

Computer cluster

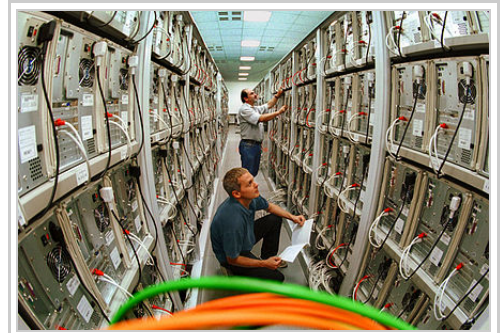
From Wikipedia, the free encyclopedia

A **computer cluster** consists of a set of loosely or tightly connected computers that work together so that, in many respects, they can be viewed as a single system. Unlike grid computers, computer clusters have each node set to perform the same task, controlled and scheduled by software.^[1]

The components of a cluster are usually connected to each other through fast local area networks ("LAN"), with each *node* (computer used as a server) running its own instance of an operating system. In most circumstances, all of the nodes use the same hardware^[2] and the same operating system, although in some setups (i.e. using Open Source Cluster Application Resources (OSCAR)), different operating systems can be used on each computer, and/or different hardware.^[3]

They are usually deployed to improve performance and availability over that of a single computer, while typically being much more cost-effective than single computers of comparable speed or availability.^[4]

Computer clusters emerged as a result of convergence of a number of computing trends including the availability of low cost microprocessors, high speed networks, and software for high-performance distributed computing. They have a wide range of applicability and deployment, ranging from small business clusters with a handful of nodes to some of the fastest supercomputers in the world such as IBM's Sequoia.^[5] The applications that can be done however, are nonetheless limited, since the software needs to be purpose-built per task. It is hence not possible to use computer clusters for casual computing tasks.^[6]



Technicians working on a large Linux cluster at the Chemnitz University of Technology, Germany



Sun Microsystems Solaris Cluster

Contents

- 1 Basic concepts
- 2 History
- 3 Attributes of clusters
- 4 Benefits
- 5 Design and Configuration
- 6 Data sharing and communication
 - 6.1 Data sharing
 - 6.2 Message passing and communication
- 7 Cluster management
 - 7.1 Task scheduling

- 7.2 Node failure management
- 8 Software development and administration
 - 8.1 Parallel programming
 - 8.2 Debugging and monitoring
- 9 Some implementations
- 10 Other approaches
- 11 See also
- 12 References
- 13 Further reading
- 14 External links

Basic concepts

The desire to get more computing power and better reliability by orchestrating a number of low cost commercial off-the-shelf computers has given rise to a variety of architectures and configurations.

The computer clustering approach usually (but not always) connects a number of readily available computing nodes (e.g. personal computers used as servers) via a fast local area network.^[7] The activities of the computing nodes are orchestrated by "clustering middleware", a software layer that sits atop the nodes and allows the users to treat the cluster as by and large one cohesive computing unit, e.g. via a single system image concept.^[7]

Computer clustering relies on a centralized management approach which makes the nodes available as orchestrated shared servers. It is distinct from other approaches such as peer to peer or grid computing which also use many nodes, but with a far more distributed nature.^[7]

A computer cluster may be a simple two-node system which just connects two personal computers, or may be a very fast supercomputer. A basic approach to building a cluster is that of a Beowulf cluster which may be built with a few personal computers to produce a cost-effective alternative to traditional high performance computing. An early project that showed the viability of the concept was the 133 nodes Stone Soupercomputer.^[8] The developers used Linux, the Parallel Virtual Machine toolkit and the Message Passing Interface library to achieve high performance at a relatively low cost.^[9]

Although a cluster may consist of just a few personal computers connected by a simple network, the cluster architecture may also be used to achieve very high levels of performance. The TOP500 organization's semiannual list of the 500 fastest supercomputers often includes many clusters, e.g. the world's fastest machine in 2011 was the K computer which has a distributed memory, cluster architecture.^{[10][11]}

History



A simple, home-built Beowulf cluster.

