

# Unit V – Data Storage and Querying

**Storage and File Structure:** Overview of physical storage media, magnetic disk and flash storage, RAID, tertiary storage, File Organization, Organization of Records in Files, Data-Dictionary Storage, Database Buffer.

**Indexing and hashing:** Ordered Indices, B+-Tree Index Files, B+-Tree Extensions, Multiple-Key Access, Static hashing, Dynamic Hashing, Comparison of Ordered Indexing and Hashing, Bitmap Indices, Index Definition in SQL.

**Query processing and Optimization:** Measures of Query Cost, selection Operations, sorting, Join operations, Other Operations, Evaluation of Expressions. Transformation of Relational Expressions, Estimating Statistics of Expression Results, choice of Evaluation plans and Materialized views

# Chapter 10: Storage and File Structure

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan  
See [www.db-book.com](http://www.db-book.com) for conditions on re-use

# Chapter 10: Storage and File Structure

- Overview of Physical Storage Media
- Magnetic Disks
- RAID
- Tertiary Storage
- Storage Access
- File Organization
- Organization of Records in Files
- Data-Dictionary Storage

# Classification of Physical Storage Media

- Speed with which data can be accessed
- Cost per unit of data
- Reliability
  - data loss on power failure or system crash
  - physical failure of the storage device
- Can differentiate storage into:
  - **volatile storage:** loses contents when power is switched off
  - **non-volatile storage:**
    - ▶ Contents persist even when power is switched off.
    - ▶ Includes secondary and tertiary storage, as well as battery-backed up main-memory.

# Physical Storage Media

- **Cache** – fastest and most costly form of storage; volatile; managed by the computer system hardware.
- **Main memory:**
  - fast access (10s to 100s of nanoseconds; 1 nanosecond =  $10^{-9}$  seconds)
  - generally too small (or too expensive) to store the entire database
    - ▶ capacities of up to a few Gigabytes widely used currently
    - ▶ Capacities have gone up and per-byte costs have decreased steadily and rapidly (roughly factor of 2 every 2 to 3 years)
  - **Volatile** — contents of main memory are usually lost if a power failure or system crash occurs.

# Physical Storage Media (Cont.)

## ■ Flash memory

- Data survives power failure
- Data can be written at a location only once, but location can be erased and written to again
  - ▶ Can support only a limited number (10K – 1M) of write/erase cycles.
  - ▶ Erasing of memory has to be done to an entire bank of memory
- Reads are roughly as fast as main memory
- But writes are slow (few microseconds), erase is slower
- Widely used in embedded devices such as digital cameras, phones, and USB keys

# Physical Storage Media (Cont.)

## ■ Magnetic-disk

- Data is stored on spinning disk, and read/written magnetically
- Primary medium for the long-term storage of data; typically stores entire database.
- Data must be moved from disk to main memory for access, and written back for storage
  - ▶ Much slower access than main memory (more on this later)
- **direct-access** – possible to read data on disk in any order, unlike magnetic tape
- Capacities range up to roughly 1.5 TB as of 2009
  - ▶ Much larger capacity and cost/byte than main memory/flash memory
  - ▶ Growing constantly and rapidly with technology improvements (factor of 2 to 3 every 2 years)
- Survives power failures and system crashes
  - ▶ disk failure can destroy data, but is rare

# Physical Storage Media (Cont.)

## ■ Optical storage

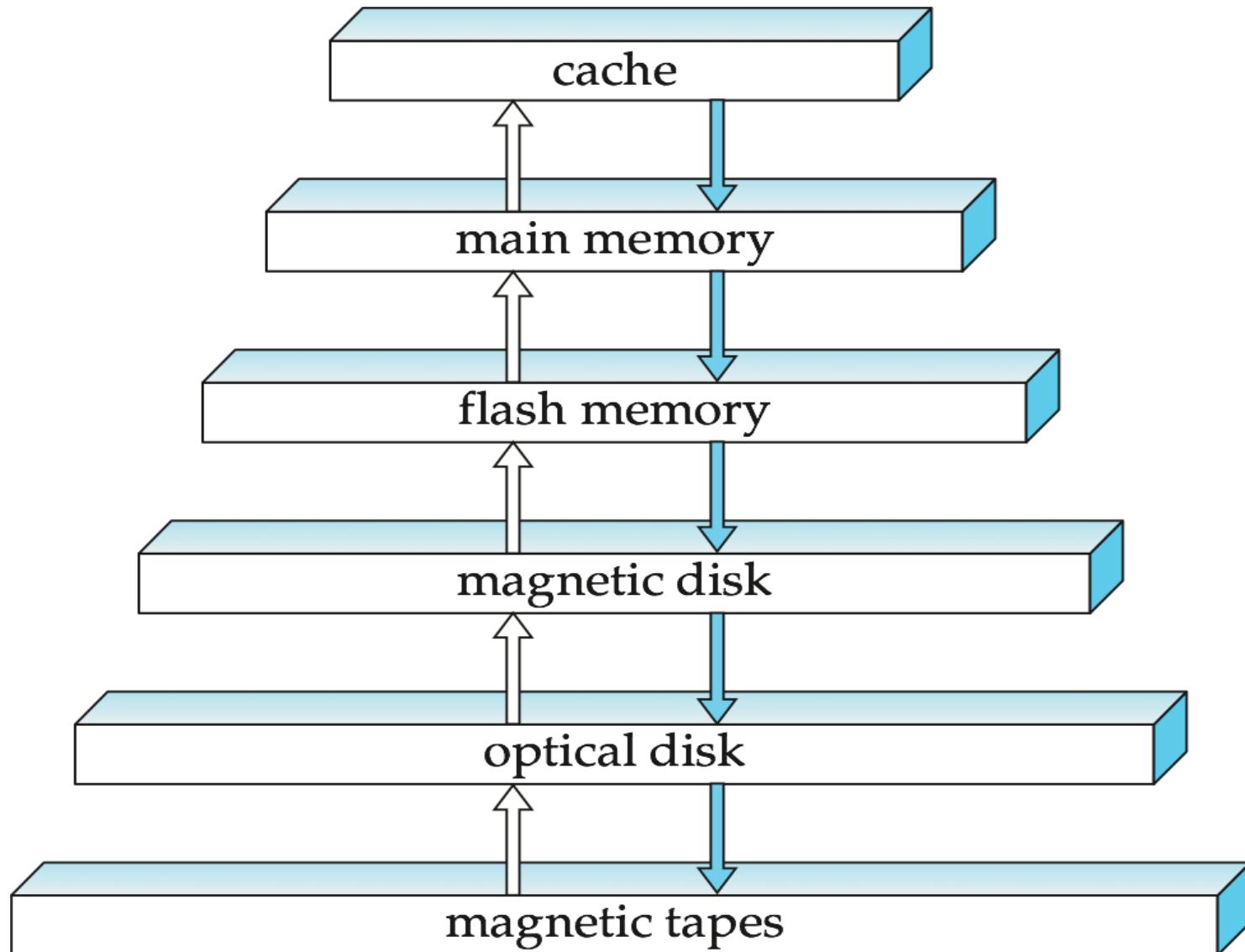
- non-volatile, data is read optically from a spinning disk using a laser
- CD-ROM (640 MB) and DVD (4.7 to 17 GB) most popular forms
- Blu-ray disks: 27 GB to 54 GB
- Write-one, read-many (WORM) optical disks used for archival storage (CD-R, DVD-R, DVD+R)
- Multiple write versions also available (CD-RW, DVD-RW, DVD+RW, and DVD-RAM)
- Reads and writes are slower than with magnetic disk
- **Juke-box** systems, with large numbers of removable disks, a few drives, and a mechanism for automatic loading/unloading of disks available for storing large volumes of data

# Physical Storage Media (Cont.)

## ■ Tape storage

- non-volatile, used primarily for backup (to recover from disk failure), and for archival data
- **sequential-access** – much slower than disk
- very high capacity (40 to 300 GB tapes available)
- tape can be removed from drive  $\Rightarrow$  storage costs much cheaper than disk, but drives are expensive
- Tape jukeboxes available for storing massive amounts of data
  - ▶ hundreds of terabytes (1 terabyte =  $10^9$  bytes) to even multiple **petabytes** (1 petabyte =  $10^{12}$  bytes)

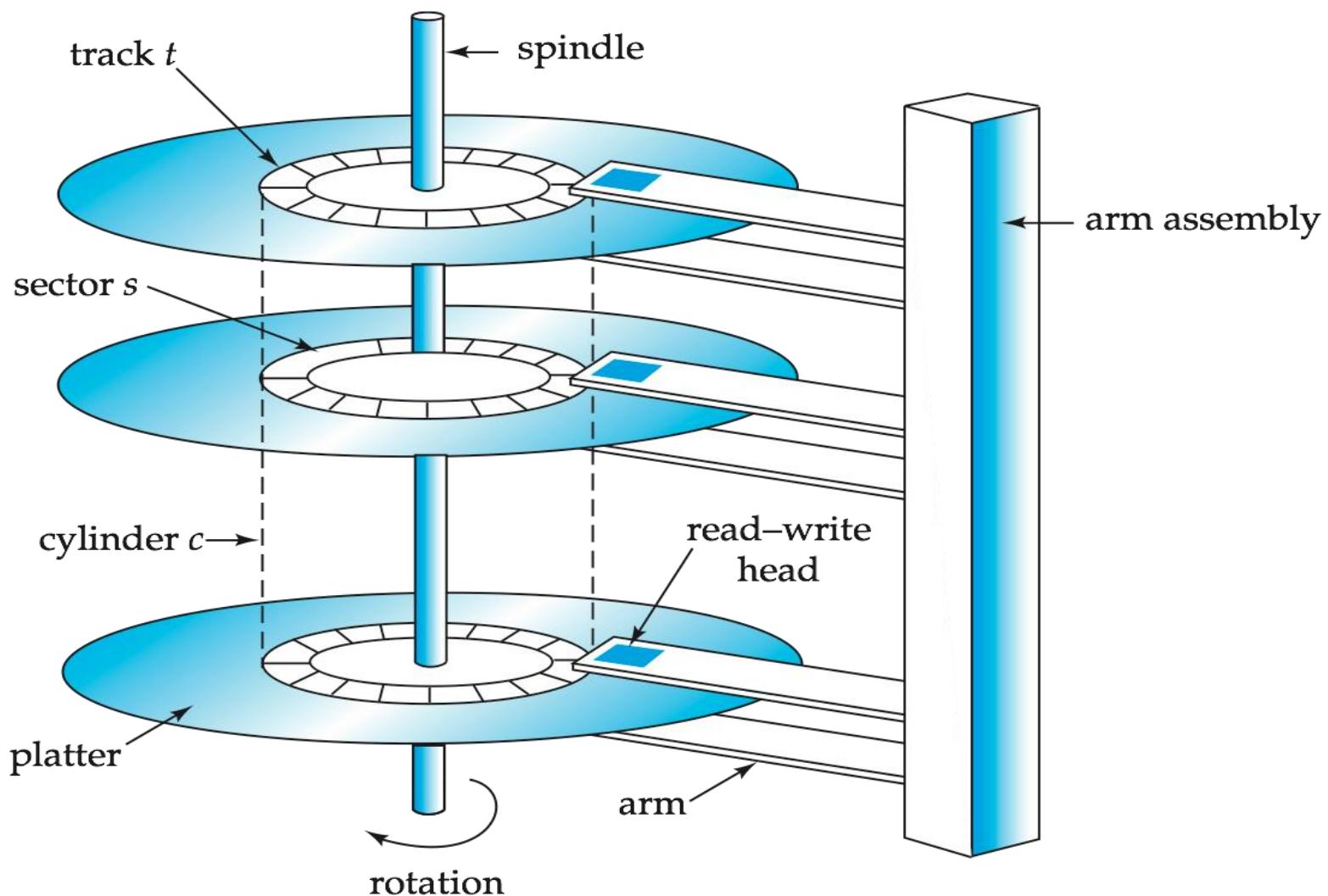
# Storage Hierarchy



# Storage Hierarchy (Cont.)

- **primary storage:** Fastest media but volatile (cache, main memory).
- **secondary storage:** next level in hierarchy, non-volatile, moderately fast access time
  - also called **on-line storage**
  - E.g. flash memory, magnetic disks
- **tertiary storage:** lowest level in hierarchy, non-volatile, slow access time
  - also called **off-line storage**
  - E.g. magnetic tape, optical storage

# Magnetic Hard Disk Mechanism



**NOTE: Diagram is schematic, and simplifies the structure of actual disk drives**

# Magnetic Disks

## ■ Read-write head

- Positioned very close to the platter surface (almost touching it)
- Reads or writes magnetically encoded information.

## ■ Surface of platter divided into circular **tracks**

- Over 50K-100K tracks per platter on typical hard disks

## ■ Each track is divided into **sectors**.

- A sector is the smallest unit of data that can be read or written.
- Sector size typically 512 bytes
- Typical sectors per track: 500 to 1000 (on inner tracks) to 1000 to 2000 (on outer tracks)

## ■ To read/write a sector

- disk arm swings to position head on right track
- platter spins continually; data is read/written as sector passes under head

## ■ Head-disk assemblies

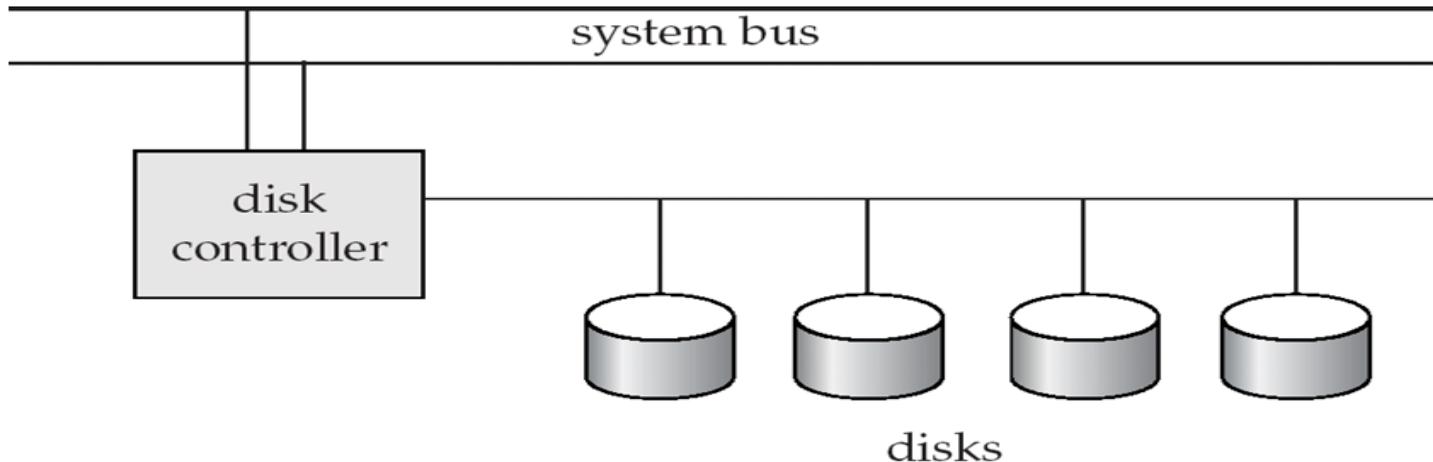
- multiple disk platters on a single spindle (1 to 5 usually)
- one head per platter, mounted on a common arm.

## ■ **Cylinder** $i$ consists of $i^{\text{th}}$ track of all the platters

# Magnetic Disks (Cont.)

- Earlier generation disks were susceptible to head-crashes
  - Surface of earlier generation disks had metal-oxide coatings which would disintegrate on head crash and damage all data on disk
  - Current generation disks are less susceptible to such disastrous failures, although individual sectors may get corrupted
- **Disk controller** – interfaces between the computer system and the disk drive hardware.
  - accepts high-level commands to read or write a sector
  - initiates actions such as moving the disk arm to the right track and actually reading or writing the data
  - Computes and attaches **checksums** to each sector to verify that data is read back correctly
    - ▶ If data is corrupted, with very high probability stored checksum won't match recomputed checksum
  - Ensures successful writing by reading back sector after writing it
  - Performs **remapping** of bad sectors

# Disk Subsystem



- Multiple disks connected to a computer system through a controller
  - Controllers functionality (checksum, bad sector remapping) often carried out by individual disks; reduces load on controller
- Disk interface standards families
  - **ATA** (AT adaptor) range of standards
  - **SATA** (Serial ATA)
  - **SCSI** (Small Computer System Interconnect) range of standards
  - **SAS** (Serial Attached SCSI)
  - Several variants of each standard (different speeds and capabilities)

# Disk Subsystem

- Disks usually connected directly to computer system
- In **Storage Area Networks (SAN)**, a large number of disks are connected by a high-speed network to a number of servers
- In **Network Attached Storage (NAS)** networked storage provides a file system interface using networked file system protocol, instead of providing a disk system interface

# Performance Measures of Disks

■ **Access time** – the time it takes from when a read or write request is issued to when data transfer begins. Consists of:

- **Seek time** – time it takes to reposition the arm over the correct track.
  - ▶ Average seek time is 1/2 the worst case seek time.
    - Would be 1/3 if all tracks had the same number of sectors, and we ignore the time to start and stop arm movement
  - ▶ 4 to 10 milliseconds on typical disks
- **Rotational latency** – time it takes for the sector to be accessed to appear under the head.
  - ▶ Average latency is 1/2 of the worst case latency.
  - ▶ 4 to 11 milliseconds on typical disks (5400 to 15000 r.p.m.)

■ **Data-transfer rate** – the rate at which data can be retrieved from or stored to the disk.

- 25 to 100 MB per second max rate, lower for inner tracks
- Multiple disks may share a controller, so rate that controller can handle is also important
  - ▶ E.g. SATA: 150 MB/sec, SATA-II 3Gb (300 MB/sec)
  - ▶ Ultra 320 SCSI: 320 MB/s, SAS (3 to 6 Gb/sec)
  - ▶ Fiber Channel (FC2Gb or 4Gb): 256 to 512 MB/s

# Performance Measures (Cont.)

■ **Mean time to failure (MTTF)** – the average time the disk is expected to run continuously without any failure.

- Typically 3 to 5 years
- Probability of failure of new disks is quite low, corresponding to a “theoretical MTTF” of 500,000 to 1,200,000 hours for a new disk
  - ▶ E.g., an MTTF of 1,200,000 hours for a new disk means that given 1000 relatively new disks, on an average one will fail every 1200 hours
- MTTF decreases as disk ages

# Optimization of Disk-Block Access

- **Block** – a contiguous sequence of sectors from a single track
  - data is transferred between disk and main memory in blocks
  - sizes range from 512 bytes to several kilobytes
    - ▶ Smaller blocks: more transfers from disk
    - ▶ Larger blocks: more space wasted due to partially filled blocks
    - ▶ Typical block sizes today range from 4 to 16 kilobytes

- **Disk-arm-scheduling** algorithms order pending accesses to tracks so that disk arm movement is minimized
  - **elevator algorithm:**



# Optimization of Disk Block Access (Cont.)

■ **File organization** – optimize block access time by organizing the blocks to correspond to how data will be accessed

- E.g. Store related information on the same or nearby cylinders.
- Files may get **fragmented** over time
  - ▶ E.g. if data is inserted to/deleted from the file
  - ▶ Or free blocks on disk are scattered, and newly created file has its blocks scattered over the disk
  - ▶ Sequential access to a fragmented file results in increased disk arm movement
- Some systems have utilities to **defragment** the file system, in order to speed up file access

# Optimization of Disk Block Access (Cont.)

- **Nonvolatile write buffers** speed up disk writes by writing blocks to a non-volatile RAM buffer immediately
  - Non-volatile RAM: battery backed up RAM or flash memory
    - ▶ Even if power fails, the data is safe and will be written to disk when power returns
  - Controller then writes to disk whenever the disk has no other requests or request has been pending for some time
  - Database operations that require data to be safely stored before continuing can continue without waiting for data to be written to disk
  - *Writes can be reordered to minimize disk arm movement*
- **Log disk** – a disk devoted to writing a sequential log of block updates
  - Used exactly like nonvolatile RAM
    - ▶ Write to log disk is very fast since no seeks are required
    - ▶ No need for special hardware (NV-RAM)
- File systems typically reorder writes to disk to improve performance
  - **Journaling file systems** write data in safe order to NV-RAM or log disk
  - Reordering without journaling: risk of corruption of file system data

# Flash Storage

## ■ NOR flash vs NAND flash

## ■ NAND flash

- used widely for storage, since it is much cheaper than NOR flash
- requires page-at-a-time read (page: 512 bytes to 4 KB)
- transfer rate around 20 MB/sec
- **solid state disks**: use multiple flash storage devices to provide higher transfer rate of 100 to 200 MB/sec
- erase is very slow (1 to 2 millisecs)
  - ▶ erase block contains multiple pages
  - ▶ **remapping** of logical page addresses to physical page addresses avoids waiting for erase
    - **translation table** tracks mapping
      - » also stored in a label field of flash page
    - remapping carried out by **flash translation layer**
  - ▶ after 100,000 to 1,000,000 erases, erase block becomes unreliable and cannot be used
    - **wear leveling**

# RAID

## ■ RAID: Redundant Arrays of Independent Disks

- disk organization techniques that manage a large numbers of disks, providing a view of a single disk of
  - ▶ **high capacity** and **high speed** by using multiple disks in parallel,
  - ▶ **high reliability** by storing data redundantly, so that data can be recovered even if a disk fails
- The chance that some disk out of a set of  $N$  disks will fail is much higher than the chance that a specific single disk will fail.
  - E.g., a system with 100 disks, each with MTTF of 100,000 hours (approx. 11 years), will have a system MTTF of 1000 hours (approx. 41 days)
  - Techniques for using redundancy to avoid data loss are critical with large numbers of disks
- Originally a cost-effective alternative to large, expensive disks
  - I in RAID originally stood for ``inexpensive''
  - Today RAIDs are used for their higher reliability and bandwidth.
    - ▶ The “I” is interpreted as independent

# Improvement of Reliability via Redundancy

- **Redundancy** – store extra information that can be used to rebuild information lost in a disk failure
- E.g., **Mirroring** (or **shadowing**)
  - Duplicate every disk. Logical disk consists of two physical disks.
  - Every write is carried out on both disks
    - ▶ Reads can take place from either disk
  - If one disk in a pair fails, data still available in the other
    - ▶ Data loss would occur only if a disk fails, and its mirror disk also fails before the system is repaired
      - Probability of combined event is very small
        - » Except for dependent failure modes such as fire or building collapse or electrical power surges
- **Mean time to data loss** depends on mean time to failure, and **mean time to repair**
  - E.g. MTTF of 100,000 hours, mean time to repair of 10 hours gives mean time to data loss of  $500 \times 10^6$  hours (or 57,000 years) for a mirrored pair of disks (ignoring dependent failure modes)

# Improvement in Performance via Parallelism

- Two main goals of parallelism in a disk system:
  1. Load balance multiple small accesses to increase throughput
  2. Parallelize large accesses to reduce response time.
- Improve transfer rate by striping data across multiple disks.
- **Bit-level striping** – split the bits of each byte across multiple disks
  - In an array of eight disks, write bit  $i$  of each byte to disk  $i$ .
  - Each access can read data at eight times the rate of a single disk.
  - But seek/access time worse than for a single disk
    - ▶ Bit level striping is not used much any more
- **Block-level striping** – with  $n$  disks, block  $i$  of a file goes to disk  $(i \bmod n) + 1$ 
  - Requests for different blocks can run in parallel if the blocks reside on different disks
  - A request for a long sequence of blocks can utilize all disks in parallel

# RAID Levels

- Schemes to provide redundancy at lower cost by using disk striping combined with parity bits
  - Different RAID organizations, or RAID levels, have differing cost, performance and reliability characteristics
- **RAID Level 0:** Block striping; non-redundant.
  - Used in high-performance applications where data loss is not critical.
- **RAID Level 1:** Mirrored disks with block striping
  - Offers best write performance.
  - Popular for applications such as storing log files in a database system.



(a) RAID 0: nonredundant striping



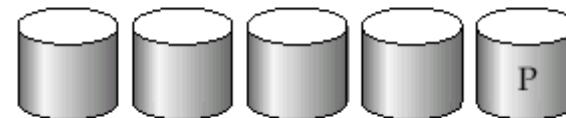
(b) RAID 1: mirrored disks

# RAID Levels (Cont.)

- **RAID Level 2:** Memory-Style Error-Correcting-Codes (ECC) with bit striping.
- **RAID Level 3:** Bit-Interleaved Parity
  - a single parity bit is enough for error correction, not just detection, since we know which disk has failed
    - ▶ When writing data, corresponding parity bits must also be computed and written to a parity bit disk
    - ▶ To recover data in a damaged disk, compute XOR of bits from other disks (including parity bit disk)



(c) RAID 2: memory-style error-correcting codes



(d) RAID 3: bit-interleaved parity

# RAID Levels (Cont.)

## ■ RAID Level 3 (Cont.)

- Faster data transfer than with a single disk, but fewer I/Os per second since every disk has to participate in every I/O.
- Subsumes Level 2 (provides all its benefits, at lower cost).

## ■ RAID Level 4: Block-Interleaved Parity; uses block-level striping, and keeps a parity block on a separate disk for corresponding blocks from $N$ other disks.

- When writing data block, corresponding block of parity bits must also be computed and written to parity disk
- To find value of a damaged block, compute XOR of bits from corresponding blocks (including parity block) from other disks.



(e) RAID 4: block-interleaved parity

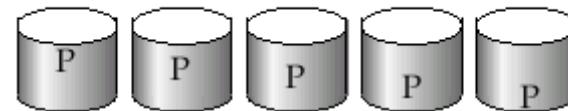
# RAID Levels (Cont.)

## ■ RAID Level 4 (Cont.)

- Provides higher I/O rates for independent block reads than Level 3
  - ▶ block read goes to a single disk, so blocks stored on different disks can be read in parallel
- Provides high transfer rates for reads of multiple blocks than no-striping
- Before writing a block, parity data must be computed
  - ▶ Can be done by using old parity block, old value of current block and new value of current block (2 block reads + 2 block writes)
  - ▶ Or by recomputing the parity value using the new values of blocks corresponding to the parity block
    - More efficient for writing large amounts of data sequentially
- Parity block becomes a bottleneck for independent block writes since every block write also writes to parity disk

# RAID Levels (Cont.)

- **RAID Level 5:** Block-Interleaved Distributed Parity; partitions data and parity among all  $N + 1$  disks, rather than storing data in  $N$  disks and parity in 1 disk.
  - E.g., with 5 disks, parity block for  $n$ th set of blocks is stored on disk  $(n \bmod 5) + 1$ , with the data blocks stored on the other 4 disks.



(f) RAID 5: block-interleaved distributed parity

P0	0	1	2	3
4	P1	5	6	7
8	9	P2	10	11
12	13	14	P3	15
16	17	18	19	P4

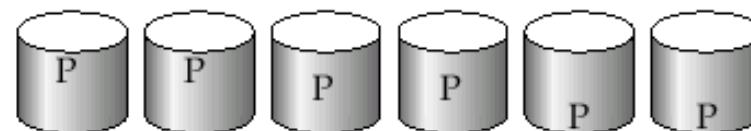
# RAID Levels (Cont.)

## ■ RAID Level 5 (Cont.)

- Higher I/O rates than Level 4.
  - ▶ Block writes occur in parallel if the blocks and their parity blocks are on different disks.
- Subsumes Level 4: provides same benefits, but avoids bottleneck of parity disk.

## ■ RAID Level 6: P+Q Redundancy scheme; similar to Level 5, but stores extra redundant information to guard against multiple disk failures.

- Better reliability than Level 5 at a higher cost; not used as widely.



