

Article development led by [acmqueue](http://queue.acm.org)
queue.acm.org

A good user experience depends on predictable performance within the data-center network.

BY DENNIS ABTS AND BOB FELDERMAN

A Guided Tour of Data-Center Networking

THE MAGIC OF the cloud is that it is always on and always available from anywhere. Users have come to expect that services are there when they need them. A data center (or warehouse-scale computer) is the nexus from which all the services flow. It is often housed in a nondescript warehouse-sized building bearing no indication of what lies inside. Amidst the whirring fans and refrigerator-sized computer racks is a tapestry of electrical cables and fiber optics weaving everything together—the data-center network. This article provides a “guided tour” through the principles and central ideas surrounding the network at the heart of a data center—the modern-day loom that weaves the digital fabric of the Internet.

Large-scale parallel computers are grounded in HPC (high-performance computing) where kilo-processor

systems were available 15 years ago. HPC systems rely on fast (low-latency) and efficient interconnection networks capable of providing both high bandwidth and efficient messaging for fine-grained (for example, cache-line size) communication. This zealous attention to performance and low latency migrated to financial enterprise systems where a fraction of a microsecond can make a difference in the value of a transaction.

In recent years, Ethernet networks have made significant inroads into bridging the performance and scalability gap between capacity-oriented clusters built using COTS (commodity-off-the-shelf) components and purpose-built custom system architectures. This is evident from the growth of Ethernet as a cluster interconnect on the Top500 list of most powerful computers (top500.org). A decade ago high-performance networks were mostly custom and proprietary interconnects, and Ethernet was used by only 2% of the Top500 systems. Today, however, more than 42% of the most powerful computers are using Gigabit Ethernet, according to the November 2011 list of Top500 computers. A close second place is InfiniBand, which is used by about 40% of the systems. These standards-based interconnects combined with economies of scale provide the genetic material of a modern data-center network.

A modern data center,^{13,17,24} as shown in Figure 1, is home to tens of thousands of hosts, each consisting of one or more processors, memory, network interface, and local high-speed I/O (disk or flash). Compute resources are packaged into racks and allocated as *clusters* consisting of thousands of hosts that are tightly connected with a high-bandwidth network. While the network plays a central role in the overall system performance, it typically represents only 10%–15% of the cluster cost. Be careful not to confuse cost with value—the network is to a cluster computer what the central nervous system is to the human body.



Figure 1. An example data center and warehouse-scale computer.

Each cluster is homogeneous in both the processor type and speed. The thousands of hosts are orchestrated to exploit thread-level parallelism central to many Internet workloads as they divide incoming requests into parallel subtasks and weave together results from many subtasks across thousands of cores. In general, in order for the request to complete, all parallel subtasks must complete. As a result, the maximum response time of any one subtask will dictate the overall response time.²⁵ Even in the presence of abundant thread-level parallelism, the communication overhead imposed by the network and protocol stack can ultimately limit application performance as the effects of Amdahl's Law² creep in.

The high-level system architecture and programming model shape both the programmer's conceptual view and application usage. The latency and bandwidth "cost" of local (DRAM) and remote (network) memory references

are often baked into the application as programming trade-offs are made to optimize code for the underlying system architecture. In this way, an application organically grows within the confines of the system architecture.

The cluster-application usage model, either dedicated or shared among multiple applications, has a significant impact on SLAs (service-level agreements) and application performance. HPC applications typically use the system in a dedicated fashion to avoid contention from multiple applications and reduce the resulting variation in application performance. On the other hand, Web applications often rely on services sourced from multiple clusters, where each cluster may have several applications simultaneously running to increase overall system utilization. As a result, a data-center cluster may use virtualization for both performance and fault isolation, and Web applications are programmed with this sharing in mind.

Web applications such as search, email, and document collaboration are scheduled resources and run within a cluster.^{4,8} User-facing applications have soft real-time latency guarantees or SLAs that the application must meet. In this model, an application has approximately tens of milliseconds to reply to the user's request, which is subdivided and dispatched to worker threads within the cluster. The worker threads generate replies that are aggregated and returned to the user. Unfortunately, if a portion of the workflow does not execute in a timely manner, then it may exceed a specified timeout delay—as a result of network congestion, for example—and consequently some portion of the coalesced results will be unavailable and thus ignored. This needlessly wastes both CPU cycles and network bandwidth, and may adversely impact the computed result.

To reduce the likelihood of congestion, the network can be overprovi-

sioned by providing ample bandwidth for even antagonistic traffic patterns. Overprovisioning within large-scale networks is prohibitively expensive. Alternatively, implementing QoS (quality of service) policies to segregate traffic into distinct classes and provide performance isolation and high-level traffic engineering is a step toward ensuring application-level SLAs are satisfied. Most QoS implementations are implemented by switch and NIC (network interface controller) hardware where traffic is segregated based on priority explicitly marked by routers and hosts or implicitly steered using port ranges. The goal is the same: a high-performance network that provides predictable latency and bandwidth characteristics across varying traffic patterns.

