

# Homework 3

for the course Parallel and distributed processing

Aleksandar Ivanovski (Student ID: 186063)

28 November 2020

## Problem 2.1

Differentiate and exemplify the following terms related to clusters:

### **a. Compact versus slack clusters**

The compact cluster has all the nodes closely attached in one or more server racks placed in one facility - server room, and the nodes are not attached to the peripherals, while in a slack cluster the nodes are attached to their usual peripherals, i.e. they are in a form of workstations or PCs, and the nodes can be located in different facilities, buildings or even on remote geographic regions.

### **b. Centralized versus decentralized clusters**

This differentiation of clusters concerns the responsibilities of the clusters. In a centralized cluster all of the nodes are owned, controlled and administered by one single entity - central operator. While in decentralized clusters each node has individual entity which controls and manages them.

### **c. Homogeneous versus heterogeneous clusters**

The primary aspect of this differentiation is the architecture of the clusters. Homogeneous clusters use nodes from the same platform, i.e. same processor architecture, same operating systems and all of the nodes are from the same vendor. In heterogeneous clusters the nodes can be from different platforms.

### **d. Enclosed versus exposed clusters**

This differentiation of clusters concerns about the security. In an exposed cluster, the communication paths among the nodes are exposed to the outside world, i.e. an outside entity can access the communication paths. On the other hand, in an enclosed cluster, the intracluster communication is shielded from the outside world - providing more security.

### **e. Dedicated versus enterprise clusters**

Dedicated clusters are usually compact centralized clusters installed in a server room, while the enterprise cluster is usually used to utilize idle resources in the nodes - they are usually slack heterogeneous clusters, i.e. workstations or PCs.

## Problem 2.2

This problem refers to the redundancy technique. Assume that when a node fails, it takes 10 seconds to diagnose the fault and another 30 seconds for the workload to be switched over.

**a. What is the availability of the cluster if planned downtime is ignored?**

If planned downtime is ignored, the availability depends on the distribution of node failure, so we need more data to determine the availability.

**b. What is the availability of the cluster if the cluster is taken down one hour per week for maintenance, but one node at a time?**

We assume that the cluster has more than one node, so the bottle neck is the time to diagnose and the time to switch the load, so we have:

$$\frac{7days}{7days+40seconds} = \frac{604800}{604840} \approx 0.999933$$

The availability of the cluster is  $\approx 99.99\%$ .

## Problem 2.4

This problem consists of two parts related to cluster computing:

**1. Define and distinguish among the following terms on scalability:**

**a. Scalability over machine size**

Machine size scalability indicates how well the performance will improve with additional processors. Scaling up in resource means gaining higher performance by investing more memory, bigger off-chip caches, bigger disks and so on.

**b. Scalability over problem size**

This indicates how well the system can handle larger problems with larger data size and workload. Apart from depending on machine size, it also depends on memory capacity, and communication capability of the machine.

**c. Resource scalability**

This refers to gaining higher performance or functionality by increasing the machine size (i.e. the number of processors), investing in more storage (cache, main memory, disks), improving the software, etc. Within this dimension, three categories have to be considered. Machine size scalability indicates how well the performance will improve with additional processors. Scaling up in resource means gaining higher performance by investing more memory, bigger off-chip caches, bigger disks and so on. Finally, software scalability indicates how the performance of a system be improved by a newer version of the OS that has more functionalities, a better compiler with more efficient optimizations, more efficient mathematical and engineering libraries, more efficient and easy-to-use applications software and more user-friendly programming environment.

**d. Generation scalability**

This refers to the capability of a system to scale up the performance using the next generation components, such as a faster processor, a faster memory, a newer version of operating system,

a more powerful compiler, etc, with the rest of the system be usable and modifiable as little as possible.