Greg Pfister has stated that clusters were not invented by any specific vendor but by customers who could not fit all their work on one computer, or needed a backup.^[12] Pfister estimates the date as some time in the 1960s. The formal engineering basis of cluster computing as a means of doing parallel work of any sort was arguably invented by Gene Amdahl of IBM, who in 1967 published what has come to be regarded as the seminal paper on parallel processing: Amdahl's Law.

The history of early computer clusters is more or less directly tied into the history of early networks, as one of the primary motivations for the development of a network was to link computing resources, creating a de facto computer cluster.

The first commercial clustering product was Datapoint Corporation's "Attached Resource Computer" (ARC) system, developed in 1977, and using ARCnet as the cluster interface. Clustering per se did not really take off until Digital Equipment Corporation released their VAXcluster product in 1984 for the VAX/VMS operating system (now named as OpenVMS). The ARC and VAXcluster products not only supported parallel computing, but also shared file systems and peripheral devices. The idea was to provide the advantages of parallel processing, while maintaining data reliability and uniqueness. Two other noteworthy early commercial clusters were the *Tandem Himalayan* (a circa 1994 high-availability product) and the *IBM S/390 Parallel Sysplex* (also circa 1994, primarily for business use).

Within the same time frame, while computer clusters used parallelism outside the computer on a commodity network, supercomputers began to use them within the same computer. Following the success of the CDC 6600 in 1964, the Cray 1 was delivered in 1976, and introduced internal parallelism via vector processing.^[13] While early supercomputers excluded clusters and relied on shared memory, in time some of the fastest supercomputers (e.g. the K computer) relied on cluster architectures.

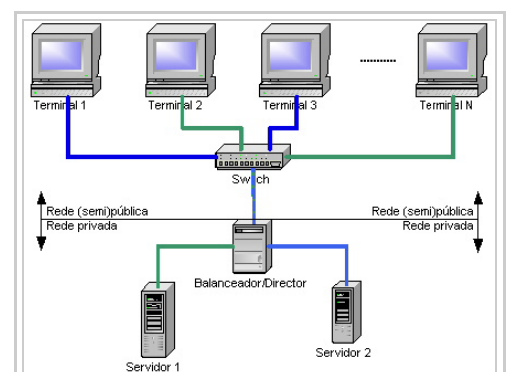
Attributes of clusters

Computer clusters may be configured for different purposes ranging from general purpose business needs such as web-service support, to computation-intensive scientific calculations. In either case, the cluster may use a high-availability approach. Note that the attributes described below are not exclusive and a "computer cluster" may also use a high-availability approach, etc.

"Load-balancing" clusters are configurations in which cluster-nodes share computational workload to provide better overall performance. For example, a web server cluster may assign different queries to different nodes, so the overall response time will be optimized.^[14] However, approaches to load-balancing may significantly differ among applications, e.g. a high-performance cluster used for scientific computations would balance load with different algorithms from a web-server cluster which may just use a simple round-robin method by assigning each new request to a different node.^[14]



A VAX 11/780, c. 1977



A load balancing cluster with two servers and N user stations (Galician).

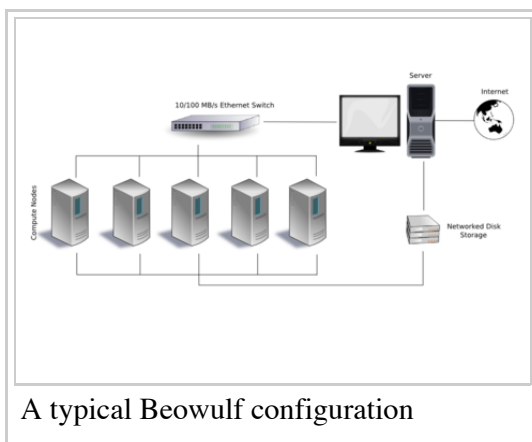
Computer clusters are used for computation-intensive purposes, rather than handling IO-oriented operations such as web service or databases.^[15] For instance, a computer cluster might support computational simulations of vehicle crashes or weather. Very tightly coupled computer clusters are designed for work that may approach "supercomputing".

