

Chapter 6

Authors: John Hennessy & David Patterson

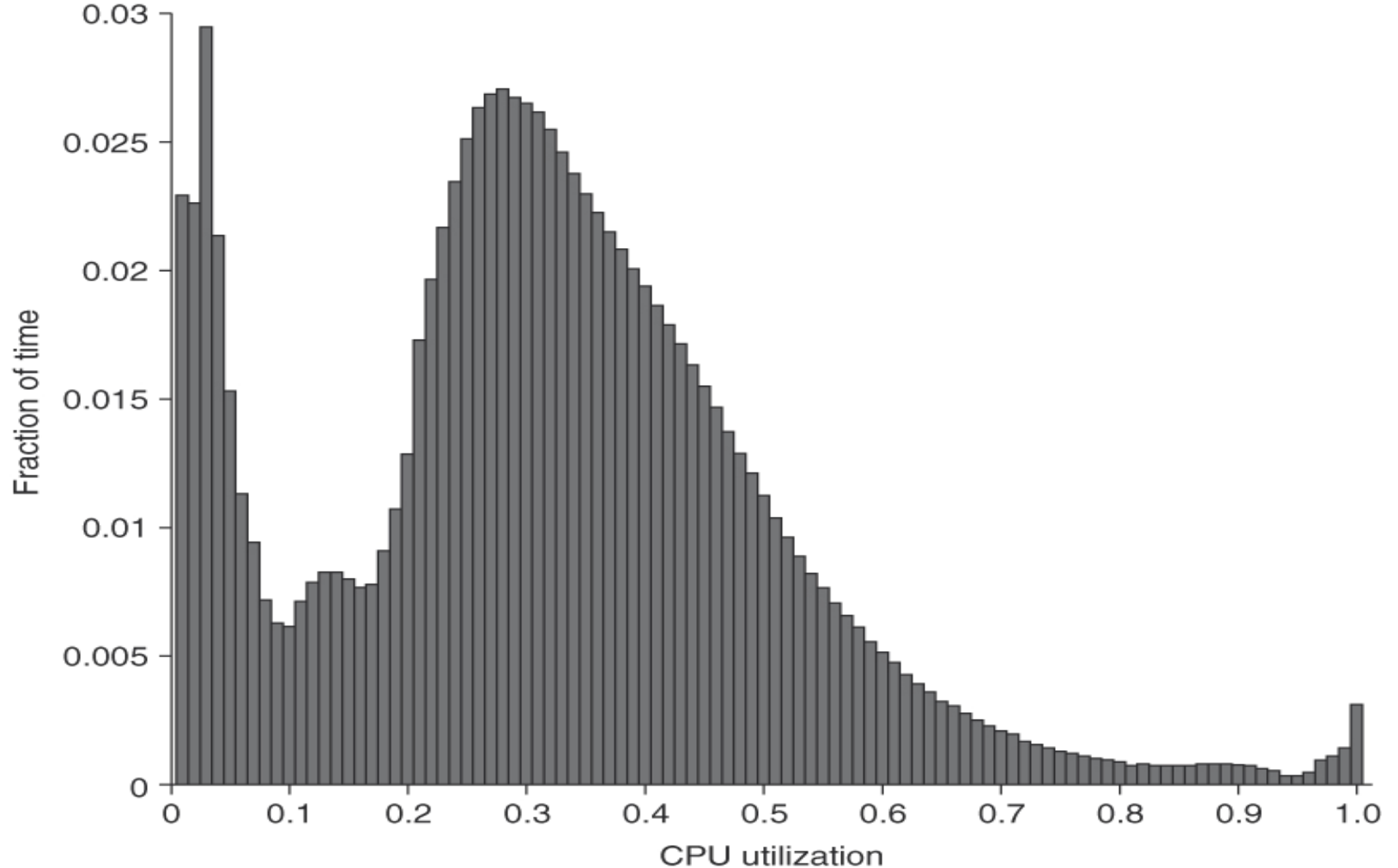


Figure 6.3 Average CPU utilization of more than 5000 servers during a 6-month period at Google. Servers are rarely completely idle or fully utilized, in-stead operating most of the time at between 10% and 50% of their maximum utilization. (From Figure 1 in Barroso and Hölzle [2007].) The column the third from the right in Figure 6.4 calculates percentages plus or minus 5% to come up with the weightings; thus, 1.2% for the 90% row means that 1.2% of servers were between 85% and 95% utilized.

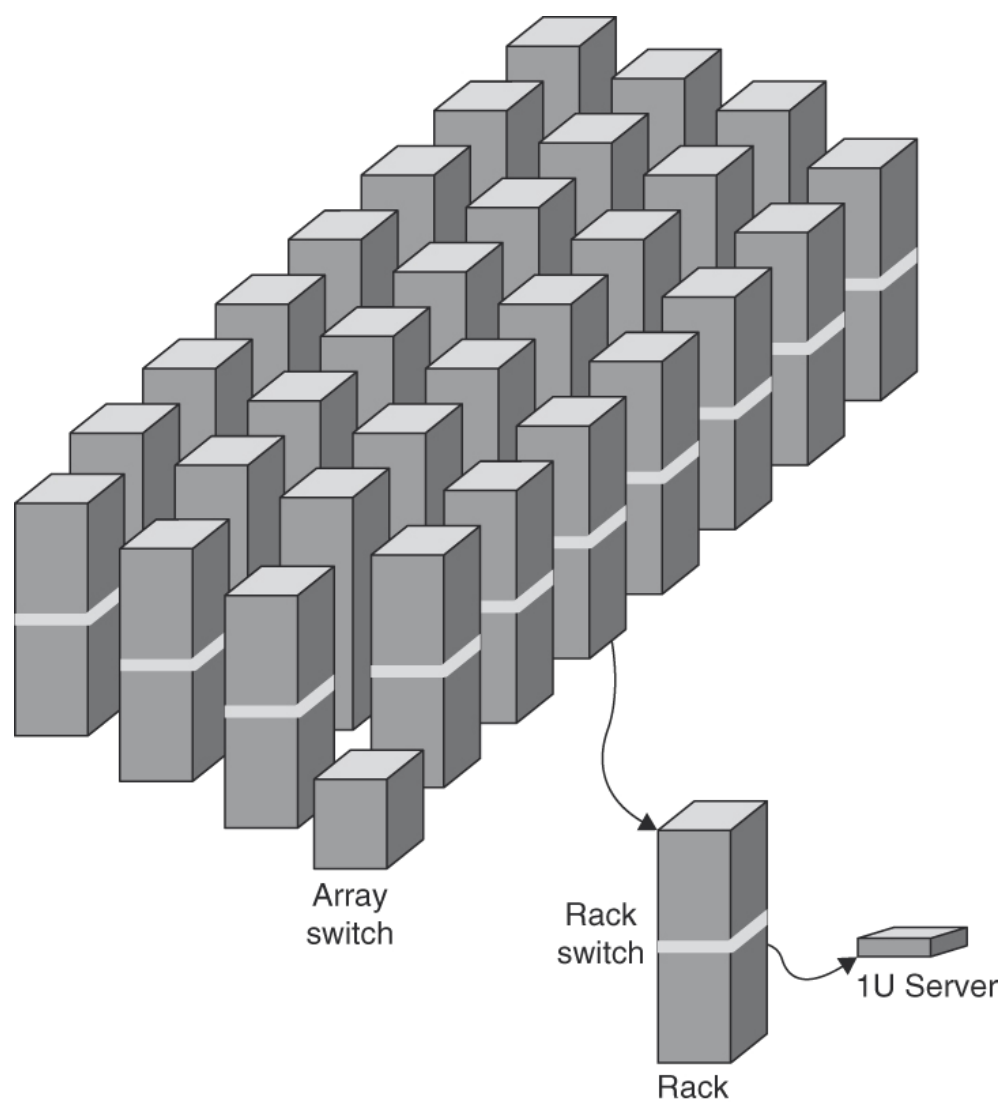


Figure 6.5 Hierarchy of switches in a WSC. (Based on Figure 1.2 of Barroso and Hölzle [2009].)

