### 2.1.1   Achieving high economies of scale with cloud data centers

Revisiting the vehicle analogy, the data center is the car's engine. A *data center*—one that you might find in any large company—is a facility (usually secure) to house a large collection of computers, networking, and communications equipment. But the

large internet-based companies, such as Amazon, Yahoo!, Google, Intuit, Apple, and others have, over the years, built up what have to be considered *mega* data centers with thousands of servers. These data centers are the starting point for what is being built out by the cloud providers.

It's useful to understand the structure and the economics of these massive data centers to gauge how much you can scale your operations, how reliable your cloud computing will be, how secure your data will be, and where the economics of public clouds are going. This is particularly important should you decide to build your own *private* cloud. You'll learn more about private clouds later in this chapter, and we've dedicated chapter 4 to the topics of security and private clouds.

**THE STRUCTURE OF A DATA CENTER**

A data center can occupy one room of a building, one or more floors, or an entire building. Most of the equipment is often in the form of servers mounted in 19-inch rack cabinets, which are usually placed in single rows with corridors between them. This allows people access to the front and rear of each cabinet. Servers differ greatly in size, from 1U servers (which occupy one of 42 slots in a standard rack) to large free-standing storage silos that occupy many tiles on the floor. Mainframe computers and storage devices may be as big as the racks themselves and are placed alongside them. Large data centers may use shipping containers packed with 1,000 or more servers each; when they need to repair or upgrade, they replace the whole container (rather than repairing individual servers).

Clean, unwavering power—and lots of it—is essential. Data centers need to keep their computers running at all times. They should be prepared to handle brownouts and even power outages. The power must be conditioned, and backup batteries and diesel generators must be available to keep power flowing no matter what.

As you can imagine, all that power generates a lot of heat. Data centers must cool their racks of equipment. The most common mode of cooling is air-conditioning; water-cooling is also an option when it's easily available, such as at some of the new data centers along the Columbia River in Washington State. Air-conditioning not only cools the environment but also controls humidity to avoid condensation or static electric buildup.

Network connectivity and ample bandwidth to and from the network backbones are vital, to handle the input and output from the entire collection of servers and storage units. All these servers will be idle if no one can access them.

Another important aspect is physical and logical security. Bigger data centers are targets for hackers all over the world. Some freestanding data centers begin with security through obscurity and disguise the fact that a data center even exists at that location. Guards, mantraps, and state-of-the-art authentication technology keep unauthorized people from physically entering. Firewalls, VPN gateways, intrusion-detection software, and so on keep unauthorized people from entering over the network. (More on all aspects of cloud security in chapter 4.)

Finally, data centers must always assume the worst and have disaster recovery contingencies in place that avoid loss of data and experience the minimum loss of service in case of disaster.

## DATA CENTERS: SCALING FOR THE CLOUD

A traditional, large data center dedicated to a single large corporation costs approximately $100-200 million.[1] Contrast that to the total cost of building the largest mega data centers that provide cloud services: $500 million or more.[2,3] What is going into that much higher cost, and what can the biggest cloud data centers do that normal companies can't do with their dedicated data centers?

The largest data-center operators like Google, Amazon, and Microsoft situate their data centers in geographic proximity to heavy usage areas to keep network latency to a minimum and to provide failover options. They also choose geographies with access to cheap power. The northwest is particularly advantageous because the available hydropower is the cheapest power in the country and air-conditioning needs are low to zero. Major data centers can use a whopping amount of wattage and cost their owners upward of $30 million a year for electricity alone, which is why data-center power consumption across the U.S. represents 1.2 percent of total power consumption in the country—and it's rising. The positive side is that cloud data centers use so much power and have so much clout that they can negotiate huge power volume discounts.

Additionally, these giant data centers tend to buy so much hardware that they can negotiate huge volume discounts far beyond the reach of even the largest company that's building a dedicated data center. For example, Amazon spent about $90 million for 50,000 servers from Rackable/SGI in 2008,[4] which, without the massive volume discounts, would have cost $215 million.