(g) RAID 6: P + Q redundancy

# Choice of RAID Level

- Factors in choosing RAID level
  - Monetary cost
  - Performance: Number of I/O operations per second, and bandwidth during normal operation
  - Performance during failure
  - Performance during rebuild of failed disk
    - ▶ Including time taken to rebuild failed disk
- RAID 0 is used only when data safety is not important
  - E.g. data can be recovered quickly from other sources
- Level 2 and 4 never used since they are subsumed by 3 and 5
- Level 3 is not used anymore since bit-striping forces single block reads to access all disks, wasting disk arm movement, which block striping (level 5) avoids
- Level 6 is rarely used since levels 1 and 5 offer adequate safety for most applications

# Choice of RAID Level (Cont.)

- Level 1 provides much better write performance than level 5
  - Level 5 requires at least 2 block reads and 2 block writes to write a single block, whereas Level 1 only requires 2 block writes
  - Level 1 preferred for high update environments such as log disks
- Level 1 had higher storage cost than level 5
  - disk drive capacities increasing rapidly (50%/year) whereas disk access times have decreased much less (x 3 in 10 years)
  - I/O requirements have increased greatly, e.g. for Web servers
  - When enough disks have been bought to satisfy required rate of I/O, they often have spare storage capacity
    - ▶ so there is often no extra monetary cost for Level 1!
- Level 5 is preferred for applications with low update rate, and large amounts of data
- Level 1 is preferred for all other applications

# Hardware Issues

- **Software RAID:** RAID implementations done entirely in software, with no special hardware support
- **Hardware RAID:** RAID implementations with special hardware
  - Use non-volatile RAM to record writes that are being executed
  - Beware: power failure during write can result in corrupted disk
    - ▶ E.g. failure after writing one block but before writing the second in a mirrored system
    - ▶ Such corrupted data must be detected when power is restored
      - Recovery from corruption is similar to recovery from failed disk
      - NV-RAM helps to efficiently detect potentially corrupted blocks
        - » Otherwise all blocks of disk must be read and compared with mirror/parity block

# Hardware Issues (Cont.)

- **Latent failures:** data successfully written earlier gets damaged
  - can result in data loss even if only one disk fails
- **Data scrubbing:**
  - continually scan for latent failures, and recover from copy/parity
- **Hot swapping:** replacement of disk while system is running, without power down
  - Supported by some hardware RAID systems,
  - reduces time to recovery, and improves availability greatly
- Many systems maintain **spare disks** which are kept online, and used as replacements for failed disks immediately on detection of failure
  - Reduces time to recovery greatly
- Many hardware RAID systems ensure that a single point of failure will not stop the functioning of the system by using
  - Redundant power supplies with battery backup
  - Multiple controllers and multiple interconnections to guard against controller/interconnection failures

# Optical Disks

- Compact disk-read only memory (CD-ROM)
  - Removable disks, 640 MB per disk
  - Seek time about 100 msec (optical read head is heavier and slower)
  - Higher latency (3000 RPM) and lower data-transfer rates (3-6 MB/s) compared to magnetic disks
- Digital Video Disk (DVD)
  - DVD-5 holds 4.7 GB , and DVD-9 holds 8.5 GB
  - DVD-10 and DVD-18 are double sided formats with capacities of 9.4 GB and 17 GB
  - Blu-ray DVD: 27 GB (54 GB for double sided disk)
  - Slow seek time, for same reasons as CD-ROM
- Record once versions (CD-R and DVD-R) are popular
  - data can only be written once, and cannot be erased.
  - high capacity and long lifetime; used for archival storage
  - Multi-write versions (CD-RW, DVD-RW, DVD+RW and DVD-RAM) also available

# Magnetic Tapes

- Hold large volumes of data and provide high transfer rates
  - Few GB for DAT (Digital Audio Tape) format, 10-40 GB with DLT (Digital Linear Tape) format, 100 GB+ with Ultrium format, and 330 GB with Ampex helical scan format
  - Transfer rates from few to 10s of MB/s
- Tapes are cheap, but cost of drives is very high
- Very slow access time in comparison to magnetic and optical disks
  - limited to sequential access.
  - Some formats (Accelis) provide faster seek (10s of seconds) at cost of lower capacity
- Used mainly for backup, for storage of infrequently used information, and as an off-line medium for transferring information from one system to another.
- Tape jukeboxes used for very large capacity storage
  - Multiple petabytes ( $10^{15}$  bytes)

# File Organization, Record Organization and Storage Access

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan  
See [www.db-book.com](http://www.db-book.com) for conditions on re-use

# File Organization

- The database is stored as a collection of *files*. Each file is a sequence of *records*. A record is a sequence of fields.
- One approach:
  - assume record size is fixed
  - each file has records of one particular type only
  - different files are used for different relations

This case is easiest to implement; will consider variable length records later.

# Fixed-Length Records

## ■ Simple approach:

- Store record  $i$  starting from byte  $n * (i - 1)$ , where  $n$  is the size of each record.
- Record access is simple but records may cross blocks
  - ▶ Modification: do not allow records to cross block boundaries

## ■ Deletion of record $i$ : alternatives:

- move records  $i + 1, \dots, n$  to  $i, \dots, n - 1$
- move record  $n$  to  $i$
- do not move records, but link all free records on a *free list*

record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	12121	Wu	Finance	90000
record 2	15151	Mozart	Music	40000
record 3	22222	Einstein	Physics	95000
record 4	32343	El Said	History	60000
record 5	33456	Gold	Physics	87000
record 6	45565	Katz	Comp. Sci.	75000
record 7	58583	Califieri	History	62000
record 8	76543	Singh	Finance	80000
record 9	76766	Crick	Biology	72000
record 10	83821	Brandt	Comp. Sci.	92000
record 11	98345	Kim	Elec. Eng.	80000

# Deleting record 3 and compacting

record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	12121	Wu	Finance	90000
record 2	15151	Mozart	Music	40000
record 4	32343	El Said	History	60000
record 5	33456	Gold	Physics	87000
record 6	45565	Katz	Comp. Sci.	75000
record 7	58583	Califieri	History	62000
record 8	76543	Singh	Finance	80000
record 9	76766	Crick	Biology	72000
record 10	83821	Brandt	Comp. Sci.	92000
record 11	98345	Kim	Elec. Eng.	80000

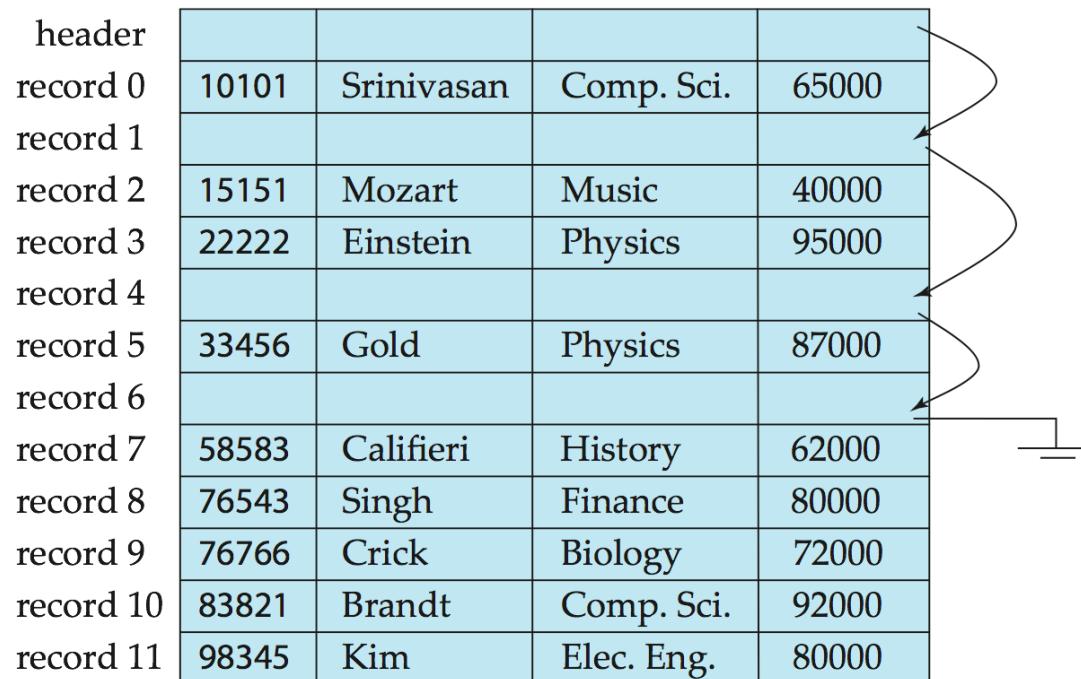
# Deleting record 3 and moving last record

record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	12121	Wu	Finance	90000
record 2	15151	Mozart	Music	40000
record 11	98345	Kim	Elec. Eng.	80000
record 4	32343	El Said	History	60000
record 5	33456	Gold	Physics	87000
record 6	45565	Katz	Comp. Sci.	75000
record 7	58583	Califieri	History	62000
record 8	76543	Singh	Finance	80000
record 9	76766	Crick	Biology	72000
record 10	83821	Brandt	Comp. Sci.	92000

# Free Lists

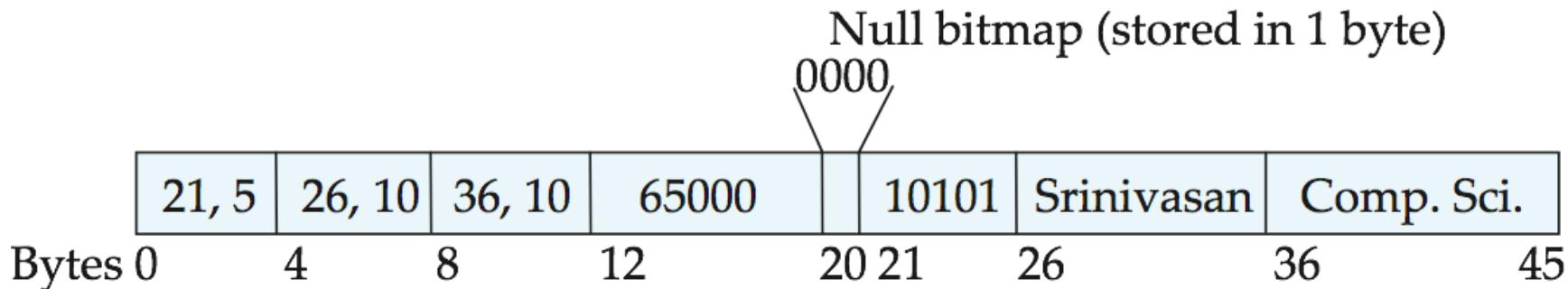
- Store the address of the first deleted record in the file header.
- Use this first record to store the address of the second deleted record, and so on
- Can think of these stored addresses as **pointers** since they “point” to the location of a record.
- More space efficient representation: reuse space for normal attributes of free records to store pointers. (No pointers stored in in-use records.)

header				
record 0	10101	Srinivasan	Comp. Sci.	65000
record 1				
record 2	15151	Mozart	Music	40000
record 3	22222	Einstein	Physics	95000
record 4				
record 5	33456	Gold	Physics	87000
record 6				
record 7	58583	Califieri	History	62000
record 8	76543	Singh	Finance	80000
record 9	76766	Crick	Biology	72000
record 10	83821	Brandt	Comp. Sci.	92000
record 11	98345	Kim	Elec. Eng.	80000

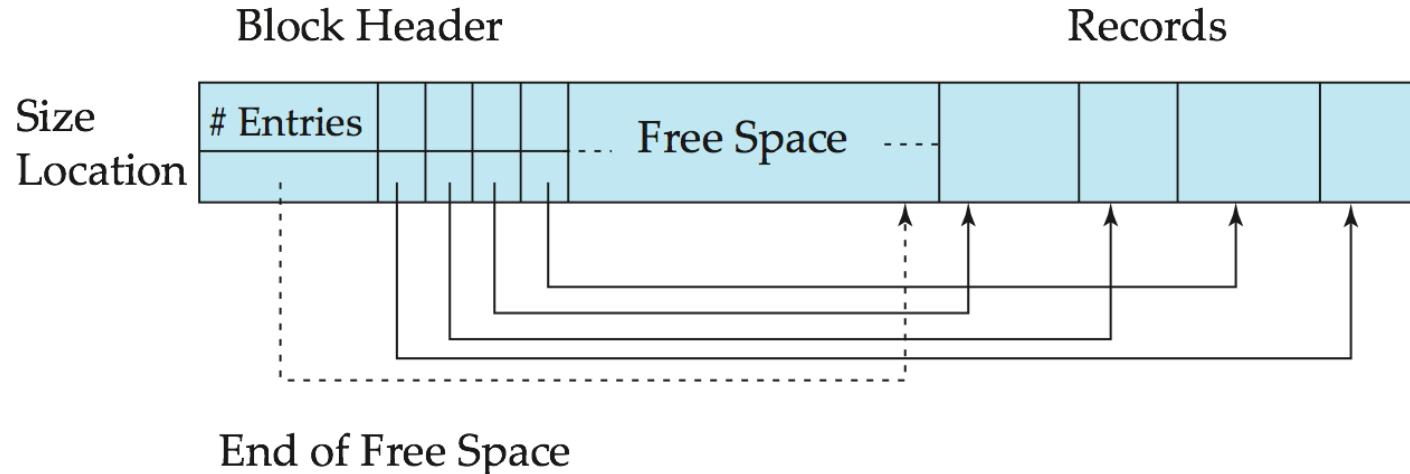


# Variable-Length Records

- Variable-length records arise in database systems in several ways:
  - Storage of multiple record types in a file.
  - Record types that allow variable lengths for one or more fields such as strings (**varchar**)
  - Record types that allow repeating fields (used in some older data models).
- Attributes are stored in order
- Variable length attributes represented by fixed size (offset, length), with actual data stored after all fixed length attributes
- Null values represented by null-value bitmap



# Variable-Length Records: Slotted Page Structure



- **Slotted page** header contains:

- number of record entries
- end of free space in the block
- location and size of each record

- Records can be moved around within a page to keep them contiguous with no empty space between them; entry in the header must be updated.
- Pointers should not point directly to record — instead they should point to the entry for the record in header.

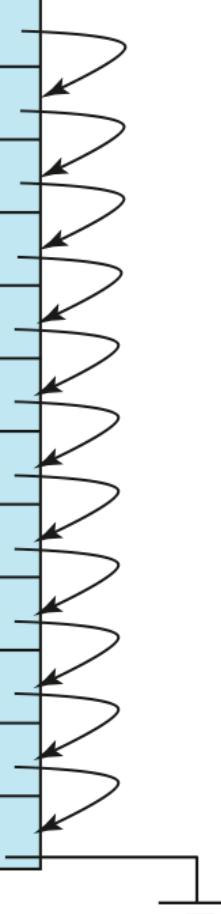
# Organization of Records in Files

- **Heap** – a record can be placed anywhere in the file where there is space
- **Sequential** – store records in sequential order, based on the value of the search key of each record
- **Hashing** – a hash function computed on some attribute of each record; the result specifies in which block of the file the record should be placed
- Records of each relation may be stored in a separate file. In a **multitable clustering file organization** records of several different relations can be stored in the same file
  - Motivation: store related records on the same block to minimize I/O

# Sequential File Organization

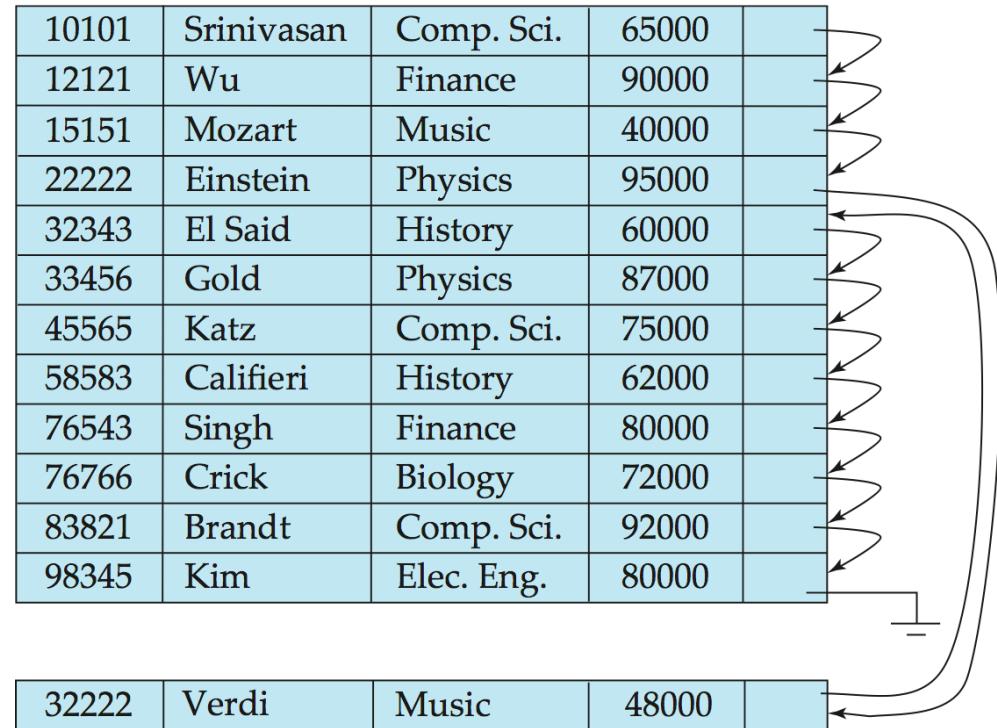
- Suitable for applications that require sequential processing of the entire file
- The records in the file are ordered by a **search-key**

10101	Srinivasan	Comp. Sci.	65000	
12121	Wu	Finance	90000	
15151	Mozart	Music	40000	
22222	Einstein	Physics	95000	
32343	El Said	History	60000	
33456	Gold	Physics	87000	
45565	Katz	Comp. Sci.	75000	
58583	Califieri	History	62000	
76543	Singh	Finance	80000	
76766	Crick	Biology	72000	
83821	Brandt	Comp. Sci.	92000	
98345	Kim	Elec. Eng.	80000	



# Sequential File Organization (Cont.)

- Deletion – use pointer chains
- Insertion – locate the position where the record is to be inserted
  - if there is free space insert there
  - if no free space, insert the record in an **overflow block**
  - In either case, pointer chain must be updated
- Need to reorganize the file from time to time to restore sequential order



# Multitable Clustering File Organization

Store several relations in one file using a **multitable clustering** file organization

*department*

<i>dept_name</i>	<i>building</i>	<i>budget</i>
Comp. Sci.	Taylor	100000
Physics	Watson	70000

*instructor*

<i>ID</i>	<i>name</i>	<i>dept_name</i>	<i>salary</i>
10101	Srinivasan	Comp. Sci.	65000
33456	Gold	Physics	87000
45565	Katz	Comp. Sci.	75000
83821	Brandt	Comp. Sci.	92000

multitable clustering  
of *department* and  
*instructor*

Comp. Sci.	Taylor	100000
45564	Katz	75000
10101	Srinivasan	65000
83821	Brandt	92000
Physics	Watson	70000
33456	Gold	87000

# Multitable Clustering File Organization (cont.)

- good for queries involving *department*   *instructor*, and for queries involving one single department and its ~~in~~structors
- bad for queries involving only *department*
- results in variable size records
- Can add pointer chains to link records of a particular relation

Comp. Sci.	Taylor	100000	
45564	Katz	75000	
10101	Srinivasan	65000	
83821	Brandt	92000	
Physics	Watson	70000	
33456	Gold	87000	



The diagram shows a curved arrow originating from the rightmost cell of the 'Physics' row and pointing to the right. At the end of this arrow is a small T-junction, from which a horizontal line extends downwards, ending in a small square symbol. This symbol is positioned directly above the rightmost cell of the 'Gold' row's corresponding cluster, indicating a pointer from the 'Physics' record to the 'Gold' record.

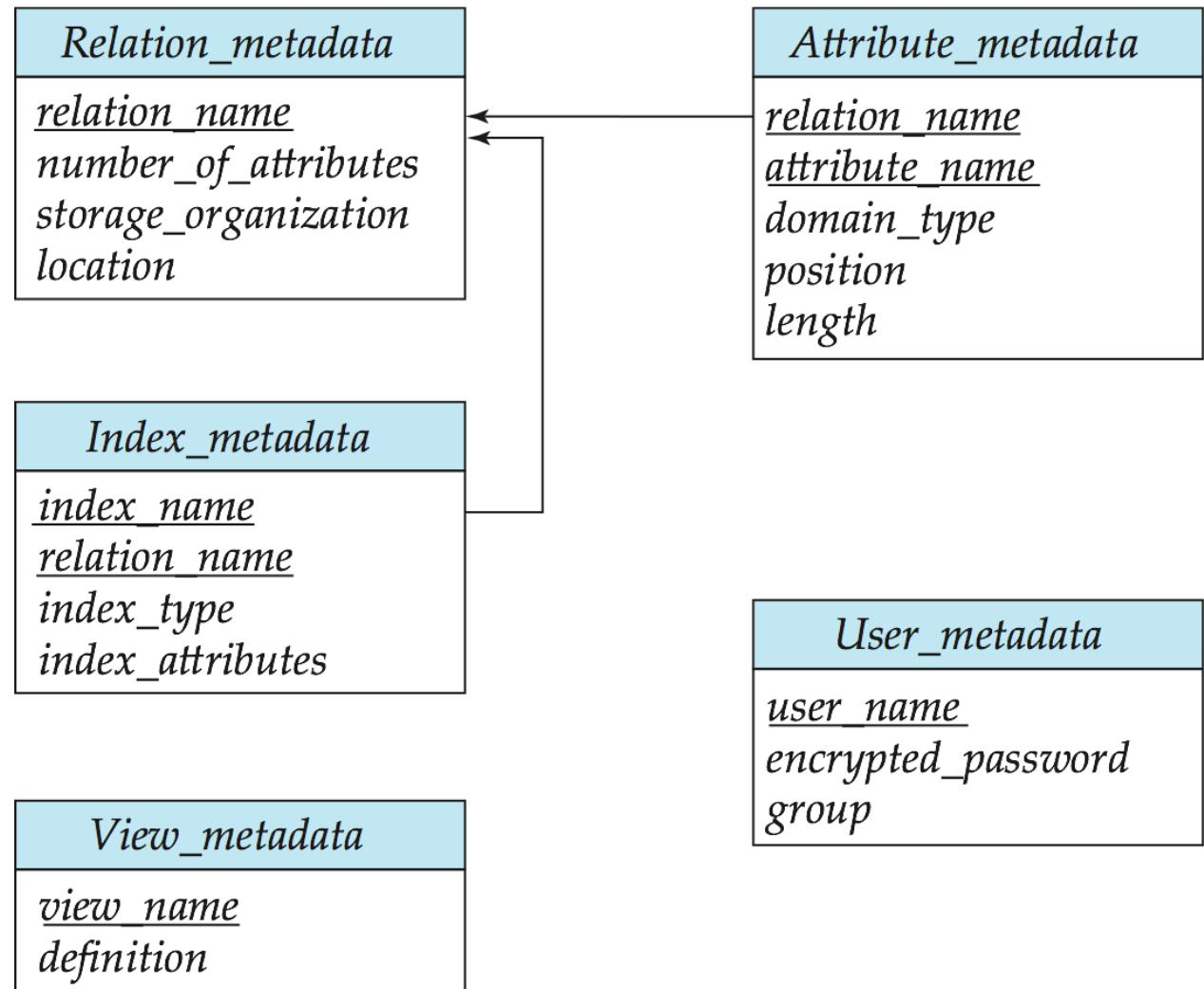
# Data Dictionary Storage

The **Data dictionary** (also called **system catalog**) stores **metadata**; that is, data about data, such as

- Information about relations
  - names of relations
  - names, types and lengths of attributes of each relation
  - names and definitions of views
  - integrity constraints
- User and accounting information, including passwords
- Statistical and descriptive data
  - number of tuples in each relation
- Physical file organization information
  - How relation is stored (sequential/hash/...)
  - Physical location of relation
- Information about indices (Chapter 11)

# Relational Representation of System Metadata

- Relational representation on disk
- Specialized data structures designed for efficient access, in memory



## Storage Access

- A database file is partitioned into fixed-length storage units called **blocks**.  
Blocks are units of both storage allocation and data transfer.
- Database system seeks to minimize the number of block transfers between the disk and memory. We can reduce the number of disk accesses by keeping as many blocks as possible in main memory.
- **Buffer** – portion of main memory available to store copies of disk blocks.
- **Buffer manager** – subsystem responsible for allocating buffer space in main memory.

# Buffer Manager

- Programs call on the buffer manager when they need a block from disk.
  1. If the block is already in the buffer, buffer manager returns the address of the block in main memory
  2. If the block is not in the buffer, the buffer manager
    1. Allocates space in the buffer for the block
      1. Replacing (throwing out) some other block, if required, to make space for the new block.
      2. Replaced block written back to disk only if it was modified since the most recent time that it was written to/fetched from the disk.
    2. Reads the block from the disk to the buffer, and returns the address of the block in main memory to requester.

# Buffer-Replacement Policies

- Most operating systems replace the block **least recently used (LRU strategy)**
- Idea behind LRU – use past pattern of block references as a predictor of future references
- Queries have well-defined access patterns (such as sequential scans), and a database system can use the information in a user's query to predict future references
  - LRU can be a bad strategy for certain access patterns involving repeated scans of data
    - ▶ For example: when computing the join of 2 relations  $r$  and  $s$  by a nested loops

for each tuple  $tr$  of  $r$  do  
    for each tuple  $ts$  of  $s$  do  
        if the tuples  $tr$  and  $ts$  match ...
  - Mixed strategy with hints on replacement strategy provided by the query optimizer is preferable

## Buffer-Replacement Policies (Cont.)

- **Pinned block** – memory block that is not allowed to be written back to disk.
- **Toss-immediate** strategy – frees the space occupied by a block as soon as the final tuple of that block has been processed
- **Most recently used (MRU) strategy** – system must pin the block currently being processed. After the final tuple of that block has been processed, the block is unpinned, and it becomes the most recently used block.
- Buffer manager can use statistical information regarding the probability that a request will reference a particular relation
  - E.g., the data dictionary is frequently accessed. Heuristic: keep data-dictionary blocks in main memory buffer
- Buffer managers also support **forced output** of blocks for the purpose of recovery (more in Chapter 16)

# End of Chapter 10

**Database System Concepts, 6<sup>th</sup> Ed.**

©Silberschatz, Korth and Sudarshan  
See [www.db-book.com](http://www.db-book.com) for conditions on re-use

# Chapter 11: Indexing and Hashing

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan  
See [www.db-book.com](http://www.db-book.com) for conditions on re-use

# Chapter 11: Indexing and Hashing

- Basic Concepts
- Ordered Indices
- B<sup>+</sup>-Tree Index Files
- B-Tree Index Files
- Static Hashing
- Dynamic Hashing
- Comparison of Ordered Indexing and Hashing
- Index Definition in SQL
- Multiple-Key Access

# Basic Concepts

- Indexing mechanisms used to speed up access to desired data.
  - E.g., author catalog in library
- **Search Key** - attribute or set of attributes used to look up records in a file.
- An **index file** consists of records (called **index entries**) of the form



- Index files are typically much smaller than the original file
- Two basic kinds of indices:
  - **Ordered indices:** search keys are stored in sorted order
  - **Hash indices:** search keys are distributed uniformly across “buckets” using a “hash function”.

# Index Evaluation Metrics

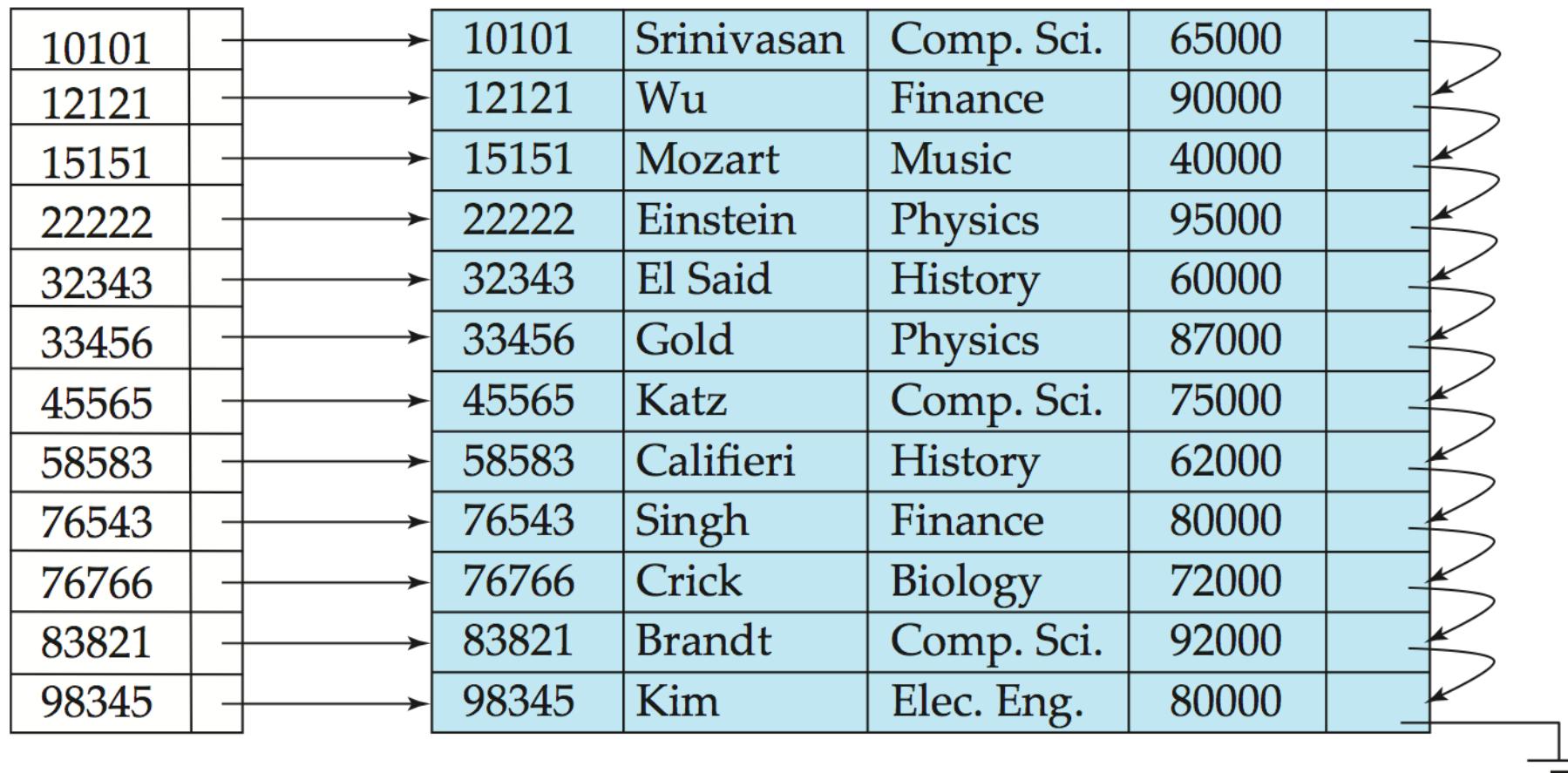
- Access types supported efficiently. E.g.,
  - records with a specified value in the attribute
  - or records with an attribute value falling in a specified range of values.
- Access time
- Insertion time
- Deletion time
- Space overhead

# Ordered Indices

- In an **ordered index**, index entries are stored sorted on the search key value. E.g., author catalog in library.
- **Primary index**: in a sequentially ordered file, the index whose search key specifies the sequential order of the file.
  - Also called **clustering index**
  - The search key of a primary index is usually but not necessarily the primary key.
- **Secondary index**: an index whose search key specifies an order different from the sequential order of the file. Also called **non-clustering index**.
- **Index-sequential file**: ordered sequential file with a primary index.

