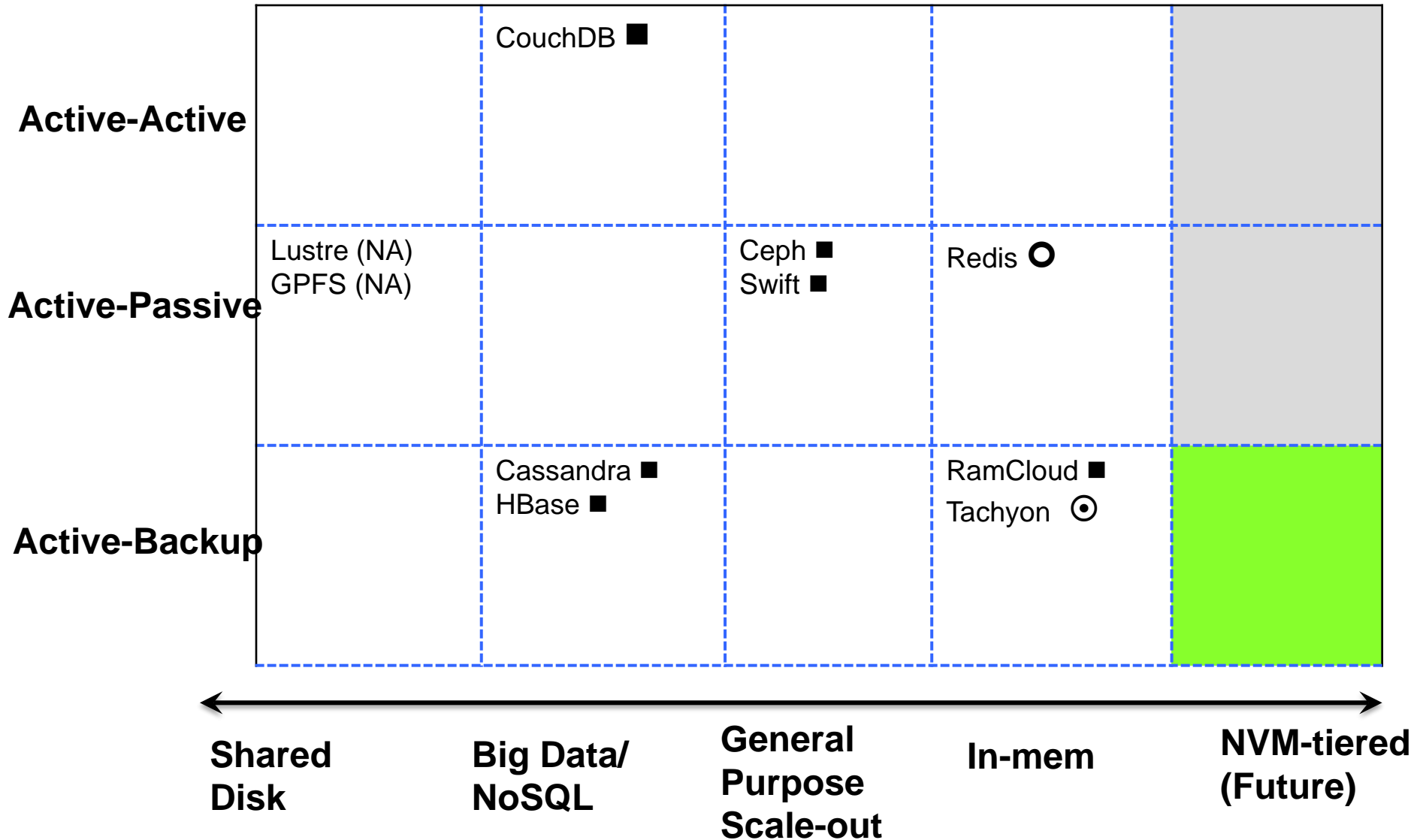




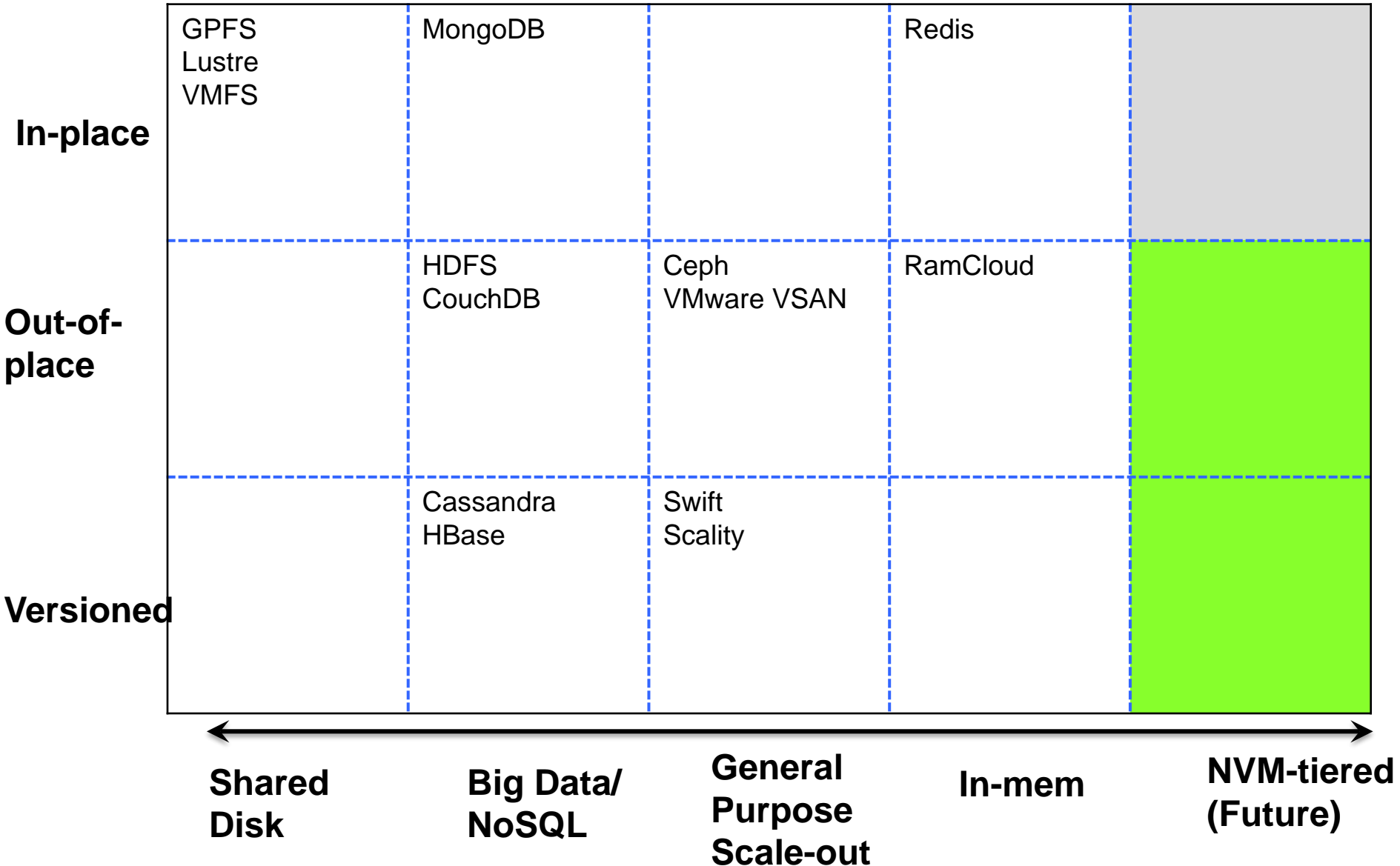


# High Availability

Data Replication: ■ State-based ○ Ops-based ⊙ Lineage



# Persistence Model





# Summary

- Why Enterprise Storage → Scale-out Architectures
  - Web 2.0 model to handle exponential data growth
- Why NVM != *<Yet Another Storage Tier>*
  - Holistic shifts across compute, network, memory, storage
- Why New Scale-out Design != Clean slate
  - Piecemeal evolution of micro-services within a Scale-out Architecture.



# Questions?

Sandeep Uttamchandani  
[suttamchandani@vmware.com](mailto:suttamchandani@vmware.com)

Blog: <https://blogs.vmware.com/cto/author/suttamchandani/>