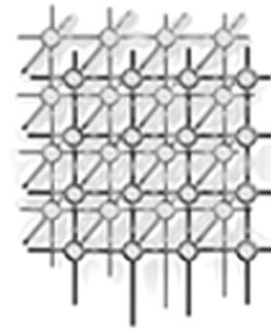


Optimizing the distribution of large data sets in theory and practice[‡]



Felix Rauch^{*,†}, Christian Kurmann and Thomas M. Stricker

*Laboratory for Computer Systems, ETH - Swiss Institute of Technology,
CH-8092 Zürich, Switzerland*

SUMMARY

Multicasting large amounts of data efficiently to all nodes of a PC cluster is an important operation. In the form of a *partition cast* it can be used to replicate entire software installations by cloning. Optimizing a partition cast for a given cluster of PCs reveals some interesting architectural tradeoffs, since the fastest solution does not only depend on the network speed and topology, but remains highly sensitive to other resources like the disk speed, the memory system performance and the processing power in the participating nodes. We present an analytical model that guides an implementation towards an optimal configuration for any given PC cluster. The model is validated by measurements on our cluster using Gigabit- and Fast-Ethernet links. The resulting simple software tool, *Dolly*, can replicate an entire 2 GB Windows NT image onto 24 machines in less than 5 min. Copyright © 2002 John Wiley & Sons, Ltd.

KEY WORDS: software installation and maintenance; data streaming; partition management; communication modelling; multicast; input output systems

1. INTRODUCTION AND RELATED WORK

The work on partition cast was motivated by our work with the Patagonia multi-purpose PC cluster. This cluster can be used for different tasks by booting different system installations [1]. The usage modes comprise traditional scientific computing workloads (Linux), research experiments in distributed data processing (data-mining) or distributed collaborative work (Linux and Windows NT)

*Correspondence to: Felix Rauch, Laboratory for Computer Systems, ETH - Swiss Institute of Technology, CH-8092 Zürich, Switzerland.

†E-mail: rauch@inf.ethz.ch

‡The original version of this article was first published as 'Rauch F, Kurmann C, Stricker TM. Optimizing the distribution of large data sets in theory and practice. *Euro-Par 2000—Parallel Processing (Lecture Notes in Computer Science, vol. 1900)*, Bode A, Ludwig T, Karl W, Wismüller R (eds.). Springer, 2000; 1118–1131', and is reproduced here by kind permission of the publisher.



and computer science education (Windows NT, Oberon). For best flexibility and maintenance, such a multiuse cluster must support the installation of new operating-system images within minutes.

The problem of copying entire partitions over a fast network leads to some interesting tradeoffs in the overall design of a PC cluster architecture. Our cluster nodes are built from advanced components such as fast microprocessors, disk drives and high speed network interfaces connected via a scalable switching fabric. Yet it is not obvious which arrangement of the network or which configuration of the software results in the fastest system to distribute large blocks of data to all the machines of the cluster.

After in-depth analytical modelling of network and cluster nodes, we create a simple, operating-system independent tool that distributes raw disk partitions. The tool can be used to clone any operating system. Most operating systems can perform automatic installation and customization at startup and a cloned partition image can therefore be used immediately after a partition cast completes.

For experimental verification of our approach we use a meta cluster at our university (ETH Zürich) that unites several PC clusters, connecting their interconnects to a dedicated cluster backbone. This cluster testbed offers a variety of topologies and networking speeds. The networks include some Gigabit networking technology like SCI [2,3] and Myrinet [4] with an emphasis on Fast and Gigabit Ethernet [5]. The evaluation work was performed on the Patagonia sub-cluster of 24 Dell 410 Desktop PCs configured as workstations with keyboards and monitors. The Intel based PC nodes are built around a dual Pentium II processor configuration (running at 400 MHz) and 256 MB SDRAM memory connected to a 100 MHz front side bus. All machines are equipped with 9 GB Ultra2 Cheetah SCSI hard-disk drives which can read and write a data stream with more than 20 MB s^{-1} .

Partition cloning is similar to general backup and restore operations. The differences between logical and physical backup are examined in [6]. We wanted our tool to remain operating-system and file-system independent and therefore we work with raw disk partitions ignoring their filesystems and their content.

Another previous study of software distribution [7] presents a protocol and a tool to distribute data to a large number of machines while putting a minimal load on the network (i.e. executing in the background). The described tool uses unicast, multicast and broadcast protocols depending on the capabilities and the location of the receivers. The different protocols drastically reduce the network usage of the tool, but also prevent the multicast from reaching near maximal speeds.

Pushing the protocols for reliable multicast over unreliable physical network towards higher speeds leads to a great variation in the perceived bandwidth, even with moderate packet loss rates, as shown in [8]. Known solutions for reliable multicast (such as [9]) require flow-control and retransmission protocols to be implemented in the application. Most of the multicast protocol work is geared to distribute audio and video streams with low delay and jitter rather than to optimize bulk data transfers at a high burst rate.

The model for partition cast is based on similar ideas presented in the throughput-oriented copy-transfer model for MPP computers [10].

A few commercial products are available for operating system installation by cloning, such as Norton Ghost [11], ImageCast [12] or DriveImagePro [13]. All these tools are capable of replicating a whole disk or individual partitions and generating compressed image files, but none of them can adapt to different networks or the different performance characteristics of the computers in PC clusters. Commercial tools also depend on the operating- and the file system, since they use knowledge of the installed operating system and file systems to provide additional services such as resizing partitions, installing individual software packages and performing customizations.



Figure 1. An active node (left) with an in-degree of 1 and an out-degree of 2 as well as a passive node (right) with an in- and out-degree of 3.