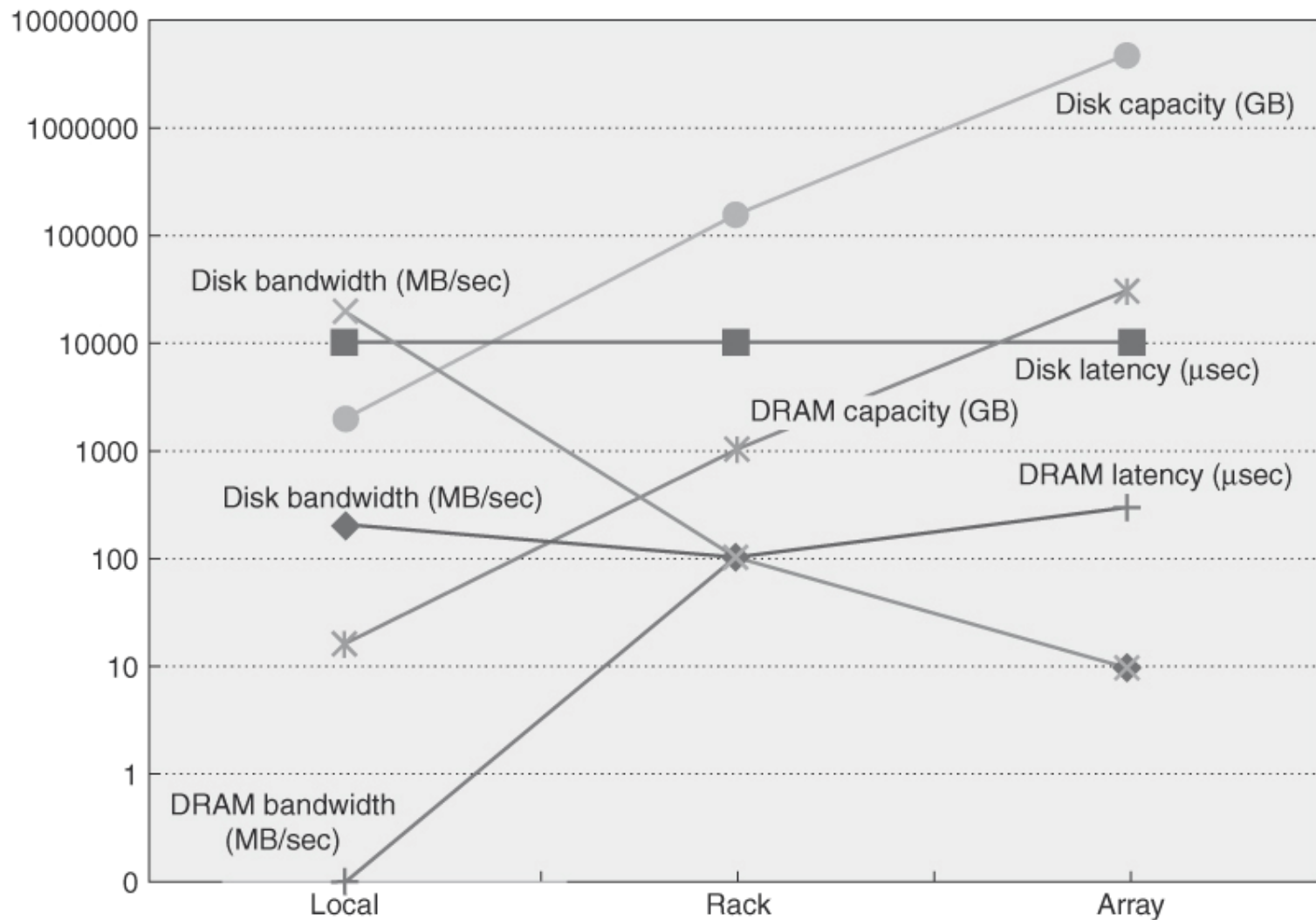


Figure 6.7 Graph of latency, bandwidth, and capacity of the memory hierarchy of a WSC for data in Figure 6.6 [Barroso and Hölzle 2009].