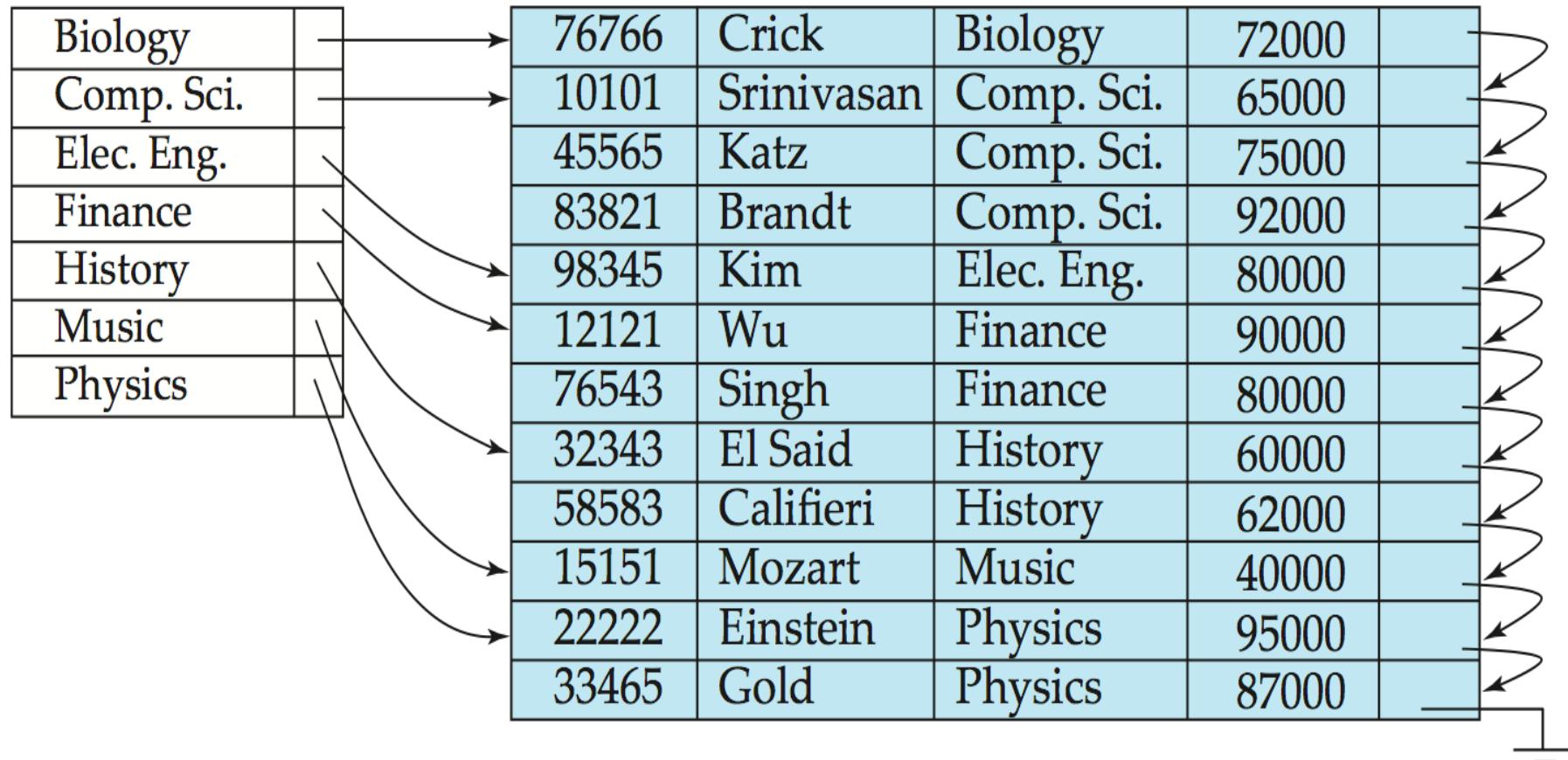
# Dense Index Files

- **Dense index** — Index record appears for every search-key value in the file.
- E.g. index on *ID* attribute of *instructor* relation



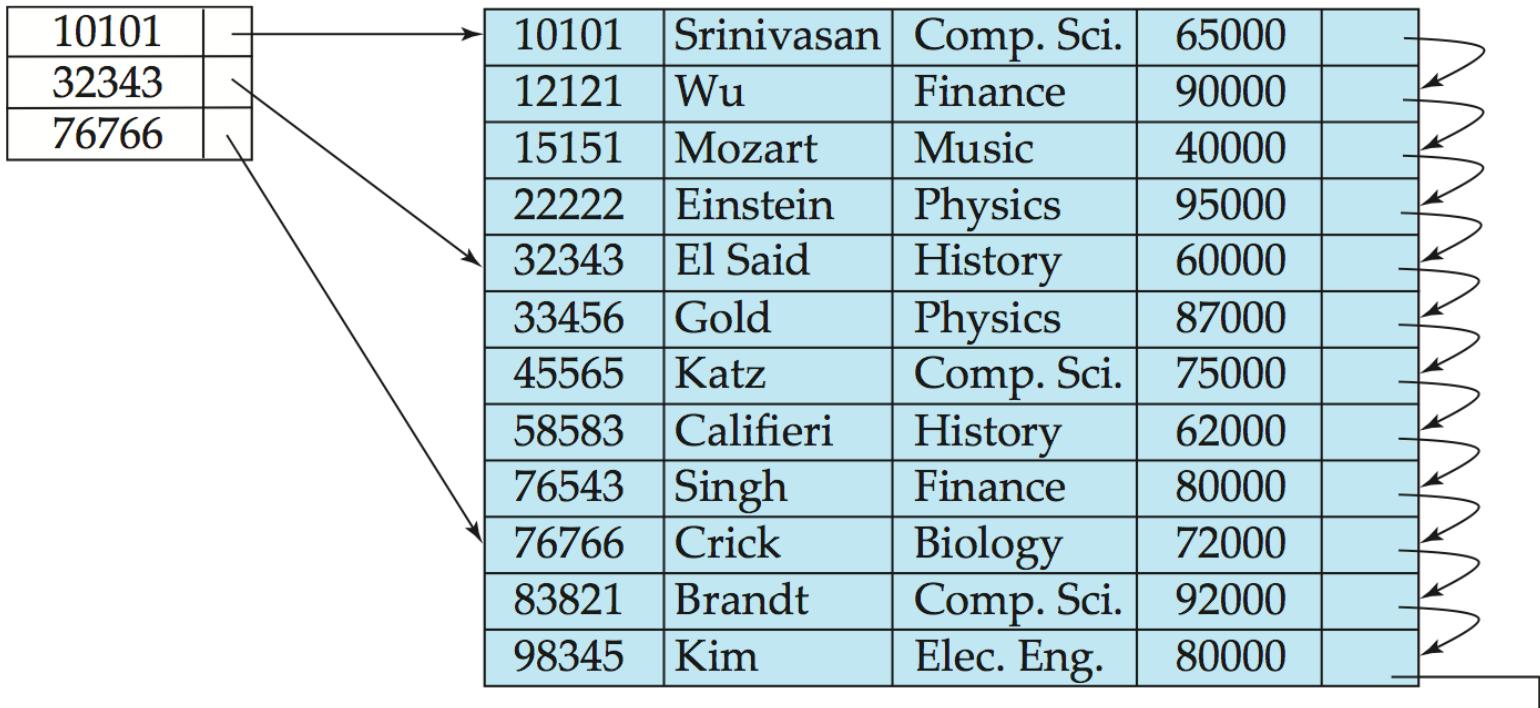
# Dense Index Files (Cont.)

- Dense index on *dept\_name*, with *instructor* file sorted on *dept\_name*



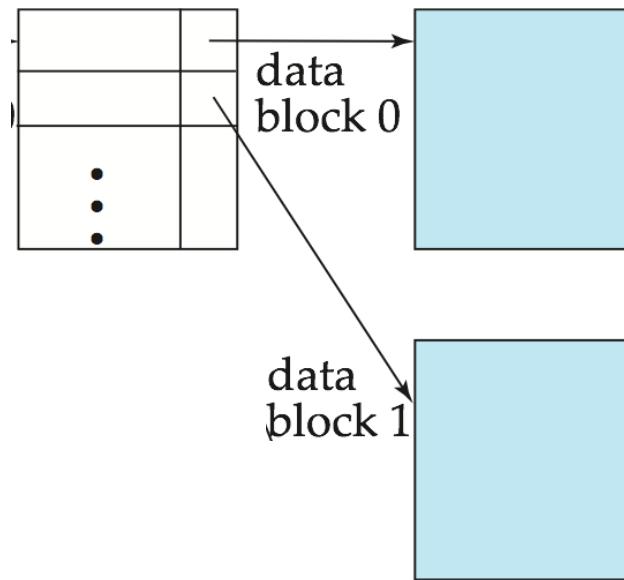
# Sparse Index Files

- **Sparse Index:** contains index records for only some search-key values.
  - Applicable when records are sequentially ordered on search-key
- To locate a record with search-key value  $K$  we:
  - Find index record with largest search-key value  $< K$
  - Search file sequentially starting at the record to which the index record points

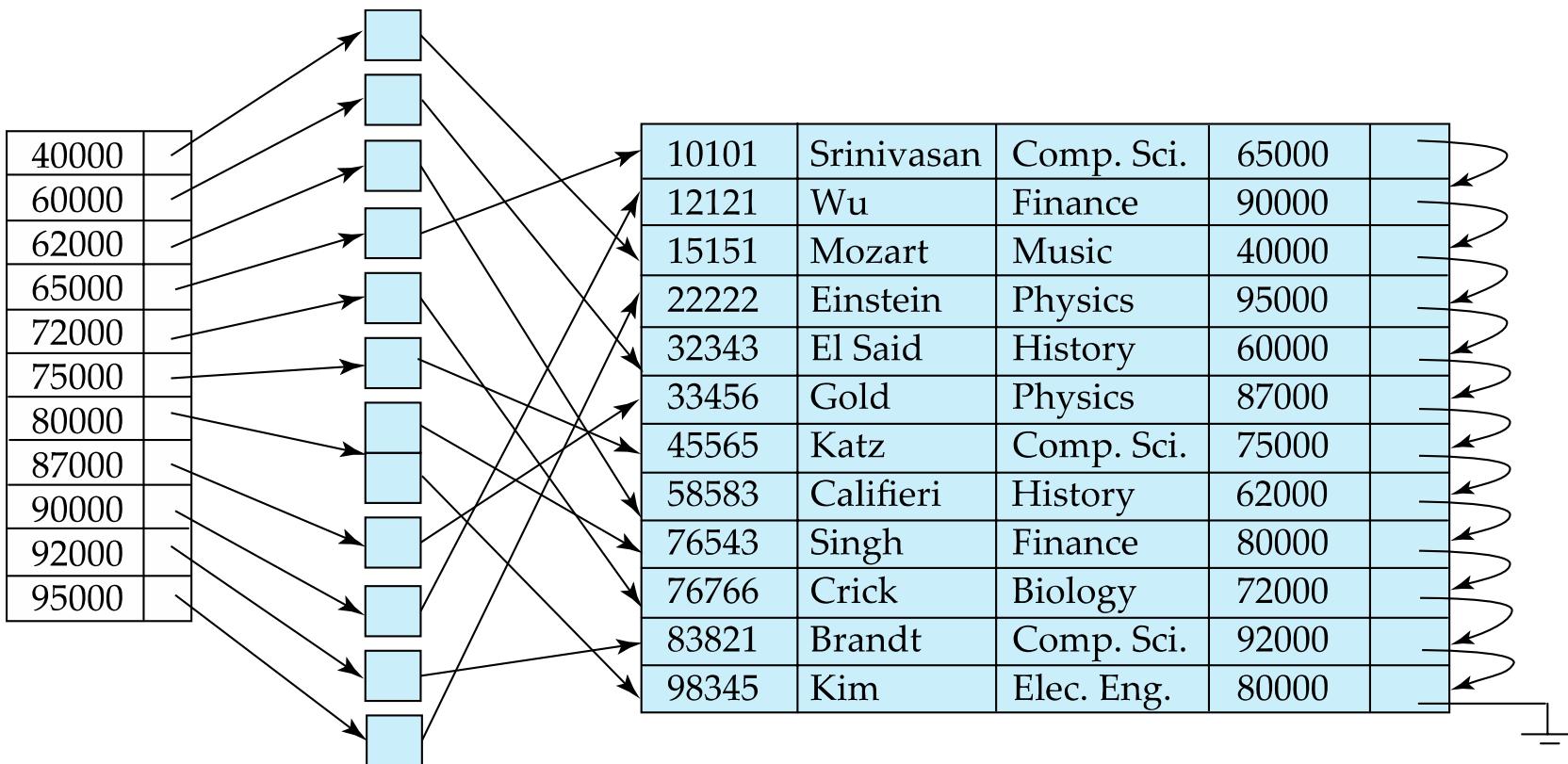


# Sparse Index Files (Cont.)

- Compared to dense indices:
  - Less space and less maintenance overhead for insertions and deletions.
  - Generally slower than dense index for locating records.
- **Good tradeoff:** sparse index with an index entry for every block in file, corresponding to least search-key value in the block.



# Secondary Indices Example



**Secondary index on salary field of *instructor***

- Index record points to a bucket that contains pointers to all the actual records with that particular search-key value.
- Secondary indices have to be dense

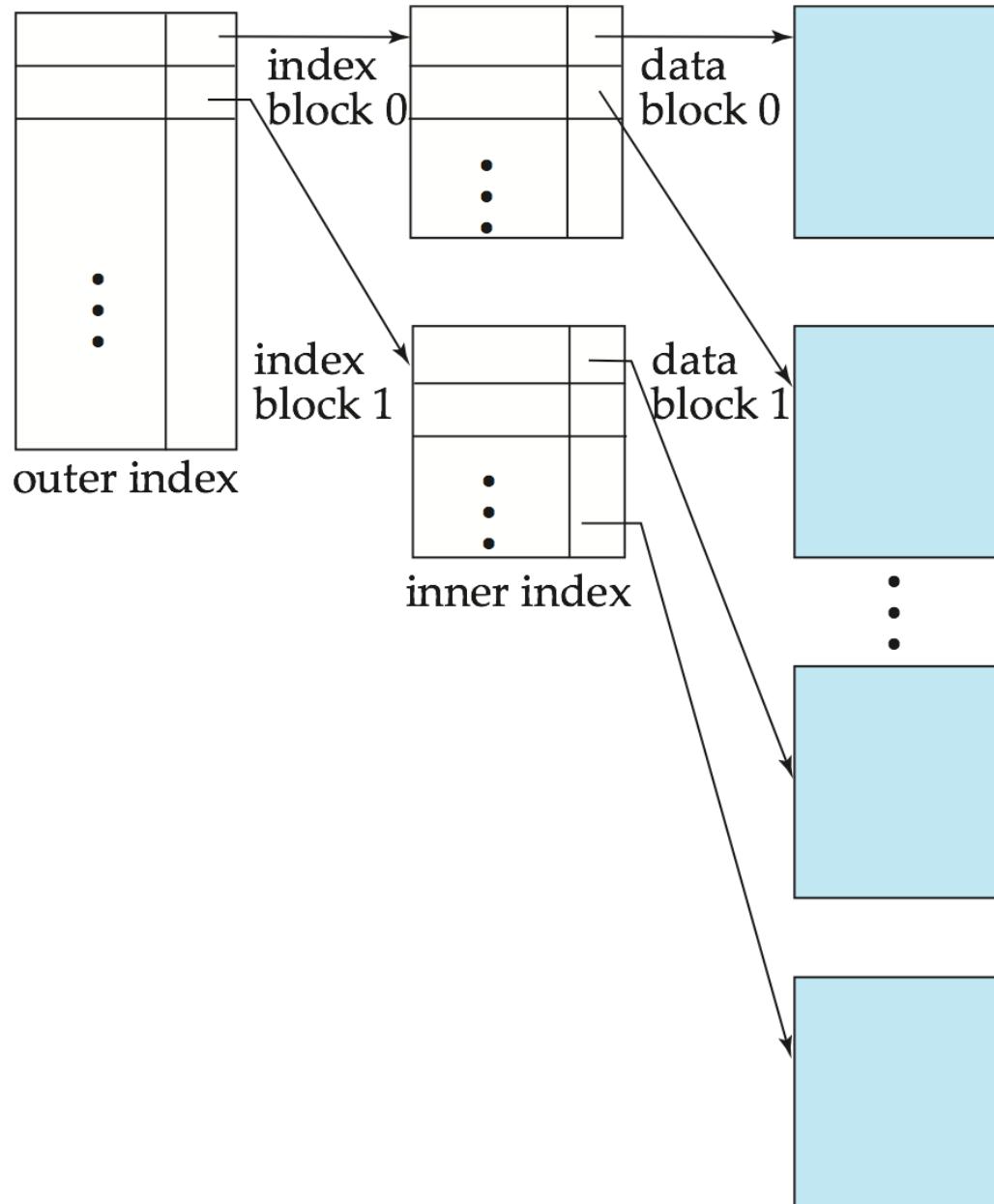
# Primary and Secondary Indices

- Indices offer substantial benefits when searching for records.
- BUT: Updating indices imposes overhead on database modification -- when a file is modified, every index on the file must be updated,
- Sequential scan using primary index is efficient, but a sequential scan using a secondary index is expensive
  - Each record access may fetch a new block from disk
  - Block fetch requires about 5 to 10 milliseconds, versus about 100 nanoseconds for memory access

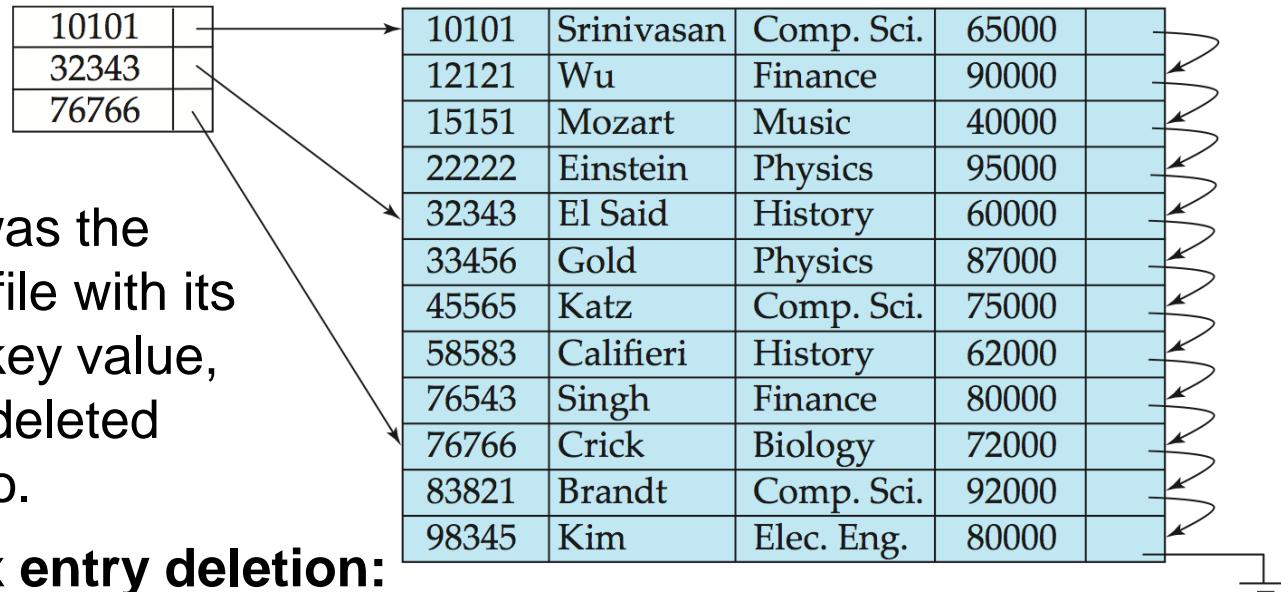
# Multilevel Index

- If primary index does not fit in memory, access becomes expensive.
- Solution: treat primary index kept on disk as a sequential file and construct a sparse index on it.
  - outer index – a sparse index of primary index
  - inner index – the primary index file
- If even outer index is too large to fit in main memory, yet another level of index can be created, and so on.
- Indices at all levels must be updated on insertion or deletion from the file.

# Multilevel Index (Cont.)



# Index Update: Deletion



- If deleted record was the only record in the file with its particular search-key value, the search-key is deleted from the index also.

## ■ Single-level index entry deletion:

- **Dense indices** – deletion of search-key is similar to file record deletion.
- **Sparse indices** –
  - ▶ if an entry for the search key exists in the index, it is deleted by replacing the entry in the index with the next search-key value in the file (in search-key order).
  - ▶ If the next search-key value already has an index entry, the entry is deleted instead of being replaced.

# Index Update: Insertion

## ■ Single-level index insertion:

- Perform a lookup using the search-key value appearing in the record to be inserted.
- **Dense indices** – if the search-key value does not appear in the index, insert it.
- **Sparse indices** – if index stores an entry for each block of the file, no change needs to be made to the index unless a new block is created.
  - ▶ If a new block is created, the first search-key value appearing in the new block is inserted into the index.

## ■ Multilevel insertion and deletion: algorithms are simple extensions of the single-level algorithms

# Secondary Indices

- Frequently, one wants to find all the records whose values in a certain field (which is not the search-key of the primary index) satisfy some condition.
  - Example 1: In the *instructor* relation stored sequentially by ID, we may want to find all instructors in a particular department
  - Example 2: as above, but where we want to find all instructors with a specified salary or with salary in a specified range of values
- We can have a secondary index with an index record for each search-key value

# B<sup>+</sup>-Tree Index Files

B<sup>+</sup>-tree indices are an alternative to indexed-sequential files.

■ Disadvantage of indexed-sequential files

- performance degrades as file grows, since many overflow blocks get created.
- Periodic reorganization of entire file is required.

■ Advantage of B<sup>+</sup>-tree index files:

- automatically reorganizes itself with small, local, changes, in the face of insertions and deletions.
- Reorganization of entire file is not required to maintain performance.

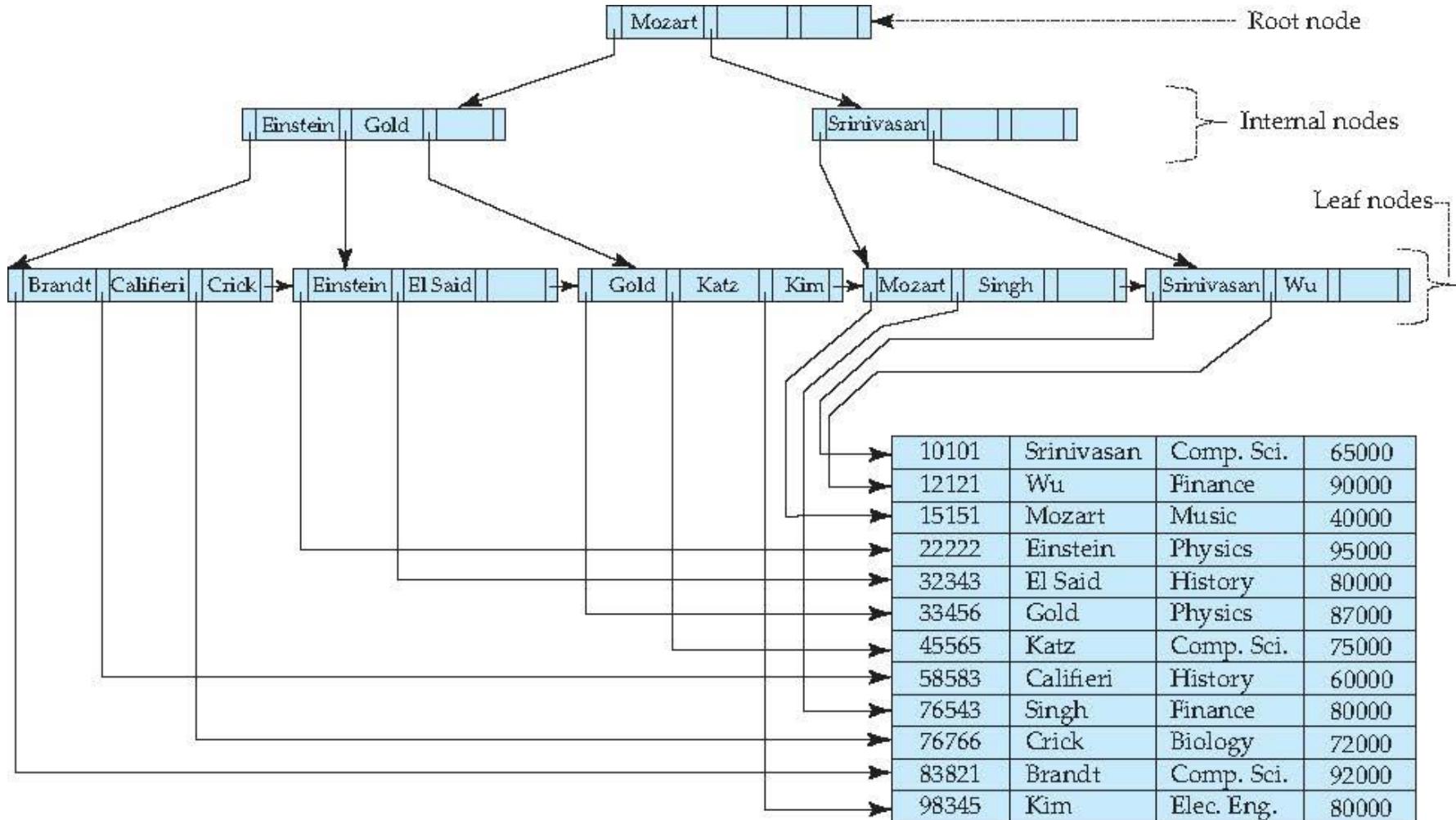
■ (Minor) disadvantage of B<sup>+</sup>-trees:

- extra insertion and deletion overhead, space overhead.

■ Advantages of B<sup>+</sup>-trees outweigh disadvantages

- B<sup>+</sup>-trees are used extensively

# Example of B+-Tree



# B<sup>+</sup>-Tree Index Files (Cont.)

**A B<sup>+</sup>-tree is a rooted tree satisfying the following properties:**

- All paths from root to leaf are of the same length
- Each node that is not a root or a leaf has between  $\lceil n/2 \rceil$  and  $n$  children.
- A leaf node has between  $\lceil (n-1)/2 \rceil$  and  $n-1$  values
- Special cases:
  - If the root is not a leaf, it has at least 2 children.
  - If the root is a leaf (that is, there are no other nodes in the tree), it can have between 0 and  $(n-1)$  values.

# B+-Tree Node Structure

## ■ Typical node

$P_1$	$K_1$	$P_2$	$\dots$	$P_{n-1}$	$K_{n-1}$	$P_n$
-------	-------	-------	---------	-----------	-----------	-------

- $K_i$  are the search-key values
- $P_i$  are pointers to children (for non-leaf nodes) or pointers to records or buckets of records (for leaf nodes).