An operating system independent open source approach is desired to support partition cast for maintenance in Beowulf installations [14]. Other applications of our tool could include *presentation-*, *database-* or *screen-image cast* for new applications in distributed data mining, collaborative work or remote tutoring on clusters of PCs. An early survey about research in that area including video-cast for clusters of PCs was done in the Tiger project [15].

2. A MODEL FOR PARTITION-CAST IN CLUSTERS

In this section we present a modelling scheme that allows to find the most efficient logical topology to distribute data streams.

2.1. Node types

We divide the nodes of a system into two categories, active nodes which duplicate a data stream and passive nodes which can only route data streams. The two node types are shown in Figure 1.

Active node. A node which is able to duplicate a data stream is called an active node. Active nodes that participate in the partition cast store the received data stream on the local disk.

An active node has at least an in-degree of 1 and is capable of passing the data stream further to one or more nodes (out-degree) by acting as a T-pipe.

Passive node. A passive node is a node in the physical network that can neither duplicate nor store a copy of the data stream. Passive nodes can pass one or more streams between active nodes in the network.

Partition cast requires *reliable* data streams with flow control. Gigabit Ethernet switches only provide *unreliable* multicast facilities and must therefore be modelled as passive switches that only route TCP/IP point-to-point connections. Incorporating intelligent network switches or genuine broadcast media (like Coax Ethernet or Hubs) could be achieved by making them active nodes and modelling them at the logical level. Such is an option for expensive Gigabit ATM switches that feature multicast capability on logical channels with separate flow control or for simple multicast enhanced switches. In the later case a special end-to-end multicast protocol is used to make multicast data transfers reliable.

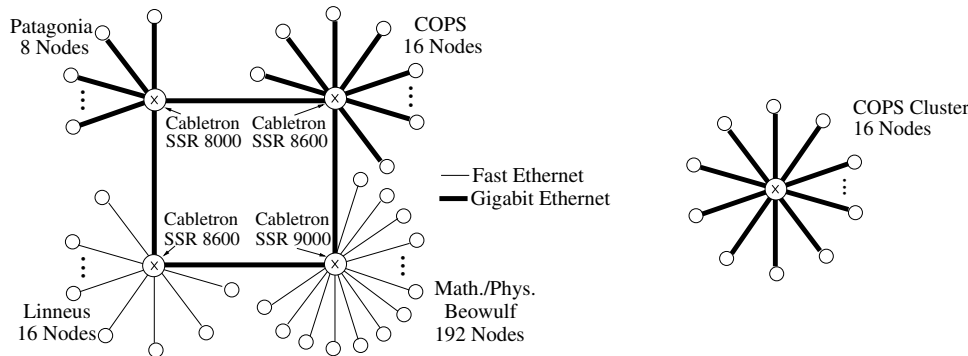


Figure 2. Physical network topologies of the ETH meta-cluster (left) and the simple sub-cluster with one central switch (right).

2.2. Network types

The different subsystems involved in a partition-cast must be specialized to transfer long data streams rather than short messages. Partitions are fairly large entities and our model is therefore purely bandwidth-oriented. We start our modelling process by investigating the topology of the physical network determining and recording the installed link and switch capacities.

Physical network. The physical network topology is a graph given by the cables, nodes and switches installed. The vertices are labeled by the maximal switching capacity of a node, the edges by the maximal link speeds.

The model itself captures a wide variety of networks including hierarchical topologies with multiple switches. Figure 2 shows the physical topology of the meta-cluster installed at ETH Zürich and the topology of our simple sub-cluster testbed. The sub-cluster testbed is built with a single central Gigabit Ethernet switch with *full duplex* point-to-point links to all the nodes. The switch has also enough Fast Ethernet ports to accommodate all cluster nodes at the low speed. Clusters of PCs are normally built with simple and fast layer-2 switches like our Cabletron Smart Switch Routers. In our case the backplane capacity for a 24 port switch is at 4 GB s^{-1} and never results in a bottleneck.

Our goal is to combine several subsystems of the participating machines in the most efficient way for an optimal partition-cast, so that the cloning of operating system images can be completed as quickly as possible. We therefore define different setups of logical networks.

Logical network. The logical network represents a connection scheme, that is embedded into a physical network. A spanning tree of TCP/IP connections routes the stream of a partition cast to all participating nodes. Unlike the physical network, the logical network must provide reliable transport and flow control over its channels.

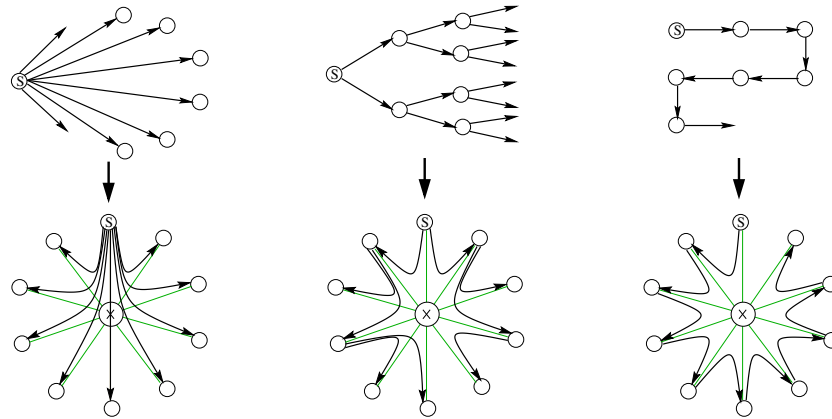


Figure 3. Logical network topologies (top) describing logical channels (star, n -ary spanning tree, multi-drop-chain) and their embedding in the physical networks.