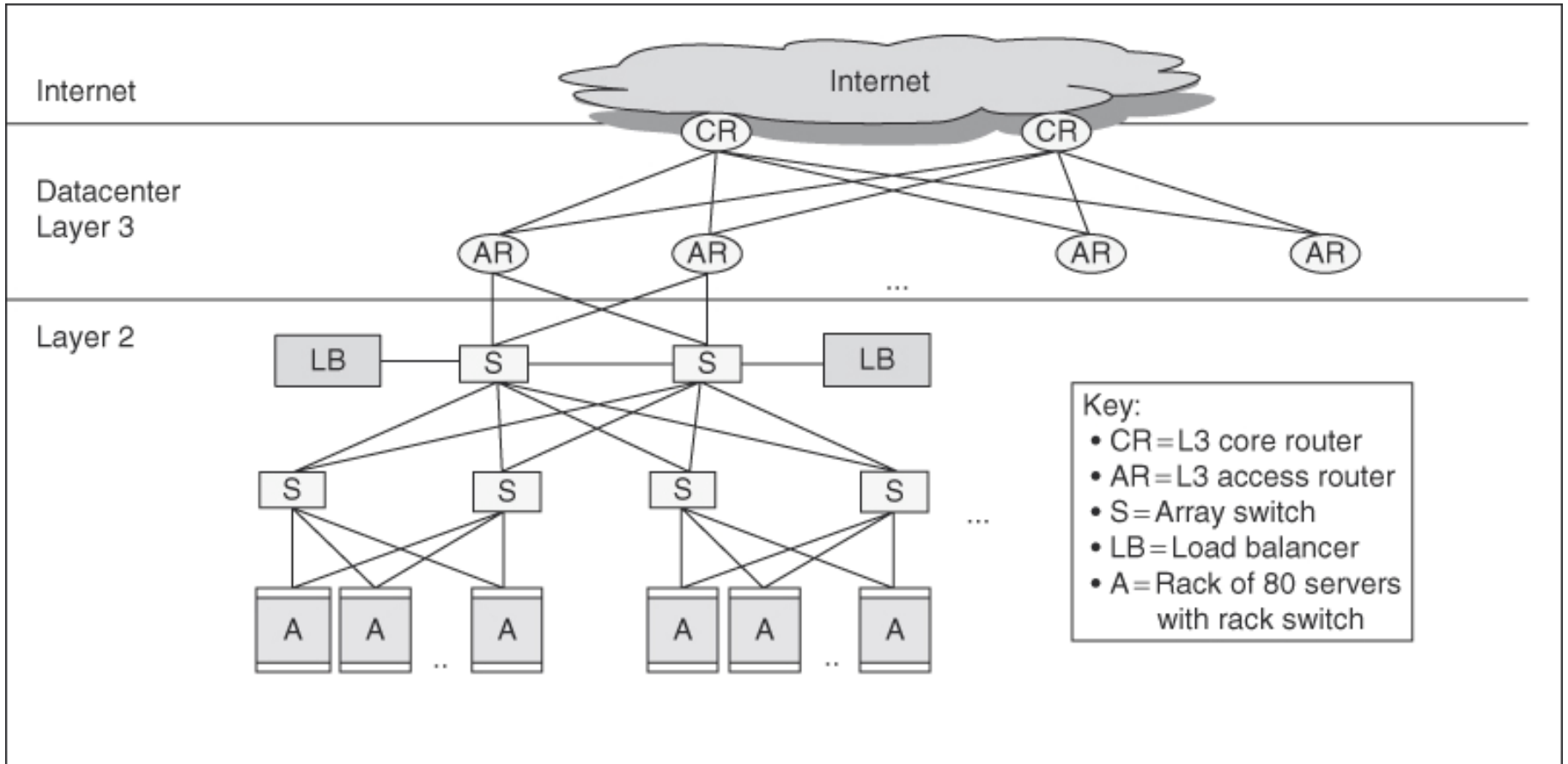


Figure 6.8 The Layer 3 network used to link arrays together and to the Internet [Greenberg et al. 2009]. Some WSCs use a separate *border router* to connect the Internet to the datacenter Layer 3 switches.

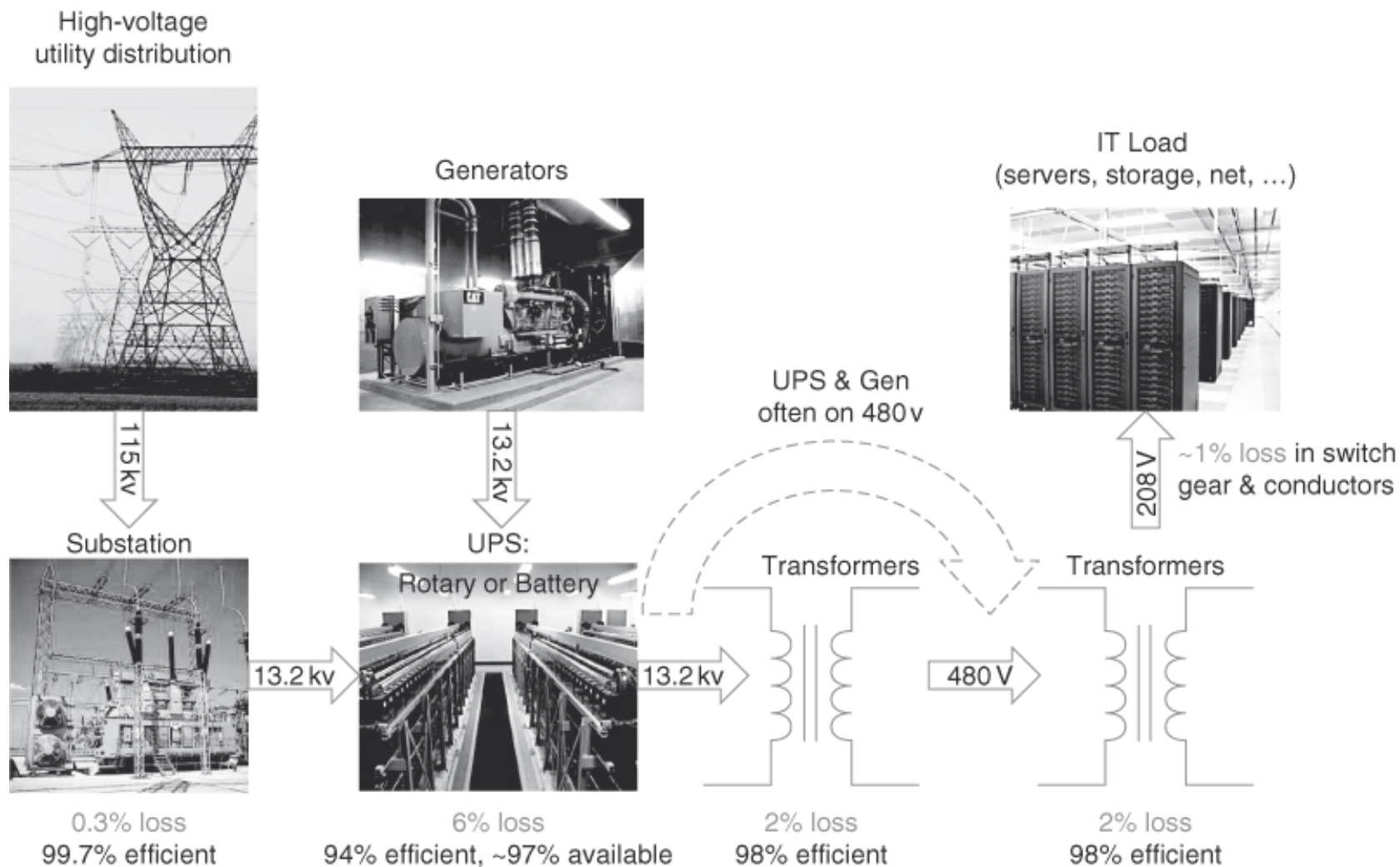


Figure 6.9 Power distribution and where losses occur. Note that the best improvement is 11%. (From Hamilton [2010].)