"High-availability clusters" (also known as failover clusters, or HA clusters) improve the availability of the cluster approach. They operate by having redundant nodes, which are then used to provide service when system components fail. HA cluster implementations attempt to use redundancy of cluster components to eliminate single points of failure. There are commercial implementations of High-Availability clusters for many operating systems. The Linux-HA project is one commonly used free software HA package for the Linux operating system.

Benefits

Clusters are primarily designed with performance in mind, but installations are based on many other factors; fault tolerance (*the ability for a system to continue working with a malfunctioning node*) also allows for simpler scalability, and in high performance situations, low frequency of maintenance routines, resource consolidation, and centralized management.^{[16][17]}

Design and Configuration



One of the issues in designing a cluster is how tightly coupled the individual nodes may be. For instance, a single computer job may require frequent communication among nodes: this implies that the cluster shares a dedicated network, is densely located, and probably has homogeneous nodes. The other extreme is where a computer job uses one or few nodes, and needs little or no inter-node communication, approaching grid computing.

In a Beowulf system, the application programs never see the computational nodes (also called slave computers) but only interact with the "Master" which is a specific computer

handling the scheduling and management of the slaves.^[15] In

a typical implementation the Master has two network interfaces, one that communicates with the private Beowulf network for the slaves, the other for the general purpose network of the organization.^[15] The slave computers typically have their own version of the same operating system, and local memory and disk space. However, the private slave network may also have a large and shared file server that stores global persistent data, accessed by the slaves as needed.^[15]

By contrast, the special purpose 144 node DEGIMA cluster is tuned to running astrophysical N-body simulations using the Multiple-Walk parallel treecode, rather than general purpose scientific computations.^[18]

Due to the increasing computing power of each generation of game consoles, a novel use has emerged where they are repurposed into High-performance computing (HPC) clusters. Some examples of game console clusters are Sony PlayStation clusters and Microsoft Xbox clusters. Another example of consumer game product is the Nvidia Tesla Personal Supercomputer workstation, which uses multiple graphics accelerator processor chips. Besides game consoles, high-end graphics cards too can be used

instead. The use of graphics cards (or rather their GPU's) to do calculations for grid computing is vastly more economical than using CPU's, despite being less precise. However, when using double-precision values, they become as precise to work with as CPU's, and still be much less costly (purchase cost).^[19]

Computer clusters have historically run on separate physical computers with the same operating system. With the advent of virtualization, the cluster nodes may run on separate physical computers with different operating systems which are painted above with a virtual layer to look similar.^[20] The cluster may also be virtualized on various configurations as maintenance takes place. An example implementation is Xen as the virtualization manager with Linux-HA.^[21]

Data sharing and communication

Data sharing

As the computer clusters were appearing during the 1980s, so were supercomputers. One of the elements that distinguished the three classes at that time was that the early supercomputers relied on shared memory. To date clusters do not typically use physically shared memory, while many supercomputer architectures have also abandoned it.

However, the use of a clustered file system is essential in modern computer clusters. Examples include the IBM General Parallel File System, Microsoft's Cluster Shared Volumes or the Oracle Cluster File System.

