

Chapter 11

Mass information storage

Introduction to Internet of Things





Introduction to Internet of Things



How to meet the requirement of mass information storage under the background of the association of things? With the development of Internet of things, **Data Center** will become the main means to solve mass data storage.

This chapter introduces the typical network storage architecture and the basic concepts of the data center.

内容提要



Review

Chapter 10 introduces the basic concepts of database management system

- Development of database models
- Basic concepts of relational databases
- Write query expressions using relational algebra
- Characteristics of Internet of things data management

This chapter mainly introduces three basic network storage architectures, and takes Google data center as an example to introduce related technologies of large-scale data center. Finally, it briefly introduces Hadoop, an open-source distributed computing framework.



Content

11.1 The Internet of things needs massive information storage

11.2 Network storage architecture

11.3 Data center

**How has information storage evolved?
What technologies have been driven by
the Internet of things' demand for
mass information storage?**





Introduction to Internet of Things

✓ History of data storage

Oracle bone inscriptions → Paper books → Digital storage





✓ The Internet of things needs massive information storage

The amount of information in the world is growing rapidly

- The amount of data generated in 2007 was 281EB (1EB=1 billion GB)
- The number of objects in the Internet of things will grow to tens of billions

The need for objects in the Internet of things to actively participate in business processes

- High intensity computing requirements
- Continuous online accessibility of data

This led to networked storage and large data centers



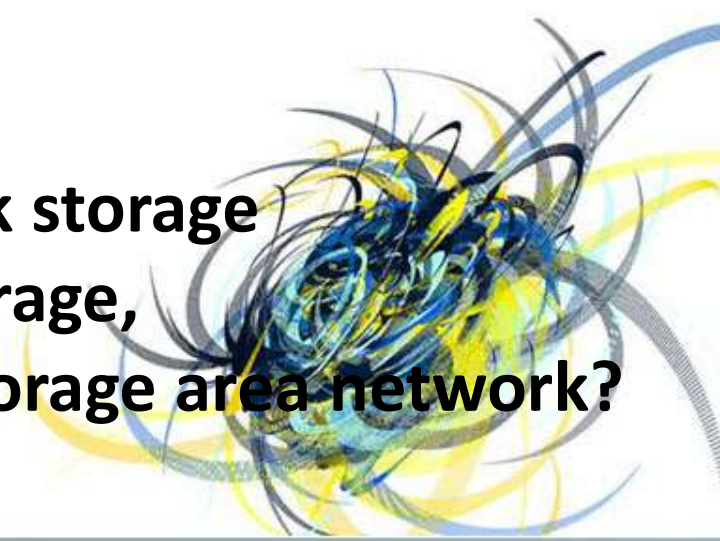
Content

11.1 The Internet of things needs massive information storage

11.2 Network storage architecture

11.3 Data center

**There are the three basic network storage architectures: direct attached storage, network attached storage, and storage area network?
What are their characteristics?**





✓ Direct-Attached Storage

Direct-attached Storage (DAS)

- Connect the storage system directly to the server or workstation via a cable
- Typically includes multiple hard disk drives, with host bus adapters via cable or fiber optic
- There are no other network devices between the storage device and the host bus adapter
- Realized the storage in the computer to the storage subsystem of the span



✓ Network Attached Storage

Network Attached Storage (NAS)

- File level computer data storage architecture
- The computer is connected to a network that provides file-based data storage services only for other devices

The difference between NAS and DAS

DAS is a simple extension of an existing server without really interconnecting the network. NAS uses the network as a storage entity, making file-level sharing easier. NAS performance is enhanced over DAS



✓ Storage Area Network

Storage Area Network (SAN)