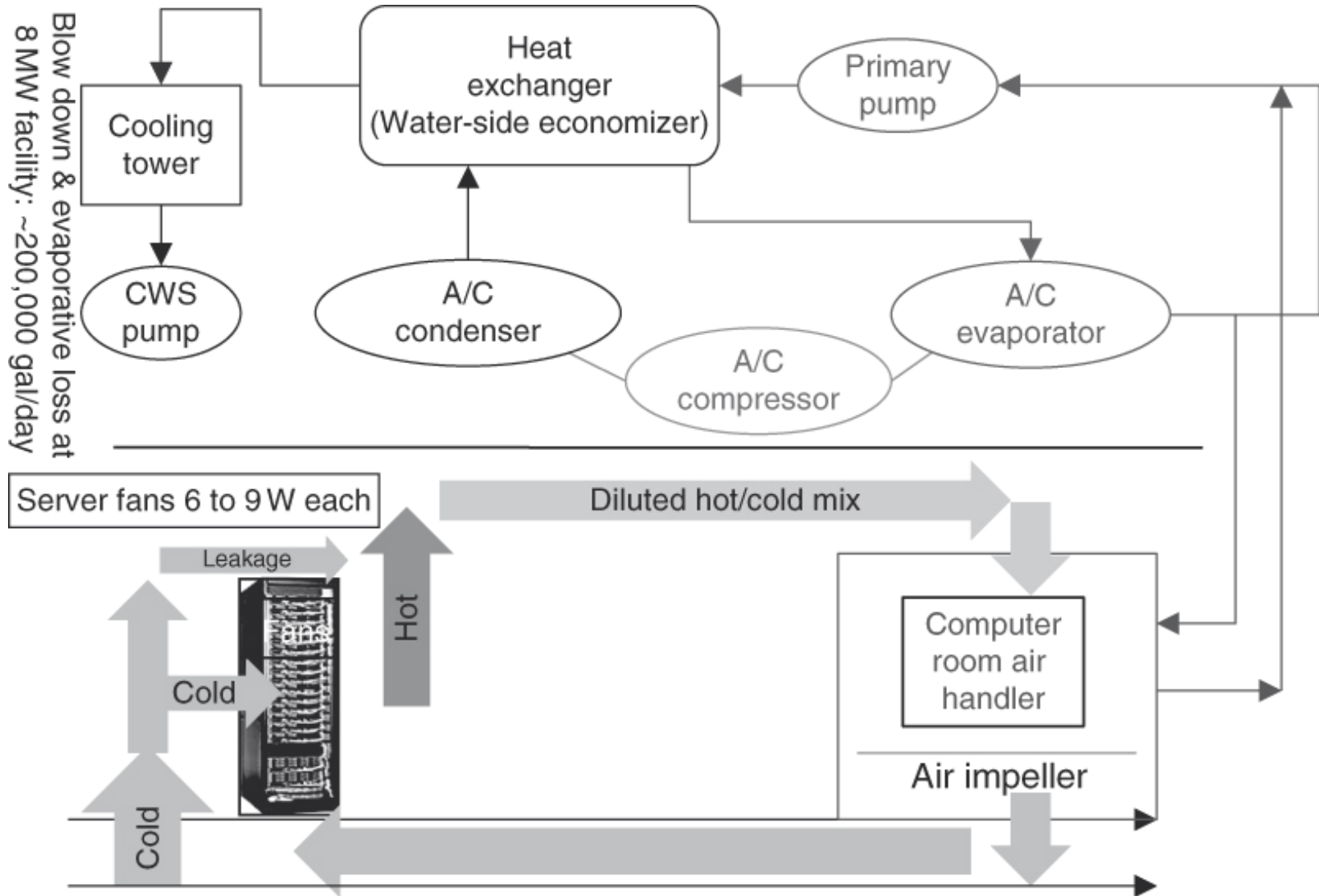


Figure 6.10 Mechanical design for cooling systems. CWS stands for circulating water system. (From Hamilton [2010].)

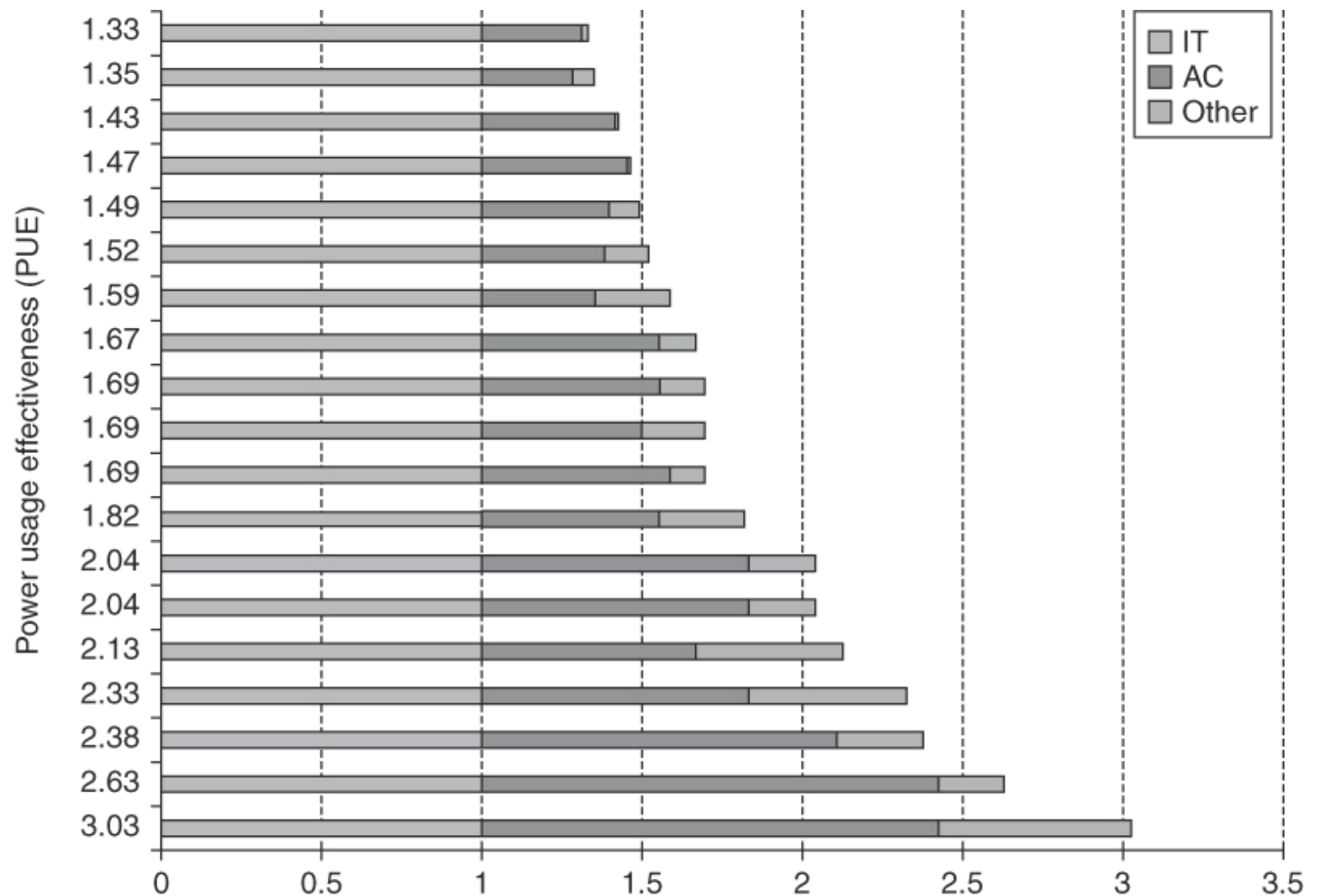


Figure 6.11 Power utilization efficiency of 19 datacenters in 2006 [Greenberg et al. 2006]. The power for air conditioning (AC) and other uses (such as power distribution) is normalized to the power for the IT equipment in calculating the PUE. Thus, power for IT equipment must be 1.0 and AC varies from about 0.30 to 1.40 times the power of the IT equipment. Power for “other” varies from about 0.05 to 0.60 of the IT equipment.

