

# Overview of Data Centers

---

This chapter presents an overview of enterprise Data Center environments, current application environment trends, the Data Center network architecture, and the services provided by the architecture. The approach to develop the architecture of the Data Center network is typically an internal process based on the requirement of the enterprise. This chapter provides the design criteria used by the authors to define the Data Center design best practices presented throughout the book.

## Data Centers Defined

Data Centers house critical computing resources in controlled environments and under centralized management, which enable enterprises to operate around the clock or according to their business needs. These computing resources include mainframes; web and application servers; file and print servers; messaging servers; application software and the operating systems that run them; storage subsystems; and the network infrastructure, whether IP or storage-area network (SAN). Applications range from internal financial and human resources to external e-commerce and business-to-business applications. Additionally, a number of servers support network operations and network-based applications. Network operation applications include Network Time Protocol (NTP); TN3270; FTP; Domain Name System (DNS); Dynamic Host Configuration Protocol (DHCP); Simple Network Management Protocol (SNMP); TFTP; Network File System (NFS); and network-based applications, including IP telephony, video streaming over IP, IP video conferencing, and so on.

According to a report from the Renewable Energy Policy Project on Energy Smart Data Centers, Data Centers are

... an essential component of the infrastructure supporting the Internet and the digital commerce and electronic communication sector. Continued growth of these sectors requires a reliable infrastructure because ... interruptions in digital services can have significant economic consequences.

Virtually, every enterprise has one or more Data Centers. Some have evolved rapidly to accommodate various enterprise application environments using distinct operating systems and hardware platforms. The evolution has resulted in complex and disparate environments that are expensive to manage and maintain. In addition to the application environment, the

supporting network infrastructure might not have changed fast enough to be flexible in accommodating ongoing redundancy, scalability, security, and management requirements.

A Data Center network design lacking in any of these areas risks not being able to sustain the expected service level agreements (SLAs). Data Center downtime, service degradation, or the inability to roll new services implies that SLAs are not met, which leads to a loss of access to critical resources and a quantifiable impact on normal business operation. The impact could be as simple as increased response time or as severe as loss of data.

## Data Center Goals

The benefits provided by a Data Center include traditional business-oriented goals such as the support for business operations around the clock (resiliency), lowering the total cost of operation and the maintenance needed to sustain the business functions (total cost of ownership), and the rapid deployment of applications and consolidation of computing resources (flexibility).

These business goals generate a number of information technology (IT) initiatives, including the following:

- Business continuance
- Increased security in the Data Center
- Application, server, and Data Center consolidation
- Integration of applications whether client/server and multitier (n-tier), or web services-related applications
- Storage consolidation

These IT initiatives are a combination of the need to address short-term problems and establishing a long-term strategic direction, all of which require an architectural approach to avoid unnecessary instability if the Data Center network is not flexible enough to accommodate future changes. The design criteria are

- Availability
- Scalability
- Security
- Performance
- Manageability

These design criteria are applied to these distinct functional areas of a Data Center network:

- **Infrastructure services**—Routing, switching, and server-farm architecture
- **Application services**—Load balancing, Secure Socket Layer (SSL) offloading, and caching

- **Security services**—Packet filtering and inspection, intrusion detection, and intrusion prevention
- **Storage services**—SAN architecture, Fibre Channel switching, backup, and archival
- **Business continuance**—SAN extension, site selection, and Data Center interconnectivity

The details of these services are discussed later in this chapter.

## Data Center Facilities

Because Data Centers house critical computing resources, enterprises must make special arrangements with respect to both the facilities that house the equipment and the personnel required for a 24-by-7 operation. These facilities are likely to support a high concentration of server resources and network infrastructure. The demands posed by these resources, coupled with the business criticality of the applications, create the need to address the following areas:

- Power capacity
- Cooling capacity
- Cabling
- Temperature and humidity controls
- Fire and smoke systems
- Physical security: restricted access and surveillance systems
- Rack space and raised floors

Discussing the facilities where the Data Center resides and the related planning functions is outside the scope of this book.

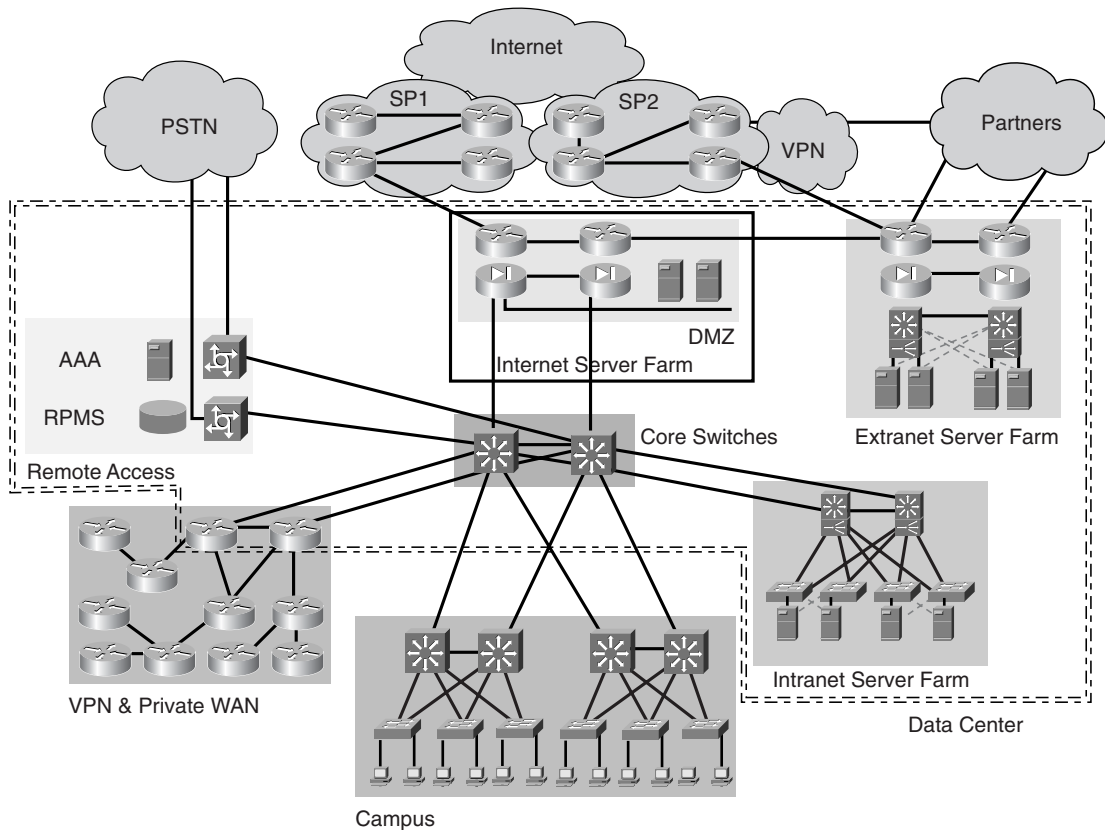
The sections that follow introduce the role of the Data Center in the enterprise network.

