

Reference Architecture and Best Practices for Virtualizing Hadoop Workloads

Justin Murray VMware







- The Hadoop Journey
- Why Virtualize Hadoop?
- Elasticity and Scalability
- Performance Tests
- Storage Reference Architectures
- Isilon Architecture and Benefits
- vSphere Big Data Extensions
- Conclusion and Q&A

The Customer Journey with Hadoop



The Hadoop Journey



Why Virtualize Hadoop?





Customer Example: Enterprise Adoption of Hadoop





What if you could...





Big Data Extensions Value Propositions BIG

Operational Simplicity with Performance

- ✓ Rapid Deployment
- ✓ Self service tools
- ✓ Performance

Maximize Resource Utilization

- ✓ Elastic scaling
- Avoid dedicated hardware
- ✓ VM-based isolation
- ✓ Increase resource utilization
- ✓ True multi-tenancy

Architect Scalable Platform

- ✓ Deployment choice
- Maintain management flexibility at scale
- ✓ Control Costs
- Leverage vSphere features

Hadoop 2.0 – Yet Another Resource Negotiator



BIG

A Virtualized Hadoop 2.0 Cluster

BIG

SUMMIT



vSphere Big Data Extensions

- Deploy Hadoop Clusters in Minutes



		Create New Big Data Cluster					
Server preparation		Big data cluster name:	foo				
		Hadoop distribution:	apache				
			Vendor: Apache	Version: 1.2.0			
		Deployment type:	Basic Hadoop Clu	ister 🗸 🗸			
OS installation		DataMaster Node Group	0				•
		Number of nodes:	1	Resource template:	Medium	-	
				2 vCPU, 7500 MB RAI Shared datastore	M, 50 GB storage on		
		ComputeMaster Node Gro	oup 🚯				
	,	Number of nodes:	1	Resource template:	Medium	-	
Network Configuration				2 vCPU, 7500 MB RAI Shared datastore	M, 50 GB storage on		
		Worker Node Group 🏾 🕕					
		Number of nodes:	3 _	Resource template:	Small	-	•
		Hadoop topology:	HOST_AS_RACK	•			
Hadoop Installation and		Network:	dhcpNetwork				
Configuration		Resources:	CI	Sele	et		
					ОК	Can	icel
From a manual process	-	To fully aut	omated,	using the	GUI		

Elastic, Multi-Tenant Hadoop with Virtualization



VM

Hadoop Node

Combined Compute and Storage

Unmodified Hadoop node in a VM

- + VM lifecycle determined by Datanode
- + Limited elasticity





Separate Compute from Storage

- + Separate compute from data
- + Stateless compute
- + Elastic compute

Separate Virtual Compute Clusters per tenant

- + Separate virtual compute
- + Compute cluster per tenant
- + Stronger VM-grade security and resource isolation

Performance and Reference Architectures



Native vs. Virtual, 32 hosts, 16 disks per automatic



2014 SNIA Analytics and Big Data Summit. © Insert Your Company Name. All Rights Reserved.

Source: http://www.vmware.com/resources/techresources/10360

Reference Architecture: 32-Server Performance Test





Up to four VMs per server vCPUs per VM fit within socket size (e.g. 4 VMs x 4 vCPUs, 2 X 8) Memory per VM - fit within NUMA node size

2013 Tests done using Hadoop 1.0

I/O Profile of a Hadoop MapReduce Job SUMMIT

(TeraSort example application)



The Combined Model – Standard Deployment





Combined Model – Two Virtual Machines on a Host Server





The Data-Compute Separation Deployment Model





Data Paths: Combined vs Data-Compute Separation





Alternative Storage for Data/Compute Separation



G

SUMMIT

VMDK Isolation for Performance





Data/Compute Separation With Isilon



BIG

Z SUMMIT

Larger Architecture with Data Compute BIG Separated



Hybrid storage model - the best of both worlds





Shared Storage



Local Storage



Master nodes

- NameNode, ResourceManager, ZooKeeper etc. on shared storage
 - Leverage vSphere vMotion, HA and FT

Worker nodes

- NodeManager/DataNode on local storage
- Lower cost, scalable bandwidth
- Temp data is written to local storage for best performance
- NFS storage for HDFS data is a very good alternative to local

vSphere Big Data Extensions and Project Serengeti



Big Data Extensions - Highlights





Contributed back to Apache Hadoop

Introducing vSphere Big Data Extensions (BDE)





Brief Tour of Big Data Extensions



One Click to Scale out the Cluster on the Fly



BIG

BDE Allows Flexible Configurations

"name": "master",

"instanceNum": 1,

"type": "SHARED",

"cpuNum": 2,

"storage": {

"sizeGB": 20

"haFlag":"on",

"rpNames": [

"rp1"

