

Solution of Exercise Sheet 2

Exercise 1 (Cluster Computing)

1. Give a short definition of Cluster Computing.

Clustering is parallel computing on systems with distributed memory.

2. What is a Cluster of Workstations?

The nodes are only available at specific times. During normal working times, the employees use the nodes of such a cluster system as workstations.

3. Name an advantage of Clusters, which follow the Glass-house installation concept.

Some advantages are:

- *Fast access to all components for maintenance and troubleshooting.*
- *Nodes can be connected via high-performance networks.*
- *Increased protection against sabotage.*

4. Name a drawback of Clusters, which follow the Glass-house installation concept.

In case of a power failure or fire in the building, the operation of the entire cluster is at risk.

5. Name an advantage of Clusters, which follow the Campus-wide installation concept.

It is hard to destroy the cluster completely.

6. Name a drawback of Clusters, which follow the Campus-wide installation concept.

Some drawbacks are:

- *It is impossible to use high-performance computer networks.*
- *Often, the nodes contain different hardware components.*

7. How can the availability of a system be calculated?

mean uptime = Mean Time Between Failures

mean downtime = Mean Time To Repair

$$\text{availability} = \frac{\text{mean uptime}}{\text{mean uptime} + \text{mean downtime}}$$

8. What is the objective of High Availability Clustering?

High availability of the cluster service(s).

9. By which approach does High Availability Clustering achieve its objective?

With redundancy of nodes and their components and by avoiding a single point of failure.

10. What is an Active/Active-Cluster?

All nodes run the same services. All nodes are in active state. If nodes fail, the remaining active nodes need to take over their tasks.

11. What is an Active/Passive-Cluster?

During normal operation, at least a single node is in passive state. Nodes in passive state do not provide services during normal operation. If a node fails, a passive node takes over its services.

12. What is the meaning of Failover?

A node takes over the services of a failed node.

13. What is the meaning of Failback?

If failed nodes are operational again, they report their status to the load balancer and get new jobs assigned in the future.

14. What causes a Split Brain?

A split brain may result, when a connection failure between nodes occur, while the computers operate properly. A tool like Heartbeat, which monitors the presence (or disappearance) of nodes, assumes that nodes are broken. Each node declares itself to be the primary node. In a Active/Passive-Cluster, this causes a failure of the cluster (offered services).

15. What is the result of a Split Brain?

If distributed storage is used, write requests cause differing data on the nodes. The data is no longer consistent. It is difficult to fix the broken consistency without losing data.

16. Explain the difference between the Shared Nothing and the Shared Disk Clusters.

In a Shared Nothing cluster, each node has its own storage resource.

In a Shared Disk cluster, all nodes have access to a shared storage.

17. Explain the main difference between a SAN (Storage Area Network) and a NAS (Network Attached Storage).

A SAN provides block-level access to storage devices via the network.

A NAS provides file system-level access to storage devices via the network.

18. What is the objective of High Performance Clustering?

High computing power

19. Name an advantage of High Performance Clusters versus supercomputers.

Some examples for advantages of High Performance Clusters are:

- Low price and vendor independence.*
- Defective components can be obtained in a quick and inexpensive way.*
- It is easy to increase the performance in a short time via additional nodes.*

20. Name a drawback of High Performance Clusters versus supercomputers.

High administrative and maintenance costs, compared with mainframes.

21. What is a Beowulf Cluster?

It is a High Performance Cluster and the nodes use a free operating system. Beowulf clusters consist of commodity PCs or workstations. The nodes of a Beowulf cluster are used only for the cluster

22. What is a Wulfpack Cluster?

It is a High Performance Cluster and the nodes use a Windows operating system.

23. What is the objective of High Throughput Clustering?

Maximize the throughput (handling of small tasks/requests).

24. Can High Throughput Clusters be used to process the same tasks as High Performance Clusters?

No. Such clusters are not used for extensive calculations. Tasks for High Throughput Clusters must not be split into sub-tasks. The individual tasks (requests) are small and a single PC could handle them. Typical fields of application of High Throughput Clustering web servers and internet search engines.