Star. A logical network with one central server, that establishes a separate logical channel to all n other nodes. This logical network suffers heavy congestion on the outgoing link of the server.

n -ary spanning tree. Eliminates the server bottleneck by using an n -ary spanning tree topology spanning all nodes to be cloned. This approach requires active T-nodes which receive the data, store it to disk and pass it further to up to n next nodes in the tree.

Multi-drop-chain. A degenerated, specialized tree (unary case) where each active node stores a copy of the stream to disk and passes the data to just one further node. The chain is spanning all nodes to be cloned.

Figure 3 shows the above described topologies as well as their embedding in the physical networks. We assume that the central switch is a passive node and that it cannot duplicate a partition cast stream.

2.3. Capacity model

Our model for maximal throughput is based on capacity constraints expressed through a number of inequalities. These inequalities exist for active nodes, passive nodes and links, i.e. the edges in the physical net. As the bandwidth will be the limiting factor, all subsystems can be characterized by the maximal bandwidth they achieve in an isolated transfer. The extended model further introduces some more constraints e.g. for the CPU and the memory system bandwidth in a node (see Section 2.5).

Reliable transfer premise. We are looking for the fastest possible bandwidth with which we can stream data to a given number of active nodes. Since there is flow control, we know that the bandwidth b of the stream is the same in the whole system.

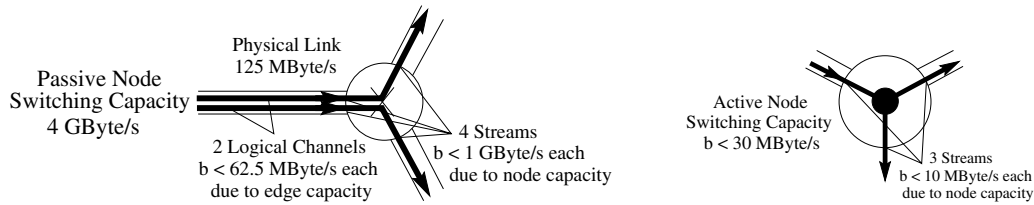


Figure 4. Edge capacities exist for the physical and logical network, node capacities for each in- and out-stream of a node.

Fair sharing of links. We assume that the flow control protocol eventually leads to a stable system and that the links or the nodes, dealing with the stream, allocate the bandwidth evenly and at a precise fraction of the capacity.

Both assumptions hold in the basic model and will be slightly extended in the refined model, which can capture raw and compressed streams at different rates simultaneously.

Edge capacity. Defines a maximum streaming capacity for each physical link and logical channel (see Figure 4).

As the physical links normally operate in *full duplex mode*, the inbound- and outbound-channels can be treated separately. If the logical-to-physical mapping suggests more than one logical channel over a single physical link, its capacity is evenly shared between them. Therefore the capacity is split in equal parts by dividing the link capacity through the number of channels that are mapped to the same physical link.

Example. For a binary tree with in-degree 1 and out-degree 2 mapped to one physical Gigabit Ethernet link the bandwidth of a stream has to comply with the following edge inequality:

$${}_1E_2 : b < 125, \quad 2b < 125 \rightarrow b < \frac{125}{2} \quad (1)$$

Node capacity. Is given by a switching capacity a node can provide, divided through the number of streams it handles.

The switching capacity of a node can be measured experimentally (by parameter fitting) or be derived directly from data of the node computer through a detailed model of critical resources. The experimental approach provides a specific limit value for each type of active node in the network, i.e. the maximal *switching capacity*. Fitting all our measurements resulted in a total switching capacity of 30 MB s^{-1} for our active nodes running on a 400 MHz Pentium II based cluster node. The switching capacity of our passive node, the 24 port Gigabit Ethernet switch is about 4 GB s^{-1} —much higher than needed for a partition cast.



2.4. Model algorithm

With the model described above we are now able to evaluate the different logical network alternatives described earlier in this section. The algorithm for evaluation of the model includes the following steps.

Algorithm basicmodel

- (1) Chose the physical network topology.
- (2) Chose the logical network topology.
- (3) Determine the mapping and the edge congestions.
- (4) **For all** edges:
 determine in-degree, out-degree of nodes attached to edge;
 evaluate channel capacity (according to logical net).
- (5) **For all** nodes:
 determine in-degree, out-degree and disk transfer of the node;
 evaluate node capacity.
- (6) Solve system of inequalities and find global minimum.
- (7) Return minimum as achievable throughput.

Example. We compare a multi-drop-chain versus the n -ary spanning tree structure for Gigabit Ethernet as well as for Fast Ethernet. The chain topology with all active nodes with in-degree i and out-degree o of exactly one (except for the source and the sink) can be considered as a special case of an unary tree (or Hamiltonian path) spanning all the active nodes receiving the partition cast.

- *Topology.* We evaluate the logical n -ary tree topology of Figure 3 with five nodes (and a streaming server) mapped on our simple physical network with a central switch of Figure 2. The out-degree shall be variable from 1 to 5, multi-drop-chain to star.
- *Edge Capacity.* The in-degree is always 1. The out-degree over one physical link varies between 1 for the multi-drop-chain and 5 for the star which leads to the following inequality:

$${}_1E_o : ob < 125 \rightarrow b < \frac{125}{o} \quad \text{for Gigabit Ethernet} \tag{2}$$

$${}_1E_o : ob < 12.5 \rightarrow b < \frac{12.5}{o} \quad \text{for Fast Ethernet} \tag{3}$$

- *Node capacity N :* For the active node we take the evaluated capacity of 30 MB s^{-1} with the given in-degree and out-degree and a disk write:

$$N_{1,o,1} : (1 + o + 1)b < 30 \rightarrow b < \frac{30}{(1 + o + 1)} \tag{4}$$

We now label all connections of the logical network with the maximal capacities and run the algorithm to find a global minimum of achievable throughput. The evaluation of the global minimum indicates that for Gigabit Ethernet the switching capacity of the active node is the bottleneck for the multi-drop-chain and for the n -ary trees. But for the slower links of a Fast Ethernet the n -ary tree the network rapidly becomes a bottleneck as we move to higher branching factors. Section 4 gives a detailed comparison of modelled and measured values for all cases considered.

