



BITS Pilani
Pilani Campus

BITS Pilani presentation

Dr. Mukesh Kumar Rohil
Department of Computer Science & Information Systems



BITS Pilani
Pilani Campus



Database Systems: Data Storage and Indexing

Database Systems: Data Storage and Indexing



Dr. Mukesh Kumar Rohil
Instructor of the course on
Database Systems

**Computer Science &
Information Systems
Department**

**Birla Institute of
Technology & Science
Pilani – 333031
(Rajasthan)**



Learning Objectives

-
- Overview of Physical Storage Media
 - Magnetic Disks
 - RAID
 - Tertiary Storage
 - Storage Access
 - File Organization
 - Organization of Records in Files
 - Data-Dictionary Storage
 - Basic Concepts of Indexing
 - Ordered Indices
 - B⁺-Tree Index Files
 - B-Tree Index Files
 - Static Hashing
 - Dynamic Hashing
 - Comparison of Ordered Indexing and Hashing
 - Index Definition in SQL
 - Multiple-Key Access



Resources

The subsequent slides are adopted and/or prepared from the material available in the following prescribed textbook (TB) for the course:

1. Silberschatz A, Korth H F, & Sudarshan S, ***Database System Concepts***, 7e, TMH, 2019..



Introduction

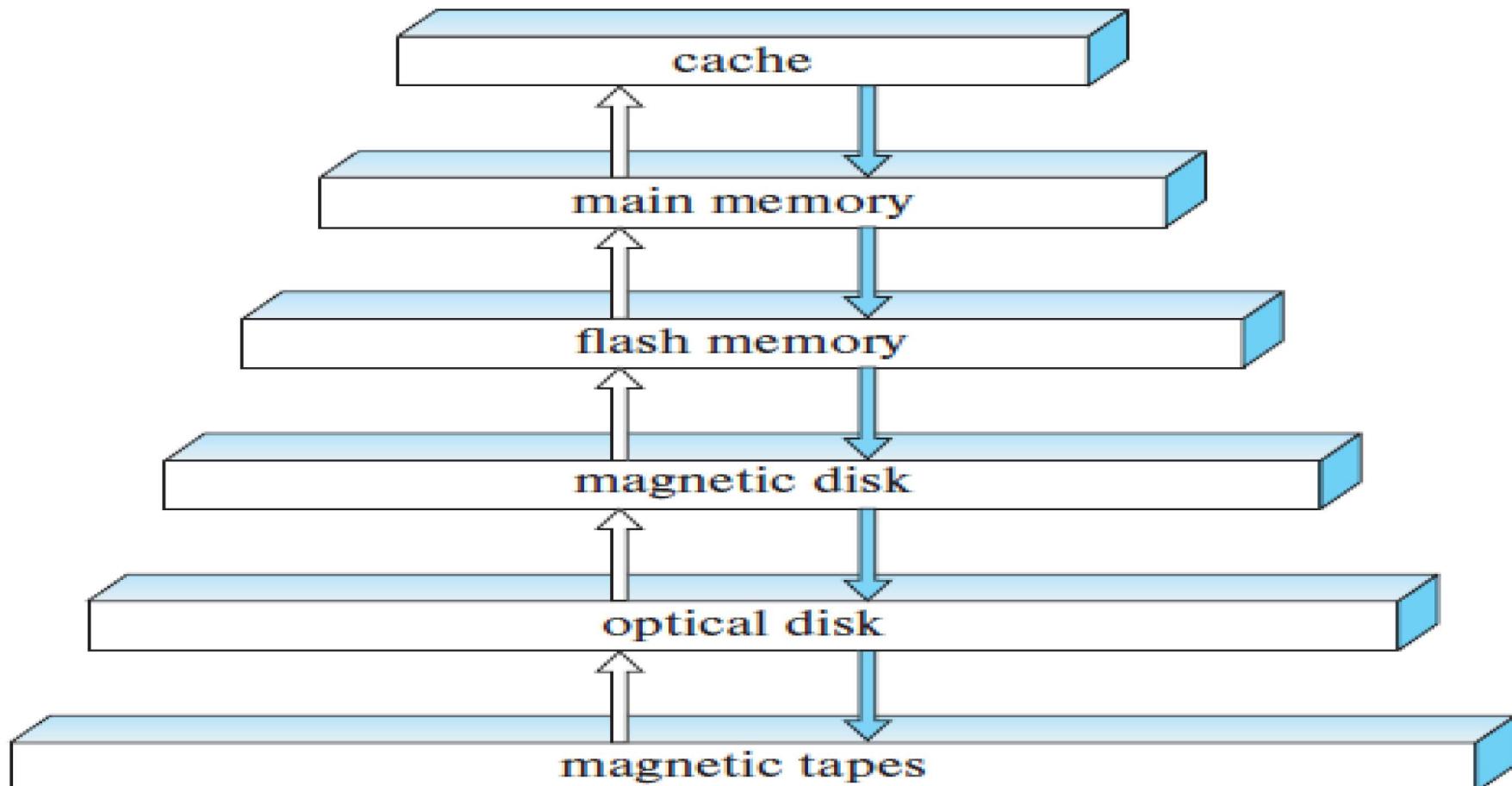
- Data have to be stored as bits on storage devices.
 - Majority of database systems store data on magnetic disks.
 - Data having higher performance requirements stored on flash-based solid-state drives.
 - Database systems fetch data into main memory for processing, and write data back to storage for persistence.
 - Data are also copied to backup devices for archival storage.
 - The physical characteristics of storage devices play a major role in the way data are stored, accessed and retrieved.
 - Overview of physical storage media, memory-hierarchy, working of magnetic storage disks.
 - Mechanisms to minimize the chance of data loss.
-