Servers dominate data-center costs. This is why Google and others are trying to get cheaper servers and have taken to building their own from components. Google relies on cheap computers with conventional multicore processors. A single Google data center has tens of thousands of these inexpensive processors and disks, held together with Velcro tape in a practice that makes for easy swapping of components.

To reduce the machines' energy appetite, Google fitted them with high-efficiency power supplies and voltage regulators, variable-speed fans, and system boards stripped of all unnecessary components, such as graphics chips. Google has also experimented with a CPU power-management feature called *dynamic voltage/frequency scaling*. It reduces a processor's voltage or frequency during certain periods (for example, when you don't need the results of a computing task right away). The server executes its work more slowly, reducing power consumption. Google engineers have reported energy savings of around 20 percent on some of their tests.

In 2006, Google built two cloud computing data centers in Dalles, Oregon, each of which has the acreage of a football field with four floors and two four-story cooling

---

[1] http://perspectives.mvdirona.com/2008/11/28/CostOfPowerInLargeScaleDataCenters.aspx
[2] http://www.datacenterknowledge.com/archives/2007/11/05/microsoft-plans-500m-illinois-data-center
[3] http://www.theregister.co.uk/2009/09/25/microsoft_chillerless_data_center
[4] http://www.datacenterknowledge.com/archives/2009/06/23/amazon-adds-cloud-data-center-in-virginia

**Figure 2.1   Photograph of Google's top-secret Dalles, OR data center, built near the Dalles Dam for access to cheap power. Note the large cooling towers on the end of each football-sized building on the left. These towers cool through evaporation rather than using more power-hungry chillers. Source: Melanie Conner,** *New York Times.*

plants (see figure 2.1). The Dalles Dam is strategic for the significant energy and cooling needs of these data centers. (Some new cloud data centers rely on cooling towers, which use evaporation to remove heat from the cooling water, instead of traditional energy-intensive chillers.)

The Dalles data center also benefits from good fiber connectivity to various locations in the U.S., Asia, and Europe, thanks to a large surplus of fiber optic networking, a legacy of the dot-com boom.

In 2007, Google built at least four new data centers at an average cost of $600 million, each adding to its Googleplex: a massive global computer network estimated to span 25 locations and 450,000 servers. Amazon also chose a Dalles location down the river for its largest data center.

Yahoo! and Microsoft chose Quincy, Washington. Microsoft's new facility there has more than 477,000 square feet of space, nearly the area of 10 football fields. The company is tight-lipped about the number of servers at the site, but it does say the facility uses 3 miles of chiller piping, 600 miles of electrical wire, 1 million square feet of drywall, and 1.6 tons of batteries for backup power. And the data center consumes 48 megawatts—enough power for 40,000 homes.

---

### World's servers surpassing Holland's emissions

The management consulting firm McKinsey & Co. reports that the world's 44 million servers consume 0.5 percent of all electricity and produce 0.2 percent of all carbon dioxide emissions, or 80 megatons a year, approaching the emissions of entire countries such as Argentina or the Netherlands.
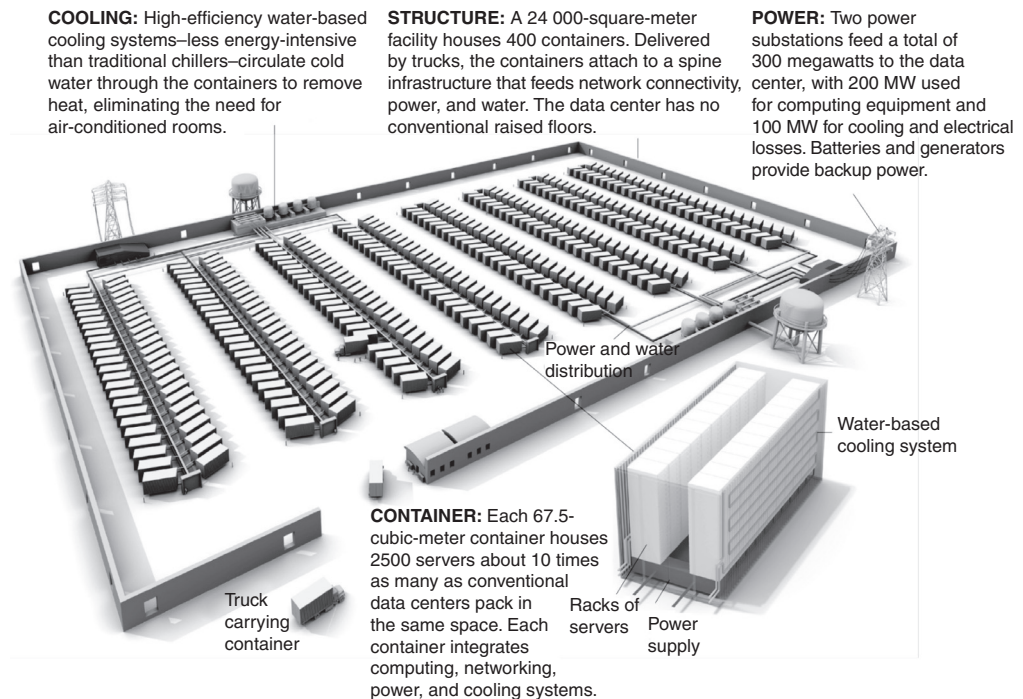
## CLOUD DATA CENTERS: BECOMING MORE EFFICIENT AND MORE FLEXIBLE THROUGH MODULARITY

