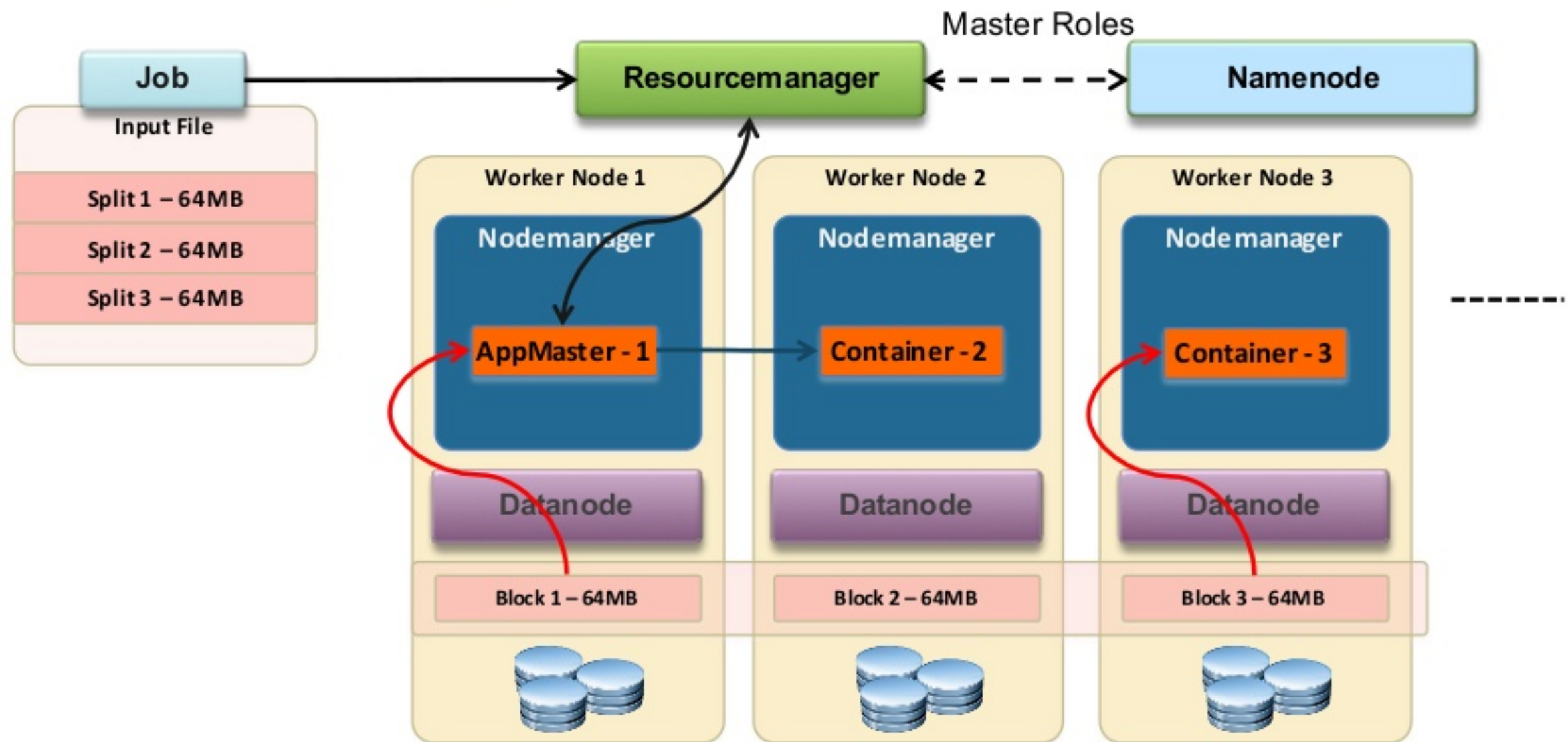# Why Virtualize Spark?