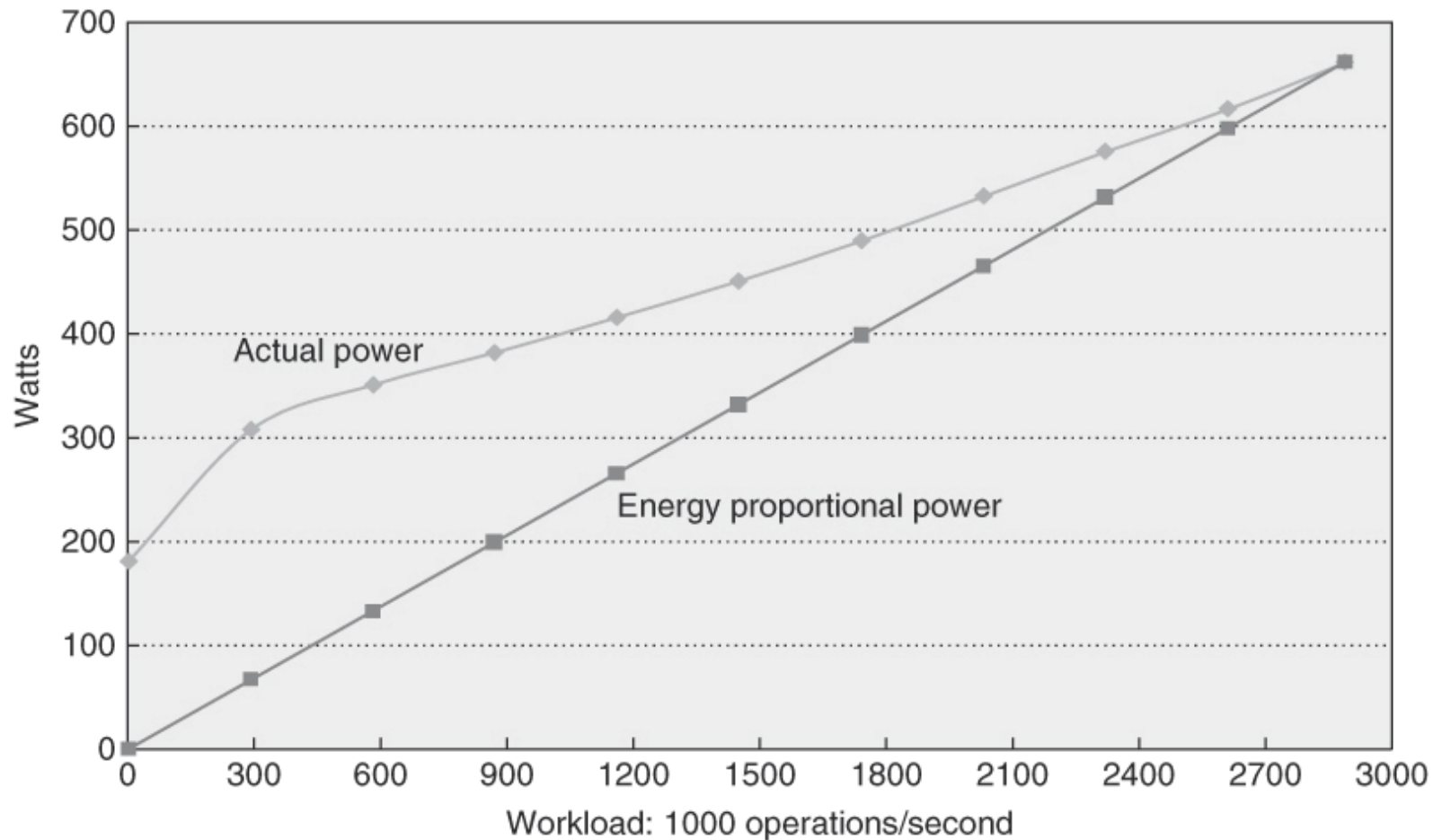


Figure 6.18 The best SPECpower results as of July 2010 versus the ideal energy proportional behavior. The system was the HP ProLiant SL2x170z G6, which uses a cluster of four dual-socket Intel Xeon L5640s with each socket having six cores running at 2.27 GHz. The system had 64 GB of DRAM and a tiny 60 GB SSD for secondary storage. (The fact that main memory is larger than disk capacity suggests that this system was tailored to this benchmark.) The software used was IBM Java Virtual Machine version 9 and Windows Server 2008, Enterprise Edition.