Data-Center Traffic

Traffic within a data-center network is often measured and characterized according to *flows*, which are sequences of packets from a source to destination host. When referring to Internet protocols, a flow is further refined to include a specific source and destination port number and transport type—UDP or TCP, for example. Traffic is asymmetric with client-to-server requests being abundant but generally small. Server-to-client responses, however, tend to be larger flows; of course, this, too,

depends on the application. From the purview of the cluster, Internet traffic becomes highly aggregated, and as a result the *mean* of traffic flows says very little because aggregated traffic exhibits a high degree of variability and is non-Gaussian.¹⁶

As a result, a network that is only 10% utilized can see lots of packet discards when running a Web search. To understand individual flow characteristics better, applications are instrumented to “sample” messages and derive a distribution of traffic flows; this knowledge allows you to infer a taxonomy of network traffic and classify individual flows. The most common classification is bimodal, using the so-called “elephant” and “mice” classes.

Elephant flows have a large number of packets and are generally long lived; they exhibit “bursty” behavior with a large number of packets injected in the network over a short time. Traffic within a flow is generally ordered, which means elephant flows can create a set of “hotspot” links that can lead to tree saturation or discarded packets in networks that use lossy flow control. The performance impact from elephant flows can be significant. Despite the relatively low number of flows—say less than 1%—they can account for more than half the data volume on the network.

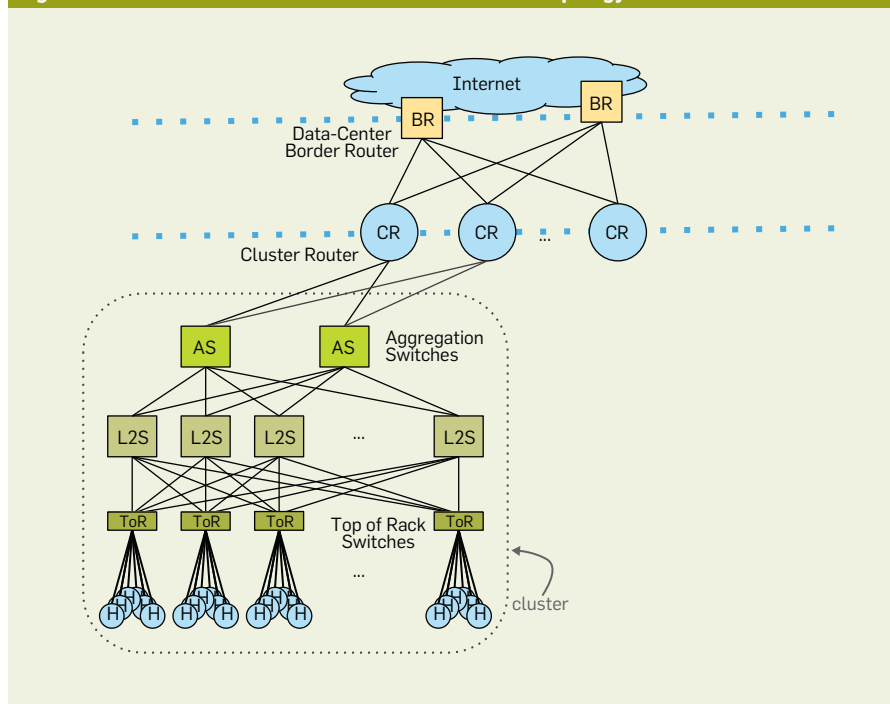
The transient load imbalance induced by elephant flows can adversely affect any innocent-bystander flows that are patiently waiting for a heavily utilized link common to both routes. For example, an elephant flow from A to B might share a common link with a flow from C to D. Any long-lived contention for the shared link increases the likelihood of discarding a packet from the C-to-D flow. Any packet discards will result in an unacknowledged packet at the sender’s transport layer and be retransmitted when the timeout period expires. Since the timeout period is generally one or two orders of magnitude more than the network’s round-trip time, this additional latency²² is a significant source of performance variation.³

Today’s typical multitiered data-center network²³ has a significant amount of *oversubscription*, where the hosts attached to the rack switch (that is, first tier) have significantly more—say an order of magnitude more—provisioned bandwidth between one another than they do with hosts in *other* racks. This *rack affinity* is necessary to reduce network cost and improve utilization. The traffic intensity emitted by each host fluctuates over time, and the transient load imbalance that results from this varying load can create contention and ultimately result in discarded packets for flow control. Traffic *between* clusters is typically less time critical and as a result can be staged and scheduled. Inter-cluster traffic is less orchestrated and consists of much larger payloads, whereas intra-cluster traffic is often fine-grained with bursty behavior. At the next level, between data centers, bandwidth is often very expensive over vast distances with highly regulated traffic streams and patterns so that expensive links are highly utilized. When congestion occurs the most important traffic gets access to the links. Understanding the granularity and distribution of network flows is essential to capacity planning and traffic engineering.