## ■ The search-keys in a node are ordered

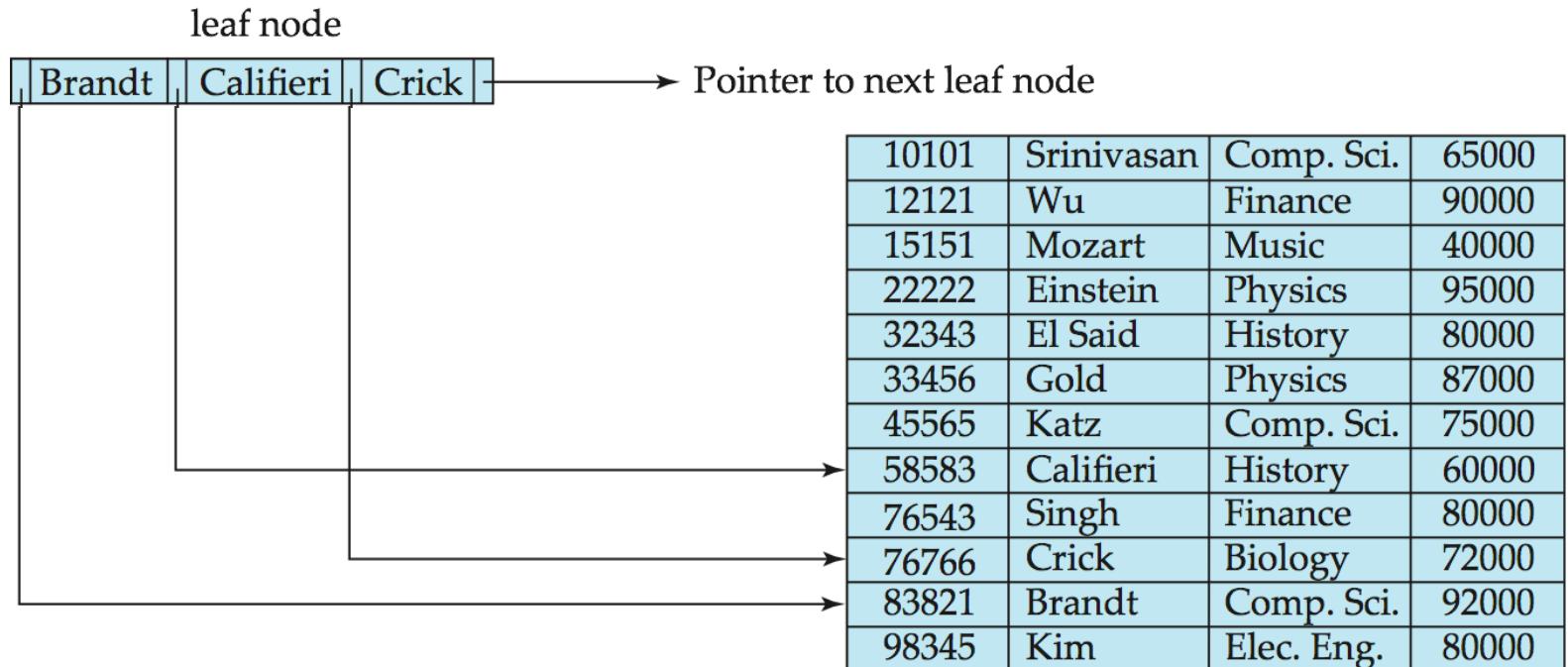
$$K_1 < K_2 < K_3 < \dots < K_{n-1}$$

(Initially assume no duplicate keys, address duplicates later)

# Leaf Nodes in B+-Trees

## Properties of a leaf node:

- For  $i = 1, 2, \dots, n-1$ , pointer  $P_i$  points to a file record with search-key value  $K_i$ ,
- If  $L_i, L_j$  are leaf nodes and  $i < j$ ,  $L_i$ 's search-key values are less than or equal to  $L_j$ 's search-key values
- $P_n$  points to next leaf node in search-key order

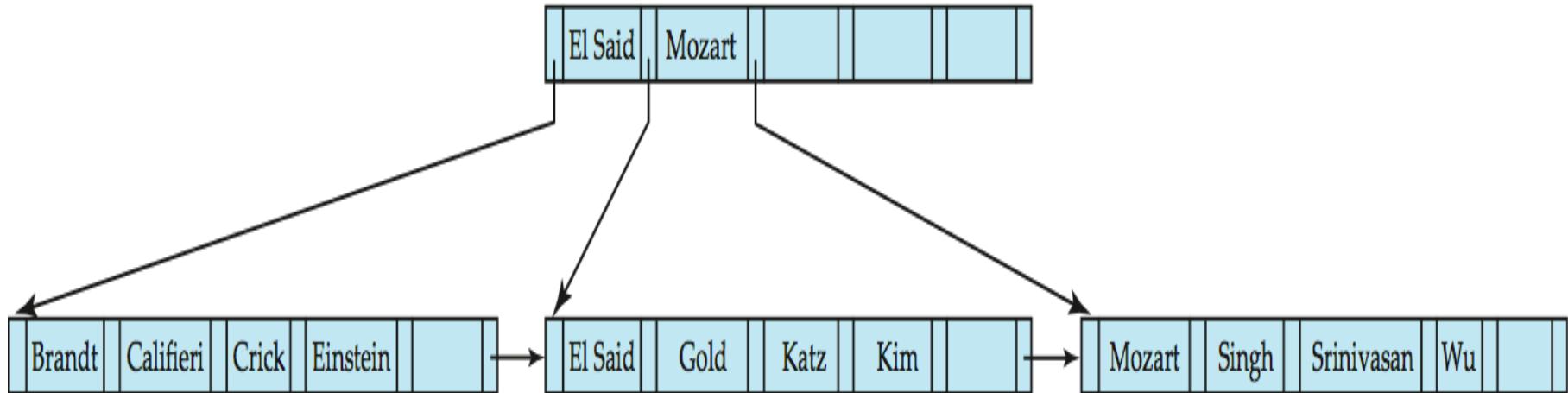


# Non-Leaf Nodes in B<sup>+</sup>-Trees

- Non leaf nodes form a multi-level sparse index on the leaf nodes. For a non-leaf node with  $m$  pointers:
  - All the search-keys in the subtree to which  $P_1$  points are less than  $K_1$
  - For  $2 \leq i \leq n - 1$ , all the search-keys in the subtree to which  $P_i$  points have values greater than or equal to  $K_{i-1}$  and less than  $K_i$
  - All the search-keys in the subtree to which  $P_n$  points have values greater than or equal to  $K_{n-1}$

$P_1$	$K_1$	$P_2$	...	$P_{n-1}$	$K_{n-1}$	$P_n$
-------	-------	-------	-----	-----------	-----------	-------

# Example of B<sup>+</sup>-tree



B<sup>+</sup>-tree for *instructor* file ( $n = 6$ )

- Leaf nodes must have between 3 and 5 values ( $\lceil (n-1)/2 \rceil$  and  $n-1$ , with  $n = 6$ ).
- Non-leaf nodes other than root must have between 3 and 6 children ( $\lceil (n/2) \rceil$  and  $n$  with  $n = 6$ ).
- Root must have at least 2 children.

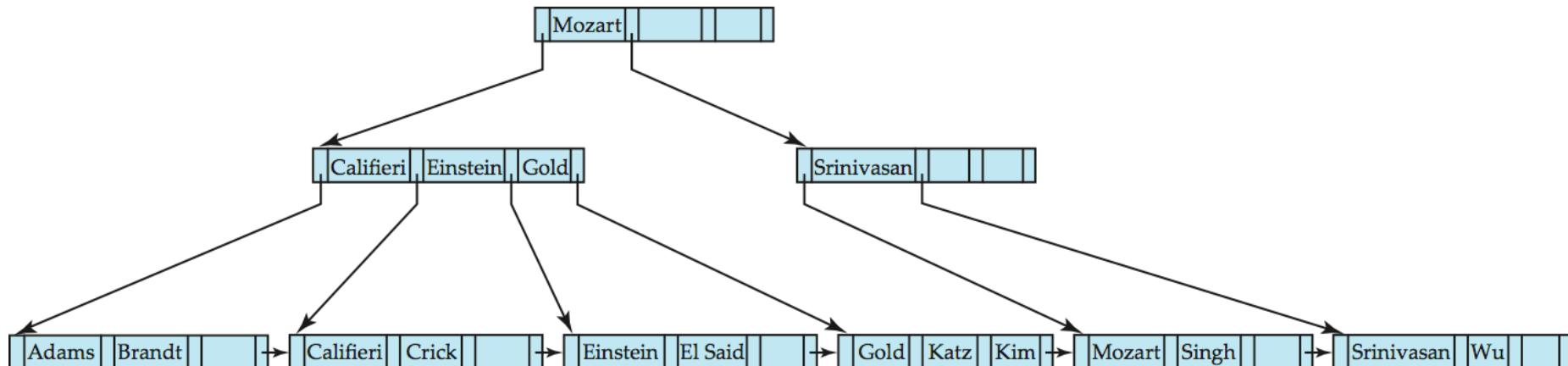
# Observations about B+-trees

- Since the inter-node connections are done by pointers, “logically” close blocks need not be “physically” close.
- The non-leaf levels of the B+-tree form a hierarchy of sparse indices.
- The B+-tree contains a relatively small number of levels
  - ▶ Level below root has at least  $2 * \lceil n/2 \rceil$  values
  - ▶ Next level has at least  $2 * \lceil n/2 \rceil * \lceil n/2 \rceil$  values
  - ▶ .. etc.
- If there are  $K$  search-key values in the file, the tree height is no more than  $\lceil \log_{\lceil n/2 \rceil}(K) \rceil$ 
  - thus searches can be conducted efficiently.
- Insertions and deletions to the main file can be handled efficiently, as the index can be restructured in logarithmic time (as we shall see).

# Queries on B<sup>+</sup>-Trees

- Find record with search-key value  $V$ .

1.  $C = \text{root}$
2. While  $C$  is not a leaf node {
  1. Let  $i$  be least value s.t.  $V \leq K_i$ .
  2. If no such exists, set  $C = \text{last non-null pointer in } C$
  3. Else { if ( $V = K_i$ ) Set  $C = P_{i+1}$  else set  $C = P_i$ }
3. Let  $i$  be least value s.t.  $K_i = V$
4. If there is such a value  $i$ , follow pointer  $P_i$  to the desired record.
5. Else no record with search-key value  $k$  exists.



# Handling Duplicates

## ■ With duplicate search keys

- In both leaf and internal nodes,
  - ▶ we cannot guarantee that  $K_1 < K_2 < K_3 < \dots < K_{n-1}$
  - ▶ but can guarantee  $K_1 \leq K_2 \leq K_3 \leq \dots \leq K_{n-1}$
- Search-keys in the subtree to which  $P_i$  points
  - ▶ are  $\leq K_i$ , but not necessarily  $< K_i$ ,
  - ▶ To see why, suppose same search key value  $V$  is present in two leaf node  $L_i$  and  $L_{i+1}$ . Then in parent node  $K_i$  must be equal to  $V$

# Handling Duplicates

## ■ We modify find procedure as follows

- traverse  $P_i$  even if  $V = K_i$
- As soon as we reach a leaf node  $C$  check if  $C$  has only search key values less than  $V$ 
  - ▶ if so set  $C =$  right sibling of  $C$  before checking whether  $C$  contains  $V$

## ■ Procedure printAll

- uses modified find procedure to find first occurrence of  $V$
- Traverse through consecutive leaves to find all occurrences of  $V$

\*\* Errata note: modified find procedure missing in first printing of 6<sup>th</sup> edition

## Queries on B+-Trees (Cont.)

- If there are  $K$  search-key values in the file, the height of the tree is no more than  $\lceil \log_{n/2}(K) \rceil$ .
- A node is generally the same size as a disk block, typically 4 kilobytes
  - and  $n$  is typically around 100 (40 bytes per index entry).
- With 1 million search key values and  $n = 100$ 
  - at most  $\log_{50}(1,000,000) = 4$  nodes are accessed in a lookup.
- Contrast this with a balanced binary tree with 1 million search key values
  - around 20 nodes are accessed in a lookup
  - above difference is significant since every node access may need a disk I/O, costing around 20 milliseconds

# Updates on B+-Trees: Insertion

1. Find the leaf node in which the search-key value would appear
2. If the search-key value is already present in the leaf node
  1. Add record to the file
  2. If necessary add a pointer to the bucket.
3. If the search-key value is not present, then
  1. add the record to the main file (and create a bucket if necessary)
  2. If there is room in the leaf node, insert (key-value, pointer) pair in the leaf node
  3. Otherwise, split the node (along with the new (key-value, pointer) entry) as discussed in the next slide.

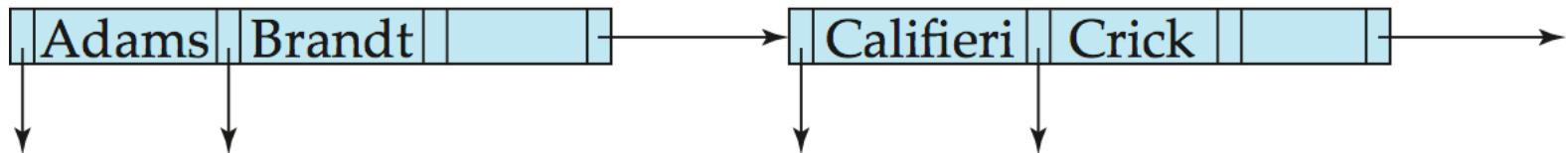
# Updates on B+-Trees: Insertion (Cont.)

## ■ Splitting a leaf node:

- take the  $n$  (search-key value, pointer) pairs (including the one being inserted) in sorted order. Place the first  $\lceil n/2 \rceil$  in the original node, and the rest in a new node.
- let the new node be  $p$ , and let  $k$  be the least key value in  $p$ . Insert  $(k,p)$  in the parent of the node being split.
- If the parent is full, split it and **propagate** the split further up.

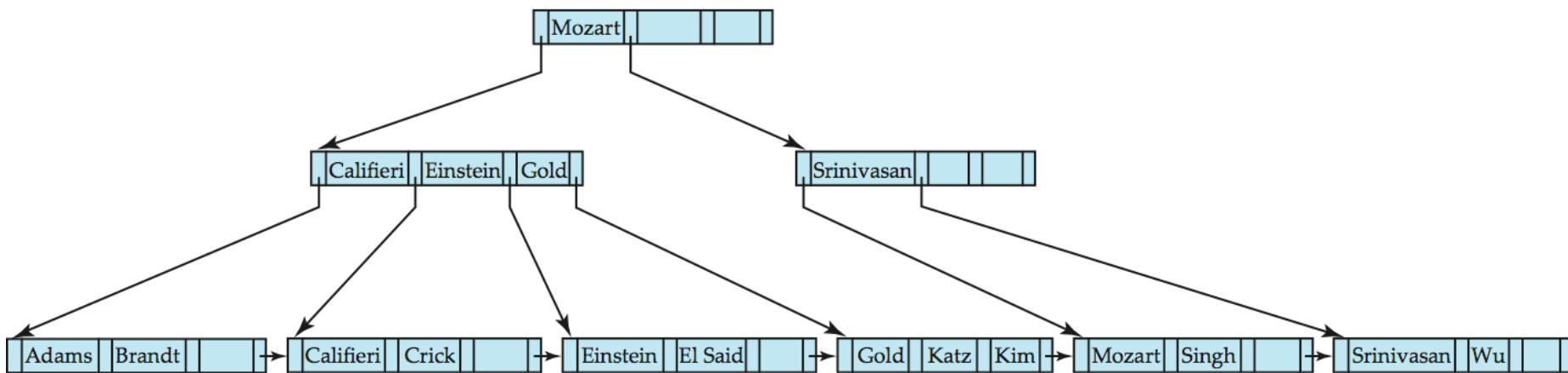
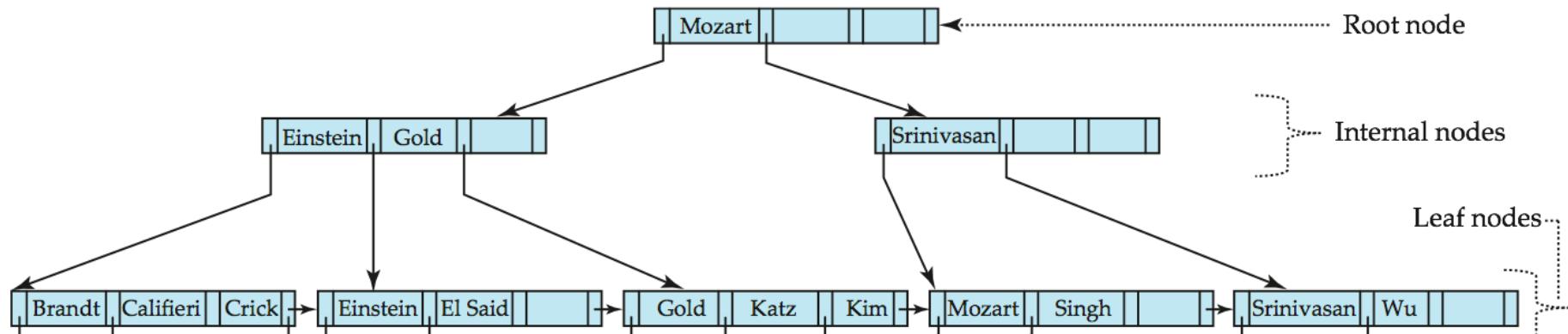
## ■ Splitting of nodes proceeds upwards till a node that is not full is found.

- In the worst case the root node may be split increasing the height of the tree by 1.



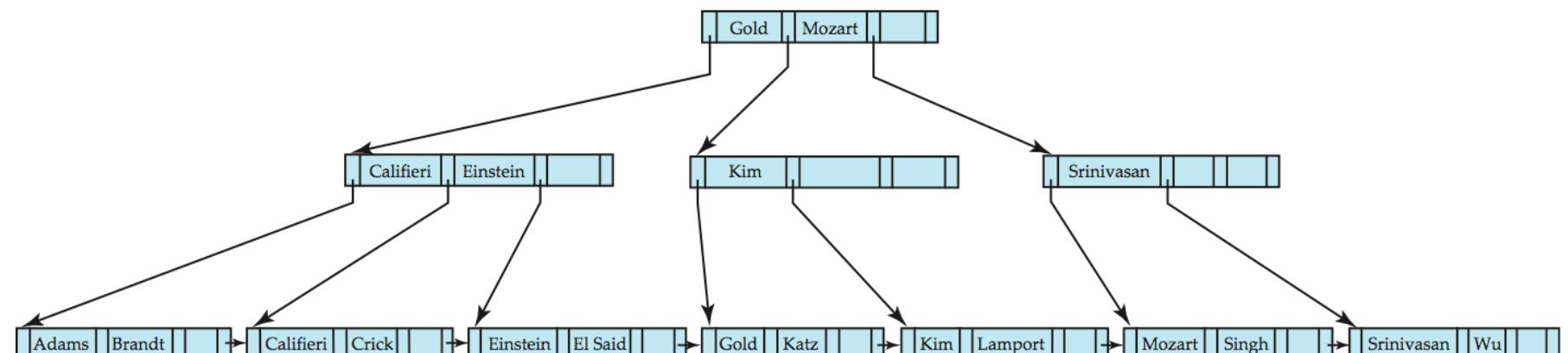
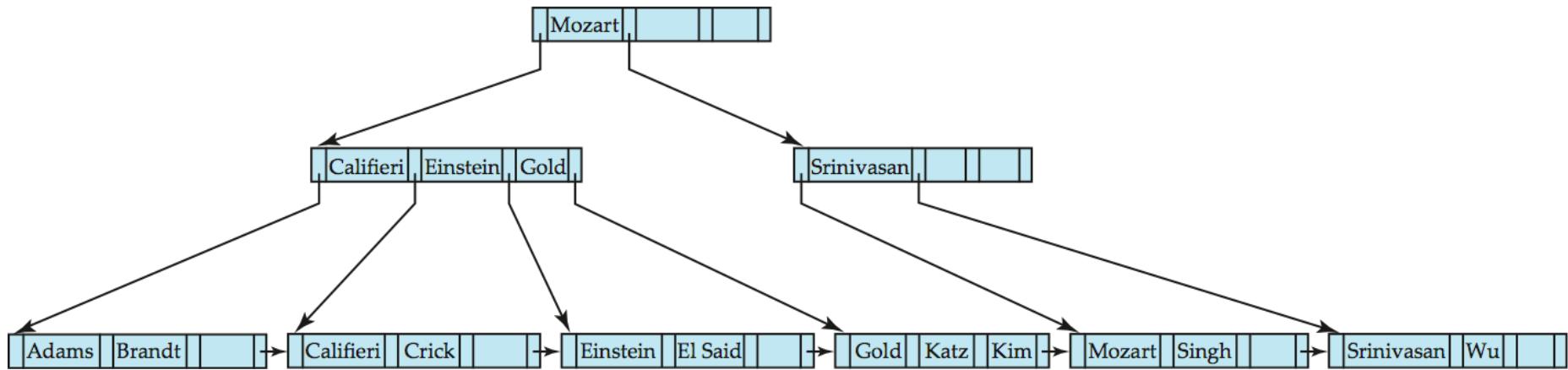
Result of splitting node containing Brandt, Califieri and Crick on inserting Adams  
Next step: insert entry with (Califieri,pointer-to-new-node) into parent

# B+-Tree Insertion



B+-Tree before and after insertion of “Adams”

# B<sup>+</sup>-Tree Insertion

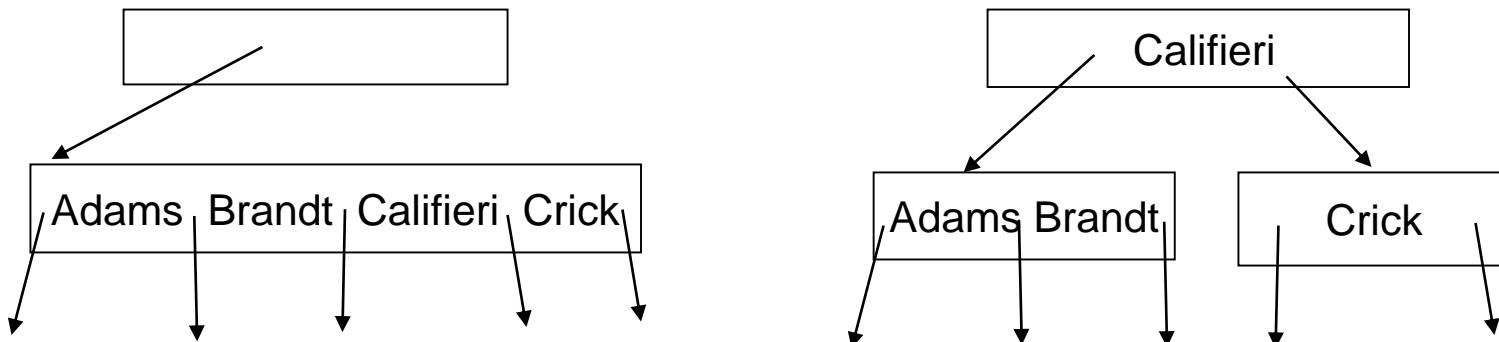


B<sup>+</sup>-Tree before and after insertion of “Lamport”

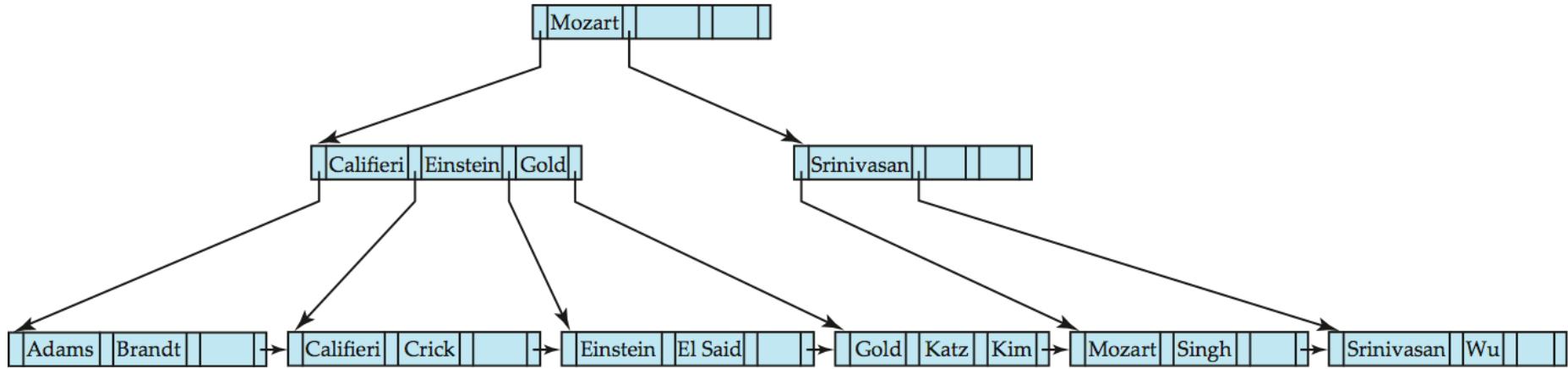
# Insertion in B+-Trees (Cont.)

- Splitting a non-leaf node: when inserting  $(k,p)$  into an already full internal node  $N$ 
  - Copy  $N$  to an in-memory area  $M$  with space for  $n+1$  pointers and  $n$  keys
  - Insert  $(k,p)$  into  $M$
  - Copy  $P_1, K_1, \dots, K_{\lceil n/2 \rceil - 1}, P_{\lceil n/2 \rceil}$  from  $M$  back into node  $N$
  - Copy  $P_{\lceil n/2 \rceil + 1}, K_{\lceil n/2 \rceil + 1}, \dots, K_n, P_{n+1}$  from  $M$  into newly allocated node  $N'$
  - Insert  $(K_{\lceil n/2 \rceil}, N')$  into parent  $N$

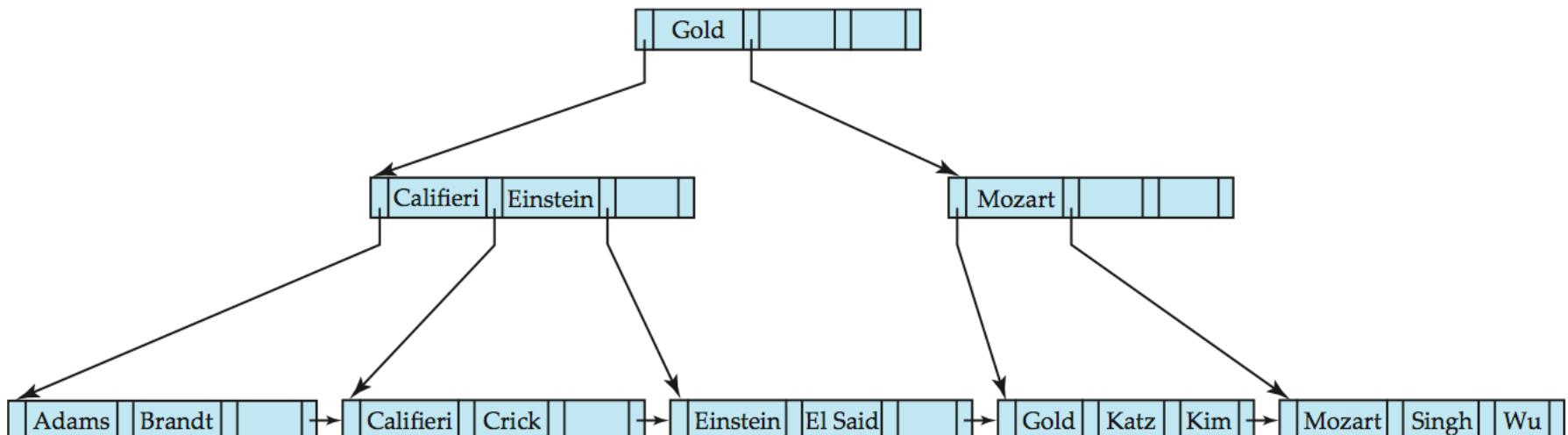
## ■ Read pseudocode in book!