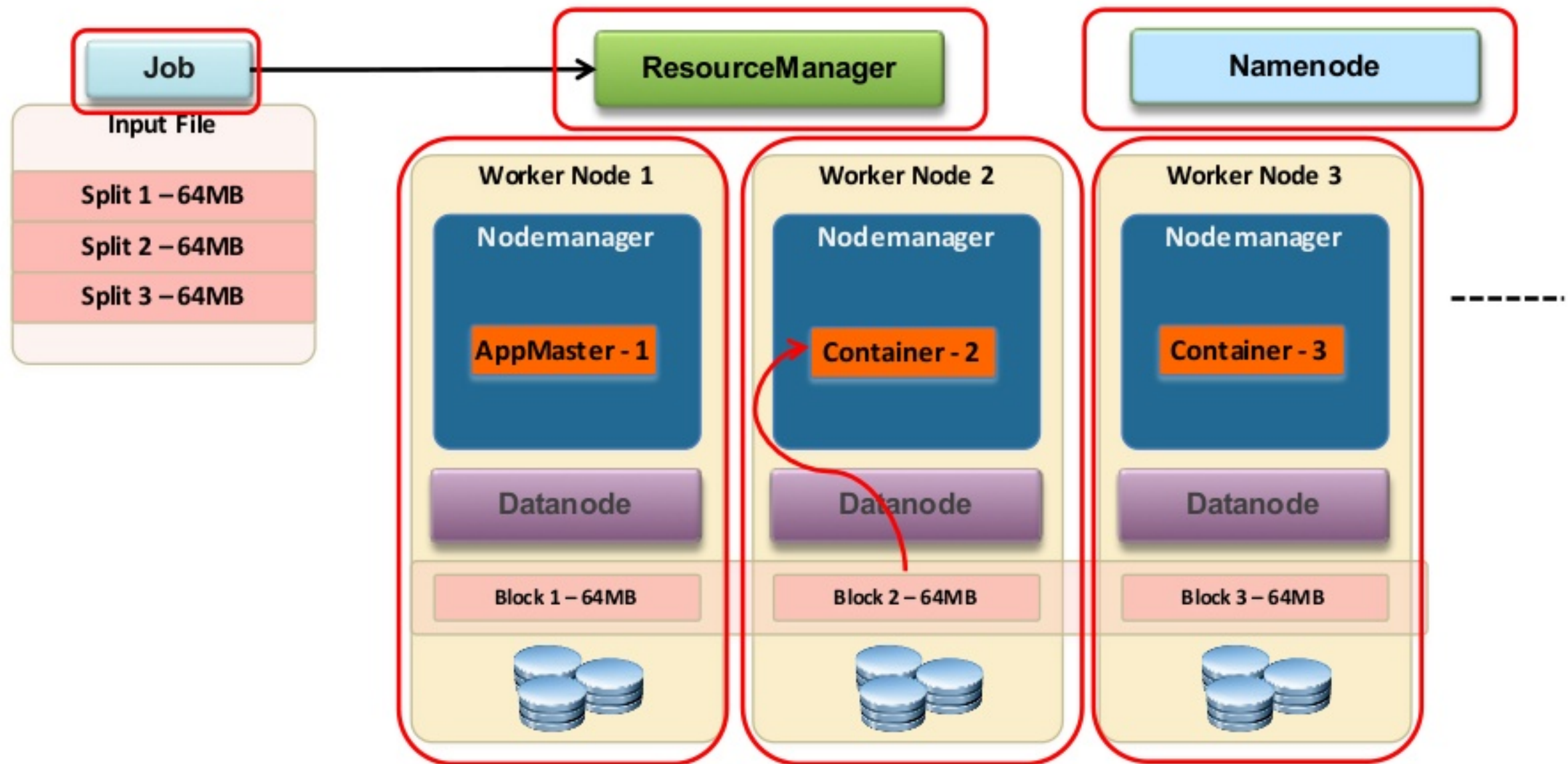# Use Cases : Virtualization of Big Data

- IT wants to provide **Spark clusters as a service** on-demand for its end users

- Enterprises have development, test, pre-prod staging and production clusters that are required to be separated from each other and provisioned independently

- Organizations need different versions of Spark to be available to different teams - with possibly different services available