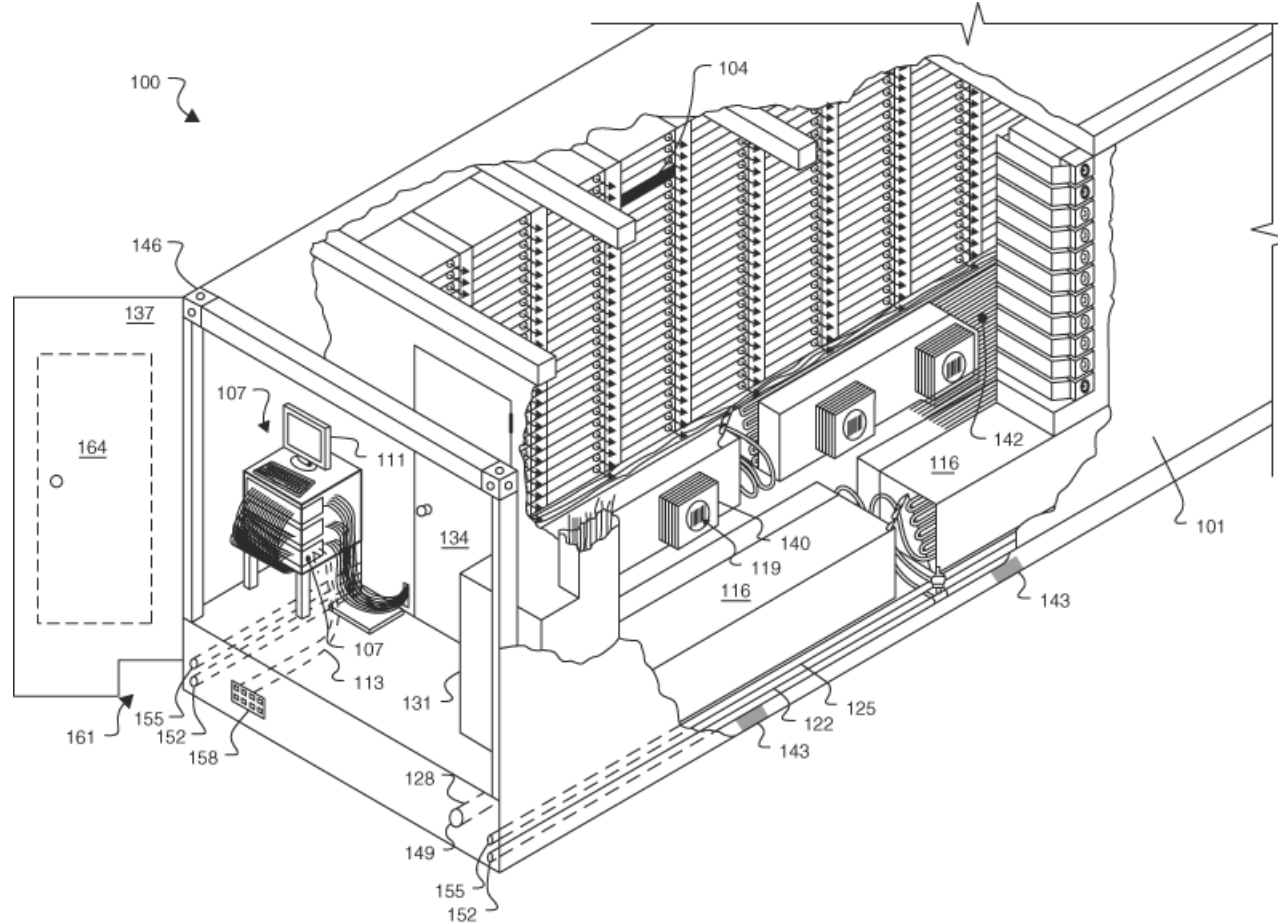


Figure 6.19 Google customizes a standard 1AAA container: 40 x 8 x 9.5 feet (12.2 x 2.4 x 2.9 meters). The servers are stacked up to 20 high in racks that form two long rows of 29 racks each, with one row on each side of the container. The cool aisle goes down the middle of the container, with the hot air return being on the outside. The hanging rack structure makes it easier to repair the cooling system without removing the servers. To allow people inside the container to repair components, it contains safety systems for fire detection and mist-based suppression, emergency egress and lighting, and emergency power shut-off. Containers also have many sensors: temperature, airflow pressure, air leak detection, and motion-sensing lighting. A video tour of the datacenter can be found at <http://www.google.com/corporate/green/datacenters/summit.html>. Microsoft, Yahoo!, and many others are now building modular datacenters based upon these ideas but they have stopped using ISO standard containers since the size is inconvenient.

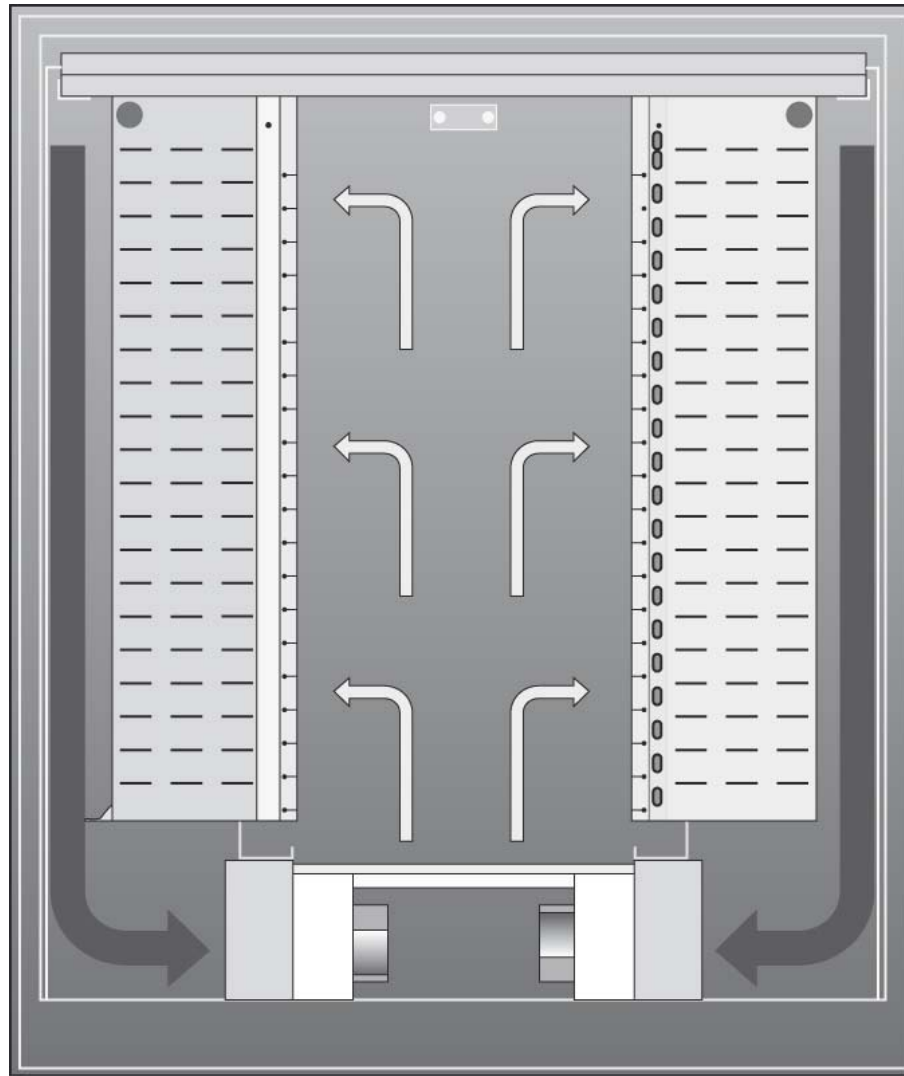


Figure 6.20 Airflow within the container shown in Figure 6.19. This cross-section diagram shows two racks on each side of the container. Cold air blows into the aisle in the middle of the container and is then sucked into the servers. Warm air returns at the edges of the container. This design isolates cold and warm airflows.