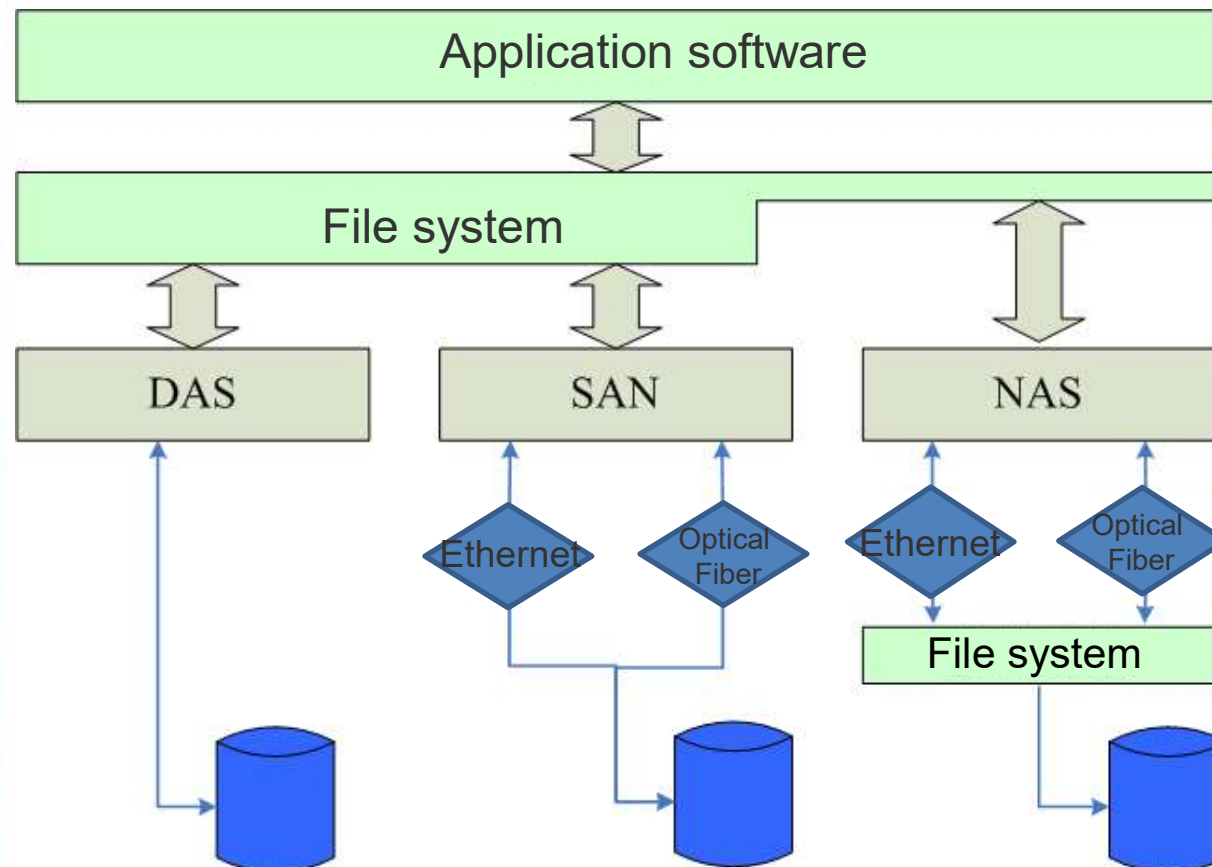
- Storage architecture that connects storage devices and application servers over a network
- It consists of servers, storage devices, and san-connected devices

The characteristics of SAN

- Shared storage
- Support server boot directly from SAN



Q Comparison of three network storage structures





Q Comparison of three network storage structures

DAS

Easy management, simple structure; Centralized architecture cannot meet the demand of large-scale data access; Storage resource utilization rate is low, resource sharing ability is poor, resulting in "information island".

NAS

Network storage entity, easy to achieve file level sharing; The performance is heavily dependent on network traffic. There are too many users, and the performance is limited when reading and writing frequently.

SAN

Storage management is simplified and storage capacity utilization is improved. No direct file level access, but file systems can be built on a SAN basis.



Content

11.1 The Internet of things needs massive information storage

11.2 Network storage architecture

11.3 Data center

What is the data center?

What are the typical data centers?

What is the research hot spot of the data center?





Q What is the data center?

Wikipedia: "a data center is a complex set of facilities. It includes not only computer systems and other associated equipment (such as communication and storage systems), but also redundant data communication connections, environmental control equipment, monitoring equipment and various security devices."

Google: "a multi-purpose building that can accommodate multiple servers and communications devices. These devices are placed together because they share the same environmental and physical security requirements and are easy to maintain."