Data-Center Network Architecture

The network topology describes precisely how switches and hosts are interconnected. This is commonly represented as a graph in which vertices represent switches or hosts, and links

Figure 2. A conventional tree-like data-center network topology.




are the edges that connect them. The topology is central to both the performance and cost of the network. The topology affects a number of design trade-offs, including performance, system packaging, path diversity, and redundancy, which, in turn, affect the network's resilience to faults, average and maximum cable length, and, of course, cost.¹² The *Cisco Data Center Infrastructure 3.0 Design Guide*⁶ describes common practices based on a tree-like topology¹⁵ resembling early telephony networks proposed by Charles Clos,⁷ with bandwidth aggregation at different levels of the network.

A fat-tree or folded-Clos topology, similar to that shown in Figure 2, has an aggregate bandwidth that grows in proportion to the number of host ports in the system. A *scalable* network is one in which increasing the number of ports in the network should linearly increase the delivered bisection bandwidth. Scalability and reliability are inseparable since growing to large system size requires a robust network.


Network addressing. A host's *address* is how endpoints are identified in the network. Endpoints are distinguished from intermediate switching elements traversed en route since messages are created by and delivered to an endpoint. In the simplest terms, the address can be thought of as the numerical equivalent of a host name similar to that reported by the Unix `hostname` command.

An address is unique and must be represented in a canonical form that can be used by the *routing function* to determine where to route a packet. The switch inspects the packet header corresponding to the layer in which routing is performed—for example, IP address from layer 3 or Ethernet address from layer 2. Switching over Ethernet involves ARP (address resolution protocol) and RARP (reverse address resolution protocol) that broadcast messages on the layer 2 network to update local caches mapping layer 2 to layer 3 addresses and vice versa. Routing at layer 3 requires each switch to maintain a subnet *mask* and assign IP addresses statically or disseminate host addresses using DHCP (dynamic host configuration protocol), for example.

The layer 2 routing tables are automatically populated when a switch is



The high-level system architecture and programming model shape both the programmer's conceptual view and application usage.



plugged in and learns its identity and exchanges route information with its peers; however, the capacity of the packet-forwarding tables is limited to, say, 64K entries. Further, each layer 2 switch will participate in an STP (spanning tree protocol) or use the TRILL (transparent interconnect of lots of links) link-state protocol to exchange routing information and avoid transient routing loops that may arise while the link state is exchanged among peers. Neither layer 2 nor layer 3 routing is perfectly suited to data-center networks, so to overcome these limitations many new routing algorithms have been proposed (for example, PortLand^{1,18} and VL2¹¹).

Routing. The *routing* algorithm determines the path a packet traverses through the network. A packet's route, or path, through the network can be asserted when the message is created, called source routing, or may be asserted hop by hop in a distributed manner as a packet visits intermediate switches. Source routing requires that every endpoint know the prescribed path to reach all other endpoints, and each source-routed packet carries the full information to determine the set of port/link traversals from source to destination endpoint. As a result of this overhead and inflexible fault handling, source-routed packets are generally used only for topology discovery and network initialization, or during fault recovery when the state of a switch is unknown. A more flexible method of routing uses distributed *lookup tables* at each switch, as shown in Figure 3.

For example, consider a typical Ethernet switch. When a packet arrives at a switch input port, it uses fields from the packet header to index into a lookup table and determine the next hop, or egress port, from the current switch.

A good topology will have abundant *path diversity* in which multiple possible egress ports may exist, with each one leading to a distinct path. Path diversity in the topology may yield ECMP (equal-cost multipath) routing; in that case the routing algorithm attempts to load-balance the traffic flowing across the links by spreading traffic uniformly. To accomplish this *uniform spreading*, the routing function in the switch will *hash* several fields of the packet header to produce a deterministic egress port.

In the event of a link or switch failure, the routing algorithm will take advantage of path diversity in the network to find another path.