# Examples of B+-Tree Deletion

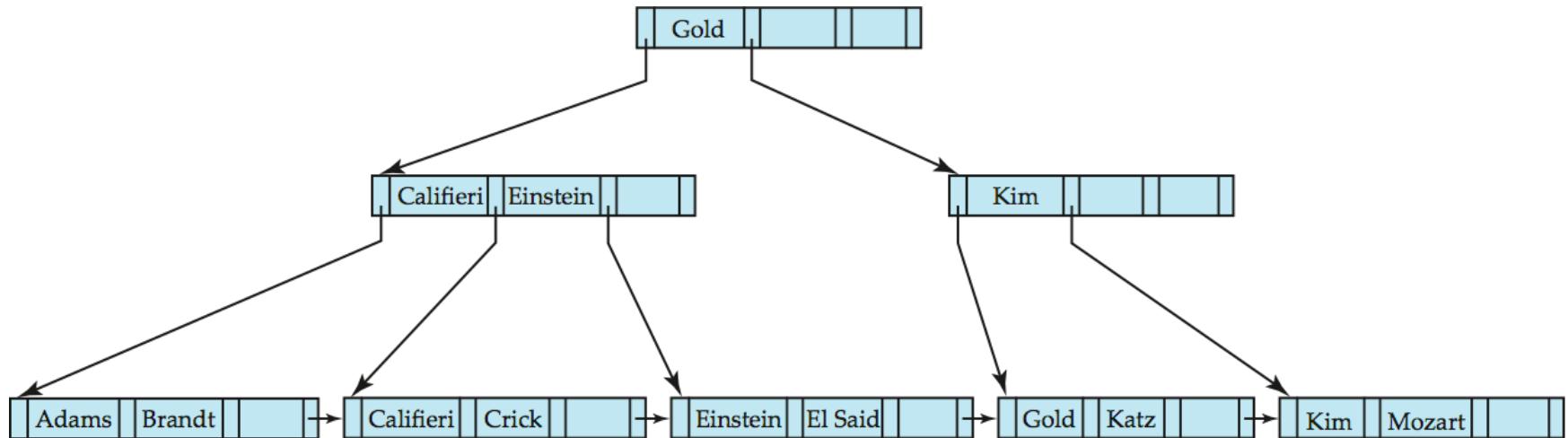


Before and after deleting “Srinivasan”



- Deleting “Srinivasan” causes merging of under-full leaves

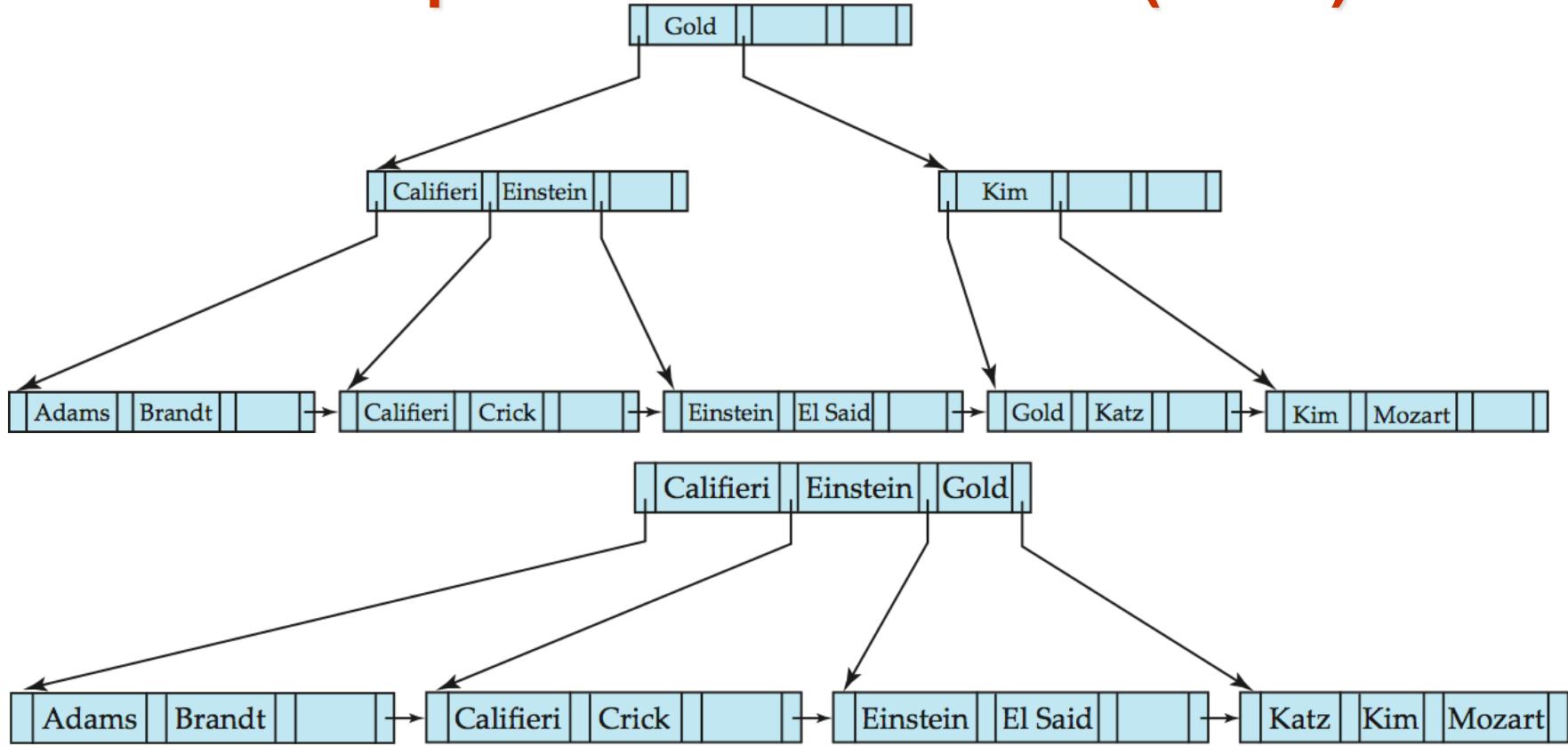
# Examples of B+-Tree Deletion (Cont.)



Deletion of “Singh” and “Wu” from result of previous example

- Leaf containing Singh and Wu became underfull, and borrowed a value Kim from its left sibling
- Search-key value in the parent changes as a result

# Example of B<sup>+</sup>-tree Deletion (Cont.)



## Before and after deletion of “Gold” from earlier example

- Node with Gold and Katz became underfull, and was merged with its sibling
  - Parent node becomes underfull, and is merged with its sibling
    - Value separating two nodes (at the parent) is pulled down when merging
  - Root node then has only one child, and is deleted

# Updates on B+-Trees: Deletion

- Find the record to be deleted, and remove it from the main file and from the bucket (if present)
- Remove (search-key value, pointer) from the leaf node if there is no bucket or if the bucket has become empty
- If the node has too few entries due to the removal, and the entries in the node and a sibling fit into a single node, then ***merge siblings***:
  - Insert all the search-key values in the two nodes into a single node (the one on the left), and delete the other node.
  - Delete the pair  $(K_{i-1}, P_i)$ , where  $P_i$  is the pointer to the deleted node, from its parent, recursively using the above procedure.

# Updates on B+-Trees: Deletion

- Otherwise, if the node has too few entries due to the removal, but the entries in the node and a sibling do not fit into a single node, then **redistribute pointers**:
  - Redistribute the pointers between the node and a sibling such that both have more than the minimum number of entries.
  - Update the corresponding search-key value in the parent of the node.
- The node deletions may cascade upwards till a node which has  $\lceil n/2 \rceil$  or more pointers is found.
- If the root node has only one pointer after deletion, it is deleted and the sole child becomes the root.

# Non-Unique Search Keys

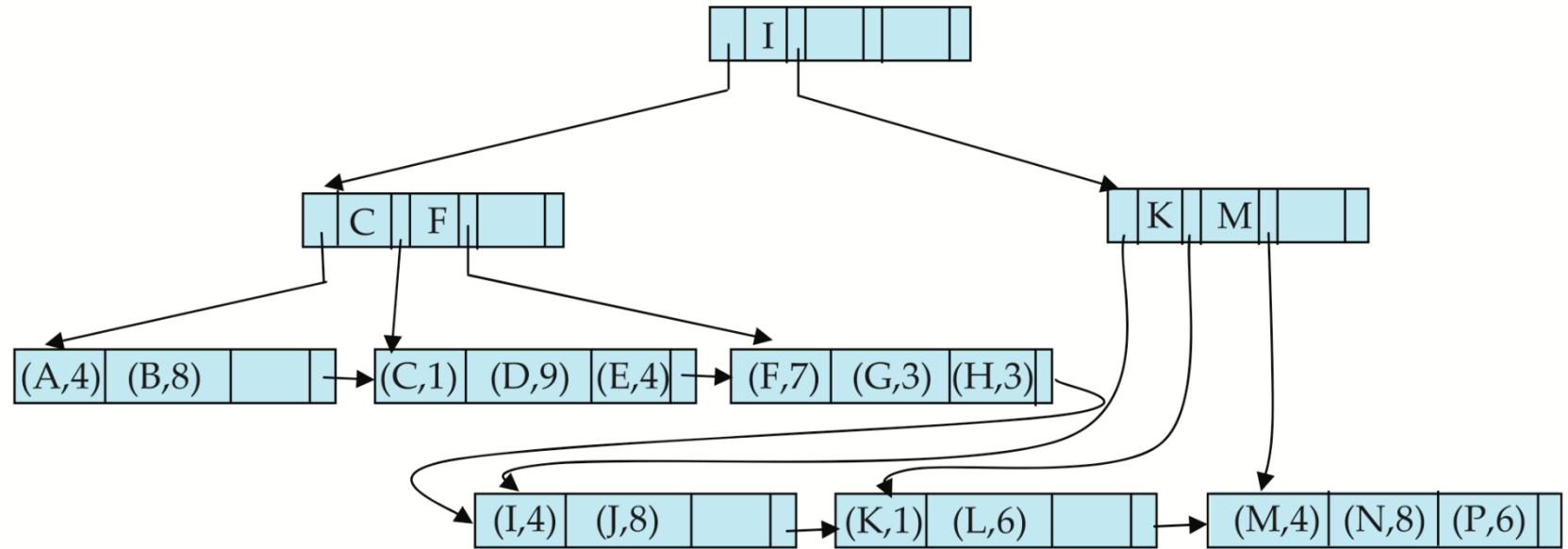
## ■ Alternatives to scheme described earlier

- Buckets on separate block (bad idea)
- List of tuple pointers with each key
  - ▶ Extra code to handle long lists
  - ▶ Deletion of a tuple can be expensive if there are many duplicates on search key (why?)
  - ▶ Low space overhead, no extra cost for queries
- Make search key unique by adding a record-identifier
  - ▶ Extra storage overhead for keys
  - ▶ Simpler code for insertion/deletion
  - ▶ Widely used

# B+-Tree File Organization

- Index file degradation problem is solved by using B+-Tree indices.
- Data file degradation problem is solved by using B+-Tree File Organization.
- The leaf nodes in a B+-tree file organization store records, instead of pointers.
- Leaf nodes are still required to be half full
  - Since records are larger than pointers, the maximum number of records that can be stored in a leaf node is less than the number of pointers in a nonleaf node.
- Insertion and deletion are handled in the same way as insertion and deletion of entries in a B+-tree index.

# B<sup>+</sup>-Tree File Organization (Cont.)



Example of B<sup>+</sup>-tree File Organization

- Good space utilization important since records use more space than pointers.
- To improve space utilization, involve more sibling nodes in redistribution during splits and merges
  - Involving 2 siblings in redistribution (to avoid split / merge where possible) results in each node having at least  $\lfloor \frac{2n}{3} \rfloor$  entries

# Other Issues in Indexing

## ■ Record relocation and secondary indices

- If a record moves, all secondary indices that store record pointers have to be updated
- Node splits in B<sup>+</sup>-tree file organizations become very expensive
- *Solution:* use primary-index search key instead of record pointer in secondary index
  - ▶ Extra traversal of primary index to locate record
    - Higher cost for queries, but node splits are cheap
  - ▶ Add record-id if primary-index search key is non-unique

# Indexing Strings

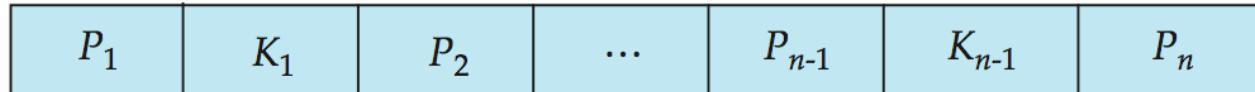
- Variable length strings as keys
  - Variable fanout
  - Use space utilization as criterion for splitting, not number of pointers
- Prefix compression
  - Key values at internal nodes can be prefixes of full key
    - ▶ Keep enough characters to distinguish entries in the subtrees separated by the key value
      - E.g. “Silas” and “Silberschatz” can be separated by “Silb”
  - Keys in leaf node can be compressed by sharing common prefixes

# Bulk Loading and Bottom-Up Build

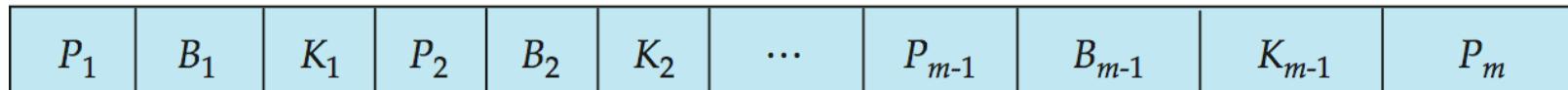
- Inserting entries one-at-a-time into a B<sup>+</sup>-tree requires  $\geq 1$  IO per entry
  - assuming leaf level does not fit in memory
  - can be very inefficient for loading a large number of entries at a time (**bulk loading**)
- Efficient alternative 1:
  - sort entries first (using efficient external-memory sort algorithms discussed later in Section 12.4)
  - insert in sorted order
    - ▶ insertion will go to existing page (or cause a split)
    - ▶ much improved IO performance, but most leaf nodes half full
- Efficient alternative 2: **Bottom-up B<sup>+</sup>-tree construction**
  - As before sort entries
  - And then create tree layer-by-layer, starting with leaf level
    - ▶ details as an exercise
  - Implemented as part of bulk-load utility by most database systems

# B-Tree Index Files

- Similar to B+-tree, but B-tree allows search-key values to appear only once; eliminates redundant storage of search keys.
- Search keys in nonleaf nodes appear nowhere else in the B-tree; an additional pointer field for each search key in a nonleaf node must be included.
- Generalized B-tree leaf node



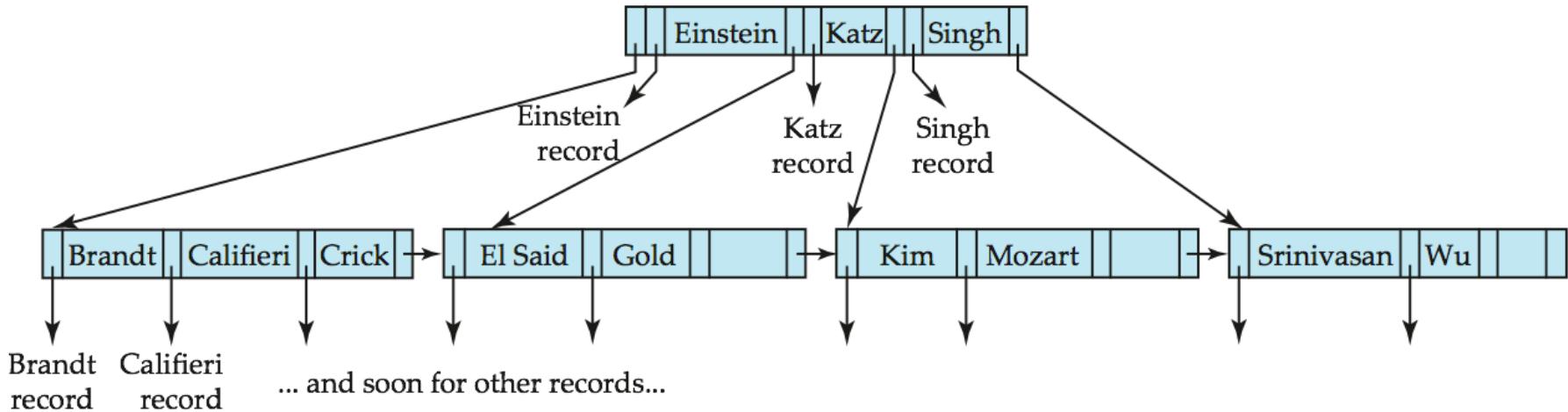
(a)



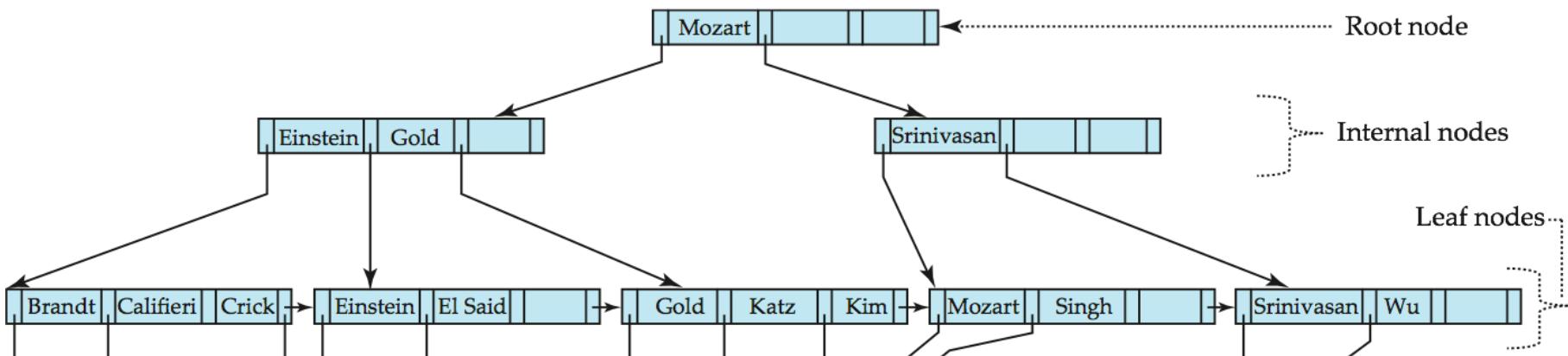
(b)

- Nonleaf node – pointers Bi are the bucket or file record pointers.

# B-Tree Index File Example



B-tree (above) and B+-tree (below) on same data



# B-Tree Index Files (Cont.)

## ■ Advantages of B-Tree indices:

- May use less tree nodes than a corresponding B<sup>+</sup>-Tree.
- Sometimes possible to find search-key value before reaching leaf node.

## ■ Disadvantages of B-Tree indices:

- Only small fraction of all search-key values are found early
- Non-leaf nodes are larger, so fan-out is reduced. Thus, B-Trees typically have greater depth than corresponding B<sup>+</sup>-Tree
- Insertion and deletion more complicated than in B<sup>+</sup>-Trees
- Implementation is harder than B<sup>+</sup>-Trees.

## ■ Typically, advantages of B-Trees do not out weigh disadvantages.

# Multiple-Key Access

- Use multiple indices for certain types of queries.
- Example:

**select** *ID*

**from** *instructor*

**where** *dept\_name* = “Finance” **and** *salary* = 80000

- Possible strategies for processing query using indices on single attributes:
  1. Use index on *dept\_name* to find instructors with department name Finance; test *salary* = 80000
  2. Use index on *salary* to find instructors with a salary of \$80000; test *dept\_name* = “Finance”.
  3. Use *dept\_name* index to find pointers to all records pertaining to the “Finance” department. Similarly use index on *salary*. Take intersection of both sets of pointers obtained.

# Indices on Multiple Keys

- **Composite search keys** are search keys containing more than one attribute
  - E.g.  $(dept\_name, salary)$
- Lexicographic ordering:  $(a_1, a_2) < (b_1, b_2)$  if either
  - $a_1 < b_1$ , or
  - $a_1 = b_1$  and  $a_2 < b_2$

# Indices on Multiple Attributes

Suppose we have an index on combined search-key  
 $(dept\_name, salary)$ .

■ With the **where** clause

**where**  $dept\_name = \text{"Finance"} \text{ and } salary = 80000$

the index on  $(dept\_name, salary)$  can be used to fetch only records that satisfy both conditions.

- Using separate indices is less efficient — we may fetch many records (or pointers) that satisfy only one of the conditions.

■ Can also efficiently handle

**where**  $dept\_name = \text{"Finance"} \text{ and } salary < 80000$

■ But cannot efficiently handle

**where**  $dept\_name < \text{"Finance"} \text{ and } balance = 80000$

- May fetch many records that satisfy the first but not the second condition

# Other Features

## ■ Covering indices

- Add extra attributes to index so (some) queries can avoid fetching the actual records
  - ▶ Particularly useful for secondary indices
    - Why?
- Can store extra attributes only at leaf

# Hashing

**Database System Concepts, 6<sup>th</sup> Ed.**

©Silberschatz, Korth and Sudarshan  
See [www.db-book.com](http://www.db-book.com) for conditions on re-use

# Static Hashing

- A **bucket** is a unit of storage containing one or more records (a bucket is typically a disk block).
- In a **hash file organization** we obtain the bucket of a record directly from its search-key value using a **hash function**.
- Hash function  $h$  is a function from the set of all search-key values  $K$  to the set of all bucket addresses  $B$ .
- Hash function is used to locate records for access, insertion as well as deletion.
- Records with different search-key values may be mapped to the same bucket; thus entire bucket has to be searched sequentially to locate a record.

# Example of Hash File Organization

Hash file organization of *instructor* file, using *dept\_name* as key  
(See figure in next slide.)

- There are 10 buckets,
- The binary representation of the  $i$ th character is assumed to be the integer  $i$ .
- The hash function returns the sum of the binary representations of the characters modulo 10
  - E.g.  $h(\text{Music}) = 1 \quad h(\text{History}) = 2$   
 $h(\text{Physics}) = 3 \quad h(\text{Elec. Eng.}) = 3$

# Example of Hash File Organization

bucket 0


bucket 1

15151	Mozart	Music	40000

bucket 2

32343	El Said	History	80000
58583	Califieri	History	60000

bucket 3

22222	Einstein	Physics	95000
33456	Gold	Physics	87000
98345	Kim	Elec. Eng.	80000

bucket 4

12121	Wu	Finance	90000
76543	Singh	Finance	80000

bucket 5

76766	Crick	Biology	72000

bucket 6

10101	Srinivasan	Comp. Sci.	65000
45565	Katz	Comp. Sci.	75000
83821	Brandt	Comp. Sci.	92000

bucket 7


Hash file organization of *instructor* file, using *dept\_name* as key (see previous slide for details).

# Hash Functions

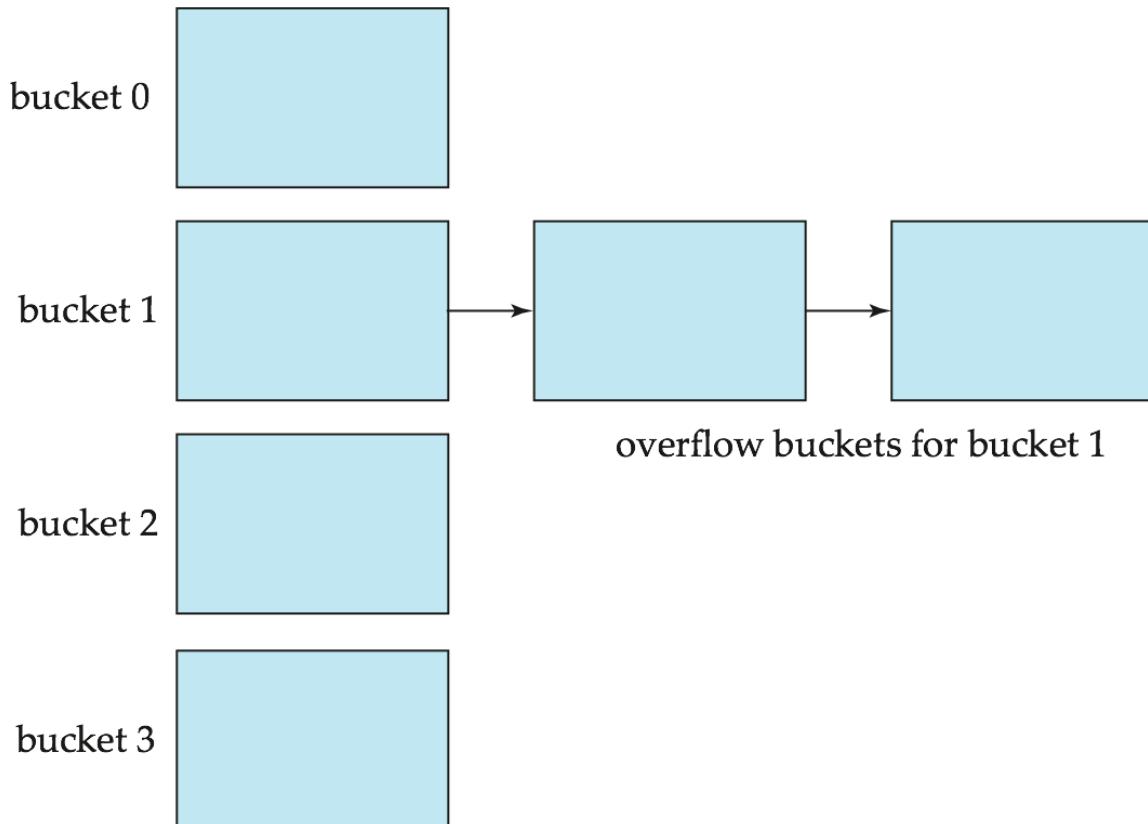
- Worst hash function maps all search-key values to the same bucket; this makes access time proportional to the number of search-key values in the file.
- An ideal hash function is **uniform**, i.e., each bucket is assigned the same number of search-key values from the set of *all* possible values.
- Ideal hash function is **random**, so each bucket will have the same number of records assigned to it irrespective of the *actual distribution* of search-key values in the file.
- Typical hash functions perform computation on the internal binary representation of the search-key.
  - For example, for a string search-key, the binary representations of all the characters in the string could be added and the sum modulo the number of buckets could be returned. .

# Handling of Bucket Overflows

- Bucket overflow can occur because of
  - Insufficient buckets
  - Skew in distribution of records. This can occur due to two reasons:
    - ▶ multiple records have same search-key value
    - ▶ chosen hash function produces non-uniform distribution of key values
- Although the probability of bucket overflow can be reduced, it cannot be eliminated; it is handled by using ***overflow buckets***.

# Handling of Bucket Overflows (Cont.)

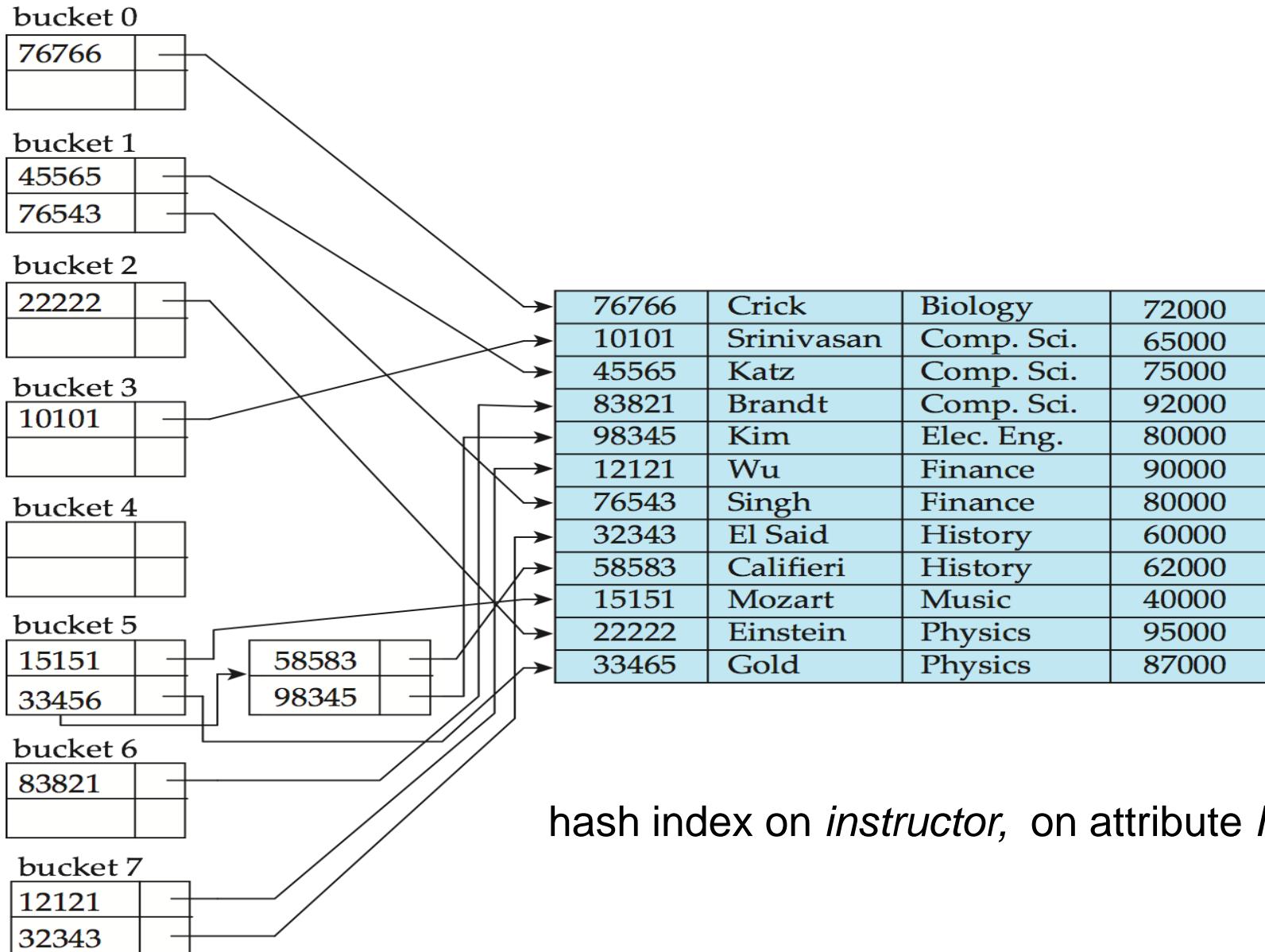
- **Overflow chaining** – the overflow buckets of a given bucket are chained together in a linked list.
- Above scheme is called **closed hashing**.
  - An alternative, called **open hashing**, which does not use overflow buckets, is not suitable for database applications.



# Hash Indices

- Hashing can be used not only for file organization, but also for index-structure creation.
- A **hash index** organizes the search keys, with their associated record pointers, into a hash file structure.
- Strictly speaking, hash indices are always secondary indices
  - if the file itself is organized using hashing, a separate primary hash index on it using the same search-key is unnecessary.
  - However, we use the term hash index to refer to both secondary index structures and hash organized files.