✓ The origin and development of data center



Mainframe



Micromachine



Mega Data Center



Introduction to Internet of Things

✓ The origin and development of data center

Large-scale data centers have been rolled out

Google
谷歌

Baidu 百度

Microsoft®

YAHOO!®

Alibaba.com
阿里巴巴

amazon.com





✓ Data center standards

Data center builders face challenges

- ✓ How to plan a new data center?
- ✓ How to upgrade the data center?

Data center standards summarize the experience

ANSI/TIA/EIA-942 (TIA-942) : Data center standard

Telecommunications industry association (TIA) present
American national standards institute (ANSI) approve



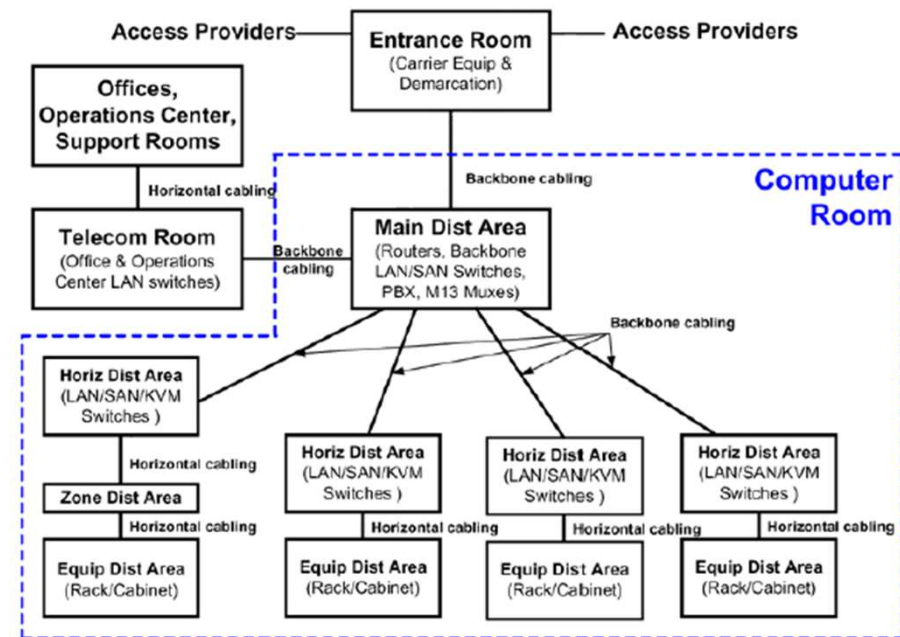
✓ Data center standard: TIA-942

Location: There are many factors

- ✓ Construction and operating costs
- ✓ The application requirements
- ✓ Preferential policy
- ✓ ...

Layout:

- ✓ Divided by functional area

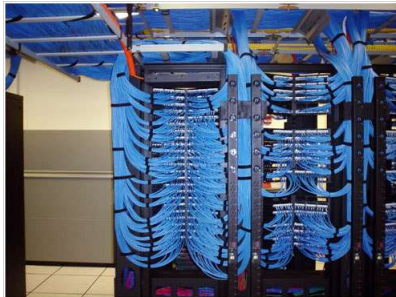


Functional Area composition



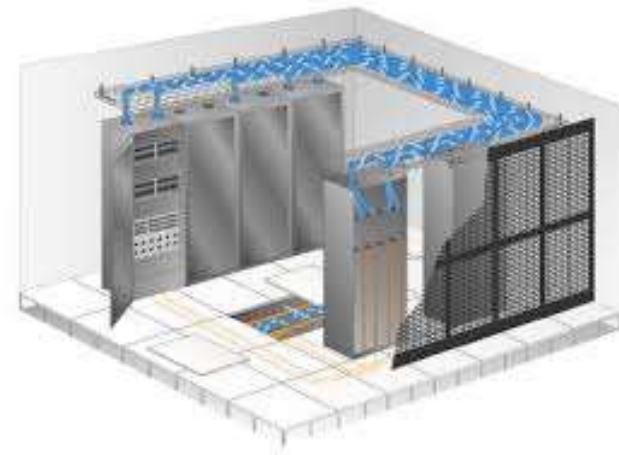
✓ Data center standard: TIA-942

TIA-942 also provides for cable systems, reliability ratings, energy systems and cooling systems.