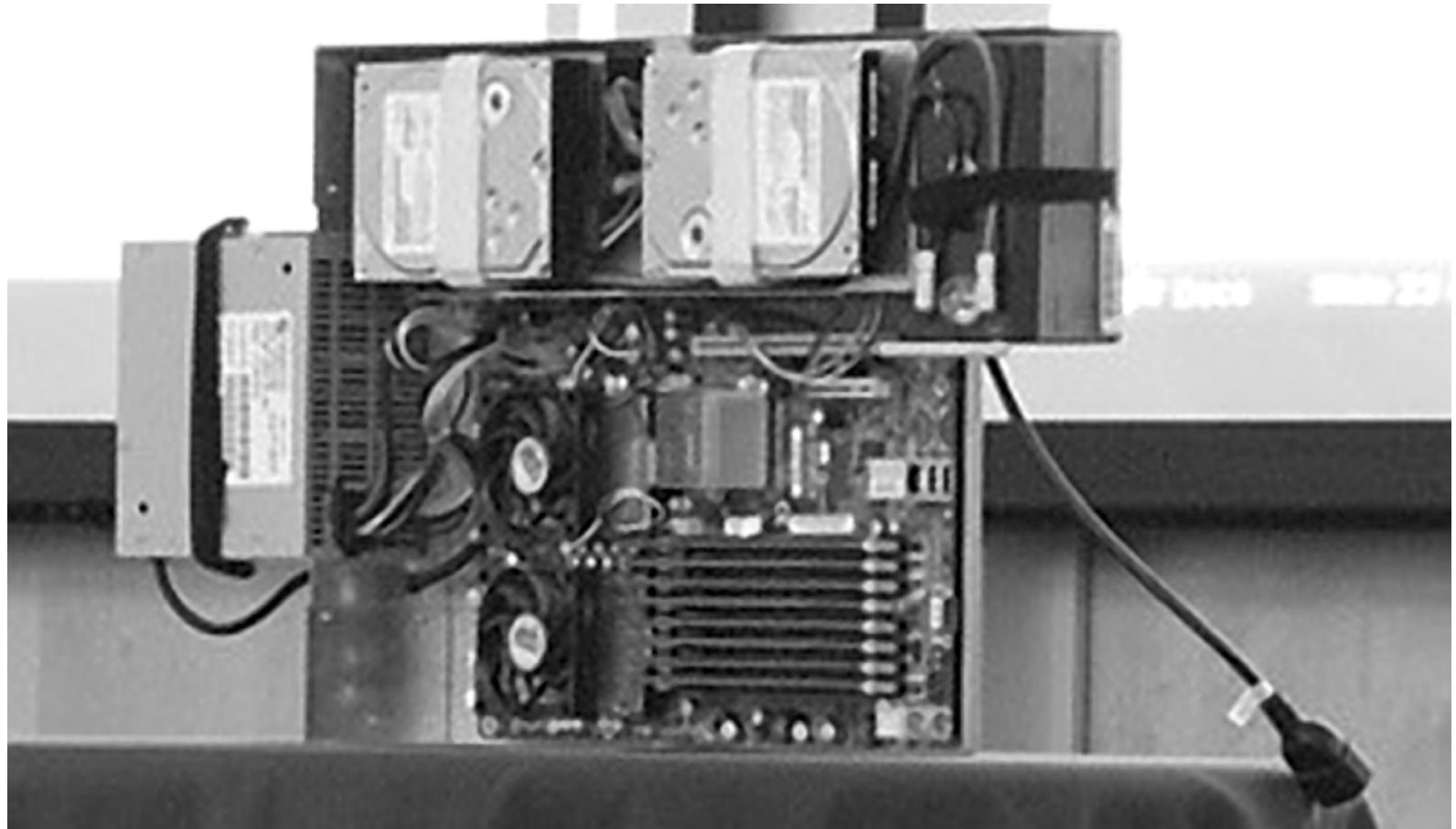


Figure 6.21 Server for Google WSC. The power supply is on the left and the two disks are on the top. The two fans below the left disk cover the two sockets of the AMD Barcelona microprocessor, each with two cores, running at 2.2 GHz. The eight DIMMs in the lower right each hold 1 GB, giving a total of 8 GB. There is no extra sheet metal, as the servers are plugged into the battery and a separate plenum is in the rack for each server to help control the airflow. In part because of the height of the batteries, 20 servers fit in a rack.

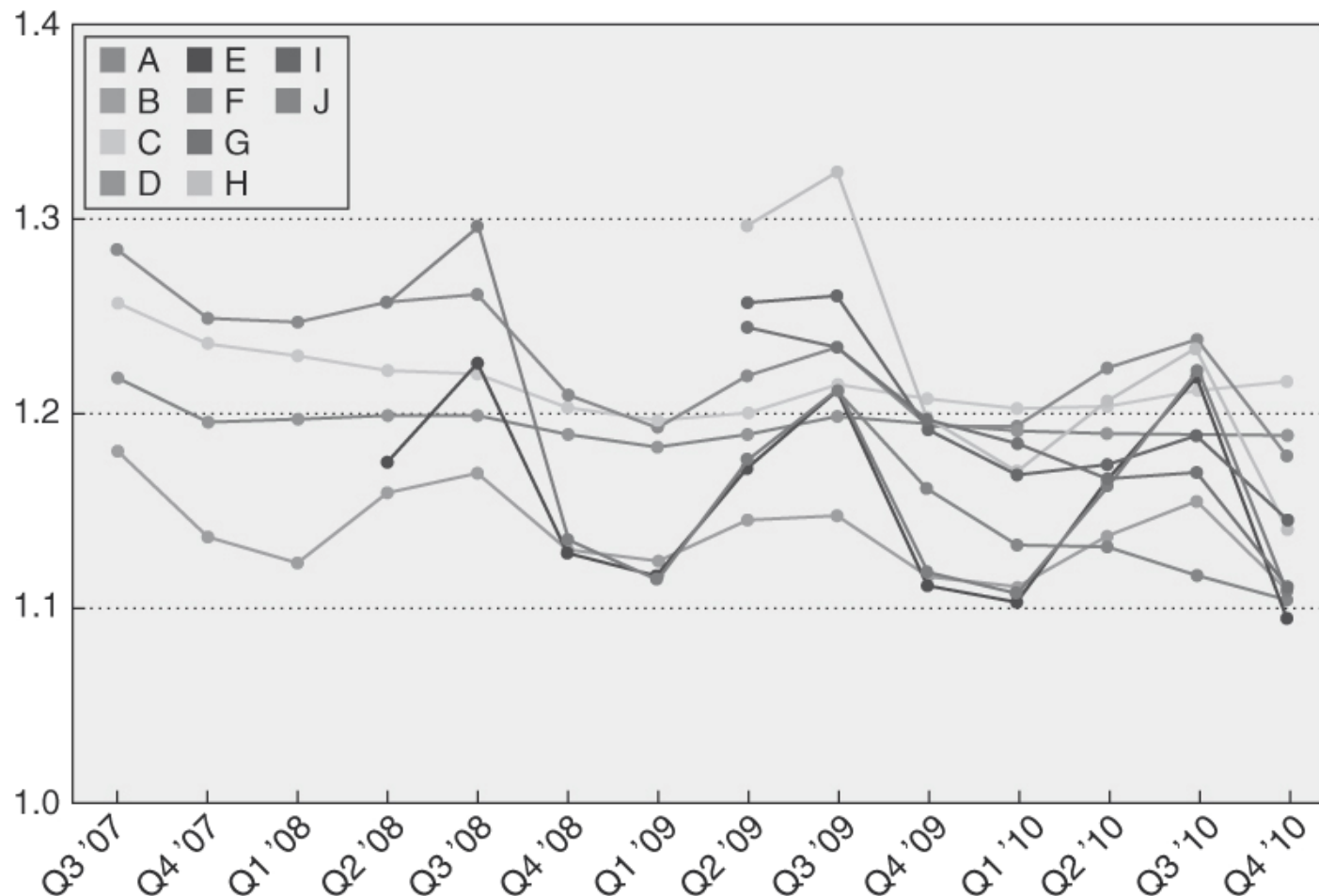


Figure 6.22 Power usage effectiveness (PUE) of 10 Google WSCs over time. Google A is the WSC described in this section. It is the highest line in Q3 '07 and Q2 '10. (From www.google.com/corporate/green/datacenters/measuring.htm.) Facebook recently announced a new datacenter that should deliver an impressive PUE of 1.07 (see <http://opencompute.org/>). The Prineville Oregon Facility has no air conditioning and no chilled water. It relies strictly on outside air, which is brought in one side of the building, filtered, cooled via misters, pumped across the IT equipment, and then sent out the building by exhaust fans. In addition, the servers use a custom power supply that allows the power distribution system to skip one of the voltage conversion steps in Figure 6.9.