A path through the network is said to be *minimal* if no shorter (that is, fewer hops) path exists; of course, there may be multiple minimal paths. A fat-tree topology,¹⁵ for example, has multiple minimal paths between any two hosts, but a butterfly topology⁹ has only a single minimal path between any two hosts. Sometimes selecting a *non-minimal* path is advantageous—for example, to avoid congestion or to route around a fault. The length of a non-minimal path can range from $\text{min}+1$ up to the length of a Hamiltonian path visiting each switch exactly once. In general, the routing algorithm might consider non-minimal paths of a length that is one more than a minimal path, since considering *all* non-minimal paths would be prohibitively expensive.

Network Performance

Here, we discuss the etiquette for sharing the network resources—specifically, the physical links and buffer spaces are resources that require *flow control* to share them efficiently. Flow control is carried out at different levels of the network stack: data-link, network, transport layer, and possibly within the application itself for explicit coordination of resources. Flow control that occurs at lower levels of the communication stack is transparent to applications.

Flow control. Network-level flow control dictates how the input buffers at each switch or NIC are managed: store-and-forward, virtual cut-through,¹⁴ or wormhole,¹⁹ for example. To understand the performance implications of flow control better, you must first understand the total delay, T , a packet incurs:

$$T = H(\tau_r + L\tau_p) + \tau_s$$

H is the number of *hops* the packet takes through the network; τ_r is the fall-through latency of the switch, measured from the time the first flit (flow-control unit) arrives to when the first flit exits; and τ_p is the propagation delay through average cable length L . For short links—say, fewer than 10

meters—electrical signaling is cost-effective. Longer links, however, require fiber optics to communicate over the longer distances. Signal propagation in electrical signaling (5 nanoseconds per meter) is faster than it is in fiber (6 nanoseconds per meter).

Propagation delay through electrical cables occurs at subluminal speeds because of a frequency-dependent component at the surface of the conductor, or “skin effect,” in the cable. This limits the signal velocity to about three-quarters the speed of light in a vacuum. Signal propagation in optical fibers is even slower because of dielectric waveguides used to alter the refractive index profile so that higher-velocity components of the signal (such as shorter wavelengths) will travel longer distances and arrive at the same time as lower-velocity components, limiting the signal velocity to about two-thirds the speed of light in a vacuum. Optical signaling must also account for the time necessary to perform electrical-to-optical signal conversion, and vice versa.

The *average* cable length, L , is largely determined by the topology and the physical placement of system racks within the data center. The packet’s serialization latency, τ_s , is the time necessary to squeeze the packet onto a narrow serial channel and is determined by the *bit rate* of the channel. For example, a 1,500-byte Ethernet packet (frame) requires more than $12\mu\text{s}$ (ignoring any interframe gap time) to be squeezed onto a 1Gb/s link. With store-and-forward flow control, as its name suggests, a packet is buffered at each hop before the switch does anything with it:

$$T_{\text{sf}} = H(\tau_r + L\tau_p + \tau_s)$$

As a result, the serialization delay, τ_s , is incurred at *each* hop, instead of just at the destination endpoint as is the case with virtual cut-through and wormhole flow control. This can potentially add on the order of $100\mu\text{s}$ to the round-trip network delay in a data-center network.

A *stable* network monotonically delivers messages as shown by a characteristic throughput-load curve in Figure 4. In the absence of end-to-end flow control, however, the network

can become unstable, as illustrated by the dotted line in the figure, when the offered load exceeds the saturation point, α . The saturation point is the offered load beyond which the network is said to be *congested*. In response to this congestion, packets may be discarded to avoid overflowing an input buffer. This *lossy* flow control is commonplace in Ethernet networks.

Discarding packets, while conceptually simple and easy to implement, puts an onus on transport-level mechanisms to detect and retransmit lost packets. Note, packets that are lost or corrupted during transmission are handled by the same transport-level reliable delivery protocol. When the offered load is low (less than α), packet loss as a result of corruption is rare, so paying the relatively large penalty for transport-level retransmission is generally tolerable. Increased traffic (greater than α) and adversarial traffic patterns will cause packet discards after the switch’s input queue is exhausted. The resulting retransmission will only further exacerbate an already congested network, yielding an unstable network that performs poorly, as shown by the dotted line in Figure 4. Alternatively, with *lossless* flow control, when congestion arises packets may be blocked or held at the source until resources are available.

A *global congestion control* mechanism prevents the network from operating in the post-saturation region. Most networks use *end-to-end flow control*, such as TCP,⁵ which uses a windowing mechanism between pairs of source and sink in an attempt to match the source’s injection rate with the sink’s acceptance rate. TCP, however, is designed for reliable packet delivery, not necessarily timely packet delivery, and as a result, requires tuning (TCP congestion-control algorithms will auto-tune to find the right rate) to balance performance and avoid unnecessary packet duplication from eagerly retransmitting packets under heavy load.

