# Chapter 1:

## Introduction to the Design and Specification of File Structure

# Why is File Structure Design Necessary?

- **Disks are very slow compared with RAM**

- **How slow is a disk?**

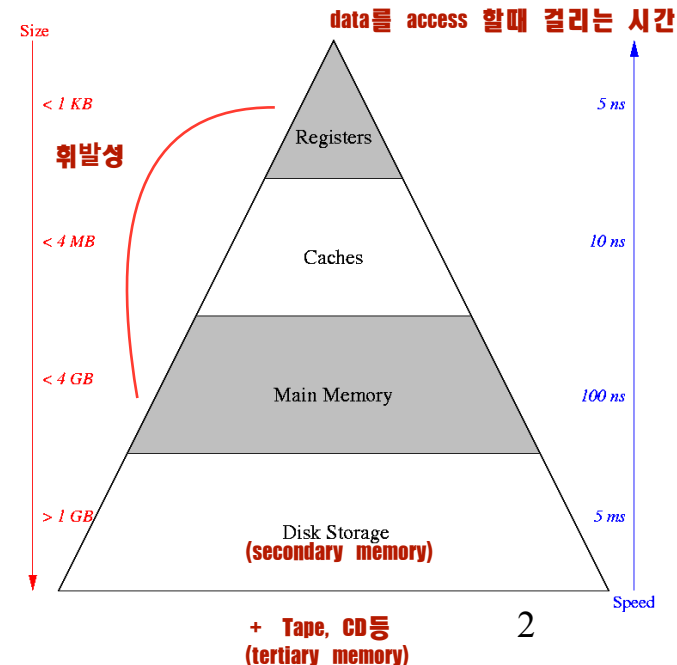목적: Cache는 레지스터와 메인메모리 사이의 속도차를 줄여주기 위해 사용
파일처리는 메인과 세컨더리사이의 속도차를 줄여줌

**RAM : Disk = 120 nanosec : 30 milisec**

$$= 1 \text{ sec} : \frac{30 \times 10^{-3}}{120 \times 10^{-9}} \text{ sec} = 1 \text{sec} : 25 \times 10^{4} \text{ sec}$$

**= 1 sec : 2 days and 22 hours**

- **However, disks provide enormous and**

  **nonvolatile capacity at much less**

cost than RAM

Performance를 측정할때
여러기준에서 볼 수 있어야함
1. Main memory사용량
2. Elapsed time & User response time
3. Network
4. Throughput(단위시간당 job처리량)
5. etc..
But 파일처리에서 가장 중요한 부분은
"2번"이다.

data를 access 할때 걸리는 시간

Size

< 1 KB    Registers    5 ns

휘발성

< 4 MB    Caches    10 ns

< 4 GB    Main Memory    100 ns

> 1 GB    Disk Storage
(secondary memory)    5 ms

+ Tape, CD등
(tertiary memory)

Speed

Dong-Joo Park

2

# What is a Good File Structure Design?

- **What is a file?** 세컨더리 메모리에 저장되는 같은종류의
  레코드의 집합
  - **<u>A set of same kind of records</u> which are stored into secondary memory like hard disk, solid state disk, CD, tape, etc.**

- **What is a file structure?** C에서 구조체,
  cpp에서 클래스처럼
  "표현"하는 것                      insert,delete,seratch,read,write등
  - **A combination of <u>representations</u> for data in files and of <u>operations</u> for accessing the data**

- **Goal of good file structure design**

  - **Allowing us to get the information with as small cost as possible**

  - **Main performance factor: <u>number of disk accesses</u>**

    디스크 엑세스의 비용이 크니까
    최대한 접근하지 않으면서 정확한 구조가
    좋은 파일구조임

# A Short History of File Structure Design (1)

- Sequential access 순차적으로 읽는 것(=:linear search)
  - Accessing records in order, looking at the first, then the next, and so on
  - Most files were on tape in early work with files

Hard disk의 경우 앞과 다르게 하나의 data를 찾기위해
전체를 보지 않아도 됨
〉일부를 읽으면서도 그사이 공간을 읽을 필요가 없음
〉〉random access

```
struct {

        char   name[20];
        char   telephone[20];
        char   address[60];
} phone_book;
```