Cable system

- ✓ Specifications
- ✓ How to place the cable



Energy systems

- ✓ External power supply
- ✓ The battery pack
- ✓ The generator

The cooling system

- ✓ Cooling equipment
- ✓ Overhead floor
- ✓ Cold channel and hot channel



Introduction to Internet of Things

✓ Typical data center: Google data center

Introduction

- There are nearly 40 large-scale data centers worldwide
- A single data center requires at least 50 megawatts of power, about the same as all the homes in a small city
- Unique hardware: custom Ethernet switches, energy systems, etc
- Self-developed software technologies: Google File System, MapReduce, BigTable, etc



✓ Google File System

Design concepts for GFS

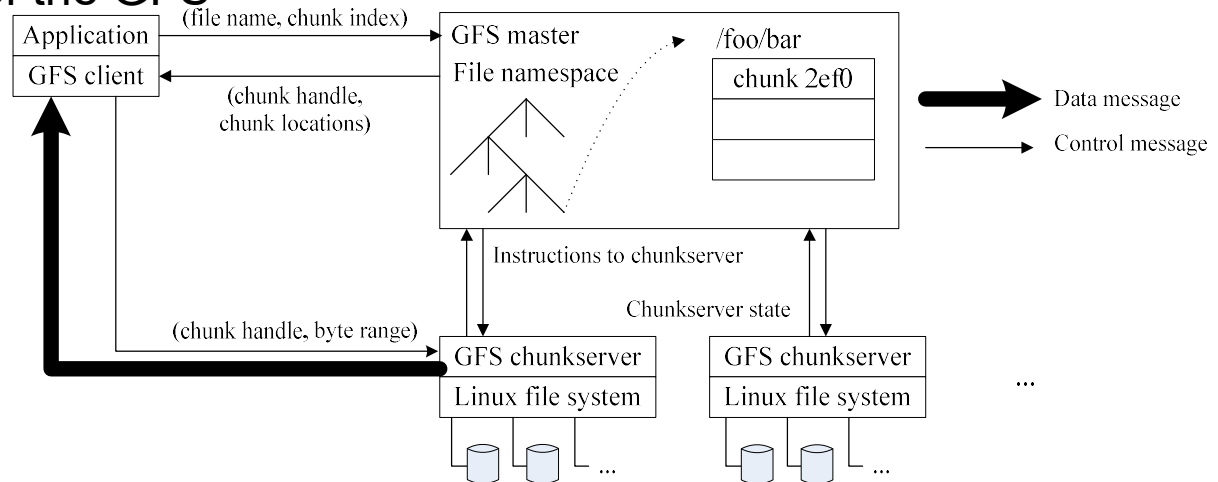
- Component failures are no longer considered accidental, but normal
- The GFS files are huge
- There are specific patterns for manipulating files
- Collaborative design of applications and file system apis increases the flexibility of the entire system



✓ Google File System

The design architecture of the GFS

A GFS cluster consists of one primary server and multiple block servers and is accessed by multiple clients.



Files are divided into fixed-size "chunks". Each block is assigned a fixed 64-bit handle unique identity by the primary server at creation time. The block server stores blocks as Linux files on the local disk and reads and writes to them based on the specified block handle and byte range.



✓ Google File System

The design architecture of the GFS

The primary server maintains metadata for all file systems, including namespaces, access control information, file to block mapping information, and the current location of blocks. In addition, the primary server controls other system-level activities. The primary server periodically communicates with the block server, following instructions and collection status.

GFS client code is embedded in each application. It implements the file system API to communicate between the primary server and the block server to implement read and write operations on behalf of the application. The client interacts with the server to implement metadata operations, but all data operations are done by interacting directly with the block server.