- Enterprises do not wish to dedicate a specific set of hardware to each different requirement above, and want to reduce overall costs

# The Traditional Hadoop Architecture

**Master Roles**

**Job** → **Resourcemanager** ← - - - - → **Namenode**

**Input File**

| |
|---|
| Split 1 – 64MB |
| Split 2 – 64MB |
| Split 3 – 64MB |

**Worker Node 1**

Nodemanager

AppMaster - 1

Datanode

Block 1 – 64MB

**Worker Node 2**

Nodemanager

Container - 2

Datanode

Block 2 – 64MB

**Worker Node 3**

Nodemanager

Container - 3

Datanode

Block 3 – 64MB

# Hadoop – in Virtual Machines

Master Roles

| Job |

**Input File**

Split 1 – 64MB

Split 2 – 64MB

Split 3 – 64MB

**ResourceManager**

**Namenode**

**Worker Node 1**

Nodemanager

AppMaster - 1

Datanode

Block 1 – 64MB

**Worker Node 2**

Nodemanager

Container - 2

Datanode

Block 2 – 64MB

**Worker Node 3**

Nodemanager

Container - 3

Datanode

Block 3 – 64MB

# The Spark Architecture – Standalone

# Spark Standalone - Virtualized

Virtual Machine

Job → Driver

**Worker Node 1**
- Executor
  - JVM
- Executor
  - JVM

**Worker Node 2**
- Executor
  - JVM
- Executor
  - JVM

**Worker Node 3**
- Executor
  - JVM
- Executor
  - JVM