Improving the network stack. Several decades ago the network designers of early workstations made trade-offs that led to a single TCP/IP/Ethernet network stack, whether communicating over a few meters or a few kilometers. As processor speed and density improved, the cost of network communication grew

relative to processor cycles, exposing the network stack as a critical latency bottleneck.²² This is, in part, the result of a user-kernel context switch in the TCP/IP/Ethernet stack—and possibly additional work to copy the message from the application buffer into the kernel buffer and back again at the receiver. A two-pronged hardware/software approach tackled this latency penalty: OS bypass, and zero copy, both of which are aimed at eliminating the user-kernel switch for every message and avoiding a redundant memory copy by allowing the network transport to grab the message payload directly from the user application buffers.

To ameliorate the performance impact of a user/kernel switch, OS bypass can be used to deposit a message directly into a user-application buffer. The application participates in the messaging protocol by spin-waiting on a doorbell memory location. Upon arrival, the NIC deposits the message contents in the user-application buffer, and then “rings” the doorbell to indicate message arrival by writing the offset into the buffer where the new message can be found. When the user thread detects the updated value, the incoming message is processed entirely from user space.

Zero-copy message-passing protocols avoid this additional memory copy from user to kernel space, and vice versa at the recipient. An interrupt signals the arrival of a message, and an interrupt handler services the new message and returns control to the user application. The interrupt latency—the time from when the interrupt is raised until control is handed to the interrupt handler—can be significant, especially if interrupt coalescing is used to amortize the latency penalty across multiple interrupts. Unfortunately, while interrupt coalescing improves message efficiency (that is, increased effective bandwidth), it does so at the cost of both increased message latency and latency variance.

Scalable, Manageable, And Flexible

In general, *cloud computing* requires two types of services: user-facing computation (for example, serving Web pages) and inward computation (for example, indexing, search, and map/reduce). Outward-facing functionality

can sometimes be done at the “border” of the Internet where commonly requested pages are cached and serviced by edge servers, while inward computation is generally carried out by a cluster in a data center with tightly coupled, orchestrated communication. User demand is diurnal for a geographic region; thus, multiple data centers are positioned around the globe to accommodate the varying demand. When possible, demand may be spread across nearby data centers to load-balance the traffic.

The sheer enormity of this computing infrastructure makes nimble deployment very challenging. Each cluster is built up rack by rack and tested as units (rack, top-of-rack switch, among others), as well as in its entirety with production-level workloads and traffic intensity.

The cluster ecosystem undergoes organic growth over its life span, propelled by the rapid evolution of soft-

ware—both applications and, to a lesser extent, the operating system. The fluid-like software demands of Web applications often consume the cluster resources that contain them, making *flexibility* a top priority in such a fluid system. For example, adding 10% additional storage capacity should mean adding no more than 10% more servers to the cluster. This linear growth function is critical to the *scalability* of the system—adding fractionally more servers results in a commensurate growth in the overall cluster capacity. Another aspect of this flexibility is the *granularity* of resource additions, which is often tied to the cluster packaging constraints. For example, adding another rack to a cluster, with, say, 100 new servers, is more manageable than adding a whole row, with tens of racks, on the data-center floor.

Even a modest-sized cluster will have several kilometers of fiber-optic cable acting as a vast highway inter-

Figure 3. Example packet routing through a switch chip.

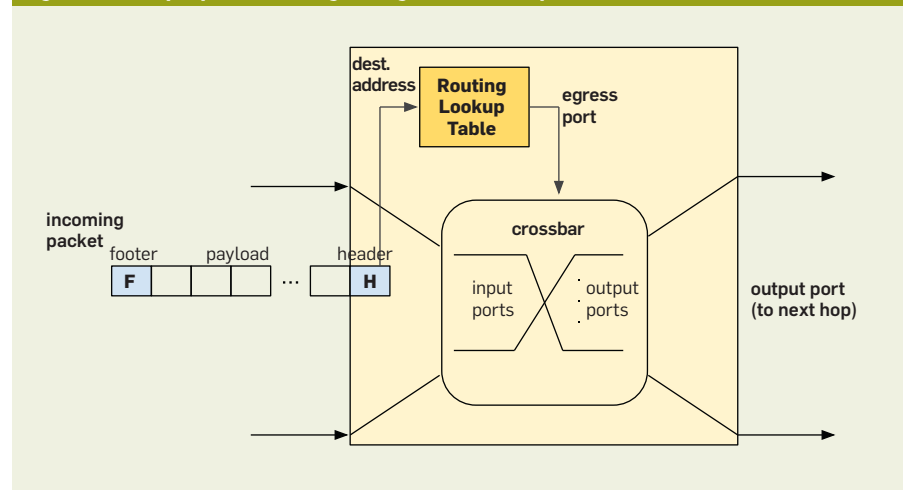
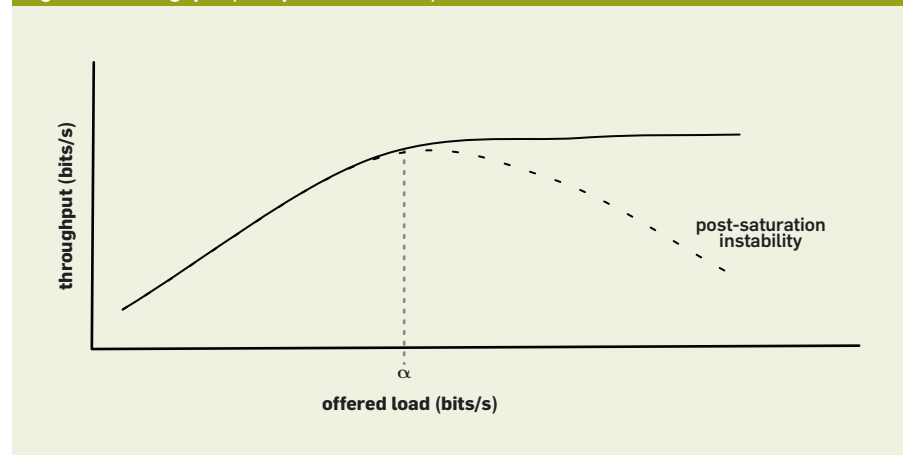