✓ MapReduce

MapReduce is a programming model and system for very large data sets

Programs developed with MapReduce can execute in parallel on large clusters of commercial computers, handle failures and schedule communication between computers

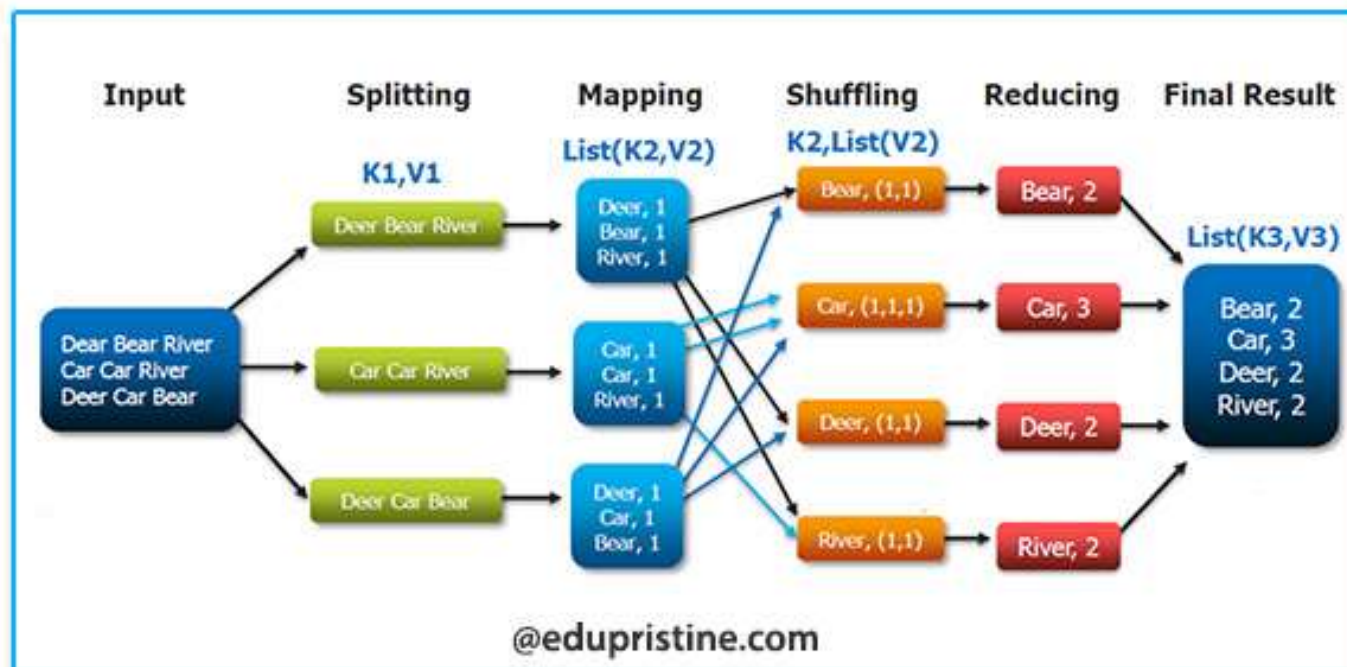
The basic idea of MapReduce

- Two programs written by users: Map and Reduce
- A framework for executing multiple program instances on a cluster of computers



✔ MapReduce

MapReduce program Execution process





✓ BigTable

BigTable is a distributed storage system designed to manage structured data at large data scales, such as petabytes of data and applications with thousands of inexpensive computers.

