

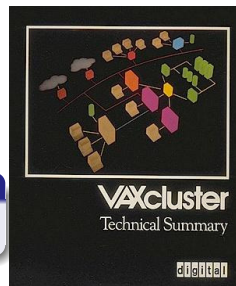


# Agenda for Today

- Cluster computing
  - History of cluster computing
  - Distinguishing criteria
    - Structure (homogeneous, heterogeneous)
    - Installation concepts (Glass-house, Campus-wide)
    - Fields of application
    - High Availability Clustering
    - High Performance Clustering
    - High Throughput Clustering
    - Behaviour in the event of failed nodes (Active/Passive, Active/Active)
  - Current situation
  - Advantages and drawbacks of clusters
  - Cluster application libraries (PVM, MPI)
  - Gearman

# History of Cluster computing

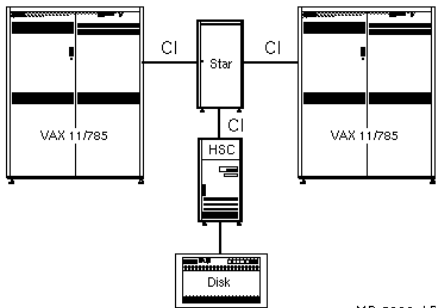
- 1983: Digital Equipment Corporation (DEC) offers for its VAX-11 system a cluster solution (VAXcluster)
  - VAXcluster allows to connect multiple computers via a serial link
  - By combining multiple VAX systems, their computing power and memory can be accessed equal to a single computer system
- 1987 DEC sells VAX 8974 and VAX 8978
  - These are clusters, which contain 4 or 8 nodes (VAX 8700 systems) and a MicroVAX II, which is used as console



## Further information

VAXcluster system. Digital Technical Journal. Number 5. September 1987  
[http://www.dtjcd.vmsresource.org.uk/pdfs/dtj\\_v01-05\\_sep1987.pdf](http://www.dtjcd.vmsresource.org.uk/pdfs/dtj_v01-05_sep1987.pdf)

# VAXcluster



CI = ComputerInterconnect  
 HSC = Hierarchical Storage Controller  
 Star = Star Coupler

MR-5062-AD

**Image sources:**

- <http://hampage.hu/oldiron/vaxen/eikvms1.jpg>
- [http://odl.sysworks.biz/disk\\$vxdocmay941/decw\\$book/d3ywaa51.p37.decw\\$book](http://odl.sysworks.biz/disk$vxdocmay941/decw$book/d3ywaa51.p37.decw$book)
- <http://www.computerhistory.org/collections/catalog/102635385>



# Definition of Cluster Computing

## Cluster computing

Clustering is parallel computing on systems with distributed memory

- A cluster consists of at least 2 nodes
  - Each node is an independent computer system
  - The nodes are connected via a computer network
    - In clusters with just a few nodes, inexpensive computer network technologies (Fast or Giga-Ethernet) are used
    - Clusters with several hundred nodes require high-speed computer networks (e.g. InfiniBand)
  - Often, the nodes are under the control of a master and are attached to a shared storage
  - Nodes can be ordinary PCs, containing commodity hardware, workstations, servers or supercomputers

From the user perspective (in a perfect world)...

- the cluster works like a single system  $\implies$  a virtual uniprocessor system
- Ideally, the users don't know, that they work with a cluster system





# Homogeneous and Heterogeneous Clusters

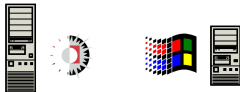
- The structure of clusters can be homogeneous and heterogeneous

## Heterogeneous structure

## Homogeneous structure



*I have never seen a heterogeneous cluster with different operating systems in practice...*



- In practice, the construction of a heterogeneous cluster is generally a bad idea
- The administration of homogeneous clusters is challenging, but the administration of heterogeneous clusters is hell (especially when commodity hardware is used)

# Installation Concepts of Clusters (1/2)



- **Glass-house**

- The cluster is located in a single room or server rack

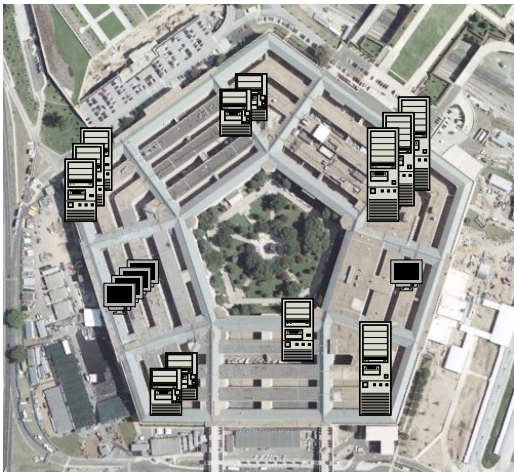
- **Advantages:**

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

- **Drawbacks:**

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

# Installation Concepts of Clusters (2/2)



- **Campus-wide**
  - The nodes are located in multiple buildings and spread across the site of the research center or company
- **Advantages:**
  - It is hard to destroy the cluster completely
- **Drawbacks:**
  - It is impossible to use high-performance computer networks
  - Often, the nodes contain different hardware components