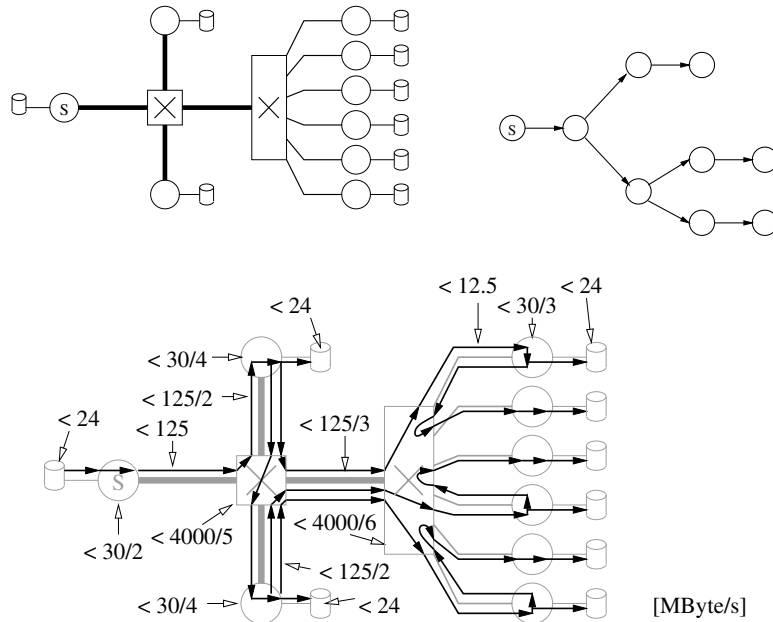


Figure 5. Example network with physical network (top left), logical network (top right) and the embedding of the logical in the physical network (bottom). The bottom figure also shows the edge and node congestions as well as the trivial limits of the hard drives. All numbers are in MB s^{-1} .

Example. In a second example we apply our model to a slightly complex network as shown in Figure 5. The data stream is sent from one server to eight clients. The server and two of the clients have direct Gigabit Ethernet links to a first switch, which in turn has a fourth Gigabit Ethernet link to a second switch. This switch connects to the remaining six clients with Fast Ethernet. We use the tree-based topology in Figure 5 as a logical network topology. One possible mapping of the logical into the physical network is shown on the bottom of the figure. In a next step, all the edge and node congestions are determined, according to the number of streams handled. This results in a number of inequalities which are also shown in the figure. Solving the inequalities of the entire model results in a maximal streaming bandwidth of 7.5 MB s^{-1} , limited by the switching capacity of the two nodes with an out-degree of two.

2.5. A more detailed model for an active node

The basic model considered two different resources: link capacity and switching capacity. The link speeds and the switch capacity of the passive node were taken from the physical data sheets of

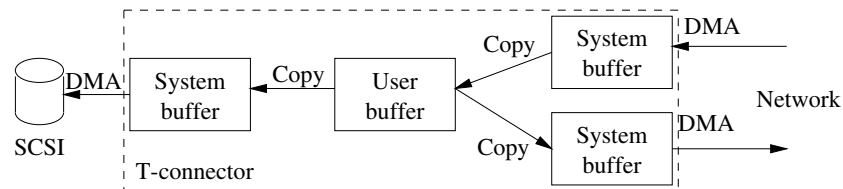


Figure 6. Schematic data flow of an active node running the Dolly [16] client.

the networking equipment, while the total switching capacity of an active node was obtained from measurements by a parameter fit. Link and switching capacity can only lead the optimization towards a graph theoretical discussion and will only be relevant to cases that have extremely low link bandwidth and high processing power or to systems that are trivially limited by disk speed. For clusters of PCs with high speed interconnects this is normally not the case and the situation is much more complex. Moving data through I/O buses and memory systems at full GB s^{-1} speed remains a major challenge in cluster computing. Most systems of today are nearly balanced between CPU performance, memory system and communication speed and some interesting tradeoffs can be observed. As indicated before, several options exist to trade off the processing power in the active node against a reduction of the load on the network. Among them are data compression or advanced protocol processing that turns some unreliable broadcast capabilities of Ethernet switches into a reliable multicast.

For a better model of an active node we consider the data streams *within* an active node and evaluate several resource constraints. For a typical client the node activity comprises receiving data from the network and writing partition images to the disk. We assume a 'one copy' TCP/IP protocol stack as provided by standard Linux. In addition to the source and sink nodes the tree and multi-drop chain topologies require active nodes that store a data stream and forward one or more copies of the data streams back into the network. Figure 6 gives a schematic data flow in an active node capable of duplicating a data stream.

2.6. Modelling the limiting resources in an active node

The switching capacity of an active node is modelled by the two limits of the network and four additional resource limits within the active node.

Link capacity. As taken from the physical specifications of the network technology (125 MB s^{-1} (Gigabit Ethernet) or 12.5 MB s^{-1} (Fast Ethernet) on current systems).

Switch capacity of passive nodes. As taken from the physical specifications of the network hardware (2 or 4 GB s^{-1} depending on the Cabletron Smart Switch Router model, 8000 or 8600).

Disk system. Similar to a link capacity in the basic model (24 MB s^{-1} for a Seagate Cheetah 10000 rpm disk).