Message passing and communication

Two widely used approaches for communication between cluster nodes are MPI, the Message Passing Interface and PVM, the Parallel Virtual Machine.^[22]

PVM was developed at the Oak Ridge National Laboratory around 1989 before MPI was available. PVM must be directly installed on every cluster node and provides a set of software libraries that paint the node as a "parallel virtual machine". PVM provides a run-time environment for message-passing, task and resource management, and fault notification. PVM can be used by user programs written in C, C++, or Fortran, etc.^{[22][23]}

MPI emerged in the early 1990s out of discussions among 40 organizations. The initial effort was supported by ARPA and National Science Foundation. Rather than starting anew, the design of MPI drew on various features available in commercial systems of the time. The MPI specifications then gave rise to specific implementations. MPI implementations typically use TCP/IP and socket connections.^[22] MPI is now a widely available communications model that enables parallel programs to be written in languages such as C, Fortran, Python, etc.^[23] Thus, unlike PVM which provides a concrete implementation, MPI is a specification which has been implemented in systems such as MPICH and Open MPI.^{[23][24]}

Cluster management

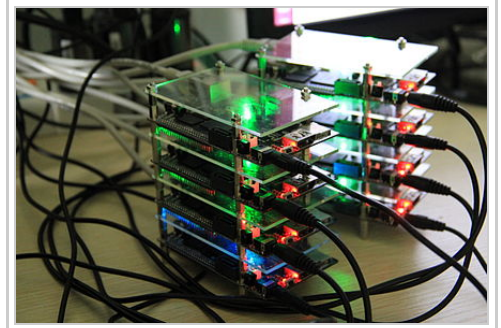


A NEC Nehalem cluster

One of the challenges in the use of a computer cluster is the cost of administrating it which can at times be as high as the cost of administrating N independent machines, if the cluster has N nodes.^[25] In some cases this provides an advantage to shared memory architectures with lower administration costs.^[25] This has also made virtual machines popular, due to the ease of administration.^[25]

Task scheduling

When a large multi-user cluster needs to access very large amounts of data, task scheduling becomes a challenge. In a heterogeneous CPU-GPU cluster with a complex application environment, the performance of each job depends on the characteristics of the underlying cluster. Therefore, mapping tasks onto CPU cores and GPU devices provides significant challenges.^[26] This is an area of ongoing research; algorithms that combine and extend MapReduce and Hadoop have been proposed and studied.^[26]



Low-cost and low energy tiny-cluster of Cubieboards, using Apache Hadoop on Lubuntu

Node failure management

When a node in a cluster fails, strategies such as "fencing" may be employed to keep the rest of the system operational.^{[27][28]} Fencing is the process of isolating a node or protecting shared resources when a node appears to be malfunctioning. There are two classes of fencing methods; one disables a node itself, and the other disallows access to resources such as shared disks.^[27]

The STONITH method stands for "Shoot The Other Node In The Head", meaning that the suspected node is disabled or powered off. For instance, *power fencing* uses a power controller to turn off an inoperable node.^[27]

The *resources fencing* approach disallows access to resources without powering off the node. This may include *persistent reservation fencing* via the SCSI3, fibre channel fencing to disable the fibre channel port, or global network block device (GNBD) fencing to disable access to the GNBD server.

Software development and administration

Parallel programming

Load balancing clusters such as web servers use cluster architectures to support a large number of users and typically each user request is routed to a specific node, achieving task parallelism without multi-node cooperation, given that the main goal of the system is providing rapid user access to shared data. However, "computer clusters" which perform complex computations for a small number of users need to take advantage of the parallel processing capabilities of the cluster and partition "the same computation" among several nodes.^[29]

Automatic parallelization of programs continues to remain a technical challenge, but parallel programming models can be used to effectuate a higher degree of parallelism via the simultaneous execution of separate portions of a program on different processors.^{[29][30]}

Debugging and monitoring

The development and debugging of parallel programs on a cluster requires parallel language primitives as well as suitable tools such as those discussed by the *High Performance Debugging Forum* (HPDF) which resulted in the HPD specifications.^{[23][31]} Tools such as TotalView were then developed to debug parallel implementations on computer clusters which use MPI or PVM for message passing.