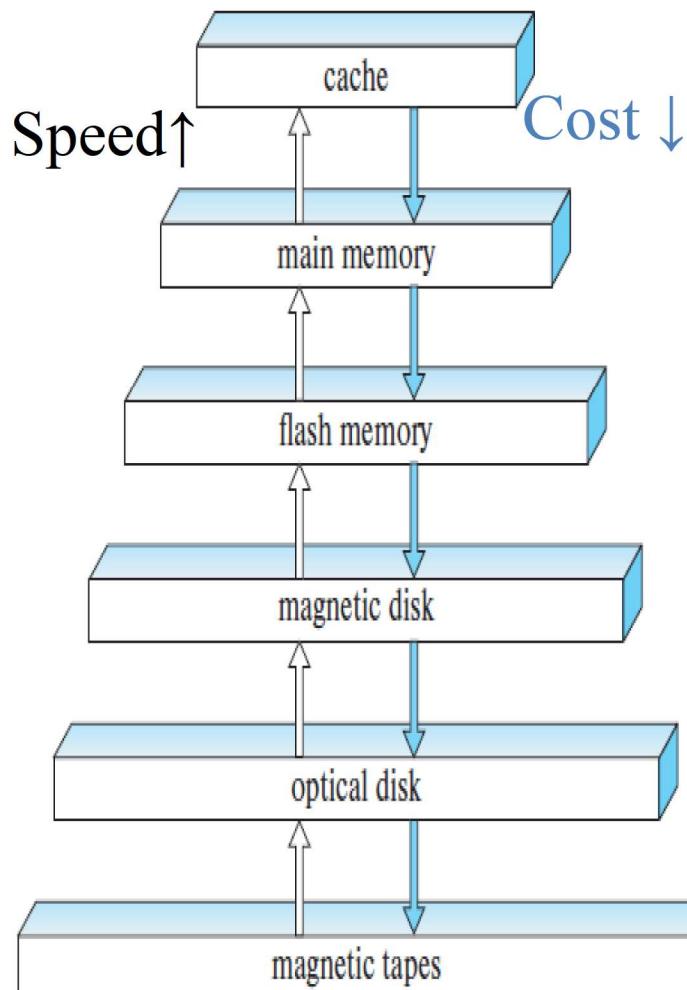
Introduction ...2

- Techniques for efficient disk-block access.
 - How records are mapped to files, which in turn are mapped to bits on the disk?
 - Techniques for the efficient management of the main-memory buffer for disk-based data.
 - How to efficiently use Column-oriented storage, used in data analytics systems?
 - Description of several types of indices (an index is a structure that helps locate desired records of a relation quickly, without examining all records) used in database systems.
-