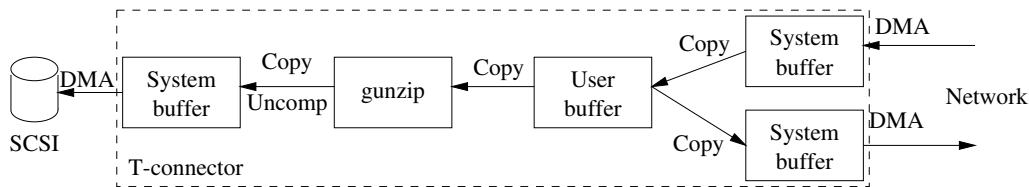


Figure 7. Schematic data flow of a *Dolly* client with data decompression.

I/O bus capacity. The sum of data streams traversing the I/O bus must be less than its capacity (132 MB s^{-1} on current, 32-bit PCI bus-based PC cluster nodes).

Memory system capacity. The sum of the data streams to and from the memory system must be less than the memory system capacity (180 MB s^{-1} on current systems with the Intel 440 BX chipset).

CPU utilization. The processing power required for the data streams at the different stages. For each operation, the fraction coefficient $1/a_1, 1/a_2, \dots, 1/a_n$ indicates a proportion corresponding to the share in CPU-time required by a subsystem to handle a stream flowing through all subsystems of a node. The coefficients a_1, a_2, \dots, a_n thereby correspond to the maximal speed of the individual subsystems handling a stream with exclusive use of the CPU. The sum of the fractions of CPU use must be < 1 ($= 100\%$) (Fractions considered: 80 MB s^{-1} SCSI transfer, 90 MB s^{-1} internal copy memory to memory, 60 MB s^{-1} send or receive over Gigabit Ethernet, 10 MB s^{-1} to decompress a data stream for a current 400 MHz single CPU cluster node).

Limitations on the four latter resources result in constraint inequalities for the maximal throughput achievable through an active node. The modelling algorithm determines and posts all constraining limits in the same manner as described in the example with a single switching capacity. The constraint over the edges of a logical network can be evaluated then into the maximum achievable throughput considering all limiting resources.

2.7. Dealing with compressed images

Partition images or multimedia presentations can be stored and distributed in compressed form. This reduces network load but puts an additional burden on the CPUs in the active nodes. Compressing and uncompressing is introduced into the model by an additional data copy to a *gunzip* process, which uncompresses data with an output data rate of about 10 MB s^{-1} (see Figure 7).

The workload is defined in raw data bytes to be distributed and the throughput rates are calculated in terms of the uncompressed data stream. For constraints inequalities involving the compressed data stream the throughput must be adjusted by the compression factor c . Hardware supported multicast could be modeled in a similar manner. For multicast, the network would be enhanced by newly introduced active switches, but a reliable multicast flow control protocol module must be added at



the endpoints and would consume a certain amount of CPU performance and memory system capacity (just like a compression module).

Example. We model the switching capacity of an active node for a multi-drop chain with Fast Ethernet and compression.

From the flow chart in Figure 7 we derive one network send and one receive stream, one disk write, three crossings of the I/O bus, eleven streams from and to buffer memory, one compression module and four internal copies of the data stream.

This leads to the following constraints for the maximal achievable throughput b :

$$\begin{aligned} \frac{b}{c} &< 12.5 \text{ MB s}^{-1} && \text{link for receive} \\ \frac{b}{c} &< 12.5 \text{ MB s}^{-1} && \text{link for send} \\ b &< 24 \text{ MB s}^{-1} && \text{SCSI Disk} \\ \frac{2b}{c} + b &< 132 \text{ MB s}^{-1} && \text{i/o bus, PCI} \\ \frac{8b}{c} + 3b &< 180 \text{ MB s}^{-1} && \text{memory system} \\ \left(\frac{2}{60c} + \frac{1}{80} + \frac{3}{90c} + \frac{1}{90} + \frac{c}{10} \right) b &< 1 \text{ (100\%)} && \text{CPU utilization} \end{aligned}$$

For a compression factor of $c = 2$, an active node in this configuration can handle 3.9 MB s^{-1} . The limiting resource is the CPU utilization.

2.8. Modelling multicast capabilities of the switches

The workstation based active nodes use TCP/IP connections to distribute the data. If IP multicast were used, it would be unreliable without flow control. It would be up to the endpoints to deal with flow control and retransmission of lost data. A multicast flow control module is modelled similar to a compression module putting load on the memory system (extra copies for retransmission which are likely to occur more often if many receivers need to be coordinated) and on the CPU (protocol processing). In a previous study one of the authors [8] implemented several well known approaches [9,17]. Unfortunately the performance reached in those implementations was not able to sustain data rates that come even close to the 10 MB s^{-1} observed on our cluster machines.

3. DIFFERENCES IN THE IMPLEMENTATIONS

Our first experimental setup for partition cast uses a simple file sharing service like Network File System (NFS) with transfers over a UDP/IP network resulting in a *star* topology. The NFS server exports partition images to all the clients in the cluster. A command-line script reads the images from a network mounted drive, possibly decompresses the data and writes it to the raw partitions of the local