The Berkeley NOW (Network of Workstations) system gathers cluster data and stores them in a database, while a system such as PARMON, developed in India, allows for the visual observation and management of large clusters.^[23]

Application checkpointing can be used to restore a given state of the system when a node fails during a long multi-node computation.^[32] This is essential in large clusters, given that as the number of nodes increases, so does the likelihood of node failure under heavy computational loads. Checkpointing can restore the system to a stable state so that processing can resume without having to recompute results.^[32]

Some implementations

The GNU/Linux world supports various cluster software; for application clustering, there is distcc, and MPICH. Linux Virtual Server, Linux-HA - director-based clusters that allow incoming requests for services to be distributed across multiple cluster nodes. MOSIX, LinuxPMI, Kerrighed, OpenSSI are full-blown clusters integrated into the kernel that provide for automatic process migration among homogeneous nodes. OpenSSI, openMosix and Kerrighed are single-system image implementations.

Microsoft Windows computer cluster Server 2003 based on the Windows Server platform provides pieces for High Performance Computing like the Job Scheduler, MSMPI library and management tools.

gLite is a set of middleware technologies created by the Enabling Grids for E-sciencE (EGEE) project.

slurm is also used to schedule and manage some of the largest supercomputer clusters (see top500 list).

Other approaches

Although most computer clusters are permanent fixtures, attempts at flash mob computing have been made to build short-lived clusters for specific computations. However, larger scale volunteer computing systems such as BOINC-based systems have had more followers.

See also

Basic concepts

- Clustered file system
- Heartbeat private network
- High-availability cluster
- Single system

Specific systems

- DEGIMA (computer cluster)
- K computer
- Microsoft Cluster Server
- Red Hat

image

- Symmetric multiprocessing
- Robotic clusters

Distributed computing

- Distributed computing
- Distributed data store
- Distributed operating system
- Distributed shared memory

Cluster Suite

- Rocks Cluster Distribution
- Solaris Cluster
- Veritas Cluster Server

Computer farms

- Compile farm
- Render farm
- Server farm

References

1. [^] Grid vs cluster computing (http://www.answers.com/Q/Comparison_Grid_computing_and_cluster_Computing)
2. [^] Cluster vs grid computing (<http://stackoverflow.com/questions/9723040/what-is-the-difference-between-cloud-grid-and-cluster>)
3. [^] Hardware of computer clusters not always needing to be the same, probably depends on software used (<http://www.pcauthority.com.au/Feature/306972,weekend-project-build-your-own-supercomputer.aspx>)
4. [^] Bader, David; Robert Pennington (June 1996). "Cluster Computing: Applications" (<http://www.cc.gatech.edu/~bader/papers/ijhpca.html>). Georgia Tech College of Computing. Retrieved 2007-07-13.
5. [^] "Nuclear weapons supercomputer reclaims world speed record for US" (<http://www.telegraph.co.uk/technology/9338651/Nuclear-weapons-supercomputer-reclaims-world-speed-record-for-US.html>). The Telegraph. 18 Jun 2012. Retrieved 18 Jun 2012.
6. [^] Grid and cluster computing, limitations (http://www.answers.com/Q/Comparison_Grid_computing_and_cluster_Computing)
7. [^] ^a ^b ^c *Network-Based Information Systems: First International Conference, NBIS 2007* ISBN 3-540-74572-6 page 375
8. [^] William W. Hargrove, Forrest M. Hoffman and Thomas Sterling (August 16, 2001). "The Do-It-Yourself Supercomputer" (<http://www.sciam.com/article.cfm?id=the-do-it-yourself-superc>). *Scientific American* **265** (2). pp. 72–79. Retrieved October 18, 2011.
9. [^] William W. Hargrove and Forrest M. Hoffman (1999). "Cluster Computing: Linux Taken to the Extreme" (<http://climate.ornl.gov/~forrest/linux-magazine-1999/>). *Linux magazine*. Retrieved October 18, 2011.
10. [^] TOP500 list (<http://i.top500.org/sublist>) To view all clusters on the TOP500 select "cluster" as architecture from the sublist menu.
11. [^] M. Yokokawa et al *The K Computer*, in "International Symposium on Low Power Electronics and Design" (ISLPED) 1–2 Aug. 2011, pp. 271–272