## Roles of Data Centers in the Enterprise

Figure 1-1 presents the different building blocks used in the typical enterprise network and illustrates the location of the Data Center within that architecture.

The building blocks of this typical enterprise network include

- Campus network
- Private WAN
- Remote access
- Internet server farm
- Extranet server farm
- Intranet server farm

**Figure 1-1** *Data Centers in the Enterprise*

Data Centers typically house many components that support the infrastructure building blocks, such as the core switches of the campus network or the edge routers of the private WAN. Data Center designs can include any or all of the building blocks in Figure 1-1, including any or all server farm types. Each type of server farm can be a separate physical entity, depending on the business requirements of the enterprise. For example, a company might build a single Data Center and share all resources, such as servers, firewalls, routers, switches, and so on. Another company might require that the three server farms be physically separated with no shared equipment. This book focuses on the details of architecting server farms in the context of a highly available and scalable Data Center. These server farms support a wide number of enterprise applications.

Enterprise applications typically focus on one of the following major business areas:

- Customer relationship management (CRM)
- Enterprise resource planning (ERP)

- Supply chain management (SCM)
- Sales force automation (SFA)
- Order processing
- E-commerce

## Roles of Data Centers in the Service Provider Environment

Data Centers in service provider (SP) environments, known as Internet Data Centers (IDCs), unlike in enterprise environments, are the source of revenue that supports collocated server farms for enterprise customers. The SP Data Center is a service-oriented environment built to house, or *host*, an enterprise customer's application environment under tightly controlled SLAs for uptime and availability. Enterprises also build IDCs when the sole reason for the Data Center is to support Internet-facing applications.

The IDCs are separated from the SP internal Data Centers that support the internal business applications environments.

Whether built for internal facing or collocated applications, application environments follow specific application architectural models such as the classic client/server or the n-tier model.

## Application Architecture Models

Application architectures are constantly evolving, adapting to new requirements, and using new technologies. The most pervasive models are the client/server and n-tier models that refer to how applications use the functional elements of communication exchange. The client/server model, in fact, has evolved to the n-tier model, which most enterprise software application vendors currently use in application architectures. This section introduces both models and the evolutionary steps from client/server to the n-tier model.

### The Client/Server Model and Its Evolution

The classic client/server model describes the communication between an application and a user through the use of a server and a client. The classic client/server model consists of the following:

- A thick client that provides a graphical user interface (GUI) on top of an application or business logic where some processing occurs
- A server where the remaining business logic resides

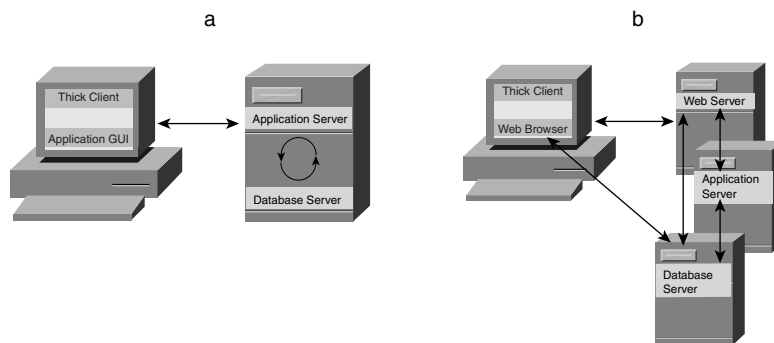
*Thick client* is an expression referring to the complexity of the business logic (software) required on the client side and the necessary hardware to support it. A thick client is then a portion of the application code running at the client's computer that has the responsibility

of retrieving data from the server and presenting it to the client. The thick client code requires a fair amount of processing capacity and resources to run in addition to the management overhead caused by loading and maintaining it on the client base.

The server side is a single server running the presentation, application, and database code that uses multiple internal processes to communicate information across these distinct functions. The exchange of information between client and server is mostly data because the thick client performs local presentation functions so that the end user can interact with the application using a local user interface.

Client/server applications are still widely used, yet the client and server use proprietary interfaces and message formats that different applications cannot easily share. Part **a** of Figure 1-2 shows the client/server model.

**Figure 1-2** *Client/Server and n-Tier Application Interaction*



The most fundamental changes to the thick client and single-server model started when web-based applications first appeared. Web-based applications rely on more standard interfaces and message formats where applications are easier to share. HTML and HTTP provide a standard framework that allows generic clients such as web browsers to communicate with generic applications as long as they use web servers for the presentation function. HTML describes how the client should render the data; HTTP is the transport protocol used to carry HTML data. Netscape Communicator and Microsoft Internet Explorer are examples of clients (web browsers); Apache, Netscape Enterprise Server, and Microsoft Internet Information Server (IIS) are examples of web servers.

The migration from the classic client/server to a web-based architecture implies the use of thin clients (web browsers), web servers, application servers, and database servers. The web browser interacts with web servers and application servers, and the web servers interact with application servers and database servers. These distinct functions supported by the servers are referred to as *tiers*, which, in addition to the client tier, refer to the *n-tier model*.