# Example of Hash Index



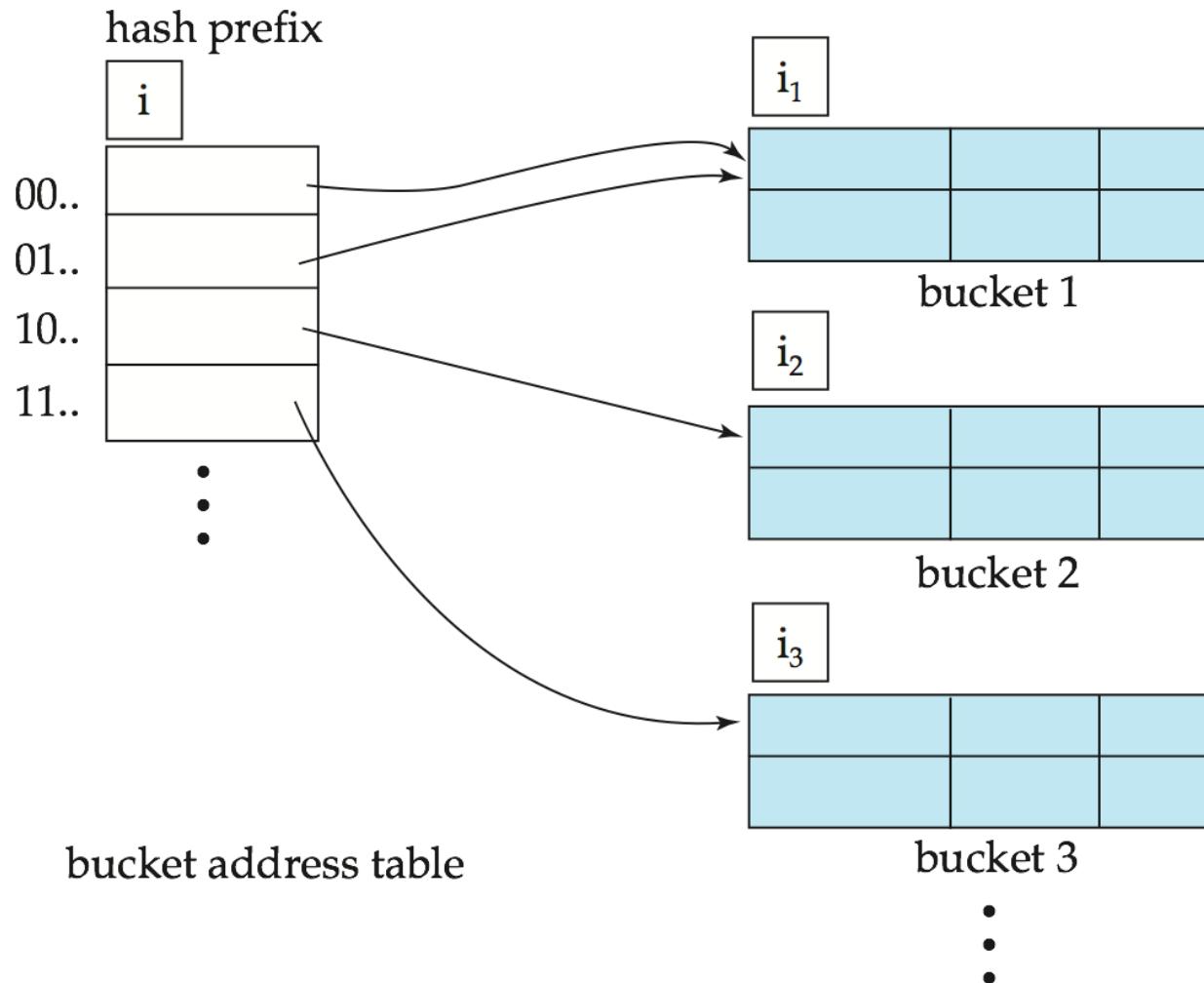
# Deficiencies of Static Hashing

- In static hashing, function  $h$  maps search-key values to a fixed set of  $B$  of bucket addresses. Databases grow or shrink with time.
  - If initial number of buckets is too small, and file grows, performance will degrade due to too much overflows.
  - If space is allocated for anticipated growth, a significant amount of space will be wasted initially (and buckets will be underfull).
  - If database shrinks, again space will be wasted.
- One solution: periodic re-organization of the file with a new hash function
  - Expensive, disrupts normal operations
- Better solution: allow the number of buckets to be modified dynamically.

# Dynamic Hashing

- Good for database that grows and shrinks in size
- Allows the hash function to be modified dynamically
- **Extendable hashing** – one form of dynamic hashing
  - Hash function generates values over a large range — typically  $b$ -bit integers, with  $b = 32$ .
  - At any time use only a prefix of the hash function to index into a table of bucket addresses.
  - Let the length of the prefix be  $i$  bits,  $0 \leq i \leq 32$ .
    - ▶ Bucket address table size =  $2^i$ . Initially  $i = 0$
    - ▶ Value of  $i$  grows and shrinks as the size of the database grows and shrinks.
  - Multiple entries in the bucket address table may point to a bucket (why?)
  - Thus, actual number of buckets is  $< 2^i$ 
    - ▶ The number of buckets also changes dynamically due to coalescing and splitting of buckets.

# General Extendable Hash Structure



In this structure,  $i_2 = i_3 = i$ , whereas  $i_1 = i - 1$  (see next slide for details)

# Use of Extendable Hash Structure

- Each bucket  $j$  stores a value  $i_j$ 
  - All the entries that point to the same bucket have the same values on the first  $i_j$  bits.
- To locate the bucket containing search-key  $K_j$ :
  1. Compute  $h(K_j) = X$
  2. Use the first  $i$  high order bits of  $X$  as a displacement into bucket address table, and follow the pointer to appropriate bucket
- To insert a record with search-key value  $K_j$ 
  - follow same procedure as look-up and locate the bucket, say  $j$ .
  - If there is room in the bucket  $j$  insert record in the bucket.
  - Else the bucket must be split and insertion re-attempted (next slide.)
    - ▶ Overflow buckets used instead in some cases (will see shortly)

# Insertion in Extendable Hash Structure (Cont)

To split a bucket  $j$  when inserting record with search-key value  $K_j$ :

- If  $i > i_j$  (more than one pointer to bucket  $j$ )
  - allocate a new bucket  $z$ , and set  $i_j = i_z = (i_j + 1)$
  - Update the second half of the bucket address table entries originally pointing to  $j$ , to point to  $z$
  - remove each record in bucket  $j$  and reinsert (in  $j$  or  $z$ )
  - recompute new bucket for  $K_j$  and insert record in the bucket (further splitting is required if the bucket is still full)
- If  $i = i_j$  (only one pointer to bucket  $j$ )
  - If  $i$  reaches some limit  $b$ , or too many splits have happened in this insertion, create an overflow bucket
  - Else
    - ▶ increment  $i$  and double the size of the bucket address table.
    - ▶ replace each entry in the table by two entries that point to the same bucket.
    - ▶ recompute new bucket address table entry for  $K_j$   
Now  $i > i_j$  so use the first case above.

# Deletion in Extendable Hash Structure

- To delete a key value,
  - locate it in its bucket and remove it.
  - The bucket itself can be removed if it becomes empty (with appropriate updates to the bucket address table).
  - Coalescing of buckets can be done (can coalesce only with a “*buddy*” bucket having same value of  $i_j$  and same  $i_j - 1$  prefix, if it is present)
  - Decreasing bucket address table size is also possible
    - ▶ Note: decreasing bucket address table size is an expensive operation and should be done only if number of buckets becomes much smaller than the size of the table

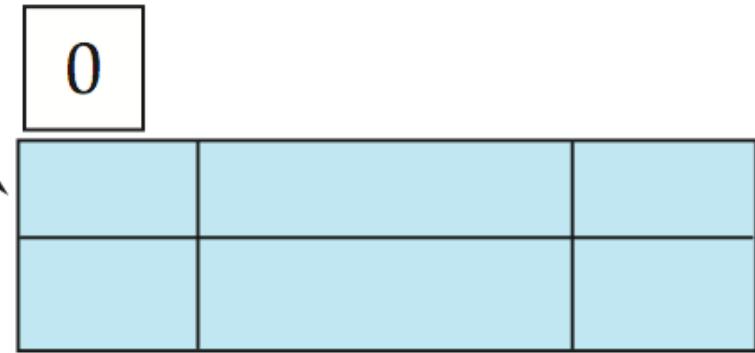
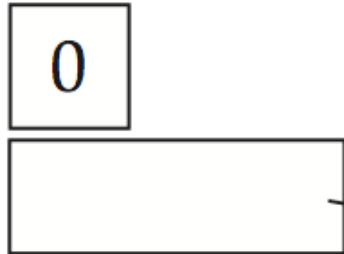
# Use of Extendable Hash Structure: Example

$dept\_name$	$h(dept\_name)$
Biology	0010 1101 1111 1011 0010 1100 0011 0000
Comp. Sci.	1111 0001 0010 0100 1001 0011 0110 1101
Elec. Eng.	0100 0011 1010 1100 1100 0110 1101 1111
Finance	1010 0011 1010 0000 1100 0110 1001 1111
History	1100 0111 1110 1101 1011 1111 0011 1010
Music	0011 0101 1010 0110 1100 1001 1110 1011
Physics	1001 1000 0011 1111 1001 1100 0000 0001

## Example (Cont.)

- Initial Hash structure; bucket size = 2

hash prefix



bucket address table

bucket 1

## Example (Cont.)

- Hash structure after insertion of “Mozart”, “Srinivasan”, and “Wu” records

hash prefix

1

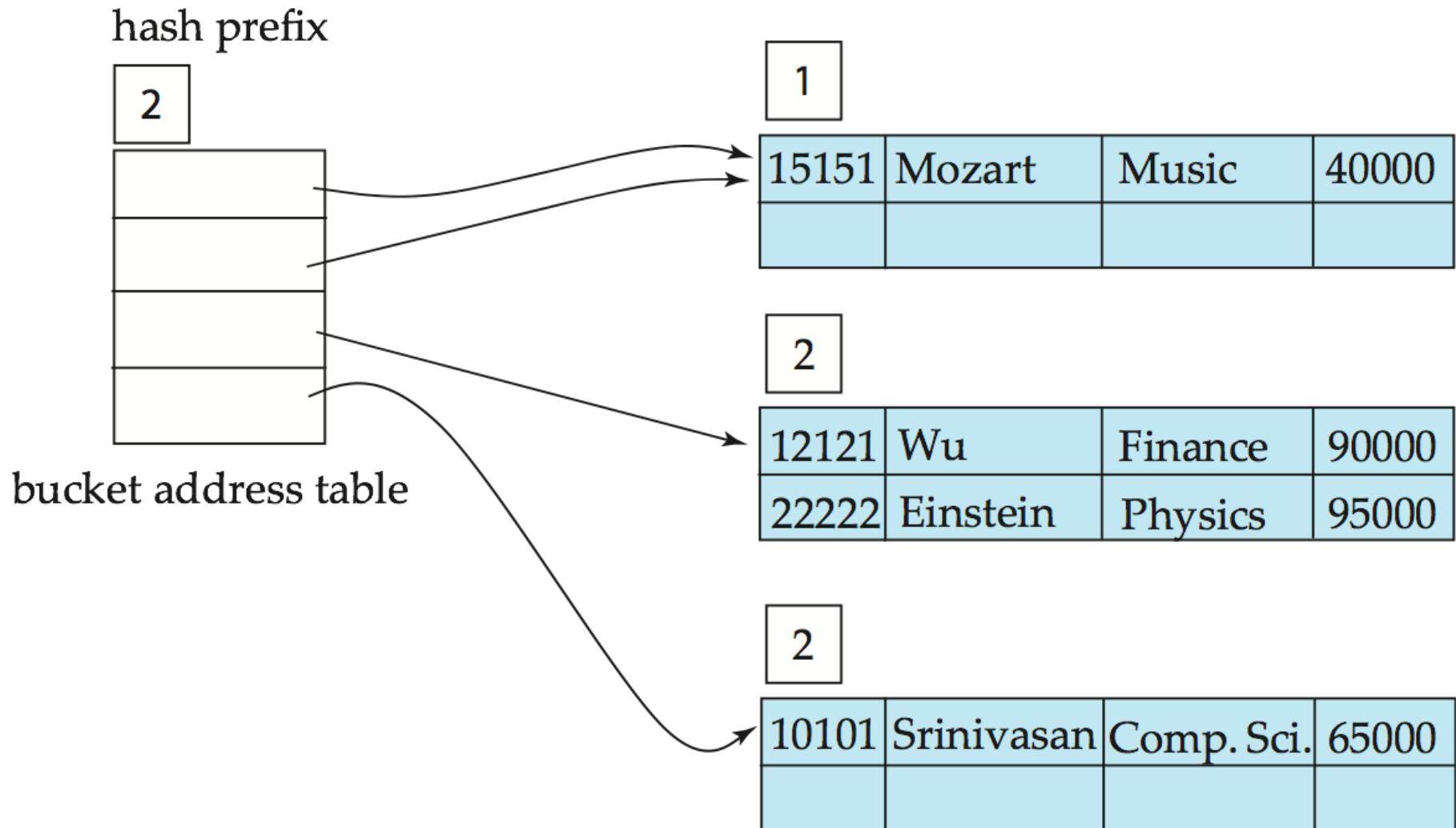
1	15151	Mozart	Music	40000

bucket address table

1	10101	Srinivasan	Comp. Sci.	90000
	12121	Wu	Finance	90000

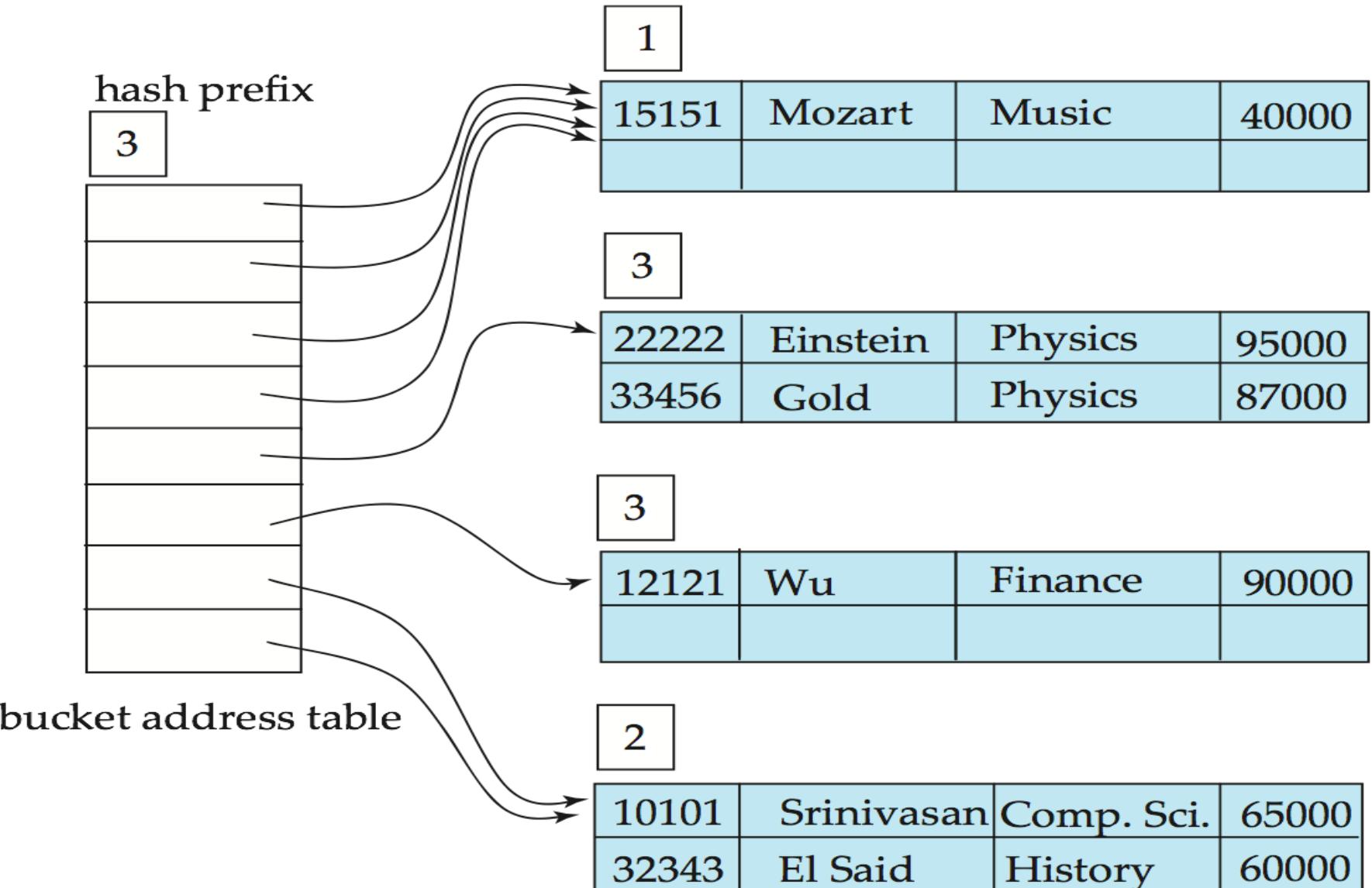
# Example (Cont.)

## ■ Hash structure after insertion of Einstein record



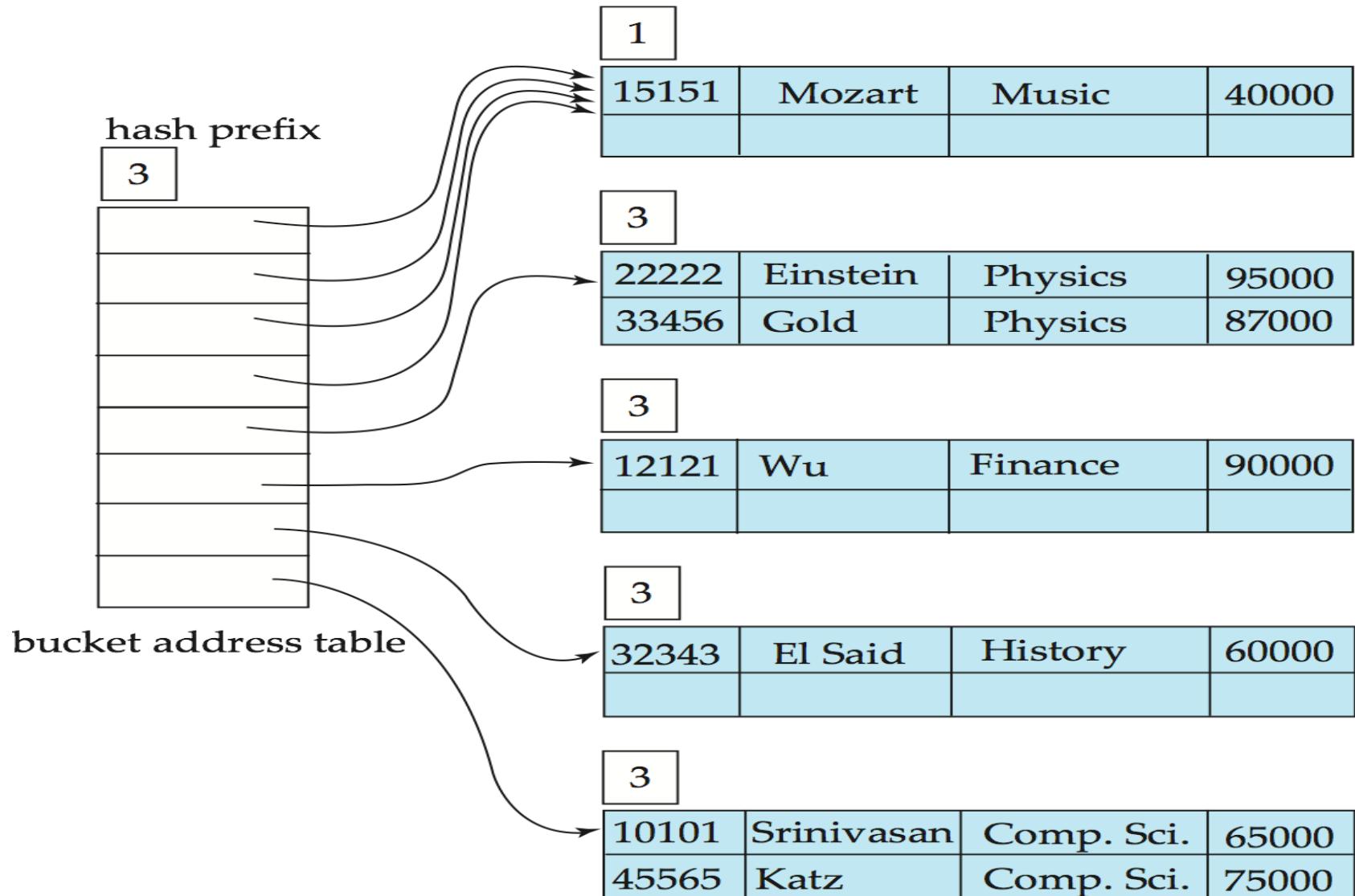
# Example (Cont.)

## ■ Hash structure after insertion of Gold and El Said records

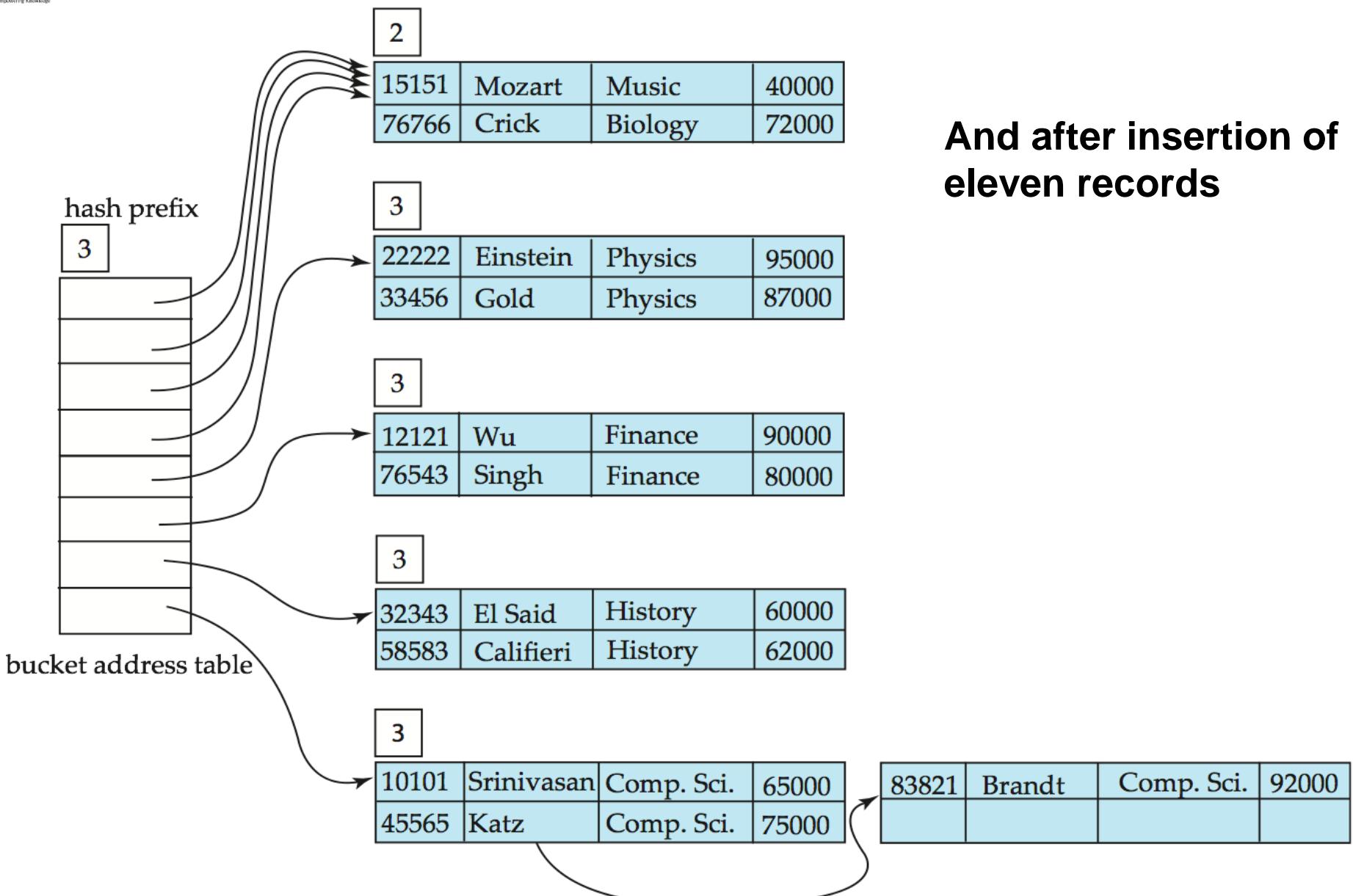


# Example (Cont.)

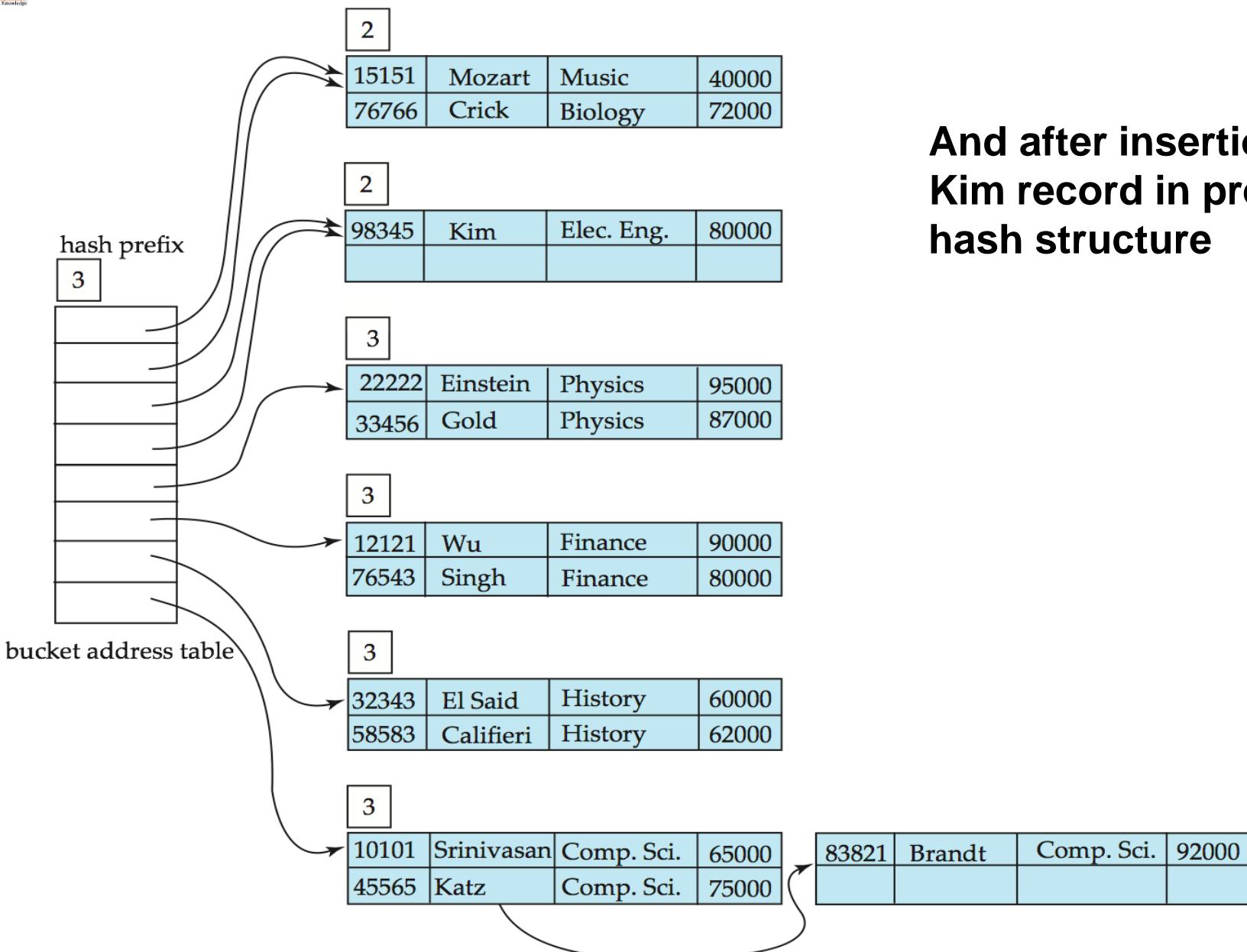
## ■ Hash structure after insertion of Katz record



# Example (Cont.)



# Example (Cont.)



# Extendable Hashing vs. Other Schemes

## ■ Benefits of extendable hashing:

- Hash performance does not degrade with growth of file
- Minimal space overhead

## ■ Disadvantages of extendable hashing

- Extra level of indirection to find desired record
  - Bucket address table may itself become very big (larger than memory)
    - ▶ Cannot allocate very large contiguous areas on disk either
    - ▶ Solution: B<sup>+</sup>-tree structure to locate desired record in bucket address table
  - Changing size of bucket address table is an expensive operation
- ## ■ Linear hashing is an alternative mechanism
- Allows incremental growth of its directory (equivalent to bucket address table)
  - At the cost of more bucket overflows

# Comparison of Ordered Indexing and Hashing

- Cost of periodic re-organization
- Relative frequency of insertions and deletions
- Is it desirable to optimize average access time at the expense of worst-case access time?
- Expected type of queries:
  - Hashing is generally better at retrieving records having a specified value of the key.
  - If range queries are common, ordered indices are to be preferred
- In practice:
  - PostgreSQL supports hash indices, but discourages use due to poor performance
  - Oracle supports static hash organization, but not hash indices
  - SQLServer supports only B<sup>+</sup>-trees

# Bitmap Indices

- Bitmap indices are a special type of index designed for efficient querying on multiple keys
- Records in a relation are assumed to be numbered sequentially from, say, 0
  - Given a number  $n$  it must be easy to retrieve record  $n$ 
    - ▶ Particularly easy if records are of fixed size
- Applicable on attributes that take on a relatively small number of distinct values
  - E.g. gender, country, state, ...
  - E.g. income-level (income broken up into a small number of levels such as 0-9999, 10000-19999, 20000-50000, 50000- infinity)
- A bitmap is simply an array of bits

# Bitmap Indices (Cont.)

- In its simplest form a bitmap index on an attribute has a bitmap for each value of the attribute
  - Bitmap has as many bits as records
  - In a bitmap for value  $v$ , the bit for a record is 1 if the record has the value  $v$  for the attribute, and is 0 otherwise

record number	<i>ID</i>	<i>gender</i>	<i>income_level</i>
0	76766	m	L1
1	22222	f	L2
2	12121	f	L1
3	15151	m	L4
4	58583	f	L3

Bitmaps for *gender*

m	10010
f	01101

Bitmaps for *income\_level*

L1	10100
L2	01000
L3	00001
L4	00010
L5	00000

# Bitmap Indices (Cont.)

- Bitmap indices are useful for queries on multiple attributes
  - not particularly useful for single attribute queries
- Queries are answered using bitmap operations
  - Intersection (and)
  - Union (or)
  - Complementation (not)
- Each operation takes two bitmaps of the same size and applies the operation on corresponding bits to get the result bitmap
  - E.g.  $100110 \text{ AND } 110011 = 100010$   
 $100110 \text{ OR } 110011 = 110111$   
 $\text{NOT } 100110 = 011001$
  - Males with income level L1:  $10010 \text{ AND } 10100 = 10000$ 
    - ▶ Can then retrieve required tuples.
    - ▶ Counting number of matching tuples is even faster

# Bitmap Indices (Cont.)

- Bitmap indices generally very small compared with relation size
  - E.g. if record is 100 bytes, space for a single bitmap is 1/800 of space used by relation.
    - ▶ If number of distinct attribute values is 8, bitmap is only 1% of relation size
- Deletion needs to be handled properly
  - **Existence bitmap** to note if there is a valid record at a record location
  - Needed for complementation
    - ▶  $\text{not}(A=v)$ :  $(\text{NOT } \text{bitmap-}A-v) \text{ AND ExistenceBitmap}$
- Should keep bitmaps for all values, even null value
  - To correctly handle SQL null semantics for  $\text{NOT}(A=v)$ :
    - ▶ intersect above result with  $(\text{NOT } \text{bitmap-}A-\text{Null})$

# Efficient Implementation of Bitmap Operations

- Bitmaps are packed into words; a single word and (a basic CPU instruction) computes and of 32 or 64 bits at once
  - E.g. 1-million-bit maps can be and-ed with just 31,250 instruction
- Counting number of 1s can be done fast by a trick:
  - Use each byte to index into a precomputed array of 256 elements each storing the count of 1s in the binary representation
    - ▶ Can use pairs of bytes to speed up further at a higher memory cost
  - Add up the retrieved counts
- Bitmaps can be used instead of Tuple-ID lists at leaf levels of B<sup>+</sup>-trees, for values that have a large number of matching records
  - Worthwhile if > 1/64 of the records have that value, assuming a tuple-id is 64 bits
  - Above technique merges benefits of bitmap and B<sup>+</sup>-tree indices

# Index Definition in SQL

## ■ Create an index

```
create index <index-name> on <relation-name>  
          (<attribute-list>)
```

E.g.: **create index b-index on branch(branch\_name)**

## ■ Use **create unique index** to indirectly specify and enforce the condition that the search key is a candidate key is a candidate key.