# The Spark Architecture (on YARN)

# Reference Architectures

# Combined Model: Two Virtual Machines on a Host

**Hadoop Node 1 Virtual Machine**

Nodemanager · Datanode

Ext4 · Ext4 · Ext4 · Ext4 · Ext4 · Ext4

**Hadoop Node 2 Virtual Machine**

Nodemanager · Datanode

Ext4 · Ext4 · Ext4 · Ext4 · Ext4 · Ext4

**Virtualization Host Server**

VMDK · VMDK · VMDK · VMDK · VMDK · VMDK · VMDK · VMDK · VMDK · VMDK · VMDK · VMDK

Six or More Local DAS disks per Virtual Machine

# #1 Reference Architecture from Cloudera

**cloudera**

CLOUDERA REFERENCE
ARCHITECTURE FOR VMWARE
vSPHERE WITH LOCALLY ATTACHED
STORAGE
VERSION CDH 5.3

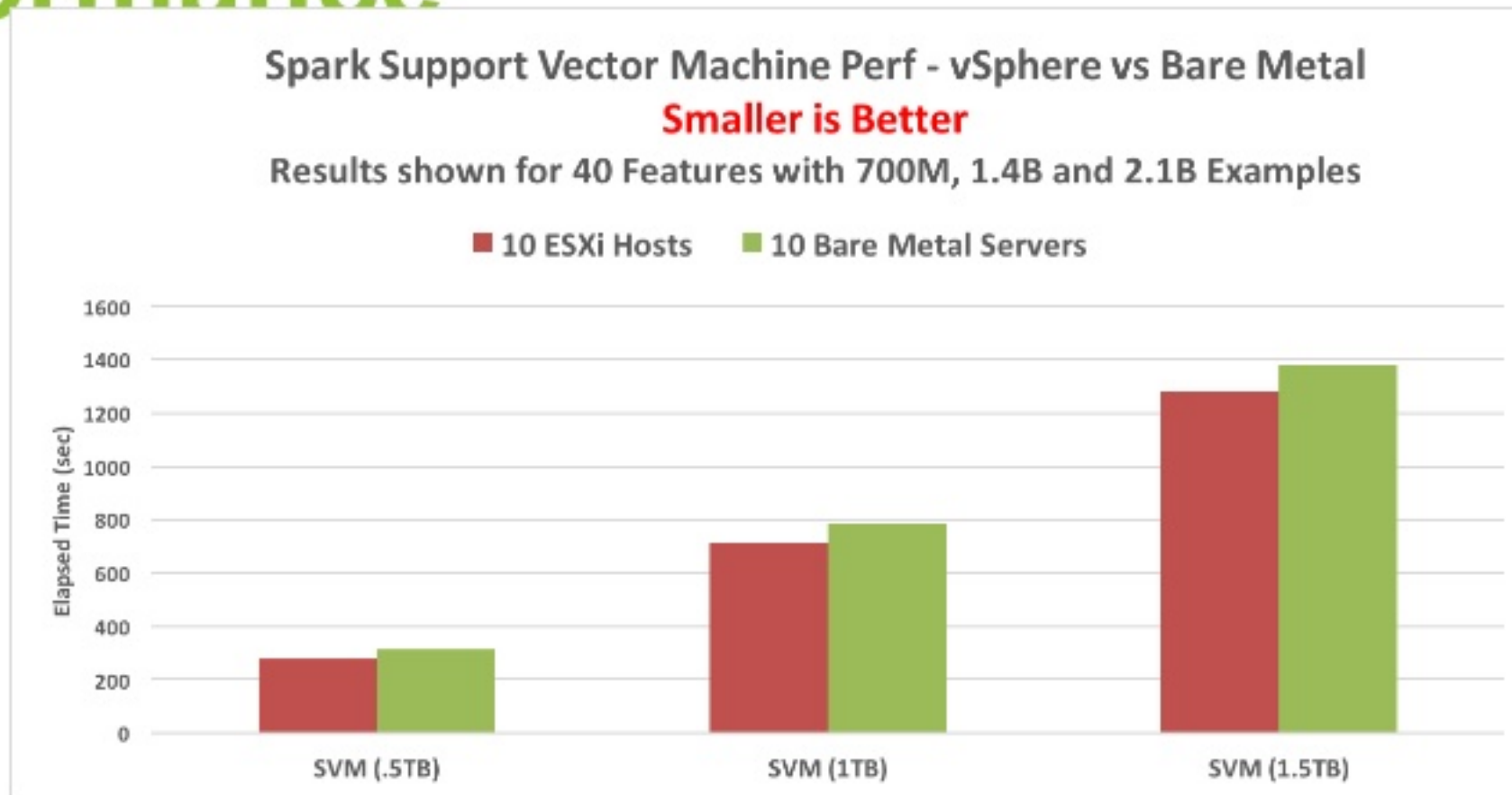# Performance

# Workloads - Spark

- Two standard analytic programs from the Spark MLLib (Machine Learning Library)

- Driven using SparkBench (https://github.com/SparkTC/spark-bench)
    - Support Vector Machine
    - Logistic Regression

# Spark Support Vector Machine Performance



Spark Support Vector Machine Perf - vSphere vs Bare Metal
**Smaller is Better**
Results shown for 40 Features with 700M, 1.4B and 2.1B Examples

# Spark Logistic Regression Performance



Spark Logistic Regression Performance - vSphere vs. Bare Metal
**Smaller is Better**
Results shown for 40 and 80 Features with 700M, 1.4B and 2.1B Examples