## The n-Tier Model

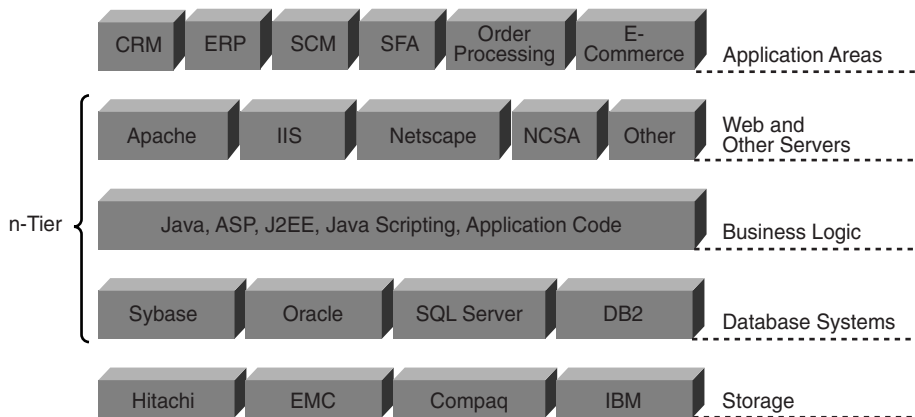
Part **b** of Figure 1-2 shows the n-tier model. Figure 1-2 presents the evolution from the classic client/server model to the n-tier model. The client/server model uses the thick client with its own business logic and GUI to interact with a server that provides the counterpart business logic and database functions on the same physical device. The n-tier model uses a thin client and a web browser to access the data in many different ways. The server side of the n-tier model is divided into distinct functional areas that include the web, application, and database servers.

The n-tier model relies on a standard web architecture where the web browser formats and presents the information received from the web server. The server side in the web architecture consists of multiple and distinct servers that are functionally separate. The n-tier model can be the client and a web server; or the client, the web server, and an application server; or the client, web, application, and database servers. This model is more scalable and manageable, and even though it is more complex than the classic client/server model, it enables application environments to evolve toward distributed computing environments.

The n-tier model marks a significant step in the evolution of distributed computing from the classic client/server model. The n-tier model provides a mechanism to increase performance and maintainability of client/server applications while the control and management of application code is simplified.

Figure 1-3 introduces the n-tier model and maps each tier to a partial list of currently available technologies at each tier.

**Figure 1-3** *n-Tier Model*



Notice that the client-facing servers provide the interface to access the business logic at the application tier. Although some applications provide a non-web-based front end, current trends indicate the process of “web-transforming” business applications is well underway.

This process implies that the front end relies on a web-based interface to face the users which interacts with a middle layer of applications that obtain data from the back-end systems.

These middle tier applications and the back-end database systems are distinct pieces of logic that perform specific functions. The logical separation of front-end application and back-end functions has enabled their physical separation. The implications are that the web and application servers, as well as application and database servers, no longer have to coexist in the same physical server. This separation increases the scalability of the services and eases the management of large-scale server farms. From a network perspective, these groups of servers performing distinct functions could also be physically separated into different network segments for security and manageability reasons.

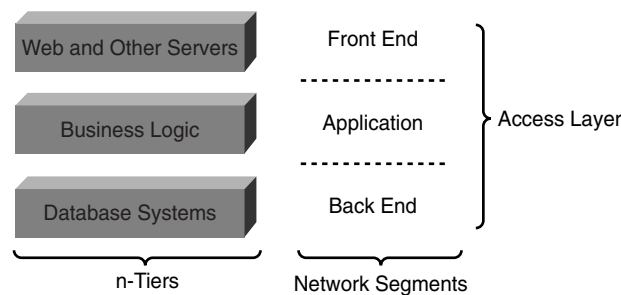
Chapter 3, “Application Architectures Overview,” discusses the details on applications that follow the n-tier model and the implications on the design of the Data Center.

## Multitier Architecture Application Environment

Multitier architectures refer to the Data Center server farms supporting applications that provide a logical and physical separation between various application functions, such as web, application, and database (n-tier model). The network architecture is then dictated by the requirements of applications in use and their specific availability, scalability, and security and management goals. For each server-side tier, there is a one-to-one mapping to a network segment that supports the specific application function and its requirements. Because the resulting network segments are closely aligned with the tiered applications, they are described in reference to the different application tiers.

Figure 1-4 presents the mapping from the n-tier model to the supporting network segments used in a multitier design.

**Figure 1-4** *Multitier Network Segments*



The web server tier is mapped to the *front-end segment*, the business logic to the *application segment*, and the database tier to the *back-end segment*. Notice that all the segments supporting the server farm connect to access layer switches, which in a multitier architecture are different access switches supporting the various server functions.



The evolution of application architectures and departing from multitier application environments still requires a network to support the interaction between the communicating entities. For example, a web service (defined as “A web service is a software system designed to support interoperable machine-to-machine interaction over a network” by the W3C web services architecture document) still refers to the network element. In this case, the network would be used for networked resources that support such interaction reliably. This layer of abstraction does not necessarily translate on to a layered network design as much as the capability of the network to support networked applications, resources, and their interaction.

The following section presents a high-level overview of the distinct network layers of the Data Center architecture.

## Data Center Architecture

The enterprise Data Center architecture is inclusive of many functional areas, as presented earlier in Figure 1-1. The focus of this section is the architecture of a generic enterprise Data Center connected to the Internet and supporting an intranet server farm. Other types of server farms, explained in Chapter 4, “Data Center Design Overview,” follow the same architecture used for intranet server farms yet with different scalability, security, and management requirements. Figure 1-5 introduces the topology of the Data Center architecture.