- Not really required if SQL **unique** integrity constraint is supported

## ■ To drop an index

```
drop index <index-name>
```

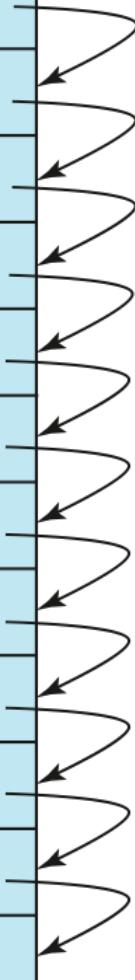
## ■ Most database systems allow specification of type of index, and clustering.

# End of Chapter

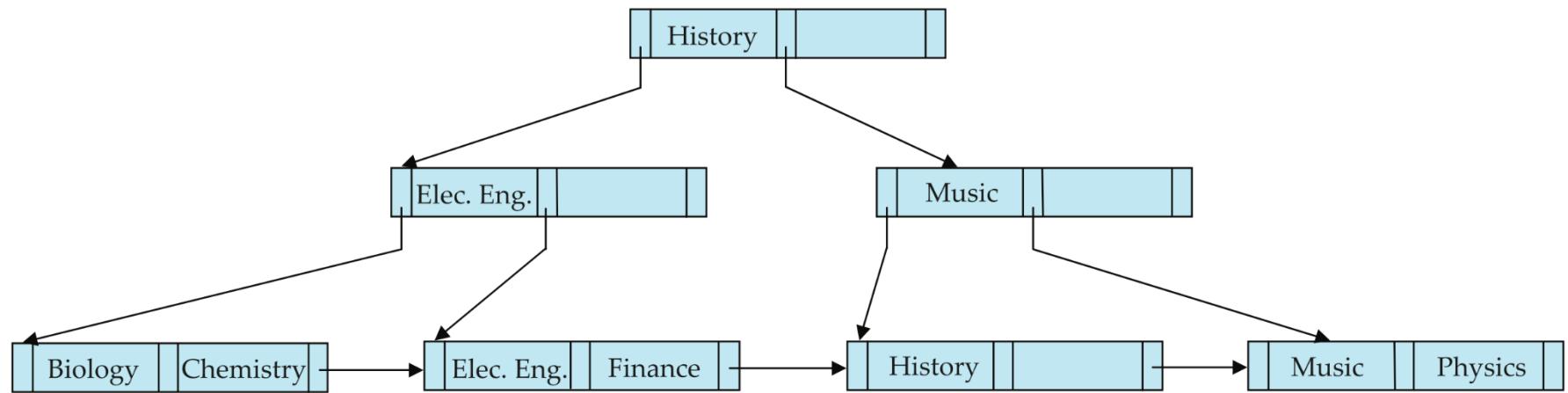
**Database System Concepts, 6<sup>th</sup> Ed.**

©Silberschatz, Korth and Sudarshan  
See [www.db-book.com](http://www.db-book.com) for conditions on re-use

# Figure 11.01

10101	Srinivasan	Comp. Sci.	65000		
12121	Wu	Finance	90000		
15151	Mozart	Music	40000		
22222	Einstein	Physics	95000		
32343	El Said	History	60000		
33456	Gold	Physics	87000		
45565	Katz	Comp. Sci.	75000		
58583	Califieri	History	62000		
76543	Singh	Finance	80000		
76766	Crick	Biology	72000		
83821	Brandt	Comp. Sci.	92000		
98345	Kim	Elec. Eng.	80000		

# Figure 11.15



# Partitioned Hashing

- Hash values are split into segments that depend on each attribute of the search-key.

$(A_1, A_2, \dots, A_n)$  for  $n$  attribute search-key

- Example:  $n = 2$ , for *customer*, search-key being (*customer-street*, *customer-city*)

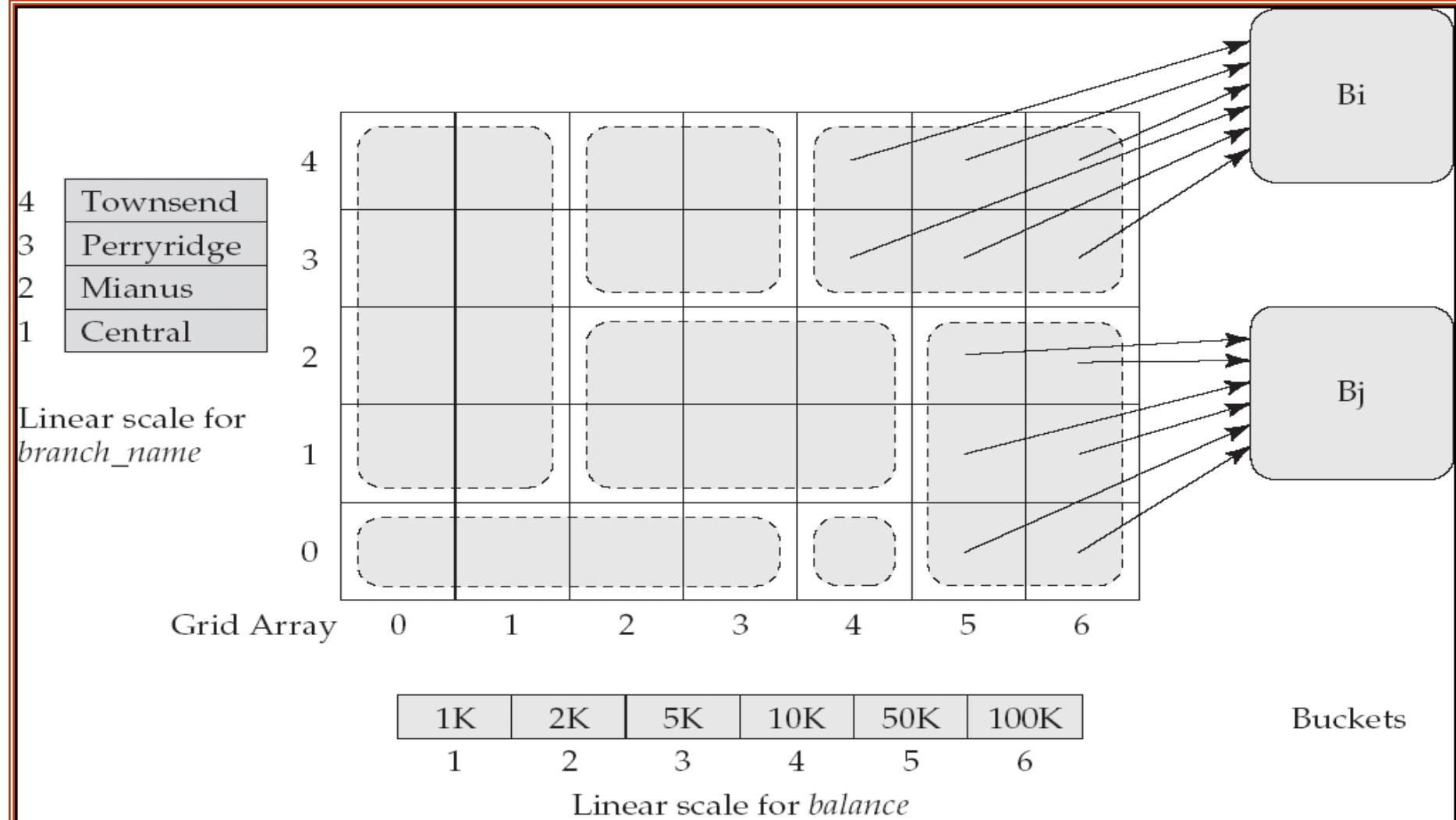
search-key value	hash value
(Main, Harrison)	101 111
(Main, Brooklyn)	101 001
(Park, Palo Alto)	010 010
(Spring, Brooklyn)	001 001
(Alma, Palo Alto)	110 010

- To answer equality query on single attribute, need to look up multiple buckets. Similar in effect to grid files.

# Grid Files

- Structure used to speed the processing of general multiple search-key queries involving one or more comparison operators.
- The **grid file** has a single grid array and one linear scale for each search-key attribute. The grid array has number of dimensions equal to number of search-key attributes.
- Multiple cells of grid array can point to same bucket
- To find the bucket for a search-key value, locate the row and column of its cell using the linear scales and follow pointer

# Example Grid File for account



# Queries on a Grid File

- A grid file on two attributes  $A$  and  $B$  can handle queries of all following forms with reasonable efficiency
  - $(a_1 \leq A \leq a_2)$
  - $(b_1 \leq B \leq b_2)$
  - $(a_1 \leq A \leq a_2 \wedge b_1 \leq B \leq b_2),..$
- E.g., to answer  $(a_1 \leq A \leq a_2 \wedge b_1 \leq B \leq b_2)$ , use linear scales to find corresponding candidate grid array cells, and look up all the buckets pointed to from those cells.

# Grid Files (Cont.)

- During insertion, if a bucket becomes full, new bucket can be created if more than one cell points to it.
  - Idea similar to extendable hashing, but on multiple dimensions
  - If only one cell points to it, either an overflow bucket must be created or the grid size must be increased
- Linear scales must be chosen to uniformly distribute records across cells.
  - Otherwise there will be too many overflow buckets.
- Periodic re-organization to increase grid size will help.
  - But reorganization can be very expensive.
- Space overhead of grid array can be high.
- R-trees (Chapter 23) are an alternative

# Chapter 12: Query Processing

Database System Concepts, 6<sup>th</sup> Ed.

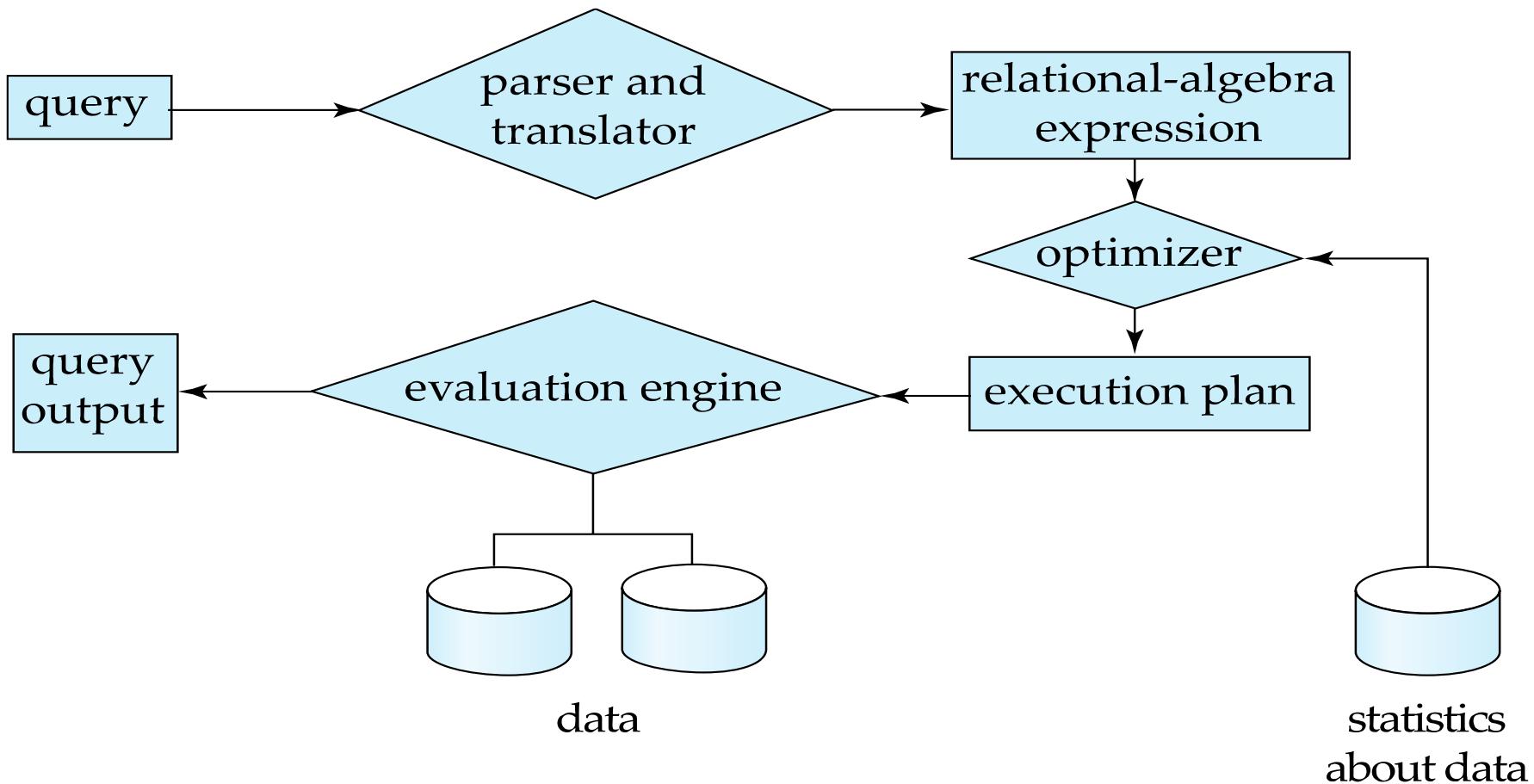
©Silberschatz, Korth and Sudarshan  
See [www.db-book.com](http://www.db-book.com) for conditions on re-use

# Chapter 12: Query Processing

- Overview
- Measures of Query Cost
- Selection Operation
- Sorting
- Join Operation
- Other Operations
- Evaluation of Expressions

# Basic Steps in Query Processing

1. Parsing and translation
2. Optimization
3. Evaluation



# Basic Steps in Query Processing (Cont.)

## ■ Parsing and translation

- translate the query into its internal form. This is then translated into relational algebra.
- Parser checks syntax, verifies relations

## ■ Evaluation

- The query-execution engine takes a query-evaluation plan, executes that plan, and returns the answers to the query.

# Basic Steps in Query Processing : Optimization

- A relational algebra expression may have many equivalent expressions
  - E.g.,  $\sigma_{\text{salary} < 75000}(\Pi_{\text{salary}}(\text{instructor}))$  is equivalent to
$$\Pi_{\text{salary}}(\sigma_{\text{salary} < 75000}(\text{instructor}))$$
- Each relational algebra operation can be evaluated using one of several different algorithms
  - Correspondingly, a relational-algebra expression can be evaluated in many ways.
- Annotated expression specifying detailed evaluation strategy is called an **evaluation-plan**.
  - E.g., can use an index on *salary* to find instructors with  $\text{salary} < 75000$ ,
  - or can perform complete relation scan and discard instructors with  $\text{salary} \geq 75000$

# Basic Steps: Optimization (Cont.)

- **Query Optimization:** Amongst all equivalent evaluation plans choose the one with lowest cost.
  - Cost is estimated using statistical information from the database catalog
    - ▶ e.g. number of tuples in each relation, size of tuples, etc.
- In this chapter we study
  - How to measure query costs
  - Algorithms for evaluating relational algebra operations
  - How to combine algorithms for individual operations in order to evaluate a complete expression
- In Chapter 14
  - We study how to optimize queries, that is, how to find an evaluation plan with lowest estimated cost

# Measures of Query Cost

- Cost is generally measured as total elapsed time for answering query
  - Many factors contribute to time cost
    - ▶ *disk accesses, CPU, or even network communication*
- Typically disk access is the predominant cost, and is also relatively easy to estimate. Measured by taking into account
  - Number of seeks \* average-seek-cost
  - Number of blocks read \* average-block-read-cost
  - Number of blocks written \* average-block-write-cost
    - ▶ Cost to write a block is greater than cost to read a block
      - data is read back after being written to ensure that the write was successful

# Measures of Query Cost (Cont.)

- For simplicity we just use the **number of block transfers** from disk and the **number of seeks** as the cost measures
  - $t_T$  – time to transfer one block
  - $t_S$  – time for one seek
  - Cost for b block transfers plus S seeks
$$b * t_T + S * t_S$$
- We ignore CPU costs for simplicity
  - Real systems do take CPU cost into account
- We do not include cost to writing output to disk in our cost formulae

## Measures of Query Cost (Cont.)

- Several algorithms can reduce disk IO by using extra buffer space
  - Amount of real memory available to buffer depends on other concurrent queries and OS processes, known only during execution
    - ▶ We often use worst case estimates, assuming only the minimum amount of memory needed for the operation is available
- Required data may be buffer resident already, avoiding disk I/O
  - But hard to take into account for cost estimation

# Selection Operation

- **File scan**
- Algorithm **A1 (linear search)**. Scan each file block and test all records to see whether they satisfy the selection condition.
  - Cost estimate =  $b_r$  block transfers + 1 seek
    - ▶  $b_r$  denotes number of blocks containing records from relation  $r$
  - If selection is on a key attribute, can stop on finding record
    - ▶ cost =  $(b_r/2)$  block transfers + 1 seek
  - Linear search can be applied regardless of
    - ▶ selection condition or
    - ▶ ordering of records in the file, or
    - ▶ availability of indices
- Note: binary search generally does not make sense since data is not stored consecutively
  - except when there is an index available,
  - and binary search requires more seeks than index search

# Selections Using Indices

- **Index scan** – search algorithms that use an index
  - selection condition must be on search-key of index.
- **A2 (primary index, equality on key)**. Retrieve a single record that satisfies the corresponding equality condition
  - $\text{Cost} = (h_i + 1) * (t_T + t_S)$
- **A3 (primary index, equality on nonkey)** Retrieve multiple records.
  - Records will be on consecutive blocks
    - ▶ Let b = number of blocks containing matching records
  - $\text{Cost} = h_i * (t_T + t_S) + t_S + t_T * b$

# Selections Using Indices

## ■ A4 (secondary index, equality on nonkey).

- Retrieve a single record if the search-key is a candidate key
  - ▶  $\text{Cost} = (h_i + 1) * (t_T + t_S)$
- Retrieve multiple records if search-key is not a candidate key
  - ▶ each of  $n$  matching records may be on a different block
  - ▶  $\text{Cost} = (h_i + n) * (t_T + t_S)$ 
    - Can be very expensive!

# Selections Involving Comparisons

- Can implement selections of the form  $\sigma_{A \leq v}(r)$  or  $\sigma_{A \geq v}(r)$  by using
  - a linear file scan,
  - or by using indices in the following ways:
- **A5 (primary index, comparison).** (Relation is sorted on A)
  - ▶ For  $\sigma_{A \geq v}(r)$  use index to find first tuple  $\geq v$  and scan relation sequentially from there
  - ▶ For  $\sigma_{A \leq v}(r)$  just scan relation sequentially till first tuple  $> v$ ; do not use index
- **A6 (secondary index, comparison).**
  - ▶ For  $\sigma_{A \geq v}(r)$  use index to find first index entry  $\geq v$  and scan index sequentially from there, to find pointers to records.
  - ▶ For  $\sigma_{A \leq v}(r)$  just scan leaf pages of index finding pointers to records, till first entry  $> v$
  - ▶ In either case, retrieve records that are pointed to
    - requires an I/O for each record
    - Linear file scan may be cheaper

# Implementation of Complex Selections

■ **Conjunction:**  $\sigma_{\theta_1} \wedge \theta_2 \wedge \dots \wedge \theta_n(r)$

■ **A7 (conjunctive selection using one index).**

- Select a combination of  $\theta_i$ , and algorithms A1 through A7 that results in the least cost for  $\sigma_{\theta_i}(r)$ .
- Test other conditions on tuple after fetching it into memory buffer.

■ **A8 (conjunctive selection using composite index).**

- Use appropriate composite (multiple-key) index if available.

■ **A9 (conjunctive selection by intersection of identifiers).**

- Requires indices with record pointers.
- Use corresponding index for each condition, and take intersection of all the obtained sets of record pointers.
- Then fetch records from file
- If some conditions do not have appropriate indices, apply test in memory.

# Algorithms for Complex Selections

- **Disjunction:**  $\sigma_{\theta_1 \vee \theta_2 \vee \dots \vee \theta_n}(r)$ .
- **A10 (disjunctive selection by union of identifiers).**
  - Applicable if *all* conditions have available indices.
    - ▶ Otherwise use linear scan.
  - Use corresponding index for each condition, and take union of all the obtained sets of record pointers.
  - Then fetch records from file
- **Negation:**  $\sigma_{\neg\theta}(r)$ 
  - Use linear scan on file
  - If very few records satisfy  $\neg\theta$ , and an index is applicable to  $\theta$ 
    - ▶ Find satisfying records using index and fetch from file

# Sorting

- We may build an index on the relation, and then use the index to read the relation in sorted order. May lead to one disk block access for each tuple.
- For relations that fit in memory, techniques like quicksort can be used. For relations that don't fit in memory, **external sort-merge** is a good choice.

# External Sort-Merge

Let  $M$  denote memory size (in pages).

1. **Create sorted runs.** Let  $i$  be 0 initially.

Repeatedly do the following till the end of the relation:

- (a) Read  $M$  blocks of relation into memory
- (b) Sort the in-memory blocks
- (c) Write sorted data to run  $R_i$ ; increment  $i$ .

Let the final value of  $i$  be  $N$

2. *Merge the runs (next slide).....*

# External Sort-Merge (Cont.)

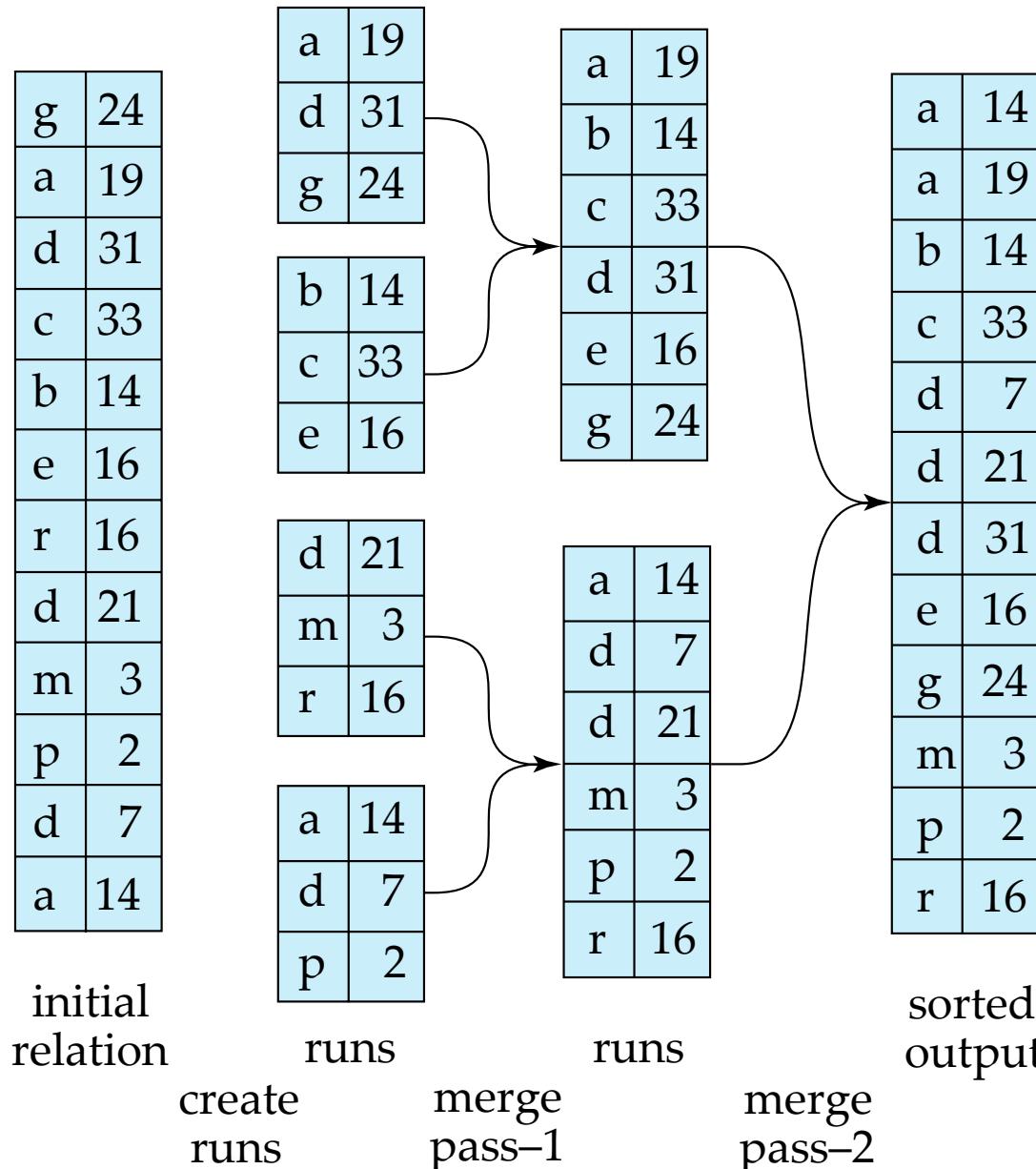
**2. Merge the runs (N-way merge).** We assume (for now) that  $N < M$ .

1. Use  $N$  blocks of memory to buffer input runs, and 1 block to buffer output. Read the first block of each run into its buffer page
2. **repeat**
  1. Select the first record (in sort order) among all buffer pages
  2. Write the record to the output buffer. If the output buffer is full write it to disk.
  3. Delete the record from its input buffer page.  
**If** the buffer page becomes empty **then**  
read the next block (if any) of the run into the buffer.
3. **until** all input buffer pages are empty:

# External Sort-Merge (Cont.)

- If  $N \geq M$ , several merge passes are required.
  - In each pass, contiguous groups of  $M - 1$  runs are merged.
  - A pass reduces the number of runs by a factor of  $M - 1$ , and creates runs longer by the same factor.
    - ▶ E.g. If  $M=11$ , and there are 90 runs, one pass reduces the number of runs to 9, each 10 times the size of the initial runs
  - Repeated passes are performed till all runs have been merged into one.