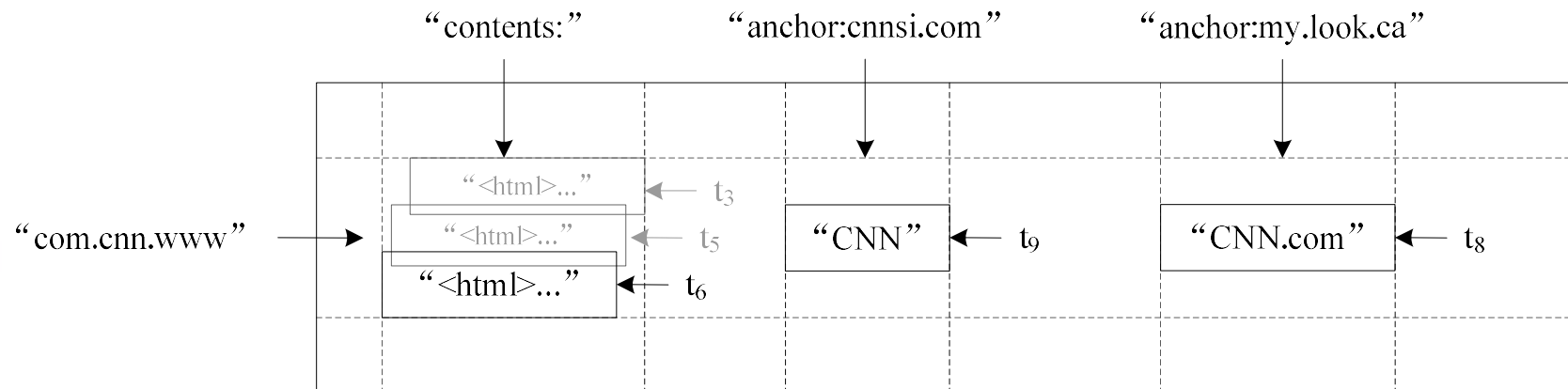
Application

- Google earth
- Web page index
- RSS reader
- ...



✓ BigTable

Each BigTable is a sparse, distributed, multidimensional ordered graph indexed by row key values, column key values, and timestamps





✓ Typical data center: Hadoop

What is Hadoop?

- Apache open source is a distributed computing framework
- Used to run data-intensive distributed applications on cheap server devices in large clusters
- In the early days, IT was actually an open source implementation of the Google file system and MapReduce distributed computing framework and related IT infrastructure services