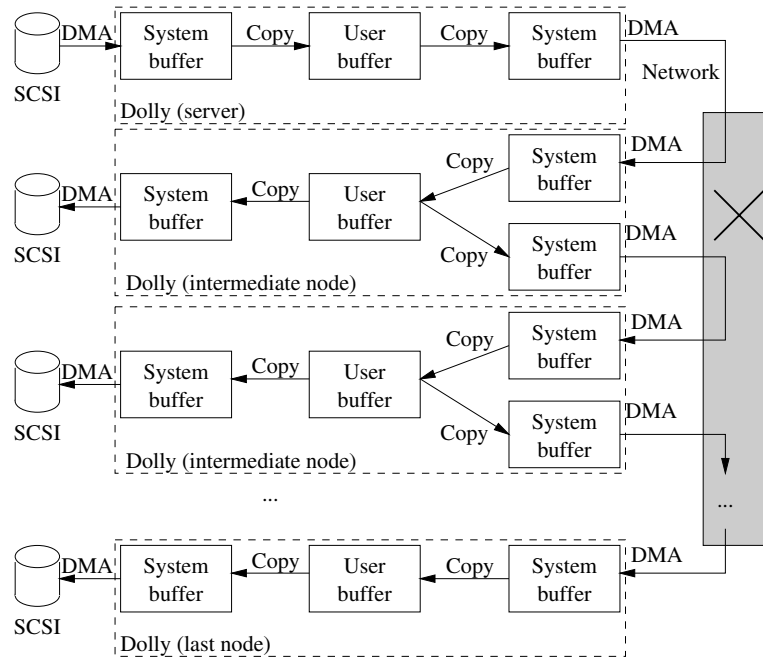


Figure 8. Schematic data flow for a number of nodes running the Dolly client.

disk. Because of the asymmetric role of one server and many clients, this approach does not scale well, since the single high speed Gigabit Ethernet can be saturated while serving a relatively small number of clients (see performance numbers in Section 4). Although this approach might look a bit naive to an experienced system architect, it is simple, highly robust and supported by every operating system. Furthermore, a single client failure or a congested network can be easily dealt with.

As a second setup, we considered putting together active clients in a n -ary spanning tree topology. This method works with the standard TCP point-to-point connections and uses the excellent switching capability of the Gigabit Ethernet switch backplane. A partition cloning program (called Dolly [16]) runs on each active node. A simple server program reads the data from disk on the image server and sends the stream over a TCP connection to the first few clients. The clients receive the stream, write the data to the local disk and send it on to the next clients in the tree. The machines are connected in an n -ary spanning tree, eliminating the bottleneck of the server link accessing the network.

Finally for the third and optimal solution the same Dolly client program can be used with a local copy to disk and just one further downstream client to serve to. The topology turns into a highly degraded unary spanning tree. We call this logical network a *multi-drop chain*. As shown in Figure 8, the flow in all the nodes does *not* depend on the number of nodes. This independency makes this system highly scalable.

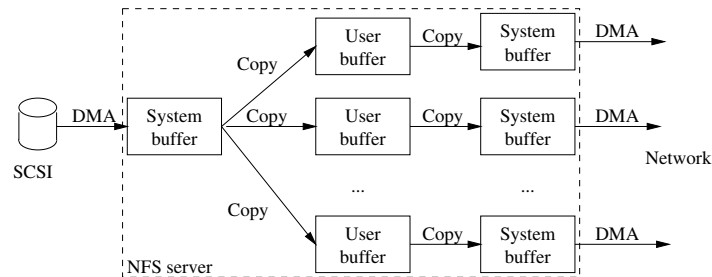


Figure 9. Schematic data flow of a NFS server.

An obvious topological alternative would be a true physical spanning tree using the multicasting feature of the networking hardware. With this option the server would only source one stream and the clients would only sink a single stream. The protocols and schemes required for reliable and robust multicast are neither trivial to implement nor included in common commodity operating systems and often depend on the multicast capabilities of the network hardware.

4. EVALUATION OF PARTITION CAST

In this section we provide measurements of partition casts in different logical topologies (as described in Section 2) with compressed and uncompressed partition images. The partition to be distributed to the target machines is a 2 GB Windows NT partition. The compressed image file is about 1 GB in size, resulting in a compression factor of 2.

The initial version of our partition-cast tool uses only existing OS services and therefore applies a simplistic star topology approach. It comprises a NFS server which exports the partition images to all the clients. The clients access the images over the network using NFS, possibly uncompressing the images and finally writing the data to the target partition. Figure 9 depicts the data flow inside the NFS server node. The flow obviously depends on the number of clients. The results of this experiment are shown in the left chart in Figure 10 (the execution time for each client machine is logged to show the variability due to congestion). The figure shows two essential results. (1) The more clients need to receive the data, the more time the distribution takes (resulting in a lower total bandwidth of the system). The bandwidth is limited by the edge capacity of the server. (2) Compression helps to increase the bandwidth for a star topology. As the edge capacity is the limiting factor, the nodes have enough CPU, memory and I/O capacity left to uncompress the incoming stream at full speed, thereby increasing the total bandwidth of the channel.

A second experiment uses an n -ary spanning tree structure. This topology was implemented in the small program Dolly [16] which acts as an active node. The program reads the partition image on the server and sends it to the first n clients. The clients write the incoming data to the target partition on disk and send the data to the next clients. The out-degree is the same for all nodes (provided there are