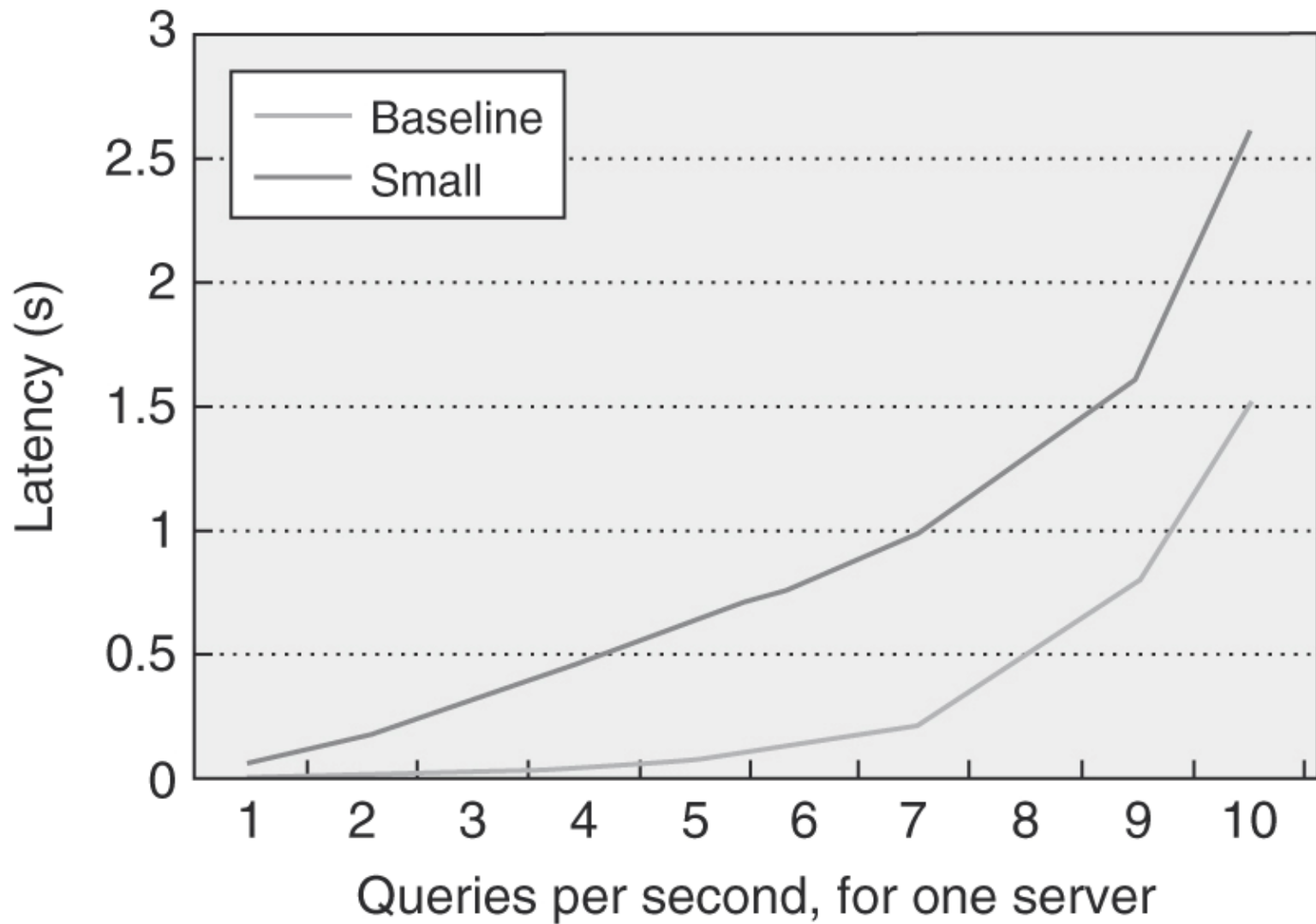
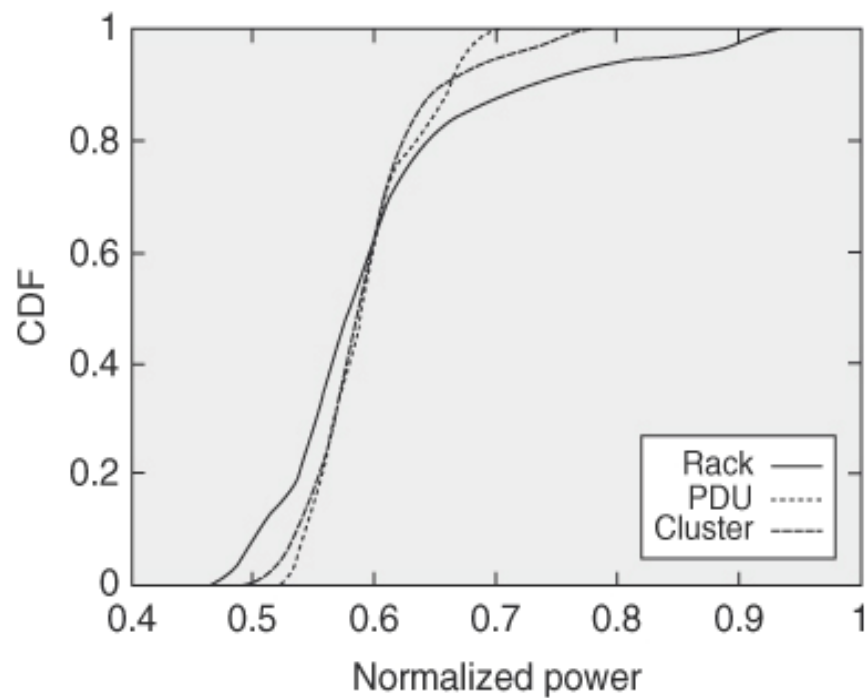
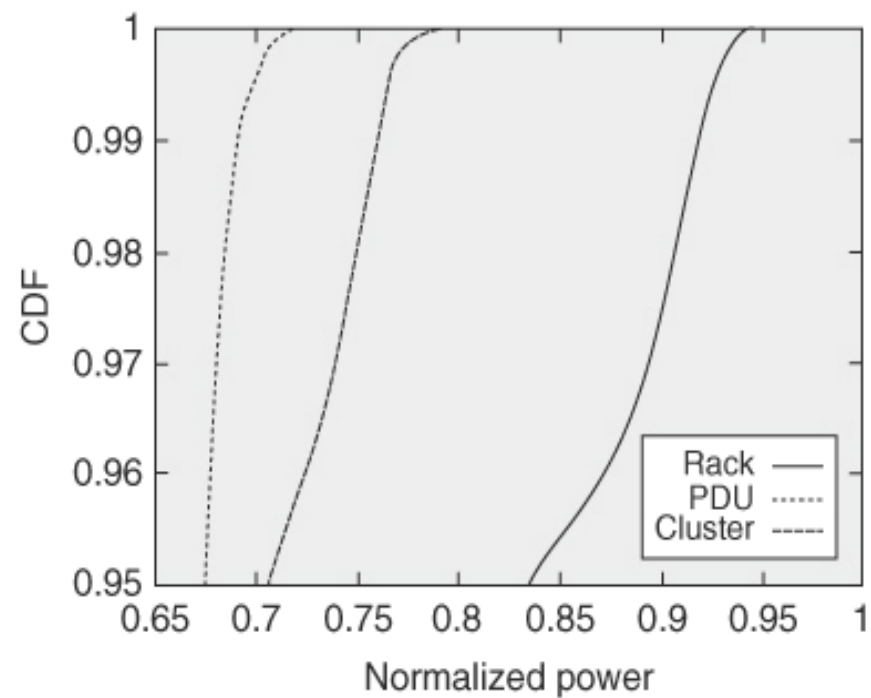


Figure 6.24 Query-response time curve.



(a) Full distribution



(b) Zoomed view

Figure 6.25 Cumulative distribution function (CDF) of a real datacenter.