# Fields of Application of Clusters

- Clusters for different applications exist
- ① **High Availability Clustering**
  - Objective: high availability
- ② **High Performance Clustering**
  - Objective: high computing power
- ③ **High Throughput Clustering**
  - Objective: high throughput







# HA-Clustering – Active/Passive and Active/Active

- **Active/Passive-Cluster** (also called: **Hot-Standby-Clusters**)
  - 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
  - **Failover** = a node takes over the services of a failed node
  - Benefit: The services must not be designed for cluster operation
  - Drawback: Much potentially available performance remains unused in normal operation
- **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
  - Advantage: Better distribution of load between nodes
  - Drawback: Services need to be designed for cluster operation, because all nodes access shared resources (data!) simultaneously

# High Availability Clustering – Failover and Failback

- **Failover:** Ability to automatically transfer the tasks of a failed node to another node for minimizing the downtime
  - The failover functionality is usually provided by the operating system
  - Example: Heartbeat for Linux

<http://www.linux-ha.org/wiki/Heartbeat>

*Heartbeat is a daemon that provides cluster infrastructure (communication and membership) services to its clients. This allows clients to know about the presence (or disappearance!) of peer processes on other machines and to easily exchange messages with them. In order to be useful to users, the Heartbeat daemon needs to be combined with a cluster resource manager (CRM) which has the task of starting and stopping the services (IP addresses, web servers, etc.) that cluster will make highly available. Pacemaker is the preferred cluster resource manager for clusters based on Heartbeat.*

- **Failback:** If failed nodes are operational again, they report their status to the load balancer and get new jobs assigned in the future
  - From that point in time, the cluster again has the same performance capability, it had before the failure of the nodes

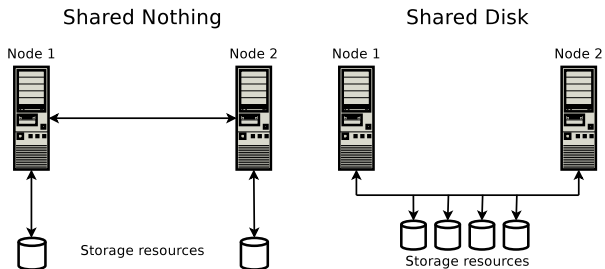
Well written articles about Heartbeat and DRBD

Andreas Sebald. *Linux-HA-Cluster mit Heartbeat und DRBD*. Linux-Magazin 07/2004.  
<http://www.linux-magazin.de/Ausgaben/2004/07/Reservespieler>



# Architectures of High Availability Clustering

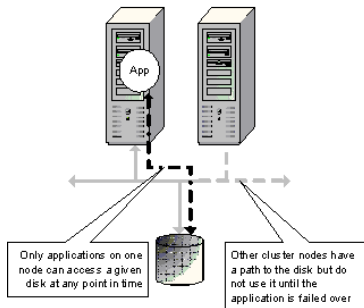
- 2 architectures of High Availability Clustering exist:
  - **Shared Nothing**  $\implies$  Distributed Storage
  - **Shared Disk**  $\implies$  Shared Storage



# Shared Nothing Architecture

Image Source: technet.microsoft.com

- In a **Shared Nothing** cluster, each node has its own storage resource
- Even, when a resource is physically connected to multiple nodes, only a single node is allowed to access it
  - Only if a node fails, the resource is acquired by another node
- Advantage: No lock management is required
  - No protocol overhead reduces the performance
  - In theory, the cluster can scale almost in a linear way
- Drawback: Higher financial effort for storage resources, because the data can not be distributed in an optimal way



# Shared Nothing with DRBD (1/3)

- Distributed Replicated Block Device (DRBD)
  - Free software to build up a network storage for Shared Nothing clusters, without an expensive Storage Area Network (SAN)
- Shared storage is always a single point of failure, since only the cluster nodes are designed in a redundant way
  - Redundant SAN solutions are expensive (> 100.000 €)

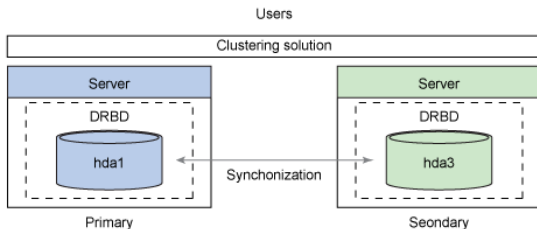


Image Source: M. Jones, <https://www.ibm.com/developerworks/library/1-drbd/index.html>

Well written articles about DRBD

iX 3/2010. Florian Haas. *Hochverfügbare Shared Nothing Cluster mit DRBD*. P.120-123  
M. Tim Jones. High availability with the Distributed Replicated Block Device. 2010.  
<https://www.ibm.com/developerworks/library/1-drbd/>

# Shared Nothing with DRBD (2/3)

- Functioning:
  - A primary server and a secondary server exist
    - Write requests are carried out by the primary server and afterwards are sent to the secondary server
    - Only if the secondary server reports the successful write operation to the primary server, the primary server reports the end of the successful write operation
  - Practically, it implements RAID 1 via TCP
  - Primary server fails  $\implies$  secondary server becomes primary server
    - If a failed system is operational again, only the data blocks, which have changed during the outage are resynchronized
  - Read access is always carried out locally ( $\implies$  better performance)