Gildong Hong | +82-2-820-0914 | 1-1 Sangdo-dong, Dongjak-gu, Seoul 156-743, Korea

Paul B. Kantor | +1-732-932-1359 | New Brunswick, NJ 08901-1071, USA

Yasushi Ogawa | +81-45-477-1569 | 3-2-3, shin-yokohama, Kohoku-ku, Yokohama-shi 222, Japan

Hans-Peter Frei | +41-1-236-5714 | Ch-8021 Zurich, Bahnhofstrasse 45, Switzerland

⋮

# A Short History of File Structure Design (2)

- Simple index
  - Storage devices like disk drives became available
  - Keeping a list of keys and pointers in a smaller file
  - Difficult to manage, especially for dynamic files

index는 각 record의 key, 위치를 저장
+ key값에 대해 sorting되어있음

| **Index File** | | | **Data File** |
|---|---|---|---|
| Gildong Hong \| 1 | | 1 | Gildong Hong \| +82-2-820-0914 \| 1-1 Sangdo- …... |
| Hans-Peter Frei \| 301 | | 101 | Paul B. Kantor \| +1-732-932-1359 \| New Br  …... |
| Paul B. Kantor \| 101 | | 201 | Yasushi Ogawa \| +81-45-477-1569 \| 3-2-3, s  …... |
| Yasushi Ogawa \| 201 | | 301 | Hans-Peter Frei \| +41-1-236-5714 \| Ch-8021 …... |

Dong-Joo Park
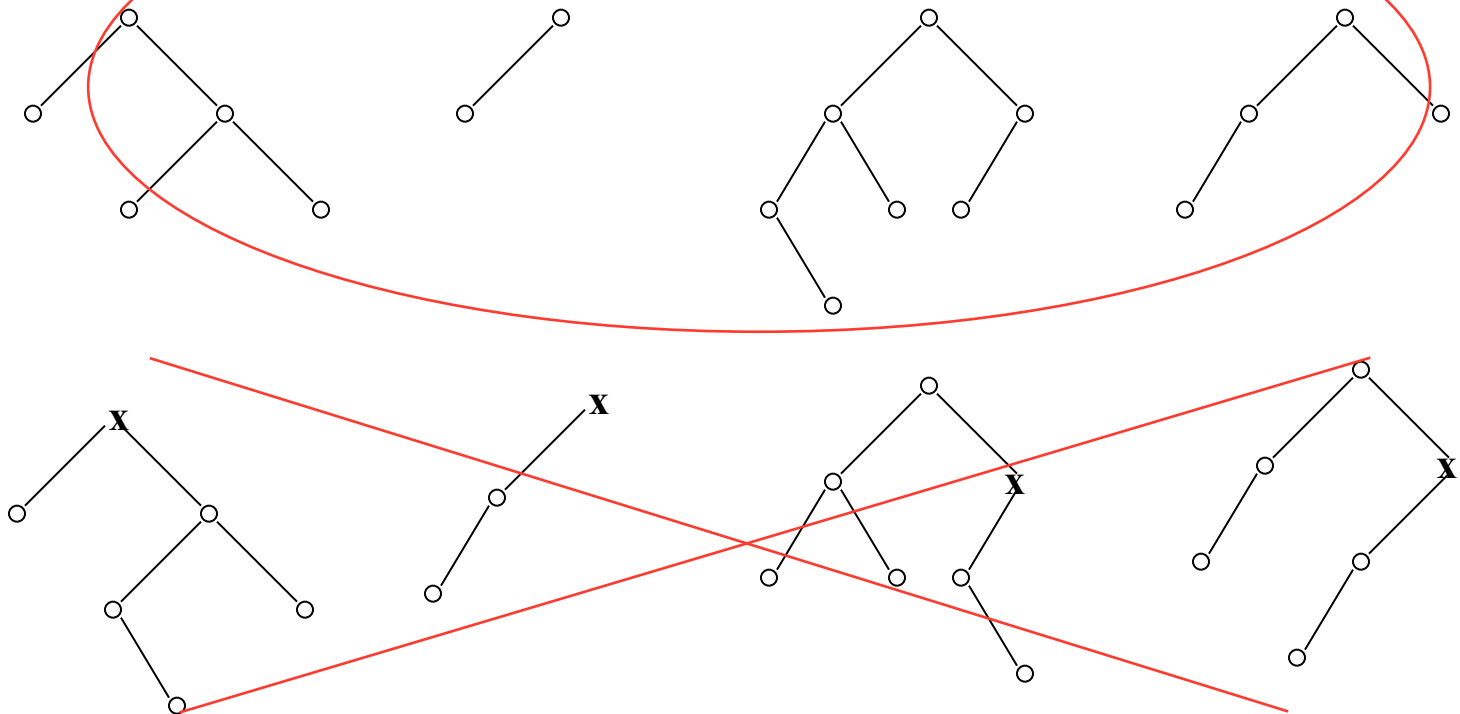
# A Short History of File Structure Design (3)