If large compute jobs need to be processed, a High Performance Cluster is required.

If multiple small compute jobs (in a short time) need to be processed, a High Throughput Cluster is required.

Exercise 2 (MPI Cluster)

Build up a MPI Cluster with at least 2 virtual nodes (e.g. with VirtualBox or any other virtualization software) or alternatively with at least 2 physical nodes.

Write down precise instructions with the steps you performed to build up the MPI cluster.

Execute a simple program to demonstrate the functioning of your cluster.

This example works for a cluster of eight physical Raspberry Pi nodes with the operating system Raspbian, based on Debian Wheezy. First, we need to install some packages.

```
1 $ sudo apt-get update && sudo apt-get -y install make gcc g++ openmpi-bin  
   openmpi-common libopenmpi-dev
```

Next, the SSH keys need to be distributed among the nodes. It is required that the master node (which is used to execute the MPI jobs) can execute commands on all worker nodes (the master is in our example a worker not too) without the need to type in a password.

The nodes have in our example the IP addresses 10.0.0.110 - 10.0.0.117

```
1 $ ssh-copy-id -i ~/.ssh/id_rsa.pub pi@10.0.0.110  
2 $ ssh-copy-id -i ~/.ssh/id_rsa.pub pi@10.0.0.111  
3 $ ssh-copy-id -i ~/.ssh/id_rsa.pub pi@10.0.0.112  
4 $ ssh-copy-id -i ~/.ssh/id_rsa.pub pi@10.0.0.113  
5 $ ssh-copy-id -i ~/.ssh/id_rsa.pub pi@10.0.0.114  
6 $ ssh-copy-id -i ~/.ssh/id_rsa.pub pi@10.0.0.115  
7 $ ssh-copy-id -i ~/.ssh/id_rsa.pub pi@10.0.0.116  
8 $ ssh-copy-id -i ~/.ssh/id_rsa.pub pi@10.0.0.117
```

Insert the hostnames (pi110 - pi117) and the IP addresses into the file /etc/hosts

```
1 10.0.0.110    pi110  
2 10.0.0.111    pi111  
3 10.0.0.112    pi112  
4 10.0.0.113    pi113  
5 10.0.0.114    pi114  
6 10.0.0.115    pi115  
7 10.0.0.116    pi116
```

8 10.0.0.117 pi117

It is helpful to have a distributed file system which can be accessed by all nodes of the cluster. This can be used to store the executable and the hosts file. Otherwise it must be distributed manually to all nodes of the cluster. In this example, a GlusterFS distributed file system is mounted on all nodes inside the folder `\glusterfs`

```
1 sudo mkdir /glusterfs/MPI_examples
2 sudo chmod 777 /glusterfs/MPI_examples/
```

Insert the following source code into the file `/glusterfs/MPI_examples/hello_world.c`

```
1 #include <mpi.h>
2 #include <stdio.h>
3 #include <unistd.h>
4
5 int main(int argc, char *argv[]) {
6     int size, rank, namelen;
7     char processor_name[MPI_MAX_PROCESSOR_NAME];
8
9     // Initialize the MPI environment
10    MPI_Init(&argc, &argv);
11
12    // How many CPU (cores) contains the MPI environment?
13    MPI_Comm_size(MPI_COMM_WORLD, &size);
14
15    // What is the number of our CPU core?
16    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
17
18    // What is the name of the CPU core?
19    MPI_Get_processor_name(processor_name, &namelen);
20
21    // Output of each process
22    printf("I am process nr. %d of %d on %s\n", rank, size, processor_name);
23
24    // Terminate the MPI environment
25    MPI_Finalize();
26
27    // Terminate the programm with exit code 0 (EXIT_SUCCESS)
28    return 0;
29 }
```

