



The City College
of New York

CSC 36000: Modern Distributed Computing *with AI Agents*

By Saptarashmi Bandyopadhyay

Email: sbandyopadhyay@ccny.cuny.edu

Assistant Professor of Computer Science

City College of New York and Graduate Center at City University of New York

November 26, 2025 CSC 36000

Today's Lecture

Distributed File Systems

- Network File System
- Google File System

Distributed Computing in the World Wide Web

- WWW Architectures
- Real-World Example: Apache

Handling Large Amounts of Traffic

- Web Server Clusters
- Content Delivery Networks

Distributed File Systems

—



Distributing Data

Distributed File Systems (DFS) allow multiple processes to share data over long periods in a secure and reliable way.

The most common way to do this is through a *Client-Server Architecture*

- **Remote Access Model:** The client is offered an interface to a file system managed by a remote server. The client is unaware of actual file locations.
- **Upload/Download Model:** A client downloads the file locally to access it, then uploads it back when finished (e.g., FTP)

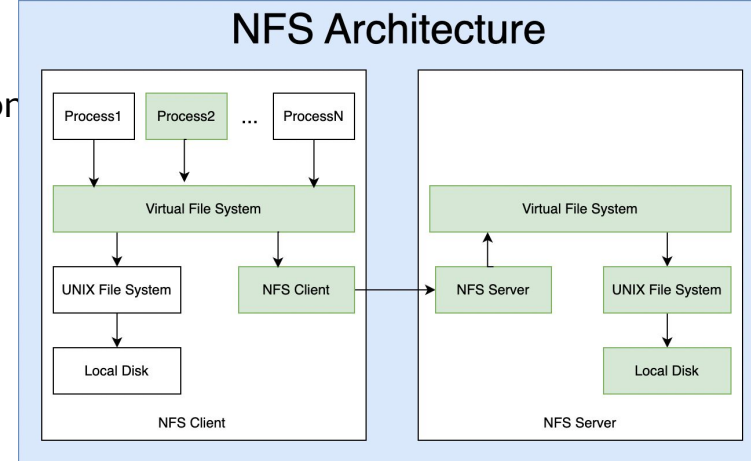
Network File System

NFS is the standard for UNIX-based systems.

- **Virtual File System (VFS):** NFS uses the VFS interface to hide differences between file systems. Operations are passed to either a local file system or the NFS client.
- **Communication:** All client-server communication is done through Remote Procedure Calls (RPCs).

Evolution

- **NFSv3 (Stateless):** The server does not maintain client state. If a server crashes, no recovery phase is needed, but locking files is difficult.
- **NFSv4 (Stateful):** Supports wide-area networks and caching. The server maintains state (e.g., file leases) to manage consistency and locking.



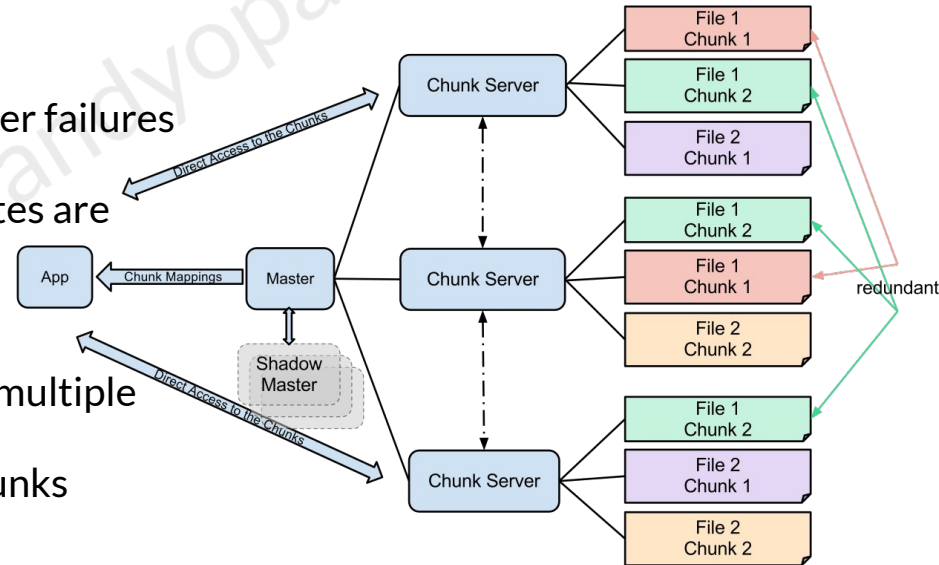
Google File System

Designed for massive data centers where server failures are the norm, not the exception.

- Files are huge (multi-gigabyte) and updates are usually appends rather than overwrites.