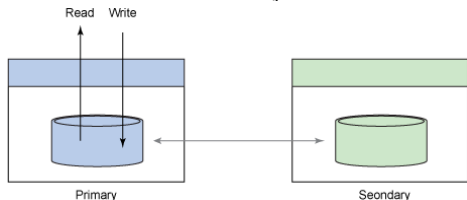


Image Source: M. Jones, <https://www.ibm.com/developerworks/library/1-drbd/index.html>

# Shared Nothing with DRBD (3/3)

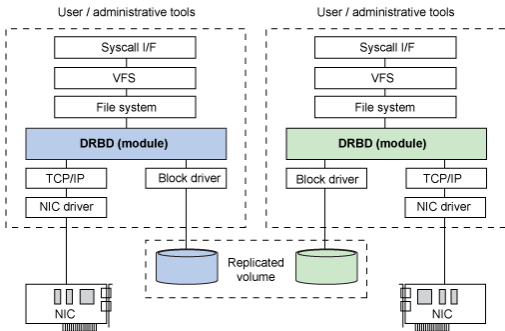


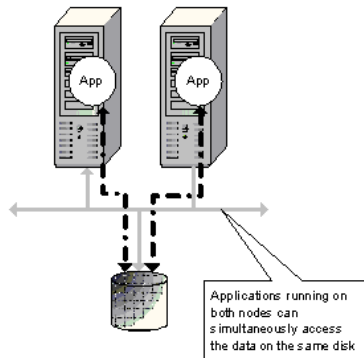
Image Source: M. Jones, <https://www.ibm.com/developerworks/library/l-drbd/index.html>

- DRBD is a part of the Linux kernel since version 2.6.33 (February 2010)
- Because DRBD operates inside the Linux kernel at block level, the system is transparent for the layers above it
- DRBD can be used as a basis for:
  - Conventional file systems, such as ext3/4 or ReiserFS
  - Shared-storage file systems, such as Oracle Cluster File System (OCFS2) and Global File System (GFS2)
    - If shared-storage file systems are used, all nodes must have direct I/O access to the device
  - Another logical block device, such as the Logical Volume Manager (LVM)

# Shared Disk Architecture

Image Source: technet.microsoft.com

- In a **Shared Disk** cluster, all nodes have access to a shared storage
- Several possible ways exist to connect the nodes to the storage:
- **SAN** (Storage Area Network) via Fibre Channel
  - Expensive, but provides high performance
  - Provides block-level access to storage devices via the network.
- **NAS** (Network Attached Storage)
  - Easy-to-use file server
  - Provides file system-level access to storage devices via the network
  - Can also be implemented as a pure software solution
    - Examples: FreeNAS and Openfiler
- **iSCSI** (Internet Small Computer System Interface)
  - SCSI protocol via TCP/IP
  - SAN-like access via the IP-network



# High Performance Clustering (1/2)

- Objective: High computing power
  - Also called: **Clustering for Scalability**
- High Performance Clusters provide the performance of mainframe computers for a much lower price
- These clusters are usually made of commodity PCs or workstations
- Typical application area:
  - Applications, which implement the Divide and Conquer principle
    - Such applications split big tasks into multiple sub-tasks, evaluates them and puts together the sub-task results to the final result
  - Applications, used for analyzing large amounts of data

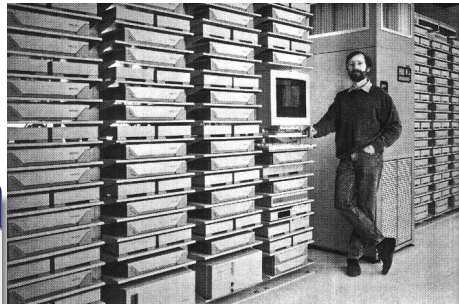
# High Performance Clustering (2/2)

Image Source: Reddit

Application examples: Crash test simulation, weather forecast, optimization of components, Monte Carlo simulation, flight path calculation, data mining, flow simulation, strength analysis, rendering of movies or clips, simulation of the night sky, variant calculating for chess, prime number computation, . . .

In 1995 Pixar rendered Toy Story on a 294 × 100MHz CPU Sun SPARCstation 20 cluster

Each SPARCstation 20 (single processor) had SunOS 5.4 installed and a HyperSPARC 100 MHz with 27.5066 MFLOPS  
⇒ The theoretical maximum performance of the setup was 294 \* 27.5066 = 8086.94 MFLOPS



## • Advantages:

- 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

## • Drawback:

- High administrative and maintenance costs, compared with mainframes

# High Performance Clustering – Beowulf Cluster

- If a free operating system is used  $\implies$  **Beowulf** cluster
- If a Windows operating system is used  $\implies$  **Wulfpack**
- A Beowulf cluster is never a cluster of workstations (COW)
  - Beowulf clusters consist of commodity PCs or workstations, but the nodes of a Beowulf cluster are used only for the cluster
- The cluster is controlled via a master node
  - The master distributes (schedules) jobs and monitors the worker nodes
- Worker nodes are only accessible via the network connection
  - They are not equipped with I/O devices like screens or keyboards
- Worker nodes contain commodity PC components and are not redundant ( $\implies$  designed for high availability)
  - A potential issue is the failure of the cooling of the system components
  - Fans in nodes and power supplies have a limited lifetime and fail without any warning
  - Modern CPUs cannot operate without adequate cooling