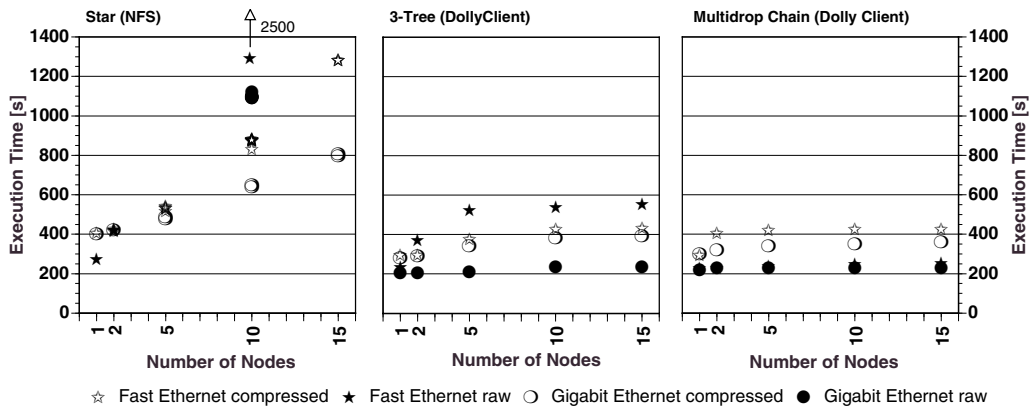


Figure 10. Total execution times for distributing a 2 GB Windows NT operating-system partition simultaneously to one, two, five, 10, and 15 machines by partition cloning with NFS based star topology, the Dolly based 3-tree and multi-drop-chain topologies on the Patagonia cluster. The star topology run with 10 clients using raw transfers over Fast Ethernet resulted in execution times around 2500 s as the NFS server's disk started thrashing.

enough successor-nodes) and can be specified at runtime. The results for a 3-ary tree are shown in the middle chart of Figure 10. For Fast Ethernet the execution time increases rapidly for a small number of clients until the number of clients (and therefore the number of successor-nodes of the server) reaches the out-degree. As soon as the number of clients is larger than the out-degree, the execution times stay roughly constant. For this network speed, the edge capacity remains a bottleneck, resulting in increased execution times for higher out-degree. In the case of Gigabit Ethernet, the link speed (the edge capacity) is high enough to satisfy an out-degree of up to 5 without the edge capacity becoming the bottleneck. The primary bottleneck in this case is given by the memory capacity of the node.

For the third experiment we use Dolly [16] to cast the partition using a multi-drop chain. The results are shown in the right chart of Figure 10. The speeds measured indicate that the execution time for this partition cast is nearly independent of the number of clients. This independence results from the fact that in a multi-drop chain configuration, the edge capacity is no longer a bottleneck as every edge carries at most one stream per direction. The new bottleneck is the nodes' memory system. The memory bottleneck also explains why compression results in a lower bandwidth for the channel (decompressing data requires more memory copy operations in UNIX pipes from/to the gunzip process).

Figure 11 shows the total, aggregate bandwidth of data transfers to all disk drives in the system with the three experiments. The figure indicates that the aggregate bandwidth of the NFS-approach increases only modestly with the number of clients while the multi-drop chain scales perfectly. The 3-ary-tree approach also scales perfectly, but increases at a lower rate. The numbers for the NFS approach clearly max out with the transfer bandwidth of the servers network interface reaching the edge capacity: our NFS-server can deliver a maximum of about 20 MB s^{-1} over Gigabit Ethernet and about 10 MB s^{-1} over Fast Ethernet (note that we are using compressed data for the NFS approach in the above figure

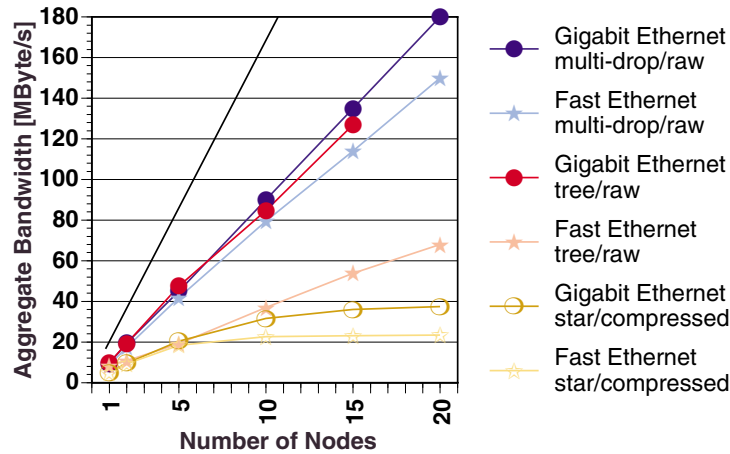


Figure 11. Total (aggregate) transfer bandwidth achieved in distributing a 2 GB Windows NT operating-system partition simultaneously to a number of hosts by partition cloning in the Patagonia cluster.

Table I. Predicted and measured bandwidths for a partition cast over a logical chain and different tree topologies for uncompressed and compressed images. All values are given in MB s⁻¹.

Topology	Network		Fast Ethernet bandwidth		Gigabit Ethernet bandwidth	
	Out-degree	Compression	Extended model	Measured	Extended model	Measured
Multi-drop-chain	1	No	12.5	8.8	12.6	9.0
Multi-drop-chain	1	Yes	3.9	4.9	3.9	6.1
2-tree	2	No	6.3	5.4	9.3	8.2
3-tree	3	No	4.2	3.8	7.4	8.0
Star	5	No	2.5	2.3	5.3	6.3
Star	5	Yes	3.0	3.6	3.0	4.1

thereby doubling the bandwidth). The predicted bandwidths are compared with measured values in our cluster in Table I.

Recently we put our next generation cluster hardware into service. The 16 dual Pentium II nodes at 400 MHz are replaced with 16 dual Pentium III nodes running at 1 GHz. A re-measurement on this new hardware demonstrated the predictive power of our model. A microbenchmark [18,19] measured the memory copy bandwidth at 200 MB s⁻¹, while the TCP transfer bandwidth with the new Linux 2.4.1 kernel is about 110 MB s⁻¹. Using these numbers our model predicts a multi-drop chain bandwidth of 25.4 MB s⁻¹. A preliminary measurement showed that the measured bandwidth is 26 MB s⁻¹ and falls within the range of a few MB s⁻¹.