Architecture

- Clusters:** Consist of a single Master and multiple Chunk Servers.
- Chunks:** Files are divided into 64 MB chunks distributed across chunk servers.
- Scalability:** The Master handles metadata (namespace), while clients communicate directly with Chunk Servers for actual data.



Distributed Computing in the World Wide Web

—

The World Wide Web

We can think of the entire internet as a giant distributed computing system.

The World Wide Web is a method of distributing documents which can have text, audio, and dynamic features, with references to these documents being made through Uniform Resource Locators (URLs).

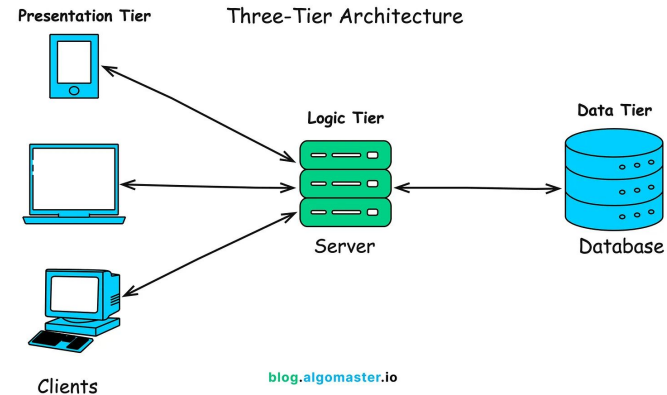
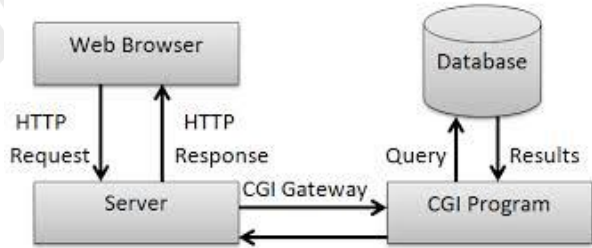
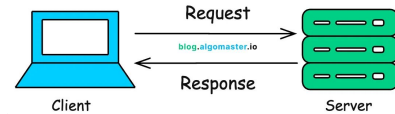


Architectural Tiers of the WWW

Traditional (2-Tier): Simple client-server interaction fetching static documents.

CGI (Common Gateway Interface): Allows the server to execute a program based on user input (e.g., a database query) to generate a document dynamically.

3-Tier Architecture: Modern sites often use a Web Server, an Application Server (business logic/servlets), and a Database.



Real-World Example: Apache

Apache is the most popular Web server, estimated to host ~70% of all sites.

- **Platform Independence:** Uses the Apache Portable Runtime (APR) to hide OS differences (file handling, networking, threa

Organization

- **Hooks:** Apache processes requests in phases (e.g., URL-to-filename translation, authentication, MIME checking). Each phase is a "hook".
- **Modules:** Separate modules provide the functions for these hooks. Developers can extend Apache by writing new modules.



Handling Large Amounts of Traffic

—



Web Server Clusters

Problem: A single Web server can become overloaded. **Solution:** Server Clusters.

Request Handling

- **Transport-Layer Switch:** Passes TCP connections to servers based on load. Simple but cannot inspect content.
- **Content-Aware Distribution:** A front end inspects the HTTP request and forwards it to a specific server (e.g., a server dedicated to video or audio).
- **TCP Handoff:** A switch hands off the connection to a server, which then communicates directly with the client without passing data back through the switch.

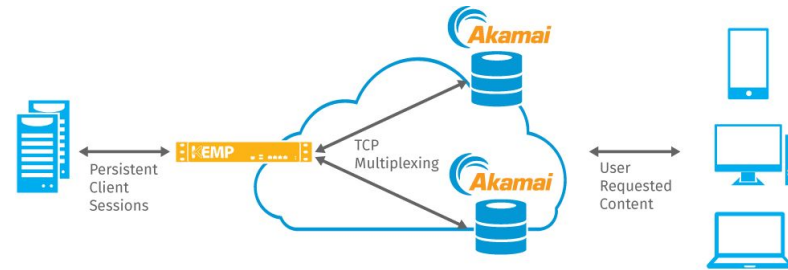
Content Delivery Networks

The Problem of Flash Crowds

- Sudden bursts of traffic (Flash Crowds) can bring down services.
- CDN Solution: Distribute and replicate documents across the Internet to offload the origin server and place content closer to clients.

Case Study: Akamai

- Uses Edge Servers to host content.
- Virtual Ghosts: Uses modified URLs that refer to "virtual ghosts" (CDN servers). The system resolves the DNS name to a CDN server close to the client.
- Consistency: If a document changes, its URL changes (embedding a unique ID), forcing the CDN server to fetch the fresh copy from the origin



Questions?

—

Saptarashmi Bandyopadhyay