# Example: External Sorting Using Sort-Merge



# External Merge Sort (Cont.)

## ■ Cost analysis:

- 1 block per run leads to too many seeks during merge
  - ▶ Instead use  $b_b$  buffer blocks per run
    - read/write  $b_b$  blocks at a time
    - ▶ Can merge  $\lfloor M/b_b \rfloor - 1$  runs in one pass
  - Total number of merge passes required:  $\lceil \log_{\lfloor M/b_b \rfloor - 1} (b_r/M) \rceil$ .
  - Block transfers for initial run creation as well as in each pass is  $2b_r$ ,
    - ▶ for final pass, we don't count write cost
      - we ignore final write cost for all operations since the output of an operation may be sent to the parent operation without being written to disk
    - ▶ Thus total number of block transfers for external sorting:  
$$b_r (2 \lceil \log_{\lfloor M/b_b \rfloor - 1} (b_r/M) \rceil + 1) \lceil$$
  - Seek: next slide

# External Merge Sort (Cont.)

## ■ Cost of seeks

- During run generation: one seek to read each run and one seek to write each run
  - ▶  $2\lceil b_r/M \rceil$
- During the merge phase
  - ▶ Need  $2\lceil b_r/b_b \rceil$  seeks for each merge pass
    - except the final one which does not require a write
  - ▶ Total number of seeks:  
$$2\lceil b_r/M \rceil + \lceil b_r/b_b \rceil (2\lceil \log_{\lfloor M/b_b \rfloor - 1}(b_r/M) \rceil - 1)$$

# Join Operation

- Several different algorithms to implement joins
  - Nested-loop join
  - Block nested-loop join
  - Indexed nested-loop join
  - Merge-join
  - Hash-join
- Choice based on cost estimate
- Examples use the following information
  - Number of records of *student*: 5,000      *takes*: 10,000
  - Number of blocks of *student*: 100      *takes*: 400

# Nested-Loop Join

- To compute the theta join  $r \bowtie_{\theta} s$   
**for each tuple  $t_r$  in  $r$  do begin**  
    **for each tuple  $t_s$  in  $s$  do begin**  
        test pair  $(t_r, t_s)$  to see if they satisfy the join condition  $\theta$   
        if they do, add  $t_r \bullet t_s$  to the result.  
    **end**  
**end**
- $r$  is called the **outer relation** and  $s$  the **inner relation** of the join.
- Requires no indices and can be used with any kind of join condition.
- Expensive since it examines every pair of tuples in the two relations.

# Nested-Loop Join (Cont.)

- In the worst case, if there is enough memory only to hold one block of each relation, the estimated cost is
  - $n_r * b_s + b_r$  block transfers, plus
  - $n_r + b_r$  seeks
- If the smaller relation fits entirely in memory, use that as the inner relation.
  - Reduces cost to  $b_r + b_s$  block transfers and 2 seeks
- Assuming worst case memory availability cost estimate is
  - with *student* as outer relation:
    - ▶  $5000 * 400 + 100 = 2,000,100$  block transfers,
    - ▶  $5000 + 100 = 5100$  seeks
  - with *takes* as the outer relation
    - ▶  $10000 * 100 + 400 = 1,000,400$  block transfers and 10,400 seeks
- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.
- Block nested-loops algorithm (next slide) is preferable.

# Block Nested-Loop Join

- Variant of nested-loop join in which every block of inner relation is paired with every block of outer relation.

```
for each block  $B_r$  of  $r$  do begin
    for each block  $B_s$  of  $s$  do begin
        for each tuple  $t_r$  in  $B_r$  do begin
            for each tuple  $t_s$  in  $B_s$  do begin
                Check if  $(t_r, t_s)$  satisfy the join condition
                if they do, add  $t_r \bullet t_s$  to the result.
            end
        end
    end
end
```

# Block Nested-Loop Join (Cont.)

- Worst case estimate:  $b_r * b_s + b_r$  block transfers +  $2 * b_r$  seeks
  - Each block in the inner relation  $s$  is read once for each *block* in the outer relation
- Best case:  $b_r + b_s$  block transfers + 2 seeks.
- Improvements to nested loop and block nested loop algorithms:
  - In block nested-loop, use  $M - 2$  disk blocks as blocking unit for outer relations, where  $M$  = memory size in blocks; use remaining two blocks to buffer inner relation and output
    - ▶ Cost =  $\lceil b_r / (M-2) \rceil * b_s + b_r$  block transfers +  
 $2 \lceil b_r / (M-2) \rceil$  seeks
  - If equi-join attribute forms a key or inner relation, stop inner loop on first match
  - Scan inner loop forward and backward alternately, to make use of the blocks remaining in buffer (with LRU replacement)
  - Use index on inner relation if available (next slide)

# Indexed Nested-Loop Join

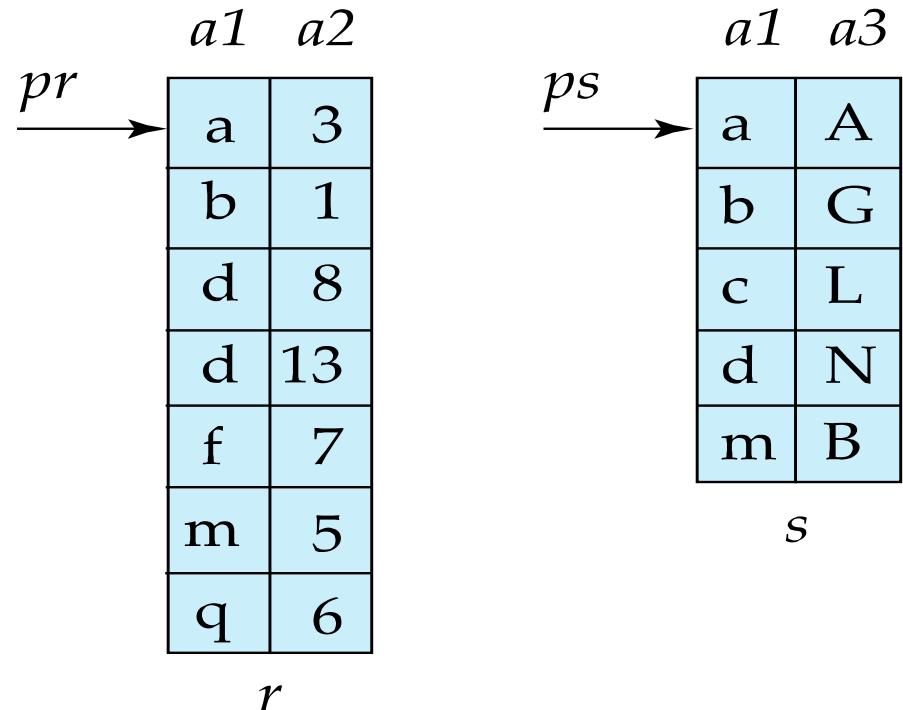
- Index lookups can replace file scans if
  - join is an equi-join or natural join and
  - an index is available on the inner relation's join attribute
    - ▶ Can construct an index just to compute a join.
- For each tuple  $t_r$  in the outer relation  $r$ , use the index to look up tuples in  $s$  that satisfy the join condition with tuple  $t_r$ .
- Worst case: buffer has space for only one page of  $r$ , and, for each tuple in  $r$ , we perform an index lookup on  $s$ .
- Cost of the join:  $b_r(t_T + t_S) + n_r * c$ 
  - Where  $c$  is the cost of traversing index and fetching all matching  $s$  tuples for one tuple of  $r$
  - $c$  can be estimated as cost of a single selection on  $s$  using the join condition.
- If indices are available on join attributes of both  $r$  and  $s$ , use the relation with fewer tuples as the outer relation.

# Example of Nested-Loop Join Costs

- Compute *student*  $\bowtie$  *takes*, with *student* as the outer relation.
- Let *takes* have a primary B<sup>+</sup>-tree index on the attribute *ID*, which contains 20 entries in each index node.
- Since *takes* has 10,000 tuples, the height of the tree is 4, and one more access is needed to find the actual data
- *student* has 5000 tuples
- Cost of block nested loops join
  - $400 * 100 + 100 = 40,100$  block transfers +  $2 * 100 = 200$  seeks
    - ▶ assuming worst case memory
    - ▶ may be significantly less with more memory
- Cost of indexed nested loops join
  - $100 + 5000 * 5 = 25,100$  block transfers and seeks.
  - CPU cost likely to be less than that for block nested loops join

# Merge-Join

1. Sort both relations on their join attribute (if not already sorted on the join attributes).
2. Merge the sorted relations to join them
  1. Join step is similar to the merge stage of the sort-merge algorithm.
  2. Main difference is handling of duplicate values in join attribute — every pair with same value on join attribute must be matched
  3. Detailed algorithm in book



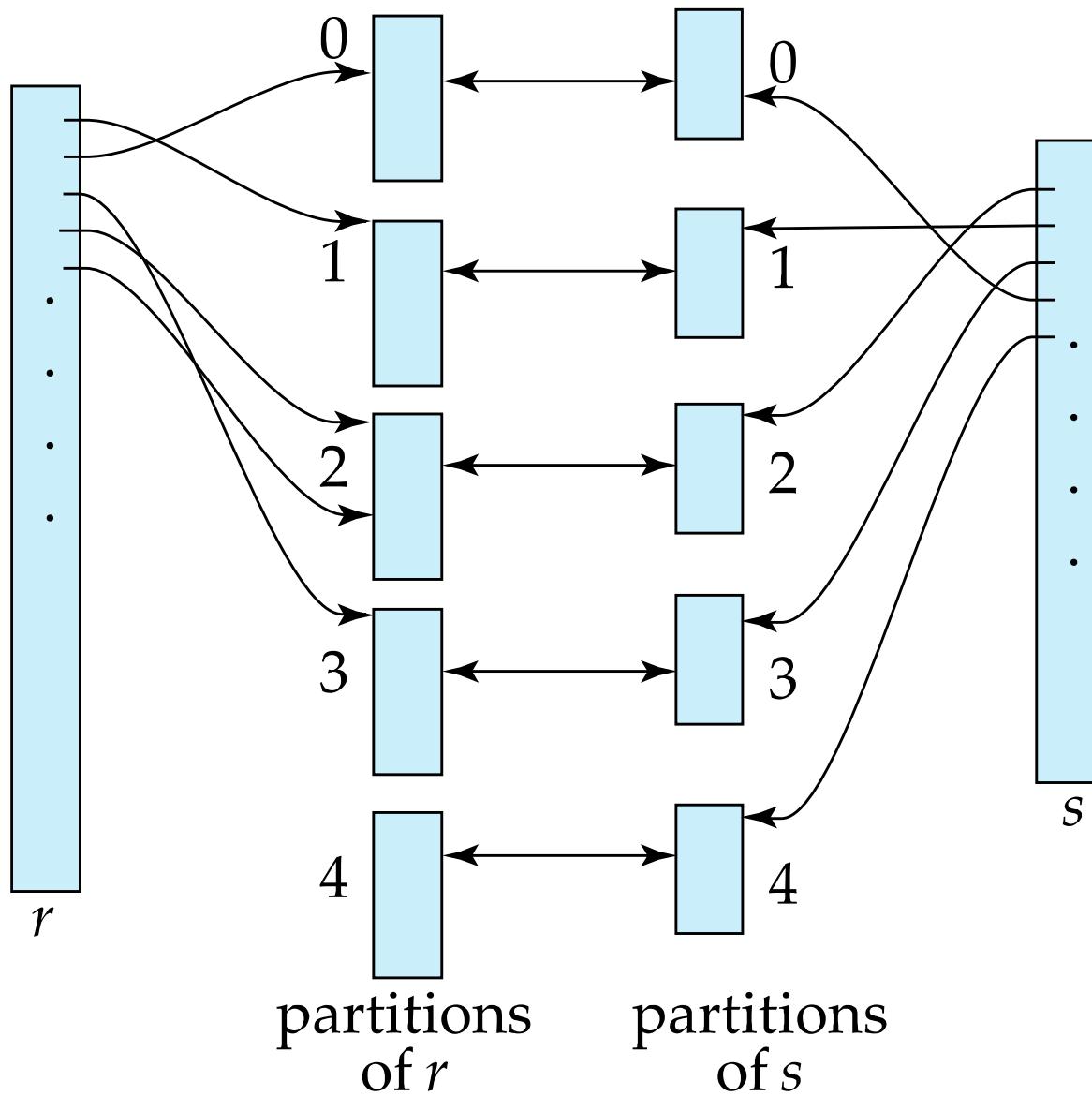
# Merge-Join (Cont.)

- Can be used only for equi-joins and natural joins
- Each block needs to be read only once (assuming all tuples for any given value of the join attributes fit in memory)
- Thus the cost of merge join is:  
 $b_r + b_s$  block transfers +  $\lceil b_r/b_b \rceil + \lceil b_s/b_b \rceil$  seeks
  - + the cost of sorting if relations are unsorted.
- **hybrid merge-join:** If one relation is sorted, and the other has a secondary B<sup>+</sup>-tree index on the join attribute
  - Merge the sorted relation with the leaf entries of the B<sup>+</sup>-tree .
  - Sort the result on the addresses of the unsorted relation's tuples
  - Scan the unsorted relation in physical address order and merge with previous result, to replace addresses by the actual tuples
    - ▶ Sequential scan more efficient than random lookup

# Hash-Join

- Applicable for equi-joins and natural joins.
- A hash function  $h$  is used to partition tuples of both relations
- $h$  maps  $JoinAttrs$  values to  $\{0, 1, \dots, n\}$ , where  $JoinAttrs$  denotes the common attributes of  $r$  and  $s$  used in the natural join.
  - $r_0, r_1, \dots, r_n$  denote partitions of  $r$  tuples
    - ▶ Each tuple  $t_r \in r$  is put in partition  $r_i$  where  $i = h(t_r[JoinAttrs])$ .
  - $r_0, r_1, \dots, r_n$  denotes partitions of  $s$  tuples
    - ▶ Each tuple  $t_s \in s$  is put in partition  $s_i$ , where  $i = h(t_s[JoinAttrs])$ .
- Note: In book,  $r_i$  is denoted as  $H_{ri}$ ,  $s_i$  is denoted as  $H_{si}$  and  $n$  is denoted as  $n_h$ .

# Hash-Join (Cont.)



# Hash-Join (Cont.)

- $r$  tuples in  $r_i$  need only to be compared with  $s$  tuples in  $s_i$ . Need not be compared with  $s$  tuples in any other partition, since:
  - an  $r$  tuple and an  $s$  tuple that satisfy the join condition will have the same value for the join attributes.
  - If that value is hashed to some value  $i$ , the  $r$  tuple has to be in  $r_i$  and the  $s$  tuple in  $s_i$ .

# Hash-Join Algorithm

The hash-join of  $r$  and  $s$  is computed as follows.

1. Partition the relation  $s$  using hashing function  $h$ . When partitioning a relation, one block of memory is reserved as the output buffer for each partition.
2. Partition  $r$  similarly.
3. For each  $i$ :
  - (a) Load  $s_i$  into memory and build an in-memory hash index on it using the join attribute. This hash index uses a different hash function than the earlier one  $h$ .
  - (b) Read the tuples in  $r_i$  from the disk one by one. For each tuple  $t_r$ , locate each matching tuple  $t_s$  in  $s_i$  using the in-memory hash index. Output the concatenation of their attributes.

Relation  $s$  is called the **build input** and  $r$  is called the **probe input**.

# Hash-Join algorithm (Cont.)

- The value  $n$  and the hash function  $h$  is chosen such that each  $s_i$  should fit in memory.
  - Typically  $n$  is chosen as  $\lceil b_s/M \rceil * f$  where  $f$  is a “**fudge factor**”, typically around 1.2
  - The probe relation partitions  $s_i$  need not fit in memory
- **Recursive partitioning** required if number of partitions  $n$  is greater than number of pages  $M$  of memory.
  - instead of partitioning  $n$  ways, use  $M - 1$  partitions for  $s$
  - Further partition the  $M - 1$  partitions using a different hash function
  - Use same partitioning method on  $r$
  - Rarely required: e.g., with block size of 4 KB, recursive partitioning not needed for relations of < 1GB with memory size of 2MB, or relations of < 36 GB with memory of 12 MB

# Handling of Overflows

- Partitioning is said to be **skewed** if some partitions have significantly more tuples than some others
- **Hash-table overflow** occurs in partition  $s_i$  if  $s_i$  does not fit in memory. Reasons could be
  - Many tuples in  $s_i$  with same value for join attributes
  - Bad hash function
- **Overflow resolution** can be done in build phase
  - Partition  $s_i$  is further partitioned using different hash function.
  - Partition  $r_i$  must be similarly partitioned.
- **Overflow avoidance** performs partitioning carefully to avoid overflows during build phase
  - E.g. partition build relation into many partitions, then combine them
- Both approaches fail with large numbers of duplicates

# Cost of Hash-Join

- If recursive partitioning is not required: cost of hash join is
$$3(b_r + b_s) + 4 * n_h \text{ block transfers} + \\ 2(\lceil b_r/b_b \rceil + \lceil b_s/b_b \rceil) \text{ seeks}$$
- If recursive partitioning required:
  - number of passes required for partitioning build relation  $s$  to less than  $M$  blocks per partition is  $\lceil \log_{\lfloor M/bb \rfloor - 1}(b_s/M) \rceil$
  - best to choose the smaller relation as the build relation.
  - Total cost estimate is:
$$2(b_r + b_s)\lceil \log_{\lfloor M/bb \rfloor - 1}(b_s/M) \rceil + b_r + b_s \text{ block transfers} + \\ 2(\lceil b_r/b_b \rceil + \lceil b_s/b_b \rceil)\lceil \log_{\lfloor M/bb \rfloor - 1}(b_s/M) \rceil \text{ seeks}$$
- If the entire build input can be kept in main memory no partitioning is required
  - Cost estimate goes down to  $b_r + b_s$ .

# Example of Cost of Hash-Join

*instructor*  $\bowtie$  *teaches*

- Assume that memory size is 20 blocks
- $b_{instructor} = 100$  and  $b_{teaches} = 400$ .
- *instructor* is to be used as build input. Partition it into five partitions, each of size 20 blocks. This partitioning can be done in one pass.
- Similarly, partition *teaches* into five partitions, each of size 80. This is also done in one pass.
- Therefore total cost, ignoring cost of writing partially filled blocks:
  - $3(100 + 400) = 1500$  block transfers +  
 $2(\lceil 100/3 \rceil + \lceil 400/3 \rceil) = 336$  seeks

# Hybrid Hash–Join

- Useful when memory sizes are relatively large, and the build input is bigger than memory.
- **Main feature of hybrid hash join:**  
**Keep the first partition of the build relation in memory.**
- E.g. With memory size of 25 blocks, *instructor* can be partitioned into five partitions, each of size 20 blocks.
  - Division of memory:
    - ▶ The first partition occupies 20 blocks of memory
    - ▶ 1 block is used for input, and 1 block each for buffering the other 4 partitions.
- *teaches* is similarly partitioned into five partitions each of size 80
  - the first is used right away for probing, instead of being written out
- Cost of  $3(80 + 320) + 20 + 80 = 1300$  block transfers for hybrid hash join, instead of 1500 with plain hash-join.
- Hybrid hash-join most useful if  $M >>$

# Complex Joins

## ■ Join with a conjunctive condition:

$$r \bowtie_{\theta_1 \wedge \theta_2 \wedge \dots \wedge \theta_n} s$$

- Either use nested loops/block nested loops, or
- Compute the result of one of the simpler joins  $r \bowtie_{\theta_i} s$ 
  - ▶ final result comprises those tuples in the intermediate result that satisfy the remaining conditions

$$\theta_1 \wedge \dots \wedge \theta_{i-1} \wedge \theta_{i+1} \wedge \dots \wedge \theta_n$$

## ■ Join with a disjunctive condition

$$r \bowtie_{\theta_1 \vee \theta_2 \vee \dots \vee \theta_n} s$$

- Either use nested loops/block nested loops, or
- Compute as the union of the records in individual joins  $r \bowtie_{\theta_i} s$ :

$$(r \bowtie_{\theta_1} s) \cup (r \bowtie_{\theta_2} s) \cup \dots \cup (r \bowtie_{\theta_n} s)$$

# Other Operations

- **Duplicate elimination** can be implemented via hashing or sorting.
  - On sorting duplicates will come adjacent to each other, and all but one set of duplicates can be deleted.
  - *Optimization:* duplicates can be deleted during run generation as well as at intermediate merge steps in external sort-merge.
  - Hashing is similar – duplicates will come into the same bucket.
- **Projection:**
  - perform projection on each tuple
  - followed by duplicate elimination.

# Other Operations : Aggregation

- **Aggregation** can be implemented in a manner similar to duplicate elimination.
  - Sorting or hashing can be used to bring tuples in the same group together, and then the aggregate functions can be applied on each group.
  - *Optimization:* combine tuples in the same group during run generation and intermediate merges, by computing partial aggregate values
    - ▶ For count, min, max, sum: keep aggregate values on tuples found so far in the group.
      - When combining partial aggregate for count, add up the aggregates
    - ▶ For avg, keep sum and count, and divide sum by count at the end

# Other Operations : Set Operations

- **Set operations** ( $\cup$ ,  $\cap$  and  $-$ ): can either use variant of merge-join after sorting, or variant of hash-join.
- E.g., Set operations using hashing:
  1. Partition both relations using the same hash function
  2. Process each partition  $i$  as follows.
    1. Using a different hashing function, build an in-memory hash index on  $r_i$ .
    2. Process  $s_i$  as follows
      - $r \cup s$ :
        1. Add tuples in  $s_i$  to the hash index if they are not already in it.
        2. At end of  $s_i$  add the tuples in the hash index to the result.

# Other Operations : Set Operations

## ■ E.g., Set operations using hashing:

1. as before partition  $r$  and  $s$ ,
2. as before, process each partition  $i$  as follows
  1. build a hash index on  $r_i$ ,
  2. Process  $s_i$  as follows
    - $r \cap s$ :
      1. output tuples in  $s_i$  to the result if they are already there in the hash index
    - $r - s$ :
      1. for each tuple in  $s_i$ , if it is there in the hash index, delete it from the index.
      2. At end of  $s_i$ , add remaining tuples in the hash index to the result.

# Other Operations : Outer Join

■ **Outer join** can be computed either as

- A join followed by addition of null-padded non-participating tuples.
- by modifying the join algorithms.

■ Modifying merge join to compute  $r \bowtie s$

- In  $r \bowtie s$ , non participating tuples are those in  $r - \Pi_R(r \bowtie s)$
- Modify merge-join to compute  $r \bowtie s$ :
  - ▶ During merging, for every tuple  $t_r$  from  $r$  that do not match any tuple in  $s$ , output  $t_r$  padded with nulls.
- Right outer-join and full outer-join can be computed similarly.

## Other Operations : Outer Join

### ■ Modifying hash join to compute $r \bowtie s$

- If  $r$  is probe relation, output non-matching  $r$  tuples padded with nulls
- If  $r$  is build relation, when probing keep track of which  $r$  tuples matched  $s$  tuples. At end of  $s$ , output non-matched  $r$  tuples padded with nulls

# Evaluation of Expressions

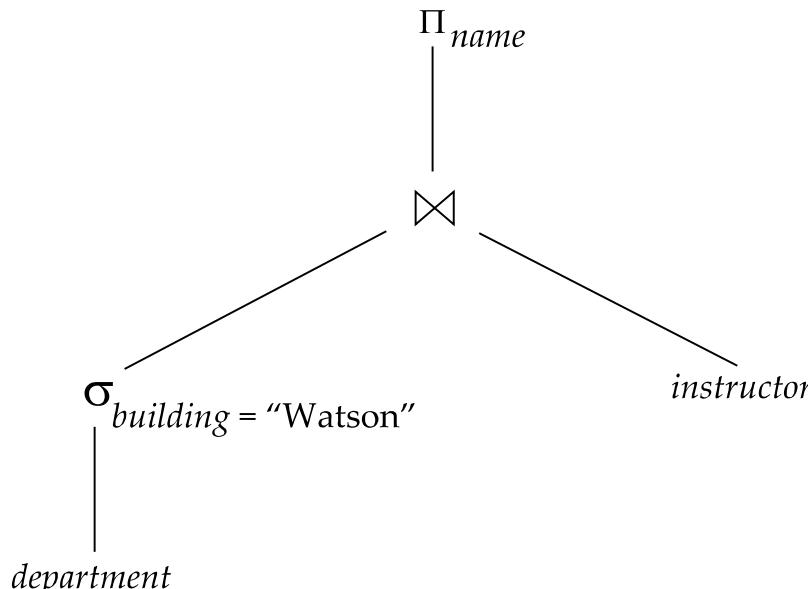
- So far: we have seen algorithms for individual operations
- Alternatives for evaluating an entire expression tree
  - **Materialization:** generate results of an expression whose inputs are relations or are already computed, **materialize** (store) it on disk. Repeat.
  - **Pipelining:** pass on tuples to parent operations even as an operation is being executed
- We study above alternatives in more detail

# Materialization

- **Materialized evaluation:** evaluate one operation at a time, starting at the lowest-level. Use intermediate results materialized into temporary relations to evaluate next-level operations.
- E.g., in figure below, compute and store

$$\sigma_{building="Watson"}(department)$$

then compute the store its join with *instructor*, and finally compute the projection on *name*.



# Materialization (Cont.)

- Materialized evaluation is always applicable
- Cost of writing results to disk and reading them back can be quite high
  - Our cost formulas for operations ignore cost of writing results to disk, so
    - ▶ Overall cost = Sum of costs of individual operations + cost of writing intermediate results to disk
- **Double buffering:** use two output buffers for each operation, when one is full write it to disk while the other is getting filled
  - Allows overlap of disk writes with computation and reduces execution time

# Pipelining

- **Pipelined evaluation** : evaluate several operations simultaneously, passing the results of one operation on to the next.
- E.g., in previous expression tree, don't store result of

$$\sigma_{building='Watson'}(department)$$

- instead, pass tuples directly to the join.. Similarly, don't store result of join, pass tuples directly to projection.
- Much cheaper than materialization: no need to store a temporary relation to disk.
- Pipelining may not always be possible – e.g., sort, hash-join.
- For pipelining to be effective, use evaluation algorithms that generate output tuples even as tuples are received for inputs to the operation.
- Pipelines can be executed in two ways: **demand driven** and **producer driven**

# Pipelining (Cont.)

- In **demand driven** or **lazy** evaluation
  - system repeatedly requests next tuple from top level operation
  - Each operation requests next tuple from children operations as required, in order to output its next tuple
  - In between calls, operation has to maintain “**state**” so it knows what to return next
- In **producer-driven** or **eager** pipelining
  - Operators produce tuples eagerly and pass them up to their parents
    - ▶ Buffer maintained between operators, child puts tuples in buffer, parent removes tuples from buffer
    - ▶ if buffer is full, child waits till there is space in the buffer, and then generates more tuples
  - System schedules operations that have space in output buffer and can process more input tuples

# Pipelining (Cont.)

- Implementation of demand-driven pipelining
  - Each operation is implemented as an **iterator** implementing the following operations
    - ▶ **open()**
      - E.g. file scan: initialize file scan
        - » state: pointer to beginning of file
      - E.g. merge join: sort relations;
        - » state: pointers to beginning of sorted relations
    - ▶ **next()**
      - E.g. for file scan: Output next tuple, and advance and store file pointer
      - E.g. for merge join: continue with merge from earlier state till next output tuple is found. Save pointers as iterator state.
    - ▶ **close()**

# Evaluation Algorithms for Pipelining

- Some algorithms are not able to output results even as they get input tuples
  - E.g. merge join, or hash join
  - intermediate results written to disk and then read back
- Algorithm variants to generate (at least some) results on the fly, as input tuples are read in
  - E.g. hybrid hash join generates output tuples even as probe relation tuples in the in-memory partition (partition 0) are read in
  - **Double-pipelined join technique:** Hybrid hash join, modified to buffer partition 0 tuples of both relations in-memory, reading them as they become available, and output results of any matches between partition 0 tuples
    - ▶ When a new  $r_0$  tuple is found, match it with existing  $s_0$  tuples, output matches, and save it in  $r_0$
    - ▶ Symmetrically for  $s_0$  tuples

# End of Chapter

**Database System Concepts, 6<sup>th</sup> Ed.**

©Silberschatz, Korth and Sudarshan  
See [www.db-book.com](http://www.db-book.com) for conditions on re-use

## Figure 12.02

$\pi_{\text{salary}}$

$\sigma_{\text{salary} < 75000; \text{use index 1}}$

*instructor*

# Selection Operation (Cont.)

- **Old-A2** (*binary search*). Applicable if selection is an equality comparison on the attribute on which file is ordered.
  - Assume that the blocks of a relation are stored contiguously
  - Cost estimate (number of disk blocks to be scanned):
    - ▶ cost of locating the first tuple by a binary search on the blocks
      - $\lceil \log_2(b_r) \rceil * (t_T + t_S)$
    - ▶ If there are multiple records satisfying selection
      - Add transfer cost of the number of blocks containing records that satisfy selection condition
      - Will see how to estimate this cost in Chapter 13