**Figure 1-5** *Topology of an Enterprise Data Center Architecture*

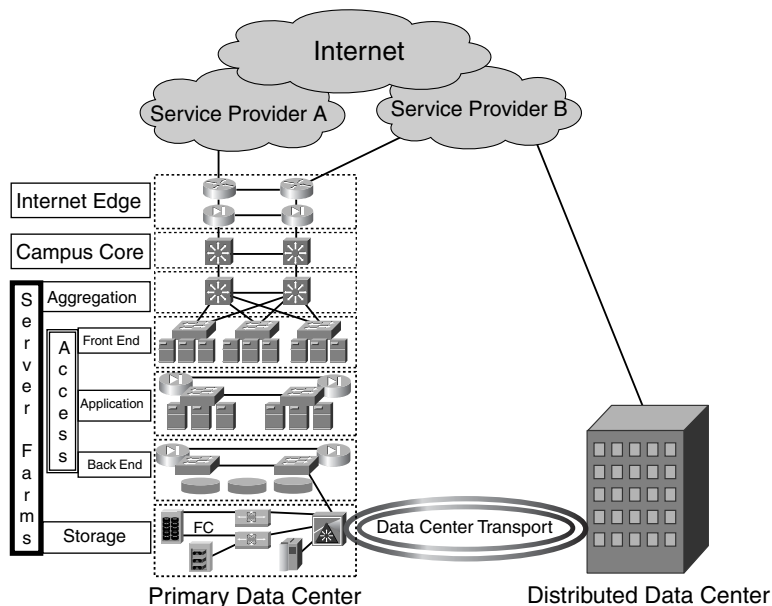


Figure 1-5 shows a fully redundant enterprise Data Center supporting the following areas:

- No single point of failure—redundant components
- Redundant Data Centers

Although the focus of this book is the architecture of an IP network that supports server farms, we include explanations pertaining to how the server farms are connected to the rest of the enterprise network for the sake of clarity and thoroughness. The core connectivity functions supported by Data Centers are Internet Edge connectivity, campus connectivity, and server-farm connectivity, as presented in Figure 1-5.

The Internet Edge provides the connectivity from the enterprise to the Internet and its associated redundancy and security functions, such as the following:

- Redundant connections to different service providers
- External and internal routing through exterior border gateway protocol (EBGP) and interior border gateway protocol (IBGP)
- Edge security to control access from the Internet
- Control for access to the Internet from the enterprise clients

The campus core switches provide connectivity between the Internet Edge, the intranet server farms, the campus network, and the private WAN. The core switches physically connect to the devices that provide access to other major network areas, such as the private WAN edge routers, the server-farm aggregation switches, and campus distribution switches.

As depicted in Figure 1-6, the following are the network layers of the server farm:

- Aggregation layer
- Access layer
  - Front-end segment
  - Application segment
  - Back-end segment
- Storage layer
- Data Center transport layer

Some of these layers depend on the specific implementation of the n-tier model or the requirements for Data Center-to-Data-Center connectivity, which implies that they might not exist in every Data Center implementation. Although some of these layers might be optional in the Data Center architecture, they represent the trend in continuing to build highly available and scalable enterprise Data Centers. This trend specifically applies to the storage and Data Center transport layers supporting storage consolidation, backup and archival consolidation, high-speed mirroring or clustering between remote server farms, and so on.

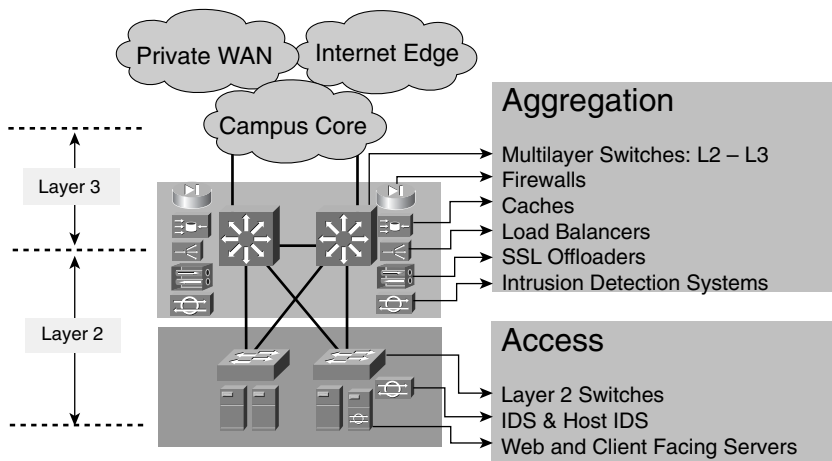
The sections that follow present the specific details of each layer.

## Aggregation Layer

The aggregation layer is the aggregation point for devices that provide services to all server farms. These devices are multilayer switches, firewalls, load balancers, and other devices that typically support services across all servers. The multilayer switches are referred to as aggregation switches because of the aggregation function they perform. Service devices are shared by all server farms. Specific server farms are likely to span multiple access switches for redundancy, thus making the aggregation switches the logical connection point for service devices, instead of the access switches.

If connected to the front-end Layer 2 switches, these service devices might not offer optimal services by creating less than optimal traffic paths between them and servers connected to different front-end switches. Additionally, if the service devices are off of the aggregation switches, the traffic paths are deterministic and predictable and simpler to manage and maintain. Figure 1-6 shows the typical devices at the aggregation layer.

**Figure 1-6** *Aggregation and Access Layers*



As depicted in Figure 1-6, the aggregation switches provide basic infrastructure services and connectivity for other service devices. The aggregation layer is analogous to the traditional distribution layer in the campus network in its Layer 3 and Layer 2 functionality.

The aggregation switches support the traditional switching of packets at Layer 3 and Layer 2 in addition to the protocols and features to support Layer 3 and Layer 2 connectivity. A more in-depth explanation on the specific services provided by the aggregation layer appears in the section, “Data Center Services.”

## Access Layer

The access layer provides Layer 2 connectivity and Layer 2 features to the server farm. Because in a multitier server farm, each server function could be located on different access switches on different segments, the following section explains the details of each segment.

### Front-End Segment

The front-end segment consists of Layer 2 switches, security devices or features, and the front-end server farms. See the section, “Data Center Services” for a detailed description of the features provided by the devices at this layer. The front-end segment is analogous to the traditional access layer of the hierarchical campus network design and provides the same functionality. The access switches are connected to the aggregation switches in the manner depicted in Figure 1-6. The front-end server farms typically include FTP, Telnet, TN3270 (mainframe terminals), Simple Mail Transfer Protocol (SMTP), web servers, DNS servers, and other business application servers, in addition to network-based application servers such as IP television (IPTV) broadcast servers and IP telephony call managers that are not placed at the aggregation layer because of port density or other design requirements.

The specific network features required in the front-end segment depend on the servers and their functions. For example, if a network supports video streaming over IP, it might require multicast, or if it supports Voice over IP (VoIP), quality of service (QoS) must be enabled. Layer 2 connectivity through VLANs is required between servers and load balancers or firewalls that segregate server farms.