Create a file `hosts.mpi` with contains the hostnames of all hosts, which shall be used for the MPI environment.

```
1 $ for i in 0 1 2 3 4 5 6 7; do echo pi11$i slots=4 >> /glusterfs/
   MPI_examples/hosts_4cores.mpi; done
2 $ cat /glusterfs/MPI_examples/hosts_4cores.mpi
3 pi110 slots=4
```

```
4 pi111 slots=4
5 pi112 slots=4
6 pi113 slots=4
7 pi114 slots=4
8 pi115 slots=4
9 pi116 slots=4
10 pi117 slots=4
```

Compile the MPI application:

```
1 $ mpicc /glusterfs/MPI_examples/hello_world.c -o /glusterfs/MPI_examples/
  hello_world
```

Execute the MPI application:

```
1 $ mpirun -np 32 --hostfile /glusterfs/MPI_examples/hosts_4cores.mpi /
  glusterfs/MPI_examples/hello_world
2 I am process nr. 0 of 32 on pi110
3 I am process nr. 2 of 32 on pi110
4 I am process nr. 3 of 32 on pi110
5 I am process nr. 7 of 32 on pi111
6 I am process nr. 6 of 32 on pi111
7 I am process nr. 5 of 32 on pi111
8 I am process nr. 27 of 32 on pi116
9 I am process nr. 17 of 32 on pi114
10 ...
```

Exercise 3 (Approximate π via Monte Carlo)

π can be approximated via Monte Carlo simulation.

$$A_S = (2r)^2 = 4r^2$$

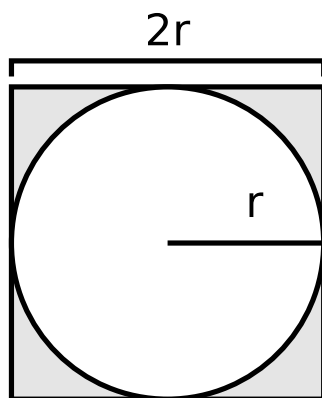
Approach: Inscribe a circle of radius r inside a square with side length $2r$.

$$A_C = \pi r^2 \implies \pi = \frac{A_C}{r^2}$$

Generate random dots in the square. The number of dots in A_C in relation to the number of dots in A_S is equal to the surface ratio.

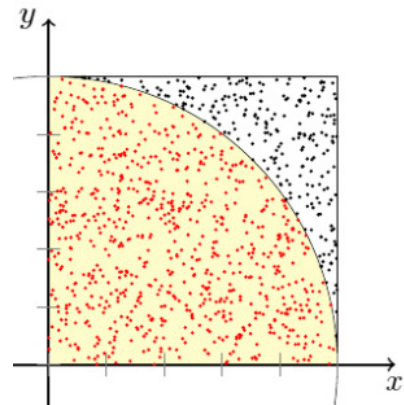
$$\frac{A_C}{A_S} = \frac{\pi r^2}{4r^2} \implies \frac{A_C}{A_S} = \frac{\pi}{4}$$

The dots can be generated in parallel by the workers. The master receives the dots and calculates π .



A = Surface ratio
 r = Radius
 C = Circle
 S = Square

Image source: Wikipedia



Develop a MPI application, which calculates π via Monte Carlo simulation.

Test your MPI application with with the MPI cluster you already created and with different numbers or worker nodes to discover if your application scales well with a growing number of worker nodes.

It is helpful to have a distribute file system which can be accessed by all nodes of the cluster. This can be used to store the executable and the hosts file. Otherwise it must be distributed manually to all nodes of the cluster.