(ISLPEU) 1-5 Aug. 2011, pages 5/1-5/2

12. ^ Pfister, Gregory (1998). *In Search of Clusters* (2nd ed.). Upper Saddle River, NJ: Prentice Hall PTR. p. 36. ISBN 0-13-899709-8.
13. ^ *Readings in computer architecture* by Mark Donald Hill, Norman Paul Jouppi, Gurindar Sohi 1999 ISBN 978-1-55860-539-8 page 41-48
14. ^ ^{a b} *High Performance Linux Clusters* by Joseph D. Sloan 2004 ISBN 0-596-00570-9 page
15. ^ ^{a b c d} *High Performance Computing for Computational Science - VECPAR 2004* by Michel Daydé, Jack Dongarra 2005 ISBN 3-540-25424-2 pages 120-121
16. ^ "IBM Cluster System : Benefits" (<http://www-03.ibm.com/systems/clusters/benefits.html>). <http://www-03.ibm.com/>. IBM. Retrieved 8 September 2014.
17. ^ "Evaluating the Benefits of Clustering" ([http://technet.microsoft.com/en-us/library/cc778629\(v=ws.10\).aspx](http://technet.microsoft.com/en-us/library/cc778629(v=ws.10).aspx)). <http://www.microsoft.com/>. Microsoft. 28 March 2003. Retrieved 8 September 2014.
18. ^ Hamada T. *et al.* (2009) A novel multiple-walk parallel algorithm for the Barnes–Hut treecode on GPUs – towards cost effective, high performance N-body simulation. *Comput. Sci. Res. Development* 24:21-31. doi:10.1007/s00450-009-0089-1 (<https://dx.doi.org/10.1007%2Fs00450-009-0089-1>)
19. ^ GPU options (<http://www.pcauthority.com.au/Feature/306972,weekend-project-build-your-own-supercomputer.aspx>)
20. ^ Using Xen (<http://www.linuxjournal.com/article/8812>)
21. ^ Maurer, Ryan: Xen Virtualization and Linux Clustering (<http://www.linuxjournal.com/article/8812>)
22. ^ ^{a b c} *Distributed services with OpenAFS: for enterprise and education* by Franco Milicchio, Wolfgang Alexander Gehrke 2007, ISBN pages 339-341 [1] (http://books.google.it/books?id=bKf4NBaIJI8C&pg=PA339&dq=%22message+passing%22+computer+cluster+MPI+PVM&hl=en&sa=X&ei=-dD7ToZCj_bhBOOxvIOI&redir_esc=y#v=onepage&q=%22message%20passing%22%20computer%20cluster%20MPI%20PVM&f=false)
23. ^ ^{a b c d e} *Grid and Cluster Computing* by Prabhu 2008 8120334280 pages 109-112 (http://books.google.it/books?id=EIVdVtGHv-0C&pg=PA112&dq=%22message+passing%22+computer+cluster+MPI+PVM&hl=en&sa=X&ei=-dD7ToZCj_bhBOOxvIOI&redir_esc=y#v=onepage&q=%22message%20passing%22%20computer%20cluster%20MPI%20PVM&f=false)
24. ^ Gropp, William; Lusk, Ewing; Skjellum, Anthony (1996). "A High-Performance, Portable Implementation of the MPI Message Passing Interface". *Parallel Computing*. CiteSeerX: 10.1.1.102.9485 (<http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.102.9485>).
25. ^ ^{a b c} *Computer Organization and Design* by David A. Patterson and John L. Hennessy 2011 ISBN 0-12-374750-3 pages 641-642
26. ^ ^{a b} K. Shirahata, et al *Hybrid Map Task Scheduling for GPU-Based Heterogeneous Clusters* in: Cloud Computing Technology and Science (CloudCom), 2010 Nov. 30 2010-Dec. 3 2010 pages 733 - 740 ISBN 978-1-4244-9405-7 [2] (http://ieeexplore.ieee.org/xpl/freeabs_all.jsp?arnumber=5708524)
27. ^ ^{a b c} Alan Robertson *Resource fencing using STONITH*. IBM Linux Research Center, 2010 [3] (ftp://ftp.telecom.uff.br/pub/linux/HA/ResourceFencing_Stonith.pdf)
28. ^ *Sun Cluster environment: Sun Cluster 2.2* by Enrique Vargas, Joseph Bianco, David Deeths 2001 ISBN page 58
29. ^ ^{a b} *Computer Science: The Hardware, Software and Heart of It* by Alfred V. Aho, Edward K. Blum 2011