The need for Layer 2 adjacency is the result of Network Address Translation (NAT) and other header rewrite functions performed by load balancers or firewalls on traffic destined to the server farm. The return traffic must be processed by the same device that performed the header rewrite operations.

Layer 2 connectivity is also required between servers that use clustering for high availability or require communicating on the same subnet. This requirement implies that multiple access switches supporting front-end servers can support the same set of VLANs to provide layer adjacency between them.

Security features include Address Resolution Protocol (ARP) inspection, broadcast suppression, private VLANs, and others that are enabled to counteract Layer 2 attacks. Security devices include network-based intrusion detection systems (IDSs) and host-based IDSs to monitor and detect intruders and prevent vulnerabilities from being exploited. In general, the infrastructure components such as the Layer 2 switches provide intelligent network services that enable front-end servers to provide their functions.

Note that the front-end servers are typically taxed in their I/O and CPU capabilities. For I/O, this strain is a direct result of serving content to the end users; for CPU, it is the connection rate and the number of concurrent connections needed to be processed.

Scaling mechanisms for front-end servers typically include adding more servers with identical content and then equally distributing the load they receive using load balancers. Load balancers distribute the load (or load balance) based on Layer 4 or Layer 5 information. Layer 4 is widely used for front-end servers to sustain a high connection rate without necessarily overwhelming the servers. See Chapter 22, “Performance Metrics of Data Center Devices,” to understand the performance of servers and load balancers under load.

Scaling mechanisms for web servers also include the use of SSL offloaders and Reverse Proxy Caching (RPC). Refer to Chapter 9, “SSL and TLS,” for more information about the use of SSL and its performance implications.

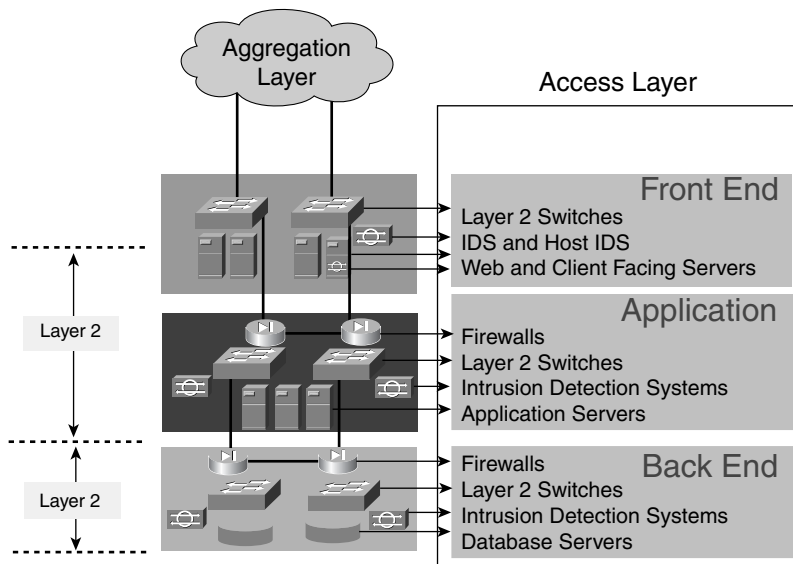
## Application Segment

The application segment has the same network infrastructure components as the front-end segment and the application servers. The features required by the application segment are almost identical to those needed in the front-end segment, albeit with additional security. This segment relies strictly on Layer 2 connectivity, yet the additional security is a direct requirement of how much protection the application servers need because they have direct access to the database systems. Depending on the security policies, this segment uses firewalls between web and application servers, IDSs, and host IDSs. Like the front-end segment, the application segment infrastructure must support intelligent network services as a direct result of the functions provided by the application services.

Application servers run a portion of the software used by business applications and provide the communication logic between the front end and the back end, which is typically referred to as the *middleware* or *business logic*. Application servers translate user requests to commands that the back-end database systems understand. Increasing the security at this segment focuses on controlling the protocols used between the front-end servers and the application servers to avoid trust exploitation and attacks that exploit known application vulnerabilities. Figure 1-7 introduces the front-end, application, and back-end segments in a logical topology.

Note that the application servers are typically CPU-stressed because they need to support the business logic. Scaling mechanisms for application servers also include load balancers. Load balancers can select the right application server based on Layer 5 information.

Deep packet inspection on load balancers allows the partitioning of application server farms by content. Some server farms could be dedicated to selecting a server farm based on the scripting language (.cgi, .jsp, and so on). This arrangement allows application administrators to control and manage the server behavior more efficiently.

**Figure 1-7** *Access Layer Segments*

## Back-End Segment

The back-end segment is the same as the previous two segments except that it supports the connectivity to database servers. The back-end segment features are almost identical to those at the application segment, yet the security considerations are more stringent and aim at protecting the data, critical or not.

The hardware supporting the database systems ranges from medium-sized servers to high-end servers, some with direct locally attached storage and others using disk arrays attached to a SAN. When the storage is separated, the database server is connected to both the Ethernet switch and the SAN. The connection to the SAN is through a Fibre Channel interface. Figure 1-8 presents the back-end segment in reference to the storage layer. Notice the connections from the database server to the back-end segment and storage layer.

Note that in other connectivity alternatives, the security requirements do not call for physical separation between the different server tiers. These alternatives are discussed in Chapter 4.

## Storage Layer

The storage layer consists of the storage infrastructure such as Fibre Channel switches and routers that support small computer system interface (SCSI) over IP (iSCSI) or Fibre Channel over IP (FCIP). Storage network devices provide the connectivity to servers, storage devices such as disk subsystems, and tape subsystems.

---

### NOTE

SAN environments in Data Centers commonly use Fibre Channel to connect servers to the storage device and to transmit SCSI commands between them. Storage networks allow the transport of SCSI commands over the network. This transport is possible over the Fibre Channel infrastructure or over IP using FCIP and iSCSI.