# Stone SouperComputer (2/2)



- Built in 1997
- Mostly 486DX-2/66 Intel CPUs
- Some Pentiums
- 10 Mbit/s Ethernet
- RedHat Linux, MPI and PVM
- Extremely heterogeneous structure
- No purchase costs
- High setup and administration effort
- Everything handmade

Image source: <http://www.climate modeling.org/~forrest/linux-magazine-1999/>

## Later Generations of Beowulf Clusters (1/2)



Image source: <http://archiv.tu-chemnitz.de/pub/2000/0089/data/clic.html>

- Vendors, such as Megaware in Chemnitz, sell complete Beowulf clusters
- Image: Chemnitzer Linux Cluster (CLIC) from 2000

# Later Generations of Beowulf Clusters (2/2)



Image source: <http://tina.nat.uni-magdeburg.de>

- Tina (Tina is no acronym) in Magdeburg from 2001

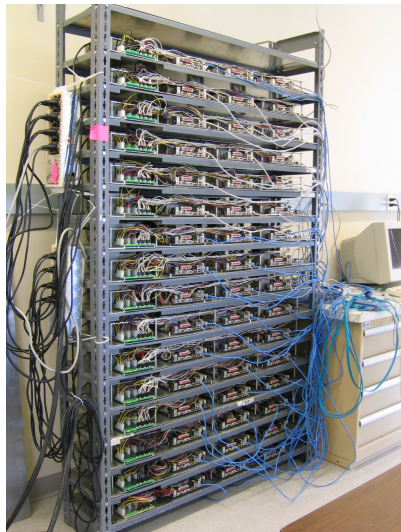
# State of the Art of Cluster Computing



Image source (right image):  
<http://physics.bu.edu/~sandvik/clusters.html>

*A Cluster of Motherboards*

The cluster in the right image has 48 nodes



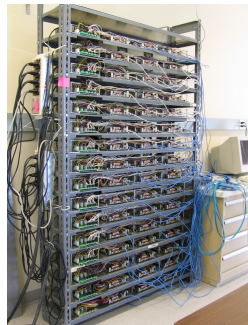
# High Throughput Clustering

- Objective: Maximize throughput
- Such clusters consist of servers, which are used to process incoming requests
- Such clusters are not used for extensive calculations
  - Tasks 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
  - Internet search engines
- Large compute jobs  $\implies$  High Performance Cluster
- Multiple small compute jobs (in a short time)  $\implies$  High Throughput Cluster

# Today: Clusters at Universities



<http://cs.boisestate.edu/~amit/research/beowulf/>



<http://physics.bu.edu/~sandvik/clusters.html>

- Beowulf clusters, built up from commodity hardware  
⇒ low acquisition cost
- High effort for administration (handmade)  
⇒ irrelevant, because students do the administration

# Today: Research and Industry (Example: HP C7000)



Image source: <http://imagehost.vendio.com/bin/imageserver.x/00000000/pdneiman/DSC04040.JPG>

- Compact blade servers or so-called *pizza boxes*
- Professional management tools (like HPE iLO) and redundant components simplify the administration

# Calculation Example about the Possible Packing Density

- A 19 inch rack contains up to 4 blade enclosures (BladeCenters)
- A HP C7000 BladeCenter provides 16 blade slots
- Blades exist, which contain 2 independent servers
  - e.g. HP Blade ProLiant BL2x220c G5
  - 2 servers per blade. Completely independent computers
  - Each server contains: 2x Intel Quad Core Xeon (2,33 GHz) and 16 GB RAM

⇒ 8 cores per server

⇒ 16 cores per blade

⇒ 256 cores per blade enclosure (BladeCenter)

⇒ 1024 cores per 19 inch rack

## The packing density increases

Intel Xeon processors with 6 cores (*Dunnington*), with 8 cores (*Nehalem-EX*), with 18 cores (*Haswell-EX*) and with 22 cores (*Broadwell*) are already available. AMD offers the Opteron (*Magny-Cours*) with 12 cores and the Ryzen (*Threadripper*) with 32 cores



# Advantages and Drawbacks of Clusters

## Advantages:

- Flexibility and extensibility
  - The number of nodes of a cluster can be dynamically increased or decreased according to the needed resources
- Lower purchase price compared with supercomputers
- Simple replacement of commodity hardware components

## Drawbacks:

- Errors occur more often compared with a single supercomputer
- Clusters consist of many independent systems  
⇒ higher administrative costs and personnel expenses compared with a single or few supercomputers
- High effort for distributing and controlling applications
  - If the number of nodes is increased, the effort increases too





# Libraries for Cluster Applications (MPI)

- **Message Passing Interface (MPI)**
  - Development started in 1993-94
  - Collection of functions (e.g. for process communication) to simplify the development of applications for parallel computers
  - The library can be used with C/C++, Fortran 77/90 and Python
    - MPI is not a programming language!
  - Contains no daemon
  - Implements message-based communication (message passing)
  - Especially suited for homogeneous environments
  - Focus: Performance and security
  - MPI implements > 100 functions and several constants
  - Implementations: LAM/MPI (obsolete), OpenMPI, MPICH2,...