- ISBN 1-4614-1167-X pages 156-166 [4] (http://books.google.it/books?id=S7QU9RRLYIYC&pg=PA156&dq=%22parallel+programming%22+computer+cluster&hl=en&sa=X&ei=C_H8TqSRCobS4QSTmY2pCA&sqi=2&redir_esc=y#v=onepage&q=%22parallel%20programming%22%20computer%20cluster&f=false)
30. ^ *Parallel Programming: For Multicore and Cluster Systems* by Thomas Rauber, Gudula Rünger 2010 ISBN 3-642-04817-X pages 94–95 [5] (http://books.google.it/books?id=wWogxOmA3wMC&pg=PA94&dq=%22parallel+programming+language%22+computer+cluster&hl=en&sa=X&ei=zfD8TpX_F5CQ4gSMvfXhBA&redir_esc=y#v=onepage&q=%22parallel%20programming%20language%22%20computer%20cluster&f=false)
31. ^ *A debugging standard for high-performance computing* by Joan M. Francioni and Cherri Pancake, in the "Journal of Scientific Programming" Volume 8 Issue 2, April 2000 [6] (<http://dl.acm.org/citation.cfm?id=1239906>)
32. ^ *^a ^b Computational Science-- ICCS 2003: International Conference* edited by Peter Sloot 2003 ISBN 3-540-40195-4 pages 291-292

Further reading

- Mark Baker, et al., *Cluster Computing White Paper* [7] (<http://arxiv.org/abs/cs/0004014>), 11 Jan 2001.
- Evan Marcus, Hal Stern: *Blueprints for High Availability: Designing Resilient Distributed Systems*, John Wiley & Sons, ISBN 0-471-35601-8
- Greg Pfister: *In Search of Clusters*, Prentice Hall, ISBN 0-13-899709-8
- Rajkumar Buyya (editor): *High Performance Cluster Computing: Architectures and Systems*, Volume 1, ISBN 0-13-013784-7, and Volume 2, ISBN 0-13-013785-5, Prentice Hall, NJ, USA, 1999.

External links

- IEEE Technical Committee on Scalable Computing (TCSC) (<https://www.ieeetcsc.org/>)
- Reliable Scalable Cluster Technology, IBM (<http://publib.boulder.ibm.com/infocenter/clresctr/vxrx/index.jsp?topic=%2Fcom.ibm.cluster.rsct.doc%2Frsctbooks.html>)
- Tivoli System Automation Wiki (<https://www.ibm.com/developerworks/wikis/display/tivoli/Tivoli+System+Automation>)



Wikimedia Commons has media related to ***Computer cluster***.

Retrieved from "http://en.wikipedia.org/w/index.php?title=Computer_cluster&oldid=640418649"

Categories: Cluster computing | Parallel computing | Concurrent computing | Supercomputers
| Local area networks | Classes of computers | Fault-tolerant computer systems

- This page was last modified on 31 December 2014, at 19:36.
- Text is available under the Creative Commons Attribution-ShareAlike License; additional terms may apply. By using this site, you agree to the Terms of Use and Privacy Policy. Wikipedia® is a registered trademark of the Wikimedia Foundation, Inc., a non-profit organization.