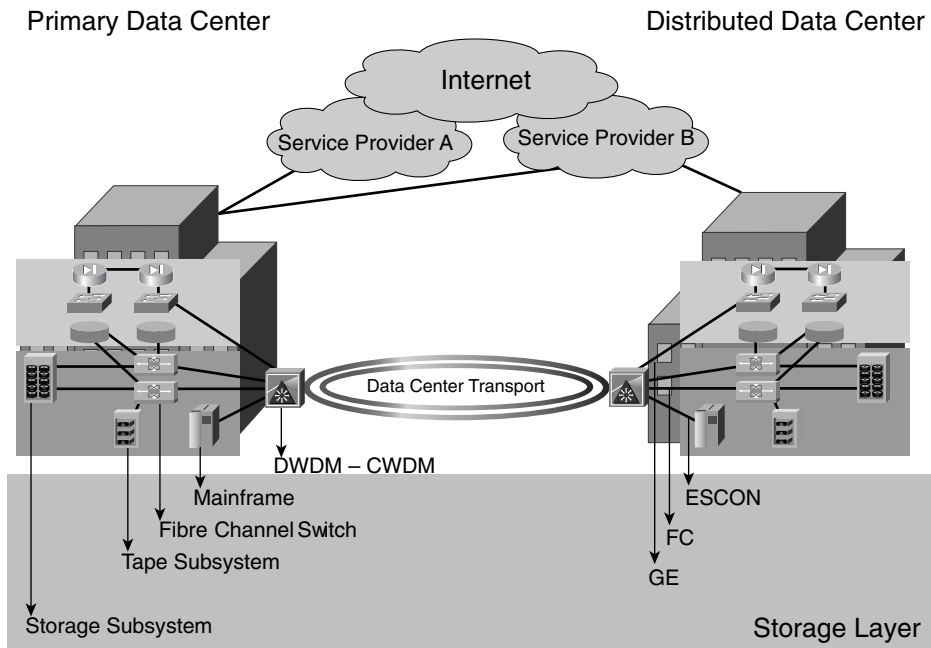
FCIP and iSCSI are the emerging Internet Engineering Task Force (IETF) standards that enable SCSI access and connectivity over IP.

---

The network used by these storage devices is referred to as a SAN. The Data Center is the location where the consolidation of applications, servers, and storage occurs and where the highest concentration of servers is likely, thus where SANs are located. The current trends in server and storage consolidation are the result of the need for increased efficiency in the application environments and for lower costs of operation.

Data Center environments are expected to support high-speed communication between servers and storage and between storage devices. These high-speed environments require block-level access to the information supported by SAN technology. There are also requirements to support file-level access specifically for applications that use Network Attached Storage (NAS) technology. Figure 1-8 introduces the storage layer and the typical elements of single and distributed Data Center environments.

Figure 1-8 shows a number of database servers as well as tape and disk arrays connected to the Fibre Channel switches. Servers connected to the Fibre Channel switches are typically critical servers and always dual-homed. Other common alternatives to increase availability include mirroring, replication, and clustering between database systems or storage devices. These alternatives typically require the data to be housed in multiple facilities, thus lowering the likelihood of a site failure preventing normal systems operation. Site failures are recovered by replicas of the data at different sites, thus creating the need for distributed Data Centers and distributed server farms and the obvious transport technologies to enable communication between them. The following section discusses Data Center transport alternatives.

**Figure 1-8** *Storage and Transport Layers*

## Data Center Transport Layer

The Data Center transport layer includes the transport technologies required for the following purposes:

- Communication between distributed Data Centers for rerouting client-to-server traffic
- Communication between distributed server farms located in distributed Data Centers for the purposes of remote mirroring, replication, or clustering

Transport technologies must support a wide range of requirements for bandwidth and latency depending on the traffic profiles, which imply a number of media types ranging from Ethernet to Fibre Channel.

For user-to-server communication, the possible technologies include Frame Relay, ATM, DS channels in the form of T1/E1 circuits, Metro Ethernet, and SONET.

For server-to-server and storage-to-storage communication, the technologies required are dictated by server media types and the transport technology that supports them transparently. For example, as depicted in Figure 1-8, storage devices use Fibre Channel and Enterprise



Systems Connectivity (ESCON), which should be supported by the metro optical transport infrastructure between the distributed server farms.

If ATM and Gigabit Ethernet (GE) are used between distributed server farms, the metro optical transport could consolidate the use of fiber more efficiently. For example, instead of having dedicated fiber for ESCON, GE, and ATM, the metro optical technology could transport them concurrently.

The likely transport technologies are dark fiber, coarse wavelength division multiplexing (CWDM), and dense wavelength division multiplexing (DWDM), which offer transparent connectivity (Layer 1 transport) between distributed Data Centers for media types such as GE, Fibre Channel, ESCON, and fiber connectivity (FICON).

Note that distributed Data Centers often exist to increase availability and redundancy in application environments. The most common driving factors are disaster recovery and business continuance, which rely on the specific application environments and the capabilities offered by the transport technologies.

- Blade servers
- Grid computing
- Web services
- Service-oriented Data Centers

All these trends influence the Data Center in one way or another. Some short-term trends force design changes, while some long-term trends force a more strategic view of the architecture.

For example, the need to lower operational costs and achieve better computing capacity at a relatively low price leads to the use of blade servers. Blade servers require a different topology when using Ethernet switches inside the blade chassis, which requires planning on port density, slot density, oversubscription, redundancy, connectivity, rack space, power consumption, heat dissipation, weight, and cabling. Blade servers can also support compute grids. Compute grids might be geographically distributed, which requires a clear understanding of the protocols used by the grid middleware for provisioning and load distribution, as well as the potential interaction between a compute grid and a data grid.

Blade servers can also be used to replace 1RU servers on web-based applications because of scalability reasons or the deployment of tiered applications. This physical separation of tiers and the ever-increased need for security leads to application layer firewalls.

An example of this is the explicit definition for application layer security (included in the Web Services Architecture [WSA] document). Security on Web Services is in reference to a secure environment for online processes from a security and privacy perspective. The development of the WSA focuses on the identification of threats to security and privacy and the architect features that are needed to respond to those threats. The infrastructure to support such security is expected to be consistently supported by applications that are expected to be distributed on the network. Past experiences suggest that some computationally

repeatable tasks would, over time, be offloaded to network devices, and that the additional network intelligence provides a more robust infrastructure to complement Web Services security (consistency- and performance-wise).