Figure 4. Throughput (accepted bandwidth) as load varies.




connecting racks of servers organized as multiple rows on the data-center floor. The data-center network topology and resulting cable complexity is “baked in” and remains a rigid fixture of the cluster. Managing cable complexity is nontrivial, which is immediately evident from the intricately woven tapestry of fiber-optic cabling laced throughout the data center. It is not uncommon to run additional fiber for redundancy, in the event of a cable failure in a “bundle” of fiber or for planned bandwidth growth. Fiber cables are carefully measured to allow some slack and to satisfy the cable’s bend radius, and they are meticulously labeled to make troubleshooting less of a needle-in-a-haystack exercise.


Reliable and Available

Abstraction is the Archimedes lever that lifts many disciplines within computer science and is used extensively in both computer system design and software engineering. Like an array of nested Russian dolls, the network-programming model provides abstraction between successive layers of the networking stack, enabling platform-independent access to both data and system management. One such example of this type of abstraction is the *protocol buffer*,²¹ which provides a structured message-passing interface for Web applications written in C++, Java, or Python.

Perhaps the most common abstraction used in networking is the notion of a *communication channel* as a virtual resource connecting two hosts. The TCP communication model provides this abstraction to the programmer in the form of a file descriptor, for example, where reads and writes performed on the socket result in the corresponding network transactions, which are hidden from the user application. In much the same way, the InfiniBand QP (queue-pair) verb model provides an abstraction for the underlying send/receive hardware queues in the NIC. Besides providing a more intuitive programming interface, abstraction also serves as a protective sheath around software when faults arise, depositing layers of software sediment to insulate it from critical faults (for example, memory corruption or, worse, host power-supply failure).



The data-center network serves as a “central nervous system” for information exchange between cooperating tasks.



Bad things happen to good software. Web applications must be designed to be fault aware and, to the extent possible, resilient in the presence of a variety of failure scenarios.¹⁰ The network is responsible for the majority of the unavailability budget for a modern cluster. Whether it is a rogue gamma ray causing a soft error in memory or an inattentive worker accidentally unearthing a fiber-optic line, the operating system and underlying hardware substrate work in concert to foster a robust ecosystem for Web applications.

The data-center network serves as a “central nervous system” for information exchange between cooperating tasks. The network’s functionality is commonly divided into *control* and *data* planes. The control plane provides an ancillary network juxtaposed with the data network and tasked with “command and control” for the primary data plane. The control plane is an autonomic system for configuration, fault detection and repair, and monitoring of the data plane. The control plane is typically implemented as an embedded system within each switch component and is tasked with fault detection, notification, and repair when possible.

For example, when a network link fails or has an uncharacteristically high number of transmission errors, the control plane will *reroute* the network to avoid the faulty link. This entails recomputing the routes according to the *routing algorithm* and emplacing new entries in the routing tables of the affected switches. Of course, the effects of this patchwork are not instantaneous. Once the routing algorithm computes new routes, taking into consideration the newfound faulty links, it must disseminate the routes to the affected switches. The time needed for this information exchange is referred to as *convergence time*, and a primary goal of the routing protocol is to ensure it is optimally confined to a small epoch.