**2. Explain the architectural and functional differences among three availability cluster configurations: hot standby, active takeover, and fault-tolerant clusters. Give two example commercial cluster systems in each availability cluster configuration. Comment on their relative strengths and weaknesses in commercial applications.**

These availability configurations are concerned about fault tolerance, so we have:

A **hot standby cluster**, in which only the primary node is actively doing all the work in normal operational mode. In the same time the standby node is powered on, hence the name hot and running some monitoring programs to communicate heartbeat signals to check the status of the primary node, but is not actively running other useful workloads. The primary node must mirror any data to shared disk storage, which is accessible by the standby node. The standby node requires a second copy of data.

In the case of **active takeover cluster**, the architecture is symmetric among multiple server nodes. Both servers are primary, doing useful work normally. Both failover and failback are often supported on both server nodes. When a node fails, the user applications fail over to the available node in the cluster. Depending on the time required to implement the failover, users may experience some delays or may lose some data that was not saved in the last checkpoint. An example of this cluster architecture is the Oracle's ZFS Storage Appliance.

## Problem 2.7

This problem is related to the use of high-end x86 processors in HPC system construction. Answer the following questions:

**a. Referring to the latest Top 500 list of supercomputing systems, list all systems that have used x86 processors. Identify the processor models and key processor characteristics such as number of cores, clock frequency, and projected performance.**

According to the latest November Top 500 list, I can't find any x86 system.

**b. Some have used GPUs to complement the x86 CPUs. Identify those systems that have procured substantial GPUs. Discuss the roles of GPUs to provide peak or sustained flops per dollar.**

These systems have used GPUs: Summit - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband with 148,600.0 TFlop/s, Sierra - IBM Power System AC922, IBM POWER9 22C 3.1GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband, IBM / NVIDIA with 94,640.0 TFlop/s, Selene - NVIDIA DGX A100, AMD EPYC 7742 64C 2.25GHz, NVIDIA A100, Mellanox HDR Infiniband 63,460.0 TFlop/s, JUWELS Booster Module - Bull Sequana XH2000 , AMD EPYC 7402 24C 2.8GHz, NVIDIA A100, Mellanox HDR InfiniBand/ParTec ParaStation ClusterSuite 44,120.0 TFlop/s, HPC5 - PowerEdge C4140, Xeon Gold 6252 24C 2.1GHz, NVIDIA Tesla V100, Mellanox HDR Infiniband, 35,450.0 TFlop/s.

The GPU has a significant role in the parallelization.

## Problem 2.9

Compare the latest Top 500 list with the Top 500 Green List of HPC systems. Discuss a few top winners and losers in terms of energy efficiency in power and cooling costs. Reveal the greenenergy winners' stories and report their special design features, packaging, cooling, and management policies that make them the winners. How different are the ranking orders in the two lists? Discuss their causes and implications based on publicly reported data.

The winners in terms of energy efficiency are Supercomputer Fugaku - Supercomputer Fugaku, A64FX 48C 2.2GHz, Tofu interconnect D, Summit - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband, Sierra - IBM Power System AC922, IBM POWER9 22C 3.1GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband, Sunway TaihuLight - Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway, Selene - NVIDIA DGX A100, AMD EPYC 7742 64C 2.25GHz, NVIDIA A100, Mellanox HDR Infiniband, Tianhe-2A - TH-IVB-FEP Cluster, Intel Xeon E5-2692v2 12C 2.2GHz, TH Express-2, Matrix-2000, JUWELS Booster Module - Bull Sequana XH2000, AMD EPYC 7402 24C 2.8GHz, NVIDIA A100, Mellanox HDR InfiniBand/ParTec ParaStation ClusterSuite. The number one winner utilizes water cooling to cut off electricity expenses, and uses smart power management to minimize the consumption of idle nodes.

But in terms of performance the most energy efficient system is ranked 170th.

## Problem 2.10

This problem is related to processor selection and system interconnects used in building the top three clustered systems with commercial interconnects in the latest Top 500 list.