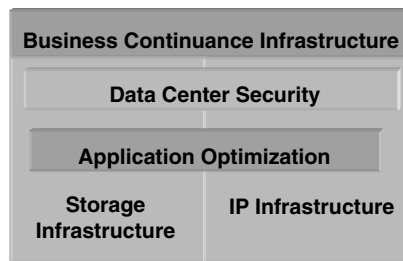
Finally, a services-oriented Data Center implies a radical change on how Data Centers are viewed by their users, which invariably requires a radical change in the integration of the likely services. In this case, interoperability and manageability of the service devices become a priority for the Data Center designers. Current trends speak to the Utility Data Center from HP and On Demand (computing) from IBM, in which both closer integration of available services and the manner in which they are managed and provisioned is adaptable to the organization. This adaptability comes from the use of standard interfaces between the integrated services, but go beyond to support virtualization and self-healing capabilities. Whatever these terms end up bringing to the Data Center, the conclusion is obvious: The Data Center is the location where users, applications, data, and the network infrastructure converge. The result of current trends will change the ways in which the Data Center is architected and managed.

Chapter 4 discusses some of these trends in more detail. The following section discusses the different services a Data Center is expected to support.

## Data Center Services

This section presents an overview of the services supported by the Data Center architecture. Related technology and features make up each service. Each service enhances the manner in which the network operates in each of the functional service areas defined earlier in the chapter. The following sections introduce each service area and its associated features. Figure 1-9 introduces the Data Center services.

**Figure 1-9** *Data Center Services*



As depicted in Figure 1-9, services in the Data Center are not only related to one another but are also, in certain cases, dependent on each other. The IP and storage infrastructure services are the pillars of all other services because they provide the fundamental building blocks of any network and thus of any service. After the infrastructure is in place, you can

build server farms to support the application environments. These environments could be optimized utilizing network technology, hence the name *application services*. Security is a service expected to leverage security features on the networking devices that support all other services in addition to using specific security technology. Finally, the business continuance infrastructure, as a service aimed at achieving the highest possible redundancy level. The highest redundancy level is possible by using both the services of the primary Data Center and their best practices on building a distributed Data Center environment.

## IP Infrastructure Services

Infrastructure services include all core features needed for the Data Center IP infrastructure to function and serve as the foundation, along with the storage infrastructure, of all other Data Center services. The IP infrastructure features are organized as follows:

- Layer 2
- Layer 3
- Intelligent Network Services

Layer 2 features support the Layer 2 adjacency between the server farms and the service devices (VLANs); enable media access; and support a fast convergence, loop-free, predictable, and scalable Layer 2 domain. Layer 2 domain features ensure the Spanning Tree Protocol (STP) convergence time for deterministic topologies is in single-digit seconds and that the failover and failback scenarios are predictable.

STP is available on Cisco network devices in three versions: Per VLAN Spanning Tree Plus (PVST+), Rapid PVST+ (which combines PVST+ and IEEE 802.1w), and Multiple Spanning Tree (IEEE 801.1s combined with IEEE 802.1w). VLANs and trunking (IEEE 802.1Q) are features that make it possible to virtualize the physical infrastructure and, as a consequence, consolidate server-farm segments. Additional features and protocols increase the availability of the Layer 2 network, such as Loopguard, Unidirectional Link Detection (UDLD), PortFast, and the Link Aggregation Control Protocol (LACP or IEEE 802.3ad).

Layer 3 features enable a fast-convergence routed network, including redundancy for basic Layer 3 services such as default gateway support. The purpose is to maintain a highly available Layer 3 environment in the Data Center where the network operation is predictable under normal and failure conditions. The list of available features includes the support for static routing, Border Gateway Protocol (BGP) and Interior Gateway Protocols (IGPs) such as Open Shortest Path First (OSPF), Enhanced Interior Gateway Routing Protocol (EIGRP), Intermediate System-to-Intermediate System (IS-IS), gateway redundancy protocols such as the Hot Standby Routing Protocol (HSRP), Multigroup HSRP (MHSRP), and Virtual Router Redundancy Protocol (VRRP) for default gateway support.

Intelligent network services encompass a number of features that enable application services network-wide. The most common features are QoS and multicast. Yet there are other important intelligent network services such as private VLANs (PVLANS) and policy-based routing

(PBR). These features enable applications, such as live or on-demand video streaming and IP telephony, in addition to the classic set of enterprise applications.

QoS in the Data Center is important for two reasons—marking at the source of application traffic and the port-based rate-limiting capabilities that enforce proper QoS service class as traffic leaves the server farms. For marked packets, it is expected that the rest of the enterprise network enforces the same QoS policies for an end-to-end QoS service over the entire network.

Multicast in the Data Center enables the capabilities needed to reach multiple users concurrently. Because the Data Center is the source of the application traffic, such as live video streaming over IP, multicast must be supported at the server-farm level (VLANs, where the source of the multicast stream is generated). As with QoS, the rest of the enterprise network must either be multicast-enabled or use tunnels to get the multicast stream to the intended destinations.

## Application Services

Application services include a number of features that enhance the network awareness of applications and use network intelligence to optimize application environments. These features are equally available to scale server-farm performance, to perform packet inspection at Layer 4 or Layer 5, and to improve server response time. The server-farm features are organized by the devices that support them. The following is a list of those features:

- Load balancing
- Caching
- SSL termination

Load balancers perform two core functions:

- Scale and distribute the load to server farms
- Track server health to ensure high availability

To perform these functions, load balancers virtualize the services offered by the server farms by front-ending and controlling the incoming requests to those services. The load balancers distribute requests across multiple servers based on Layer 4 or Layer 5 information. The mechanisms for tracking server health include both in-band monitoring and out-of-band probing with the intent of not forwarding traffic to servers that are not operational. You also can add new servers, thus scaling the capacity of a server farm, without any disruption to existing services.

Layer 5 capabilities on a load balancer allow you to segment server farms by the content they serve. For example, you can separate a group of servers dedicated to serve streaming video (running multiple video servers) from other groups of servers running scripts and application code. The load balancer can determine that a request for an .mpg file (a video

file using MPEG) goes to the first group and that a request for a .cgi file (a script file) goes to the second group.