Insert the following source code into the file `/glusterfs/MPI_examples/pi_mpi2.c`

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include "mpi.h"
4
5 int main(int argc, char *argv[]) {
6
7     int myid,nprocs;
8
9     long int npts = 1e8;
10
11     long int i,mynpts;
12
13     double f,sum,mysum;
14     double xmin,xmax,x;
15
16     MPI_Init(&argc,&argv);
17     MPI_Comm_size(MPI_COMM_WORLD,&nprocs);
18     MPI_Comm_rank(MPI_COMM_WORLD,&myid);
19
20     if (myid == 0) {
21         mynpts = npts - (nprocs-1)*(npts/nprocs);
22     } else {
23         mynpts = npts/nprocs;
24     }
25
26     mysum = 0.0;
```

```
27 xmin = 0.0;
28 xmax = 1.0;
29
30 srand(time(0));
31
32 for (i=0; i<my_npts; i++) {
33     x = (double) rand()/RAND_MAX*(xmax-xmin) + xmin;
34     mysum += 4.0/(1.0 + x*x);
35 }
36
37 MPI_Reduce(&mysum,&sum,1,MPI_DOUBLE,MPI_SUM,0,MPI_COMM_WORLD);
38
39 if (myid == 0) {
40     f = sum/npts;
41     printf("PI calculated with %ld points = %f \n",npts,f);
42 }
43
44 MPI_Finalize();
45 }
```

Compile the MPI application:

```
1 $ mpicc glusterfs/MPI_examples/pi_mpi2.c -o /glusterfs/MPI_examples/
   pi_mpi2
```

Switch to the appropriate directory:

```
1 $ cd /glusterfs/MPI_examples/
2 $ pwd
3 /glusterfs/MPI_examples/
```

In this example, the MPI application was executed on a cluster with 128 CPU cores (4 cores per physical node) with different numbers of CPU cores:

```
1 $ time mpirun -np 128 --hostfile hosts_4cores.mpi pi_mpi2
2 PI calculated with 1000000000 points = 3.141601
3
4 real 0m6.276s
5 user 0m14.050s
6 sys 0m1.520s
7
8 $ time mpirun -np 64 --hostfile hosts_4cores.mpi pi_mpi2
9 PI calculated with 1000000000 points = 3.141592
10
11 real 0m7.088s
12 user 0m20.120s
13 sys 0m1.080s
14
15 $ time mpirun -np 32 --hostfile hosts_4cores.mpi pi_mpi2
16 PI calculated with 1000000000 points = 3.141580
17
```

```
18 real 0m10.928s
19 user 0m35.800s
20 sys 0m0.850s
21
22 $ time mpirun -np 16 --hostfile hosts_4cores.mpi pi_mpi2
23 PI calculated with 1000000000 points = 3.141604
24
25 real 0m19.280s
26 user 1m9.210s
27 sys 0m0.670s
28
29 $ time mpirun -np 8 --hostfile hosts_4cores.mpi pi_mpi2
30 PI calculated with 1000000000 points = 3.141601
31
32 real 0m36.172s
33 user 2m17.150s
34 sys 0m0.500s
35
36 $ time mpirun -np 4 --hostfile hosts_4cores.mpi pi_mpi2
37 PI calculated with 1000000000 points = 3.141612
38
39 real 1m9.820s
40 user 4m33.240s
41 sys 0m0.420s
```

Explanation of the output of the time command.

Real is wall clock time - time from start to finish of the call. This is all elapsed time including time slices used by other processes and time the process spends blocked (for example if it is waiting for I/O to complete).

User is the amount of CPU time spent in user-mode code (outside the kernel) within the process. This is only actual CPU time used in executing the process. Other processes and time the process spends blocked do not count towards this figure.

Sys is the amount of CPU time spent in the kernel within the process. This means executing CPU time spent in system calls within the kernel, as opposed to library code, which is still running in user-space. Like 'user', this is only CPU time used by the process. See below for a brief description of kernel mode (also known as 'supervisor' mode) and the system call mechanism.

Real represents the actual elapsed time, while user and sys values represent CPU execution time. As a result, on a multicore system, the user and/or sys time (as well as their sum) can actually exceed the real time.

Source of the explanation: <http://stackoverflow.com/questions/556405/what-do-real-user-and-sys-mean-in-the-output-of-time1>