Already, through volume purchasing, custom server construction, and careful geographic locality, the world's largest data-center owners can build data centers at a fraction of the cost per CPU operation of private corporations. They relentlessly work to widen that gap. The economies-of-scale trend will continue in the cloud providers' favor as they become dramatically more efficient through modular data centers. These highly modular, scalable, efficient, just-in-time data centers can provide capacity that can be delivered anywhere in the world quickly and cheaply.

Figure 2.2 is an artist's rendering of a modular data center (because photographs of such facilities are highly guarded). Corporate data centers can't compete with the myriad economic efficiencies that these mega data centers can achieve today and will fall further and further behind as time goes by.

The goal behind modular data centers is to standardize them and move away from custom designs, enabling a commoditized manufacturing approach. The most striking feature is that such data centers are roofless.

Like Google, Microsoft is driven by energy costs and environmental pressures to reduce emissions and increase efficiency. The company's goal is a power usage effectiveness (PUE) at or below 1.125 by 2012 across all its data centers.

**COOLING:** High-efficiency water-based cooling systems–less energy-intensive than traditional chillers–circulate cold water through the containers to remove heat, eliminating the need for air-conditioned rooms.

**STRUCTURE:** A 24 000-square-meter facility houses 400 containers. Delivered by trucks, the containers attach to a spine infrastructure that feeds network connectivity, power, and water. The data center has no conventional raised floors.

**POWER:** Two power substations feed a total of 300 megawatts to the data center, with 200 MW used for computing equipment and 100 MW for cooling and electrical losses. Batteries and generators provide backup power.

Power and water distribution

Water-based cooling system

**CONTAINER:** Each 67.5-cubic-meter container houses 2500 servers about 10 times as many as conventional data centers pack in the same space. Each container integrates computing, networking, power, and cooling systems.

Truck carrying container

Racks of servers

Power supply

**Figure 2.2   Expandable, modular cloud data center. Notice there is no roof. New containers with servers, power, cooling and network taps can be swapped in and out as needed. Source: *IEEE Spectrum* magazine.**

**Power usage effectiveness (PUE)**

Power usage effectiveness (PUE) is a metric used to determine the energy efficiency of a data center. PUE is determined by dividing the amount of power entering a data center by the power used to run the computer infrastructure within it. PUE is therefore expressed as a ratio, with overall efficiency improving as the quotient decreases toward 1.

According to the Uptime Institute, the typical data center has an average PUE of 2.5. This means that for every 2.5 watts in at the utility meter, only 1 watt is delivered out to the IT load. Uptime estimates that most facilities could achieve 1.6 PUE using the most efficient equipment and best practices. Google and Microsoft are both approaching 1.125, far exceeding what any corporate or cohost data center can achieve.