search

- Binary tree : early 1960s

  – binary trees can grow very unevenly as records are added and deleted

  – ex) insert (A,1), (B,101), (C,201), (K,301), (L,401), (T,501), (Z,601)

A,1

B,101

C,201

K,301

L,401

T,501

Z,601

# A Short History of File Structure Design (4)

- AVL tree  **왼쪽서브트리와 오른쪽서브트리의 깊이차가 1이하라는 제한조건을 둠**
  - there is a limit on the amount of difference that is allowed between the heights of any two subtrees sharing a common root
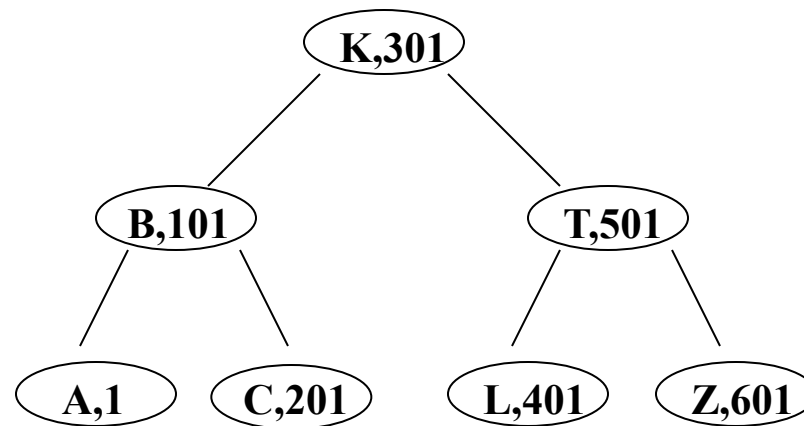
# A Short History of File Structure Design (5)

- Balanced binary tree
  - Given $N$ keys, looks at $\lfloor \log_2 N \rfloor + 1$ levels of the tree
  - ex) insert (A,1), (B,101), (C,201), (K,301), (L,401), (T,501), (Z,601)

```
                    K,301
                   /      \
              B,101        T,501
             /     \      /     \
          A,1    C,201  L,401   Z,601
```

Dong-Joo Park                                8

# A Short History of File Structure Design (6)

- B-tree
  - A tree structure that provides fast access to data stored in files
  - Unlike binary trees, in which the branching factor from a node of the tree is two, the descendents from a node of a B-tree can be a much larger number

- B+ tree
  - A variation on the B-tree structure that provides sequential access to the data as well as fast-indexed access

- Hashing
  - An access mechanism that transforms the search key into a storage address, there by providing very fast access to stored data

- Extendible hashing
  - Hashing does not work well with dynamic files
  - Dynamic hashing that could retrieve information with one or, at most, two disk accesses no matter how big the file becomes

# C++ Language: Using Objects(1)

```
class Person {
private:
    char name[20];
    int age;
    char address[50];
public:
    Person(char *pname, int page, char *paddress);
    ~Person();
    char* getName();
    int getAge();
    char* getAddress();
    void setName(char *personname);
    void setAge(int personage);
    void setAddress(char *personaddress);
};
```

# C++ Language: Using Objects(2)

```cpp
class Student: public Person {
private:
    int year;
    char major[30];
public:
    Student(int syear, char *smajor);
    ~Student();
    char* getMajor();
    int getYear();
    void setMajor(char *smajor);
    void setYear(int syear);
};
```