MPI tutorial of Wes Kendall, Dwaraka Nath, Wesley Bland: <https://mpitutorial.com/tutorials/>  
 MPI tutorial of Stefan Schaefer and Holger Blaar: <http://www2.informatik.uni-halle.de/lehre/mpi-tutorial/index.htm>

# MPI Functions – Selection of important Functions (1/5)

- `MPI_Init(&argc, &argv)`
  - Initialization routine  $\implies$  starts the MPI environment
  - Defines the communicator `MPI_COMM_WORLD`
    - A communicator contains a group of processes and a communication context
    - `MPI_COMM_WORLD` contains all processes
  - The arguments `argc` and `argv` are pointers to the parameters of the main function `main`
    - The `main` function always receives 2 parameters from the operating system
    - `argc` (**argument count**) contains the number of parameters passed
    - `argv[]` (**argument values**) contains the parameters itself
    - The names of the variables can be freely selected, but they are usually named `argc` and `argv`
    - Not command-line parameters passed  $\implies$  `argc = 1`

<https://web.archive.org/web/20070909174518/http://www2.informatik.uni-jena.de/cmc/racluster/mpi-leitfaden>

# MPI Functions – Selection of important Functions (2/5)

- `MPI_Comm_Size(MPI_Comm comm, int size)`
  - Determines the number of processes in a communicator
  - `size` is the output

```
1 #include "mpi.h"
2
3 int      size;
4 MPI_Comm comm;
5 ...
6 MPI_Comm_size(comm, &size);
7 ...
```

# MPI Functions – Selection of important Functions (3/5)

- `MPI_Comm_Rank(MPI_Comm comm, int rank)`
  - Determines the rank (identification number) of the calling process in the communicator
  - rank is the output
  - The rank is used by MPI for process identification
  - The rank number is unique within a communicator
  - Processes are numbered sequentially, starting from zero

```
1 #include "mpi.h"
2
3 int rank;
4 MPI_Comm comm;
5
6 ...
7 MPI_Comm_rank(comm, &rank);
8 if (rank==0) {
9     ... code for process 0 ...
10 }
11 else {
12     ... code for the other processes ...
13 }
```

<https://web.archive.org/web/20070909174518/http://www2.informatik.uni-jena.de/cmc/racluster/mpi-leitfaden>

# MPI Functions – Selection of important Functions (4/5)

- `MPI_Get_processor_name(char *name, int *resultlen)`
  - Determines the name of the processor
  - `name` is the output
  - The length (number of characters) of the name is returned in `resultlen`
  - The name identifies the hardware, where MPI runs
    - The exact output format is implementation-dependent and may by equal with the output of `gethostname`

```
1 #include "mpi.h"
2 int MPI_Get_processor_name(
3     char *name,
4     int *resultlen)
```







# Simple MPI Example (3/3)

- Compile program:

```
$ mpicc hello_world.c -o hello_world
```

- Distribute the program in the cluster:
  - The program must be stored on each node in the same directory!

```
$ scp hello_world node1:~  
$ scp hello_world node2:~
```

- Program execution (6 processes) in the cluster:

```
$ mpirun -np 6 --hostfile hosts.mpi hello_world  
Ich bin Prozess Nr. 0 von 6 auf domU-12-31-38-00-20-38  
Ich bin Prozess Nr. 1 von 6 auf ip-10-126-43-6  
Ich bin Prozess Nr. 2 von 6 auf domU-12-31-38-00-AD-95  
Ich bin Prozess Nr. 4 von 6 auf ip-10-126-43-6  
Ich bin Prozess Nr. 3 von 6 auf domU-12-31-38-00-20-38  
Ich bin Prozess Nr. 5 von 6 auf domU-12-31-38-00-AD-95
```

- The CPUs respond in random order
  - What is the reason?

# MPI Functions – Send-/Receive (1/3)

- `MPI_Send(int buffer, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm)`
  - Sends a message (blocking) to another process in the communicator
    - `buffer` = first address of the transmit buffer
    - `count` = number of elements in the transmit buffer (not negative)
    - `datatype` = MPI data type of the elements in the transmit buffer
    - `dest` = rank of the receiver process in the communicator
    - `tag` = ID for distinguishing the messages
    - `comm` = communicator
  - All parameters are input parameters
  - The function sends `count` data objects of type `datatype` from address `buffer` (⇒ transmit buffer) with the ID `tag` to the process with rank `dest` in communicator `comm`