**a. Compare the processors used in these clusters and identify their strengths and weaknesses in terms of potential peak floating-point performance.**

The top 3 systems from the list are:

Supercomputer Fugaku: This system is two consecutive times on the first place, it is equipped with 5,087,232 GB of memory, leverages the 7,630,848 cores from its processor - A64FX 48C 2.2GHz.

Summit: This system is manufactured by IBM, it has 2,801,664 GB of memory, utilizes 2,414,592 cores from the IBM POWER9 22C 3.07GHz processor, it can do 200 peta-floating point operations per second.

Sierra is the third ranked system, it is a joint project between IBM, NVIDIA and Mellanox. It utilizes 2,414,592 cores on the IBM POWER9 22C 3.1GHz processor, and it has 1,382,400 GB of memory.

**b. Compare the commercial interconnects of these three clusters. Discuss their potential performance in terms of their topological properties, network latency, bisection bandwidth, and hardware used.**

They use the following interconnects Tofu interconnect D, Dual-rail Mellanox EDR Infiniband (used by Summit and Sierra). Summit and Sierra utilize the same interconnect for both storage and inter-process communications traffic which delivers both 200Gb/s bandwidth between

nodes and in-network computing acceleration for communications frameworks such as MPI and SHMEM/PGAS.

While the Tofu interconnect D, uses high density node configuration which lowers the latency between nodes on the same board significantly. (from 0.91s to 0.54s)

## **Problem 2.16**

Study various SSI features and HA support for clusters in Section 2.3 and answer the following questions, providing reasons for your answers. Identify some example cluster systems that are equipped with these features. Comment on their implementation requirements and discuss the operational obstacles to establish each SSI feature in a cluster system.

### **a. Single entry point in a cluster environment**

This feature concerns the appearance of the system to the end-user. The user can not connect to a single node, but instead connects to the whole system, and the cluster uses its mechanisms to determine which node or nodes are going to be assigned for the given task.

### **b. Single memory space in a cluster system**

This feature is concerned with the node's point of view for the memory. Every node has access to the same memory, so more mechanisms need to be provided to ensure consistency.

### **c. Single file hierarchy in a cluster system**

A system with this feature has all the nodes sharing the same file hierarchy, i.e. the file system and the actual structure of the files. Also this feature requires additional mechanisms for ensuring consistency.

### **d. Single I/O space in a cluster system**

There are systems that utilize this feature to allow all nodes to access the I/O devices of other nodes, i.e. tapes, disks, serial lines, etc... There may be rules and priorities attached to the access privileges. Commercial systems that utilize this feature are: HP NSK Guardian, Inferno, LOCUS, OpenSSI, Plan9, Sprite, and many more.

### **e. Single network space in a cluster system**

The whole system exists in only one LAN network, or maybe WAN depending on the size and the type of the cluster system - centralized or decentralized.

### **f. Single networking in a cluster system**

The nodes are interconnected together using a common infrastructure which by itself is a bottleneck to the whole system.

### **g. Single point of control in a cluster system**

The whole system is being controlled by one central authority, that can be a single node, a group of nodes, or maybe a separate system. Extra measures need to be taken to ensure that the single point of control is always available, because it is the system's bottleneck and if it fails, the whole system is not operational.

### **h. Single job management in a cluster system**

There are nodes which are assigned the job management role, and all of the other nodes rely on their services for the job management. The job management algorithms are usually the same in all of the nodes which have the role of job managers.

**i. Single user interface in a cluster system**

The whole system and all of its separate modules share the same user interface, which improves the end user's experience and ease of use, and also makes the system more maintainable.

**j. Single process space in a cluster system**

Some systems use this SSI feature to provide the illusion that all processes are running on the same machine. All of the process management tools, for example "ps" on the Unix like systems operate on all the processes in the cluster. Commercial systems that utilize this feature are: Amoeba, HP NSK Guardian, Kerrighed, LOCUS, OpenSSI, TidalScale, VMScluster, z/VM, UnixWare NonStop Clusters, and many more.