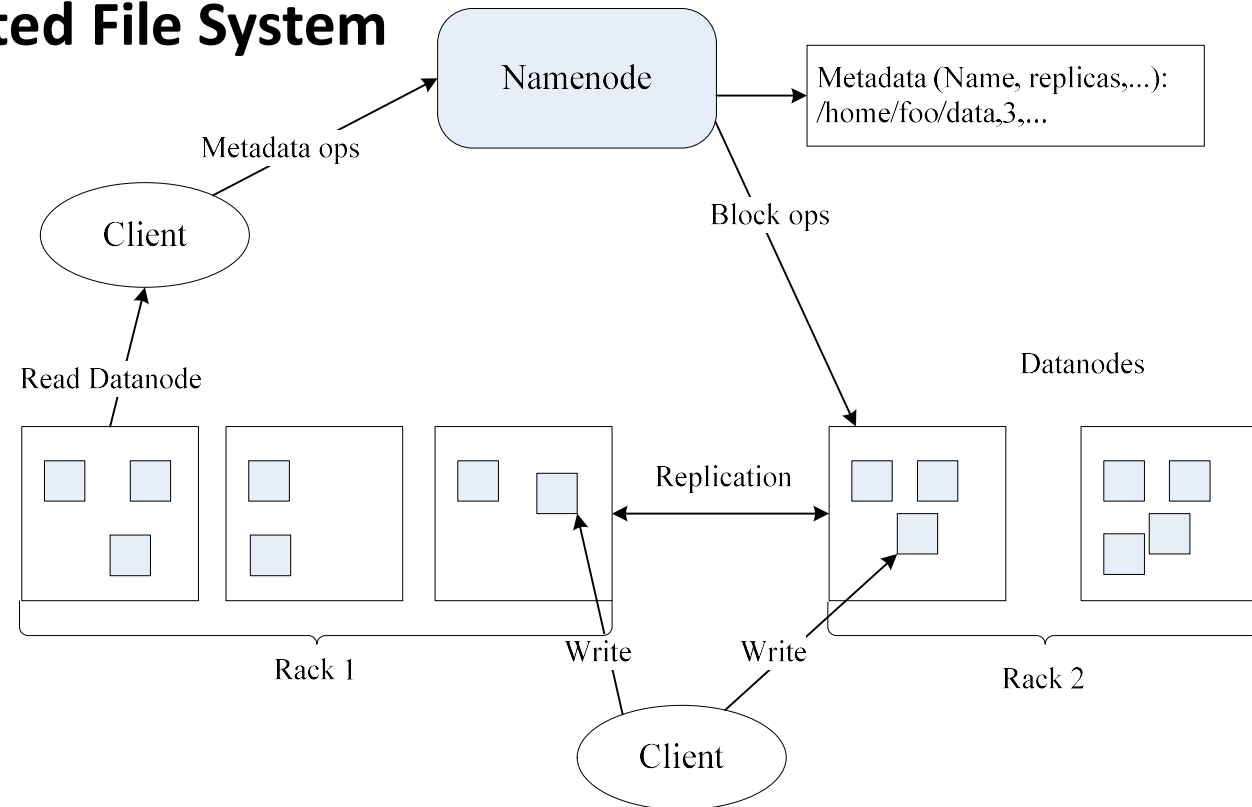
Hadoop consists of multiple subprojects

HDFS, MapReduce, HBase, Chukwa, Pig, ZooKeeper, etc



✓ HDFS

Hadoop Distributed File System





Introduction to Internet of Things

✓ Research hotspot of data center

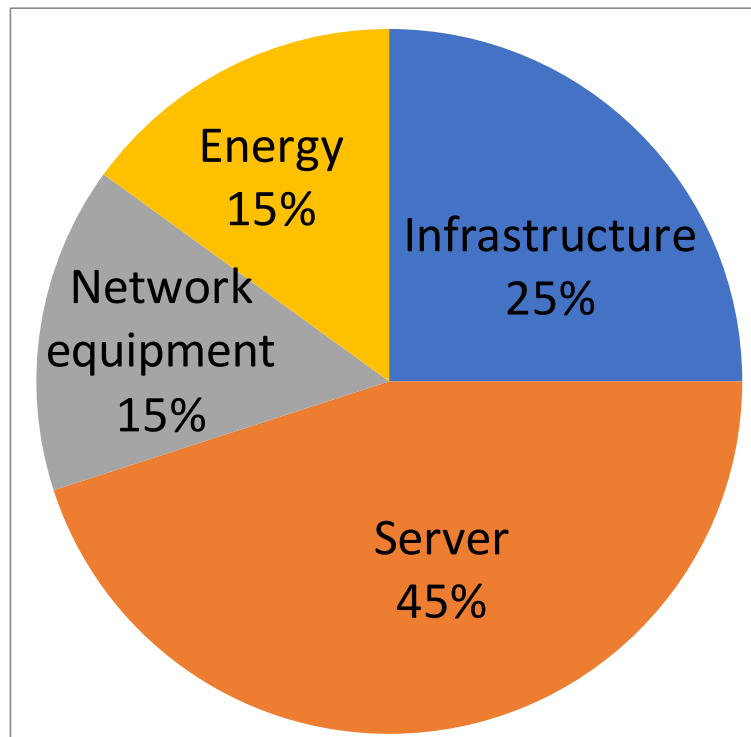
Google spent \$1.9 billion on data centers in 2006 and \$2.4 billion in 2007.

Google's data center in Oregon has nearly 100 megawatts of power, which when run at full capacity consumes roughly as much power as all the homes in a single city in Newcastle combined.

Research focus: How to reduce costs while ensuring service quality?



✓ Cost composition of the data center



The infrastructure includes energy system, cooling system, various fire prevention equipment, security equipment and so on. Reducing this part of cost often involves factors such as mechanical equipment manufacturing technology or policy preference, and has a relatively low degree of association with computer science.

We briefly introduce the reasons for the high cost and the current solutions from the aspects of server, network equipment and energy.

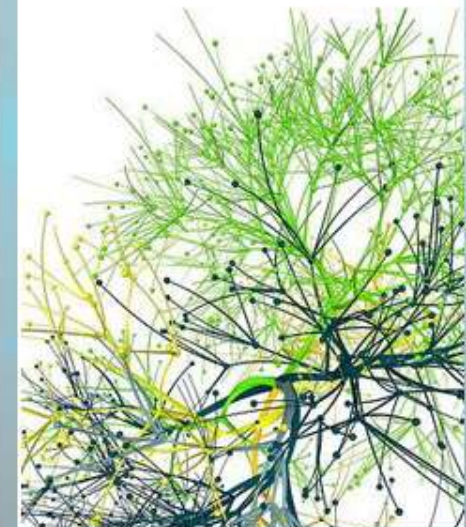


✓ Server cost

The actual utilization efficiency of the server is low

- Applications assigned to servers cannot fully utilize certain components
- It is difficult to predict the application demand, so it cannot be distributed according to demand
- In order to improve the reliability of the system, redundant equipment is usually left

The **key** to improving server utilization is to respond to dynamic changes in demand