# MPI Data Types

MPI data type	C data type	Used for...	Size	Value range
MPI_CHAR	signed char	Chars	1 Byte	-127 ... +127
MPI_UNSIGNED_CHAR	unsigned char	Chars	1 Byte	0 ... 255
MPI_SHORT	signed short int	Integers	2 Bytes	-32,768 ... 32,767
MPI_UNSIGNED_SHORT	unsigned short int	Integers	2 Bytes	0 ... 65,535
MPI_INT	signed int	Integers	2-8 Bytes	Depends on the architecture
MPI_UNSIGNED	unsigned int	Integers	2-8 Bytes	Depends on the architecture
MPI_LONG	signed long int	Integers	4 Bytes	-2.147.483.648 ... 2.147.483.647
MPI_UNSIGNED_LONG	unsigned long int	Integers	4 Bytes	0 ... 4.294.967.295
MPI_FLOAT	float	Floating point numbers	4 Bytes	Single precision
MPI_DOUBLE	double	Floating point numbers	8 Bytes	Double precision
MPI_LONG_DOUBLE	long double	Floating point numbers	16 Bytes	Quadruple precision
MPI_BYTE	—	Floating point numbers	1 Byte	0 ... 255

- The integer value range depends on the used C compiler used and architecture (2, 4 or 8 Bytes)

## MPI Functions – Send-/Receive (2/3)

- `MPI_Recv(int buffer, int count, MPI_Datatype datatype, int source, int tag, MPI_Comm comm, MPI_Status status)`
  - Receive a message (blocking)
    - `buffer` = first address of the receive buffer  $\Leftarrow$  *output parameter*
    - `count` = number of elements in the receive buffer (not negative)
    - `datatype` = MPI data type of the elements in the receive buffer
    - `source` = rank of the sender process in the communicator or `MPI_ANY_SOURCE`
    - `tag` = ID for distinguishing the messages. For receiving messages with any identifier, the constant `MPI_ANY_TAG` is used
    - `comm` = communicator
    - `status` = contains the rank of the sender process `source` and the message identifier `tag`  $\Leftarrow$  *output parameter*

# MPI Functions – Send-/Receive (3/3)

- `MPI_Get_count(status, datatype, count)`
  - Determines the number of received elements
    - `count` = number of received elements (not negative)  $\Leftarrow$  *output parameter*
    - `status` = status upon the return of the receive operation
    - `datatype` = MPI data type of the elements in the receive buffer

```
1 #include "mpi.h"
2 #define MAXBUF 1024
3
4     int         i, count;
5     void        *recvbuf;
6     MPI_Status  status;
7     MPI_Comm    comm;
8     MPI_Datatype datatype;
9
10    ...
11    MPI_Recv(recvbuf, MAXBUF, datatype, 0, 0, comm, &status);
12    MPI_Get_count(&status, datatype, &count);
13    for (i=0; i<&count; i++) {
14        ...
15    }
16    ...
```



## Simple MPI Example (2/2) – Send and Receive

### Source of the example

[http://coewww.rutgers.edu/www1/linuxclass2005/lessons/lesson13/sec\\_8.html](http://coewww.rutgers.edu/www1/linuxclass2005/lessons/lesson13/sec_8.html)

- Compile program:

```
$ mpicc sendrecv.c -o sendrecv
```

- Distribute the program in the cluster:

- The program must be stored on each node in the same directory!

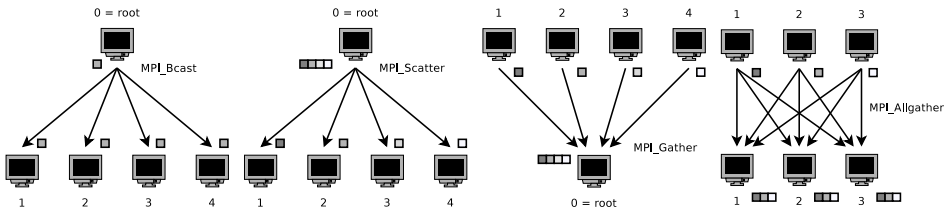
```
$ scp sendrecv node1:-  
$ scp sendrecv node2:-
```

- Program execution (2 processes) in the cluster:

```
$ mpirun -np 2 --hostfile hosts.mpi sendrecv  
Task 0: Received 1 char(s) from task 1 with tag 1  
Task 1: Received 1 char(s) from task 0 with tag 1
```

# More Send and Receive Operations

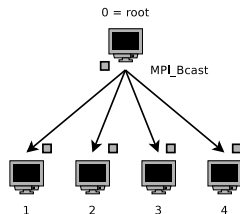
- MPI\_Send and MPI\_Recv implement *one-to-one* communication
- MPI\_Bcast and MPI\_Scatter implement *one-to-many* communication
- MPI\_Gather implements *many-to-one* communication
- MPI\_Allgather implements *many-to-many* communication



The functions MPI\_Bcast, MPI\_Scatter, MPI\_Gather and MPI\_Allgather must be called by all processes in the communicator!