Server farms benefit from caching features, specifically working in RPC mode. Caches operating in RPC mode are placed near the server farm to intercept requests sent to the server farm, thus offloading the serving of static content from the servers. The cache keeps a copy of the content, which is available to any subsequent request for the same content, so that the server farm does not have to process the requests. The process of offloading occurs transparently for both the user and the server farm.

SSL offloading features use an SSL device to offload the processing of SSL sessions from server farms. The key advantage to this approach is that the SSL termination device offloads SSL key negotiation and the encryption/decryption process away from the server farm. An additional advantage is the capability to process packets based on information in the payload that would otherwise be encrypted. Being able to see the payload allows the load balancer to distribute the load based on Layer 4 or Layer 5 information before re-encrypting the packets and sending them off to the proper server.

## Security Services

Security services include the features and technologies used to secure the Data Center infrastructure and application environments. Given the variety of likely targets in the Data Center, it is important to use a systems approach to securing the Data Center. This comprehensive approach considers the use of all possible security tools in addition to hardening every network device and using a secure management infrastructure. The security tools and features are as follows:

- Access control lists (ACLs)
- Firewalls
- IDSs and host IDSs
- Secure management
- Layer 2 and Layer 3 security features

ACLs filter packets. Packet filtering through ACLs can prevent unwanted access to network infrastructure devices and, to a lesser extent, protect server-farm application services. ACLs are applied on routers (RACLs) to filter routed packets and to VLANs (VACLs) to filter intra-VLAN traffic. Other features that use ACLs are QoS and security, which are enabled for specific ACLs.

An important feature of ACLs is the capability to perform packet inspection and classification without causing performance bottlenecks. You can perform this lookup process in hardware, in which case the ACLs operate at the speed of the media (wire speed).

The placement of firewalls marks a clear delineation between highly secured and loosely secured network perimeters. Although the typical location for firewalls remains the Internet

Edge and the edge of the Data Center, they are also used in multitier server-farm environments to increase security between the different tiers.

IDSs proactively address security issues. Intruder detection and the subsequent notification are fundamental steps for highly secure Data Centers where the goal is to protect the data. Host IDSs enable real-time analysis and reaction to hacking attempts on database, application, and Web servers. The host IDS can identify the attack and prevent access to server resources before any unauthorized transactions occur.

Secure management include the use of SNMP version 3; Secure Shell (SSH); authentication, authorization, and accounting (AAA) services; and an isolated LAN housing the management systems. SNMPv3 and SSH support secure monitoring and access to manage network devices. AAA provides one more layer of security by preventing users access unless they are authorized and by ensuring controlled user access to the network and network devices with a predefined profile. The transactions of all authorized and authenticated users are logged for accounting purposes, for billing, or for postmortem analysis.

## Storage Services

Storage services include the capability of consolidating direct attached disks by using disk arrays that are connected to the network. This setup provides a more effective disk utilization mechanism and allows the centralization of storage management. Two additional services are the capability of consolidating multiple isolated SANs on to the same larger SAN and the virtualization of storage so that multiple servers concurrently use the same set of disk arrays.

Consolidating isolated SANs on to one SAN requires the use of virtual SAN (VSAN) technology available on the SAN switches. VSANs are equivalent to VLANs yet are supported by SAN switches instead of Ethernet switches. The concurrent use of disk arrays by multiple servers is possible through various network-based mechanisms supported by the SAN switch to build logical paths from servers to storage arrays.

Other storage services include the support for FCIP and iSCSI on the same storage network infrastructure. FCIP connects SANs that are geographically distributed, and iSCSI is a lower-cost alternative to Fibre Channel. These services are used both in local SANs and SANs that might be extended beyond a single Data Center. The SAN extension subject is discussed in the next section.

## Business Continuity Infrastructure Services

Business continuity infrastructure services support the highest levels of application availability through the use of networking technology in the three major areas described next.

- Site selection
- SAN extension
- Data Center interconnectivity

Site selection refers to the features that allow the automatic detection of the failure of a Data Center on the application level and the subsequent reroute of all requests to an available site and server farm. You can use the technology for site selection over the Internet or the intranet. The mechanisms for site selection vary from the use of DNS to the host routes and the routed network.

SAN extension refers to the process of stretching an existing SAN to a secondary location, which could be located in the same Data Center or on a different geographical location. You make this extension to allow the replication of the data from the primary to the secondary SAN. Depending on the distance, the application transaction rate, and the latency between the distributed Data Centers, the replication is synchronous or asynchronous. For more information on replication technologies consult Chapter 3.

Data Center interconnectivity services are connectivity alternatives provided by various technologies. These connectivity alternatives support the communication requirements for site selection and SAN extension. The section, “Data Center Transport Layer” earlier in this chapter discussed the available technologies.

## Summary

Data Centers are strategic components of an enterprise that house the critical assets of the business: applications, data, and the computing infrastructure. The Data Center network is vital to sustaining the normal operations of the business. The Data Center network architecture is driven by business requirements.

The criteria that guide the design of a Data Center are availability, scalability, security, performance, and manageability. The Data Center designs described in this book are based on these principles.

The distinct services likely offered by the Data Center network include IP infrastructure connectivity, SAN infrastructure connectivity, application optimizations, security, and business continuance.

The IP infrastructure connectivity function refers to routing and switching. The SAN function refers to the Fibre Channel fabric switching. The application optimization functions include load balancing, caching, and SSL offloading. The security function refers to the use of ACLs, firewalls, IDSs, and secure management. The business continuance function refers to the use of site selection (IP based or DNS based), SAN extension, and Data Center Interconnectivity.

The design process consists of choosing among the available options for each function (IP connectivity, application optimization, security, and business continuance), based on how it meets the high availability, scalability, security, performance, and manageability requirements.

Additionally, the design process must take into account the current trends in application environments you have or are likely to deploy—such as the n-tier model, the adoption of blade servers, or the use of grid computing—and the Data Center network layers to support the aforementioned services. Once the application requirements are clear, the Data Center architecture needs to be qualified to ensure it meets its objectives and satisfies such requirements.

This book primarily focuses on IP-related functions, including the infrastructure design, application optimization, and security.