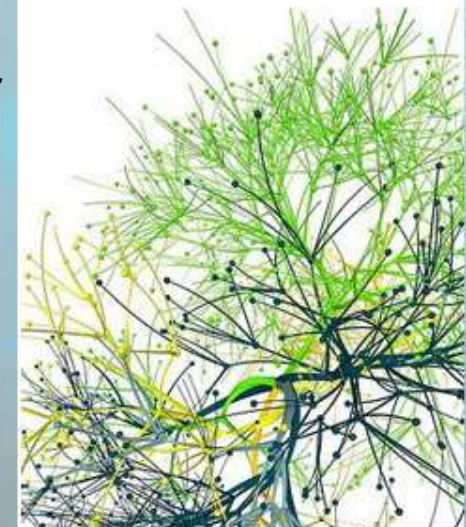
✓ Network equipment cost

The main source

- Switch, router, load balancing equipment
- Traditional data centers use tree structures, and core switches and routers are costly bottlenecks

Research hotspot: New data center network structure

- ✓ Multi-layer Tree structure centered on switch: for example, Fat-tree
- ✓ Server-centric interconnection structures: for example, Dcell

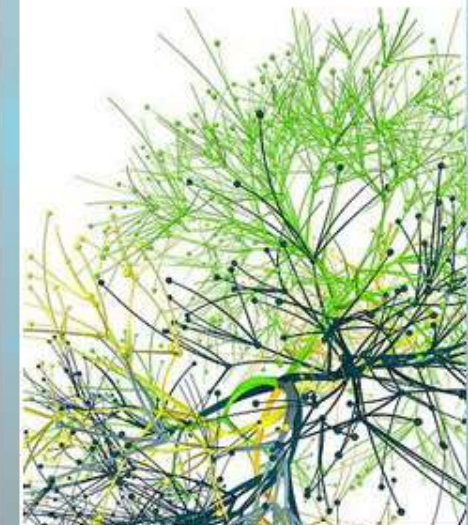
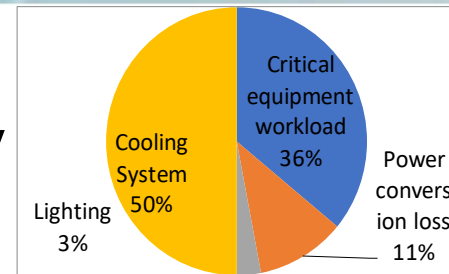




✓ Energy cost

Research hot spot

- Reduce server operating energy consumption
- ✓ Reduce energy consumption of equipment with the same performance
- ✓ Improve the performance of equipment with the same energy consumption
- ✓ Server that can adjust load
- Reduce energy consumption of cooling system
- ✓ Fine and accurate temperature control
- ✓ Containerized modular data center





Conclusion

Review

This chapter introduces three basic network storage architectures, and discusses the basic concepts of data center. Taking Google data center and Hadoop as examples, it briefly introduces the related technologies of data center, and finally points out the research hotspots of data center.

Key Points

- Understand the urgent need for mass data storage in the Internet of things.
- Focus on the basic concepts and advantages and disadvantages of three basic network storage architectures (DAS, NAS, SAN).



Conclusion

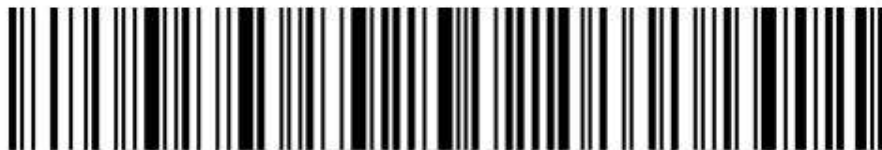
Key Points

- Understand the concept of data center. Take Google data center as an example to understand the basic concepts and characteristics of GFS, MapReduce, BigTable and other technologies. Understand the characteristics of Hadoop distributed computing open source framework.
- Know how to reduce data center cost (server cost, network equipment cost, energy cost) without sacrificing performance.

GreenOrbs
Pervasive Computing
to IoT
of
Internet
Introduction
OceanSense
Things
Smart Planet
CDMA
SQL
Smart Grid
CPS
RFID
Database
TinyOS
ITS
ZigBee
Web
ITU
nesC
ETC
BlueTooth



Thank you!



Internet of Things