# MPI Functions – Broadcast Sending (1/2)

- `MPI_Bcast(int buffer, int count, MPI_Datatype datatype, int root, MPI_Comm comm)`

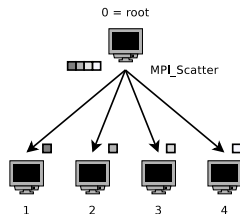


- Send a message of process root to all other processes in the communicator
  - `buffer` = first address of the transmit buffer
  - `count` = number of elements in the transmit buffer (not negative)
  - `datatype` = MPI data type of the elements in the transmit buffer
  - `root` = rank of the sender process in the communicator
  - `comm` = communicator
- All processes in the communicator must call the function



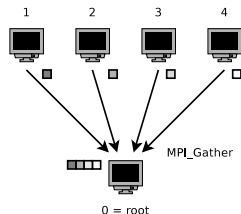
## MPI Functions – Scatter

- `MPI_Scatter(int sendbuf, int sendcount, MPI_Datatype sendtype, int recvbuf, int recvcount, MPI_Datatype recvtype, int root, MPI_Comm comm)`
- While `MPI_Bcast` sends the same piece of data to all processes, `MPI_Scatter` sends chunks of an array to different processes
  - `sendbuf + recvbuf` = first address of the transmit/receive buffer
  - `sendcount + recvcount` = number of elements in the transmit/receive buffer (not negative and typically equal size)
  - `sendtype + recvtype` = MPI data type of the elements in the transmit/receive buffer
  - `root` = rank of the sender process in the communicator
  - `comm` = communicator
- All processes in the communicator must call the function



## MPI Functions – Gather

- `MPI_Gather(int sendbuf, int sendcount, MPI_Datatype sendtype, int recvbuf, int recvcount, MPI_Datatype recvttype, int root, MPI_Comm comm)`
- `MPI_Gather` is the inverse of `MPI_Scatter`
- Instead of distributing elements from one process to many processes, `MPI_Gather` takes elements from many processes and gathers them to one single process.
- All parameters are equal to `MPI_Scatter`
- All processes in the communicator must call the function



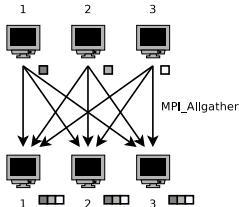
Source: <http://mpitutorial.com/tutorials/mpi-scatter-gather-and-allgather/>

# MPI Functions – Allgather

```

● MPI_Allgather(int sendbuf, int sendcount,
               MPI_Datatype sendtype,
               int recvbuf, int recvcount,
               MPI_Datatype recvtype,
               MPI_Comm comm)

```



- Given a set of elements distributed across all processes, MPI\_Allgather will gather all of the elements to all the processes.
- MPI\_Allgather is basically an MPI\_Gather followed by an MPI\_Bcast
- All parameters are equal to MPI\_Scatter and MPI\_Gather with the difference that there is no root process in MPI\_Allgather
- All processes in the communicator must call the function

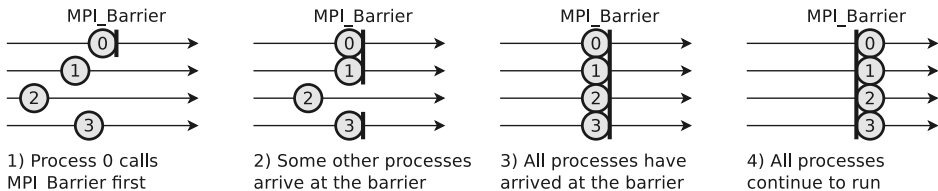
Source: <http://mpitutorial.com/tutorials/mpi-scatter-gather-and-allgather/>

# MPI Functions – Barrier

- `MPI_Barrier(MPI_Comm comm)`
  - Blocks the execution of the calling process, until all processes in the communicator `comm` have called the barrier function
  - `comm` = communicator

```
1 #include "mpi.h"
2
3 MPI_Comm comm;
4
5 ...
6 MPI_Barrier(comm);
7 ...
```

- `MPI_Barrier` can be used to synchronize a program



<https://mpitutorial.com/tutorials/mpi-broadcast-and-collective-communication/>



## Reduces Values on all Processes to a single Value

- `MPI_Reduce(const void *sendbuf, void *recvbuf, int count, MPI_Datatype datatype, MPI_Op op, int root, MPI_Comm comm)`
- Reduces values on all processes to a single value on process root
  - `sendbuf` = address of send buffer (input parameter)
  - `recvbuf` = address of receive buffer on root (output parameter)
  - `count` = number of elements in the transmit buffer (not negative)
  - `datatype` = MPI data type of the elements in the transmit buffer
  - `op` = reduce operation
  - `root` = rank of the root process in the communicator
  - `comm` = communicator (all processes in the communicator must call the function)

The reduction operations defined by MPI include:

`MPI_MAX` (Returns the maximum element)

`MPI_MIN` (Returns the minimum element)

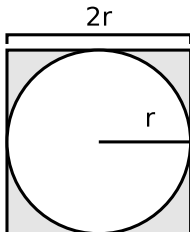
`MPI_SUM` (Sums the elements)

`MPI_PROD` (Multiplies all elements)

`MPI_MAXLOC` (Returns the maximum value and the rank of the process that owns it)

`MPI_MINLOC` (Returns the minimum value and the rank of the process that owns it)

# Example: Calculation of $\pi$ via Monte Carlo Simulation



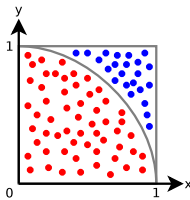
- $r$  = Radius
- $A$  = Surface ratio
- $C$  = Circle
- $S$  = Square