5. CONCLUSION

In this paper we investigated the problem of a partition-cast in clusters of PCs. We showed that optimizing a partition-cast or any distribution of a large block of raw data leads to some interesting tradeoffs between network parameters and node parameters.

In a simple analytical model we captured the network parameters (link speed and topology) as well as the basic processor resources (memory system, CPU, I/O bus bandwidth) at the intermediate nodes that are forwarding our multicast streams. The calculation of the model for our sample PC cluster pointed towards an optimal solution using uncompressed streams of raw data, forwarded along a linear multi-drop chain embedded into the Gigabit Ethernet. The optimal configuration was limited by the CPU performance in the nodes and its performance was correctly predicted at about one third of the maximal disk speed. The alternative of a star topology with one server and 24 clients suffered from heavy link congestion at the server link, while the different n -ary spanning tree solutions were slower due to the resource limitations in the intermediate nodes, that could not replicate the data into multiple streams efficiently enough. Compression resulted in a lower network utilization but was slower due to the higher CPU utilization. The existing protocols for reliable multicast on top of unreliable best-effort hardware broadcast in the Ethernet switch were not fast enough to keep up with our multi-drop solution using simple, reliable TCP/IP connections.

The resulting partition casting tool is capable of transferring a 2 GByte Windows NT operating system installation to 24 workstations in less than 5 min while transferring data at a sustained rate of about 9 MB s^{-1} per node. Fast partition cast permits the distribution of entire installations in a short time, adding flexibility to a cluster of PCs to do different tasks at different times. A setup for efficient multicast also results in easier maintenance and enhances the robustness against slowly degrading software installations in a PC cluster.

REFERENCES

1. Rauch F, Kurmann C, Stricker T, Müller BM. Patagonia—A dual use cluster of PCS for computation and education. *Proceedings of the 2nd Workshop on Cluster Computing*, Karlsruhe, Germany, March 1999.
2. Hellwagner H, Reinefeld A (eds.) *SCI: Scalable Coherent Interface—Architecture and Software for High-Performance Computer Clusters (Lecture Notes in Computer Science, vol. 1734)*. Springer: Berlin, 1999.
3. Dolphin Interconnect Solutions. The Dolphin SCI Interconnect, white paper, February 1996. <http://www.dolphinics.com/>.
4. Boden NJ, Felderman RE, Kulawik AE, Seitz CL, Seizovic JN, Su W-K. Myrinet—A gigabit per second local area network. *IEEE-Micro* 1995; **15**(1):29–36.
5. Seifert R. *Gigabit Ethernet: Technology and Applications for High-Speed LANs*. Addison-Wesley, 1998.
6. Hutchinson NC, Manley S, Federwisch M, Harris G, Hitz D, Kleiman S, Malley SO. Logical vs. physical file system backup. *Proceedings of the 3rd Symposium on Operating Systems Design and Implementation*, New Orleans, Louisiana, February 1999. The USENIX Association, 1999; 239–249.
7. Kotsopoulos S, Cooperstock J. Why use a fishing line when you have a net? an adaptive multicast data distribution protocol. *Proceedings of the USENIX 1996 Annual Technical Conference*, San Diego, California, January 1996. The USENIX Association, 1996.
8. Rauch F. Zuverlässiges Multicastprotokoll. *Master's Thesis*, ETH Zürich, 1997. English title: Reliable multicast protocol. <http://www.cs.inf.ethz.ch/>.
9. Floyd S, Jacobson V, McCanne S, Zhang L, Liu C-G. A reliable multicast framework for lightweight sessions and application level framing. *Proceedings of ACM SIGCOMM '95*, August 1995; 342–356.
10. Stricker T, Gross T. Optimizing memory system performance for communication in parallel computers. *Proceedings of the 22nd International Symposium on Computer Architecture*, Santa Margherita di Ligure, June 1995. ACM, 1995; 308–319.
11. Symantec Ghost by Symantec Corp., Cupertino, CA. <http://www.symantec.com/>.
12. ImageCast by StorageSoft, Inc., Louisville, CO. <http://www.storage-soft.com/>.



13. DriveImage Pro by PowerQuest Corp., Orem, UT. <http://www.powerquest.com/>.
14. Becker DJ, Sterling T, Savarese D, Dorband JE, Ranawake UA, Packer CV. Beowulf: a parallel workstation for scientific computation. *Proceedings of the 1995 ICPP Workshop on Challenges for Parallel Processing*, Oconomowoc, WI, August 1995. CRC Press, 1995.
15. Bolosky WJ, Barrera III JS, Draves RP, Fitzgerald RP, Gibson GA, Jones MB, Levi SP, Myhrvold NP, Rashid RF. The Tiger video fileserver. *Proceedings of the 6th International Workshop on Network and Operating System Support for Digital Audio and Video*, Zushi, Japan, April 1996. IEEE Computer Society, 1996.
16. Rauch F. Dolly—a program to clone harddisks or partitions over a fast switched network. <http://www.cs.inf.ethz.ch/CoPs/patagonia/>.
17. Paul S, Sabnani KK, Kristol DM. Multicast transport protocols for high speed networks. *Proceedings of International Conference on Network Protocols*. IEEE Computer Society Press, 1994; 4–14.
18. Kurmann C, Stricker T. Characterizing memory system performance for local and remote accesses in high-end SMPs, low-end SMPs and clusters of SMPs. *Proceedings of the 7th Workshop on Scalable Memory Multiprocessors ACM Transactions on Computer Systems, held in conjunction with ISCA98*, Barcelona, Spain, June 1998.
19. Kurmann C, Stricker T. ECT—Extended copy transfer characterization. <http://www.cs.inf.ethz.ch/CoPs/ECT/>.