Fault recovery is a very complicated subject and confounds all but the simplest of data-center networks. Among the complicating factors are *marginal* links that cause “flapping” by transitioning between active and inactive (that is, up and down), repeatedly creat-

ing a deluge of error notifications and resulting route recomputation based on fluctuating and inconsistent link status. Some link-layer protocols allow the link speed to be adjusted downward in hopes of improving the link quality. Of course, lowering the link speed results in a reduced bandwidth link, which in turn may limit the overall bandwidth of the network or at the very least will create load imbalance as a result of increased contention across the slow link. Because of these complicating factors, it is often better to logically excise the faulty link from the routing algorithm until the physical link can be replaced and validated.

Conclusion

The data-center network is generally regarded as a critical design element in the system architecture and the skeletal structure upon which processor, memory, and I/O devices are dynamically shared. The evolution from 1G to 10G Ethernet and the emerging 40G Ethernet has exposed performance bottlenecks in the communication stack that require better hardware-software coordination for efficient communication. Other approaches by Solarflare, Myricom, and InfiniBand, among others, have sought to reshape the conventional sockets programming model with more efficient abstractions. Internet sockets, however, remain the dominant programming interface for data-center networks.

Network performance and reliability are key design goals, but they are tempered by cost and serviceability constraints. Deploying a large cluster computer is done incrementally and is often limited by the power capacity of the building, with power being distributed across the cluster network so that a power failure impacts only a small fraction—say, less than 10%—of the hosts in the cluster. When hardware fails, as is to be expected, the operating system running on the host coordinates with a higher-level hypervisor or cluster operating system to allow failures to be replaced in situ without draining traffic in the cluster. Scalable Web applications are designed to expect occasional hardware failures, and the resulting software is by necessity resilient.

A good user experience relies on *predictable* performance, with the data-center network delivering predictable latency and bandwidth characteristics across varying traffic patterns. With single-thread performance plateauing, microprocessors are providing more cores to keep pace with the relentless march of Moore's Law. As a result, applications are looking for increasing thread-level parallelism and scaling to a large core count with a commensurate increase in communication among cores. This trend is motivating *communication-centric* cluster computing with tens of thousands of cores in unison, reminiscent of a flock darting seamlessly amidst the clouds. **□**

References

- Al-Fares, M., Loukissas, A. and Vahdat, A. A scalable, commodity data-center network architecture. In *Proceedings of the ACM SIGCOMM 2008 Conference on Data Communication* (2008), 63–74; <http://doi.acm.org/10.1145/1402958.1402967>.
- Amdahl's Law; http://en.wikipedia.org/wiki/Amdahl's_Law.
- Ballani, H., Costa, P., Karagiannis, T. and Rowstron, A. Towards predictable data-center networks. In *Proceedings of the ACM SIGCOMM 2011 Conference* (2011), 242–253; <http://doi.acm.org/10.1145/2018436.2018465>.
- Barroso, L.A., Dean, J. and Holze, U. Web search for a planet: The Google cluster architecture. *IEEE Micro* 23, 2 (2003), 22–28; <http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=1196112&isnumber=26907>.
- Cerf, V. and Icahn R.E. A protocol for packet network intercommunication. *SIGCOMM Computer Communication Review* 35, 2 (2005), 71–82; <http://doi.acm.org/10.1145/1064413.1064423>.
- Cisco Data Center Infrastructure 3.0 Design Guide. Data Center Design—IP Network Infrastructure; http://www.cisco.com/en/US/docs/solutions/Enterprise/Data_Center/DC_3_0/DC-3_0_IPInfra.html.
- Clos, C. A Study of non-blocking switching networks. *The Bell System Technical Journal* 32, 2 (1953), 406–424.
- Fitzpatrick, B. Distributed caching with Memcached. *Linux Journal* 2004; <http://www.linuxjournal.com/article/7451>.
- Dally, W. and Towles, B. *Principles and Practices of Interconnection Networks*. Morgan Kaufmann Publishers, San Francisco, CA, 2003.
- Gill, P., Jain, N. and Nagappan, N. Understanding network failures in data centers: measurement, analysis, and implications. In *Proceedings of the ACM SIGCOMM 2011 Conference* (2011), 350–361; <http://doi.acm.org/10.1145/2018436.2018477>.
- Greenberg, A., Hamilton, J.R., Jain, N., Kandula, S., Kim, C., Lahiri, P., Maltz, D.A., Patel, P. and Sengupta, S. VL2: A scalable and flexible data center network. In *Proceedings of the ACM SIGCOMM 2009 Conference on Data Communication* (2009), 51–62; <http://doi.acm.org/10.1145/1592568.1592576>.
- Greenberg, A., Hamilton, J., Maltz, D.A. and Patel, P. The cost of a cloud: Research problems in data center networks. *SIGCOMM Computer Communications Review* 39, 1 (2008), 68–73; <http://doi.acm.org/10.1145/1496091.1496103>.
- Hoelzle, U. and Barroso, L. A. *The Datacenter as a Computer: An Introduction to the Design of Warehouse-Scale Machines* (1st ed.). Morgan & Claypool Publishers, 2009.
- Kermani, P. and Kleinrock, L. Virtual cut-through: A new computer communication switching technique. *Computer Networks* 3, 4 (1976), 267–286; <http://www.sciencedirect.com/science/article/pii/0376507579900321>.
- Leiserson, C.E. Fat-trees: Universal networks for hardware-efficient supercomputing. *IEEE Transactions on Computers* 34, 10 (1985), 892–901.