Legend: ■ 10 ESXi hosts  ■ 10 Bare Metal Servers

Y-axis: Elapsed Time (sec), 0 to 700

Categories: LR-40F (.5TB), LR -40F (1TB), LR-40F (1.5TB), LR-80F (1TB), LR -80F (2TB), LR-8F (3TB)
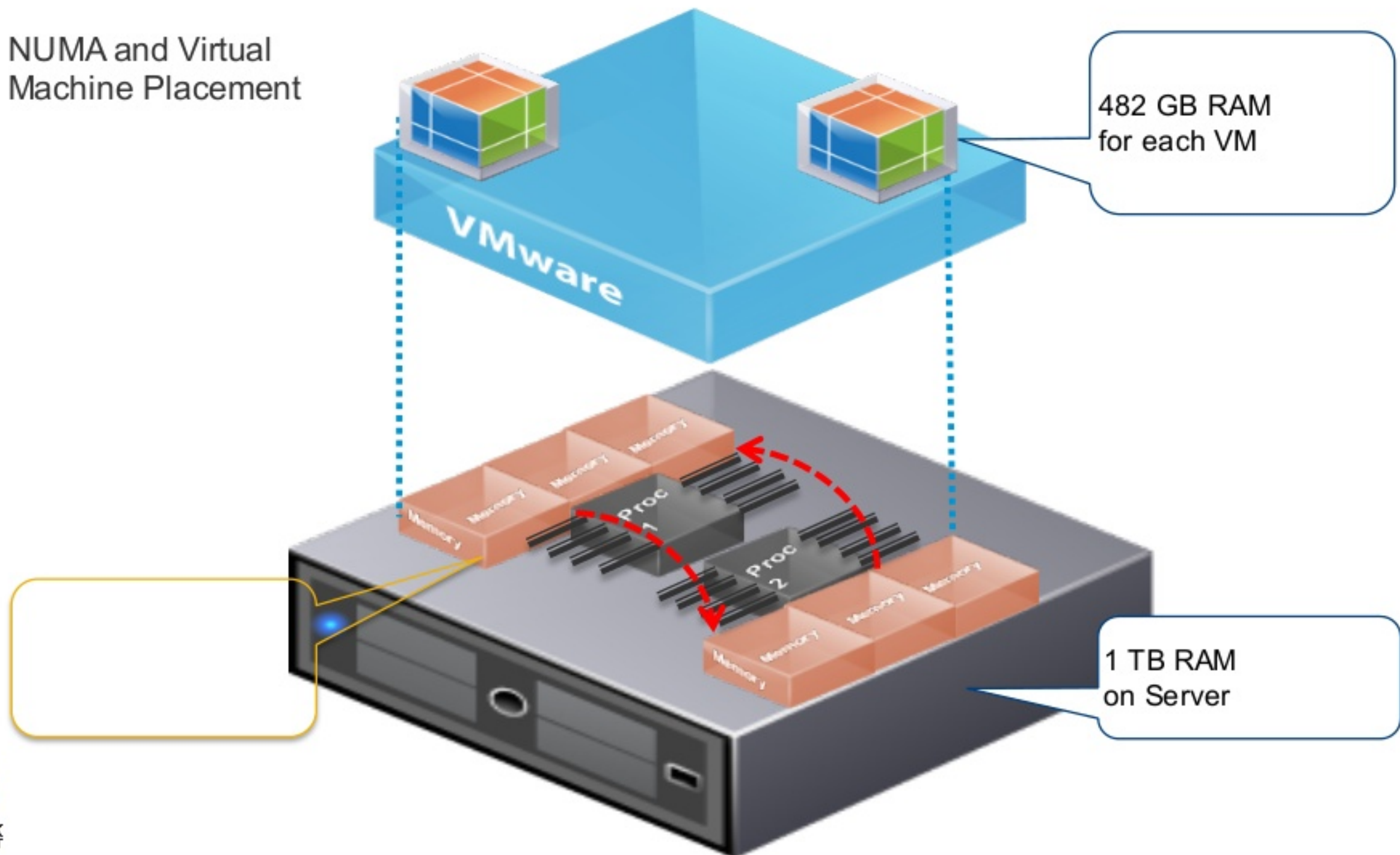
# Results - Spark

- Support Vector Machines workload, which stayed in memory, ran about 10% faster in virtualized form than on bare metal
- Logistic Regression workload, which was written to disk at the larger dataset sizes, showed a slight advantage to bare metal
    - part of the dataset was cached to disk,
    - larger memory of the bare metal Spark executors may help

- Both workloads showed linear scaling from 5 to 10 hosts and as dataset size increased

NUMA and Virtual Machine Placement

482 GB RAM for each VM

1 TB RAM on Server

# Conclusions

- Spark 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
  - Design patterns such as data-compute separation can be used to provide elasticity of your Spark cluster.

# Thank You.

Contact jmurray@vmware.com or
bigdata@vmware.com

# Add Slides as Necessary

- Supporting points go here.