- Inscribe a circle of radius  $r$  inside a square with side length  $2r$
- 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

$$\frac{A_C}{A_S} = \frac{\pi \cdot r^2}{(2 \cdot r)^2} = \frac{\pi \cdot r^2}{4 \cdot r^2} = \frac{\pi}{4}$$

- The dots can be generated (X/Y axis values via random) in parallel by the workers
- The master receives from each worker the number of calculated dots in  $A_C$  and calculates:

$$\frac{4 \cdot \text{dots in } A_C}{\text{dots in } A_S} = \pi$$



# MPI Example – Calculate $\pi$ (1/3)

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <math.h>
4 #include "mpi.h"
5
6 int main(int argc, char *argv[]) {
7
8     int myid, numprocs;
9
10    double PI25DT = 3.141592653589793238462643;
11    double t1, t2;
12
13    long long npts = 1e11;
14    long long i, mynpts;
15
16    long double f, sum, mysum;
17    long double xmin, xmax, x;
18
19    // Initialization routine => starts the MPI environment
20    // Defines the communicator MPI_COMM_WORLD
21    MPI_Init(&argc, &argv);
22    // Determines the number of processes in a communicator
23    MPI_Comm_size(MPI_COMM_WORLD, &numprocs);
24    // Determines the rank (id) of the calling process in the communicator
25    MPI_Comm_rank(MPI_COMM_WORLD, &myid);
```

This Source Code is influenced a lot by this Source...

<https://web.archive.org/web/20160812014841/http://chpc.wustl.edu/mpi-c.html>

# MPI Example – Calculate $\pi$ (2/3)

```
1 // Data decomposition. Each process gets a part of the work
2 mynpts = npts/numprocs;
3
4 if (myid == 0) {
5 // Returns the time in seconds since an arbitrary time in the past
6 t1 = MPI_Wtime();
7 }
8
9 mysum = 0.0;
10 xmin = 0.0;
11 xmax = 1.0;
12
13 // Seed the pseudo random number generator
14 srand(time(0));
15
16 for (i=0; i<mynpts; i++) {
17 // (long double)rand()/(long double)RAND_MAX
18 // returns a random number between 0 and 1.
19 // (long double)rand()/(long double)RAND_MAX*(xmax-xmin)
20 // returns a random number between 0 and max - min.
21 // the whole expression will return a random number between 0+min and min+(max-min)
22 // => between min and max.
23 x = (long double)rand()/(long double)RAND_MAX*(xmax-xmin) + xmin;
24 // Each process does a partial sum over its own points.
25 mysum += 4.0/(1.0 + x*x);
26 }
```

How to generate a random number between 0 and 1?

<https://stackoverflow.com/questions/6218399/how-to-generate-a-random-number-between-0-and-1/6219525>

MPI Example – Calculate  $\pi$  (3/3)

```

1  // Take all the processes values of mysum and add them up into sum on process 0.
2  MPI_Reduce(&mysum,&sum,1,MPI_LONG_DOUBLE,MPI_SUM,0,MPI_COMM_WORLD);
3
4  if (myid == 0) {
5      // Returns the time in seconds since an arbitrary time in the past
6      t2 = MPI_Wtime();
7
8      f = sum/npts;
9
10     printf("Pi calculated with %lld points. \n",npts);
11     printf("Pi calculated:      %.16f \n",f);
12     printf("Correct value of Pi: %.16f \n",PI25DT);
13     printf("Error is:          %.16f \n",fabs(f-PI25DT));
14     printf("Elapsed time [s] for the relevant part of the program: %f\n", t2 - t1);
15 }
16
17 // Stop the MPI environment
18 MPI_Finalize();
19 }

```

```

$ time mpirun -np 512 --hostfile hosts_4cores_128.mpi /mnt/cluster_128/pi
Pi calculated with 10000000000 points.
Pi calculated:      3.1415785751520118
Correct value of Pi: 3.1415926535897931
Error is:          0.0000140784377813
Elapsed time [s] for the relevant part of the program: 37.651207

real    0m46.394s
user    0m18.860s
sys     0m3.020s

```



# Gearman

- The name *Gearman* is an anagram for *manager*
  - Gearman only distributes jobs
- Gearman should only be used in secure private networks
  - The communication is not encrypted and uses port 4730
  - No mechanism for the authentication of the systems is implemented
- Clients and workers access shared data
  - Cluster file systems like GlusterFS or protocols such as NFS or Samba can be used

Helpful article about Gearman (in German language)

Gearman verteilt Arbeit auf Rechner im LAN, *Reiko Kaps*, c't 24/2010, P.192

- The next slides contain an application example from the article

## Gearman – Example of a Worker Script

- Client and worker both, access via `/src/media` a shared file system
  - The shared file system contains images that need to be resized
- The workers scale via ImageMagick `convert`
- Shell script `resizer-worker.sh`

```
#!/bin/bash
INFILE="$1"

echo "Converting ${INFILE} on $HOSTNAME" >> /src/media/g.log

convert "${INFILE}" -resize 1024 "${INFILE}"-small.jpg
```

- Register the worker script (`-w`) at the Job Server „`gman-jserver`“ (`-h`) with the function name „`Resizer`“ (`-f`):
  - `gearman -h gman-jserver -w -f Resizer xargs resizer-worker.sh`