- Mori, T., Uchida, M., Kawahara, R., Pan, J., and Goto, S. Identifying elephant flows through periodically sampled packets. In *Proceedings of the 4th ACM SIGCOMM Conference on Internet Measurement* (2004), 115–120; <http://doi.acm.org/10.1145/1028788.1028803>.
- Mudigonda, J., Yalagandula, P., Mogul, J., Stiekes, B., Pouffary, Y. NetLord: A scalable multi-tenant network architecture for virtualized datacenters. *SIGCOMM Computer Communication Review* 41, 4 (2011), 62–73; <http://doi.acm.org/10.1145/2043164.2018444>.
- Mysore, R.N., Pamboris, A., Farrington, N., Huang, N., Miri, P., Radhakrishnan, S., Subramanya, V., Vahdat, A. PortLand: A scalable fault-tolerant layer 2 data center network fabric. *SIGCOMM Computer Communication Review* 39, 4 (2009), 39–50; <http://doi.acm.org/10.1145/1594977.1592575>.
- Ni, L. M., McKinley, P. K. A survey of wormhole routing techniques in direct networks. *Computer* 26, 2 (1993), 62–76; <http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=191995&isnumber=4947>.
- Ousterhout, J., Agrawal, P., Erickson, D., Kozyrak, C., Leverich, J., Mazieres, D., Mitra, S., Narayanan, A., Parulkar, G., Rosenblum, M., Rumble, S.M., Stratmann, E. and Stutsman, R. The case for RAMClouds: Scalable high-performance storage entirely in DRAM. *SIGOPS Operating Systems Review* 43, 4 (2010), 92–105; <http://doi.acm.org/10.1145/1713254.171327>.
- Protocol buffers; <http://code.google.com/apis/protocolbuffers/>.
- Rumble, S.M., Ongaro, D., Stutsman, R., Rosenblum, M., and Ousterhout, J.K. It's time for low latency. In *Proceedings of the 13th Usenix Conference on Hot Topics in Operating Systems* (2011).
- Vahdat, A., Al-Fares, M., Farrington, N., Mysore, R.N., Porter, G., and Radhakrishnan, S. Scale-out networking in the data center. *IEEE Micro* 30, 4 (2010), 29–41; <http://dx.doi.org/10.1109/MM.2010.72>.
- Vahdat, A., Liu, H., Zhao, X. and Johnson, C. The emerging optical data center. Presented at the *Optical Fiber Communication Conference*. OSA Technical Digest (2011); <http://www.opticsinfobase.org/abstract.cfm?URI=OFC-2011-OTuH2>.
- Wilson, C., Ballani, H., Karagiannis, T., and Rowstron, A. Better never than late: Meeting deadlines in datacenter networks. In *Proceedings of the ACM SIGCOMM 2011 Conference* (2011), 50–61; <http://doi.acm.org/10.1145/2018436.2018443>.

Related articles on queue.acm.org

Enterprise Grid Computing

Paul Strang

<http://queue.acm.org/detail.cfm?id=1080877>

Cooling the Data Center

Andy Woods

<http://queue.acm.org/detail.cfm?id=1737963>

Improving Performance on the Internet

Tom Leighton

<http://queue.acm.org/detail.cfm?id=1466449>

Dennis Abts is a member of the technical staff at Google, where he is involved in the architecture and design of next-generation large-scale clusters. Prior to joining Google, Abts was a senior principal engineer and system architect at Cray Inc. He has numerous technical publications and patents in areas of interconnection networks, data-center networking, cache-coherence protocols, high-bandwidth memory systems, and supercomputing.

Bob Felderman spent time at both Princeton and UCLA before starting a short stint at Information Sciences Institute. He then helped found Myricom, which became a leader in cluster-computing networking technology. After seven years there, he moved to Packet Design where he applied high-performance computing ideas to the IP and Ethernet space. He later was a founder of Precision I/O. All of that experience eventually led him to Google where he is a principal engineer working on issues in data-center networking and general platforms system architecture.

© 2012 ACM 0001-0782/12/06 \$10.00