"roles": [

],

},

ł



External HDFS : Simple to Set Up



Firefox T							
ங http://10.111.89.51:8080/serengeti/# 🛛 🗙 🙋	VSphere We	eb Client 🔅	🗵 🧽 Image Mining	Demo × -	+		
vm ware ⁻							
44	🕋 Clus	Create New Cluster			×		
⊂ Clusters	+ 🚳	Cluster name:	Tier2_demo				
i ieri	Cluster Nam	Hadoop distro:	PivotalHD	-		Cluster Information	P
√ ∭ Resources	Tier1		Vendor: GPHD V	'ersion: 1.2.0.0		1 master, 2 data, 24 worker	
🕞 Resource Pools	11612	Deployment type:	Compute-only Hado	oop Cluster 🛛 👻	J	T master, To worker, T chent	
Datastores		DataMaster					
m Distros	-	DataMaster URL:	hdfs://isilon1.vmv	/are.com			
			Please input exter	mal hdfs RPC URL here.			
		ComputeMaster Node Group 🕕					
		Number of nodes:	1	Resource template: Medi	um 🚽		
				2 vCPU, 7500 MB RAM, 50 G Shared datastore	38 storage on		
		Worker Node Group 🕕					
		Number of nodes:	16 🌲	Resource template: Custo	omize 🔻		
				1 vCPU, 3748 MB RAM, 25 G	B storage on		
		Client Node Group(Optionz	a) (i)	Snared batastore			
		Number of nodes:	1	Resource template: Smal	•		
				1 vCPU, 3748 MB RAM, 50 G	B storage on		
				Shared datastore			
		Notwork					
		NEWUIK.	defaultNetwork	•			
					Cancel		

How BDE works





vSphere Configuration



Provision the virtual machines at the right size

- Reserve 6% of physical memory on the ESXi Server for vSphere usage
- Avoid over-commitment
- Enable NUMA and keep the virtual machine memory and cpu size within the NUMA node
- NUMA scheduler is important for virtualized Hadoop performance
 - Poor configuration can result in performance degradation
 - Data VM preferably should be distributed across NUMA nodes

VMware vSphere BDE and Hadoop Resources



- VMware vSphere BDE web site
 - http://www.vmware.com/bde



- Virtualized Hadoop Performance with VMware vSphere 5.1
 - http://www.vmware.com/resources/techresources/10220
- Benchmarking Case Study of Virtualized Hadoop Performance on vSphere 5
 - http://vmware.com/files/pdf/VMW-Hadoop-Performance-vSphere5.pdf
- Hadoop Virtualization Extensions (HVE) :
 - http://www.vmware.com/files/pdf/Hadoop-Virtualization-Extensions-on-VMware-vSphere-5.pdf
- Apache Hadoop High Availability Solution on VMware vSphere 5.1 <u>http://vmware.com/files/pdf/Apache-Hadoop-VMware-HA-solution.pdf</u>





- Hadoop workloads work very well on VMware vSphere
 - Various performance studies have shown that any difference between virtualized performance and native performance is minimal
 - Follow the general best practice guidelines that VMware has published
- vSphere Big Data Extensions enhances your Hadoop experience on the VMware virtualization platform
 - Rapid provisioning tool for deployment of Hadoop components in virtual machines
 - Algorithms for best layout of your Hadoop data and cluster components are built into the BDE HVE components
 - Design patterns such as data-compute separation can be used to provide elasticity of your Hadoop cluster on demand.
 - User self service available with Hadoop using tools such as vCloud Automation Center integrated with BDE

Thank You

jmurray@vmware.com Justin Murray









Backup Slides

Today's Challenges on Hadoop Infrastructure

- Fixed compute and storage coupling^{UM} leads to low utilization and inflexibility

Compute Node Storage Node

Server

- Compute and storage linked together with fixed ratio based on the hardware specification
- Not all jobs are created equal (data vs. compute intensive)
- Inflexible infrastructure leads to waste
 - Too little compute power \rightarrow slow processing
 - Too much compute power \rightarrow sitting idle
- Problem compounds with larger clusters

Getting more out of your infrastructure



Compute layer

Storage layer

Decouple the linkage between

compute and storage

for other worklands

Stateless compute can grow and shrink elastically

Z SUMM

- Data locality is preserved, place the compute where data resides
- Extra compute capacity can be used



Elasticity and Scalability





Elastic Scalability & Multiple Workloads BIG

- Deploy separate compute clusters for different tenants sharing HDFS.
- Commission/decommission compute nodes according to priority and available resources



Hadoop 1.0