Overview of Physical Storage Media and Storage device hierarchy



Overview of Physical Storage Media and Storage device hierarchy



Cache (Static RAM or SRAM) are volatile, the fastest but most costly. **Concern:** designing query processing data structures and algorithms.

Main memory (Dynamic RAM or DRAM) are volatile. Up to 10 GB for personal computers and 100 GB for servers. **Concern:** Volatile.

Flash memory does not contain any moving parts, so it delivers the faster access time, noiseless operation, higher reliability, & lower power consumption. **Concern:** Costly.

Magnetic disk storage survives power failures and system crashes. Storage capacity: up to 14 TB. **Concern:** Lower performance in terms of number of data access operations that they can support per second

Optical disk comes as ROM, WORM, RW, DVD (4.7 GB to 17 GB), Blue-Ray DVD (27 GB to 128 GB). Useful for backups of database. **Concern:** Not suitable for storing active database data. **Optical disk jukebox** systems contain a few drives and numerous disks.

Magnetic Tape storage (1TB to 12 TB) is used primarily for backup and archival data. Tape libraries (jukeboxes) are used to hold large collections of tapes. **Concern:** data must be accessed sequentially.

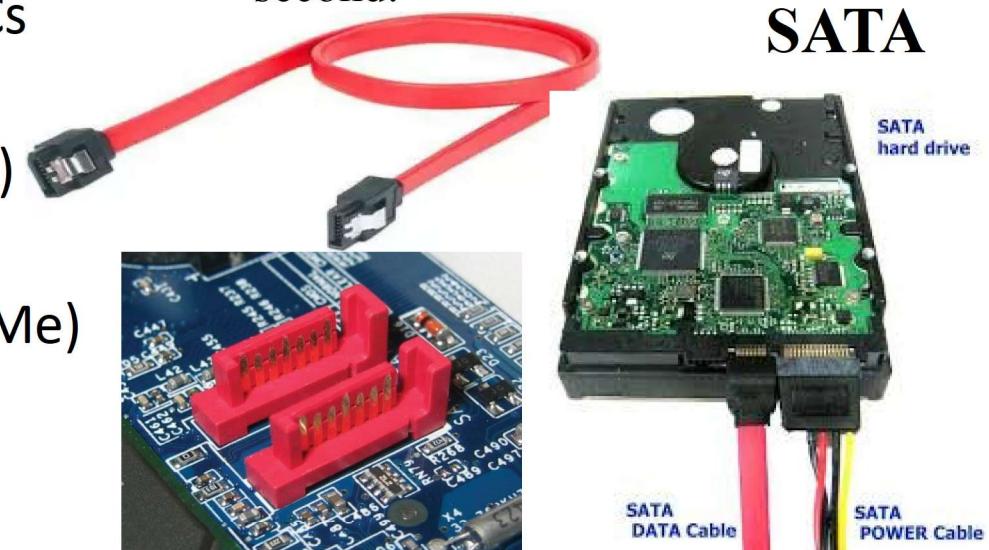
Storage Interfaces

- Magnetic disks as well as flash-based solid state disks are connected to a computer system through a high-speed interconnection.

- 1. SATA (Serial Advanced Technology Attachment)** – For hard disks in PCs
- 2. Serial Attached SCSI (Small Computer System Interface) (SAS)** – For servers
- 3. Non-Volatile Memory Express (NVMe)** – For flash drives
- 4. iSCSI** – For Storage area networks

SATA or Serial ATA

- SATA is an interface that connects various storage devices such as hard disks, optical drives, SSD's, etc. to the motherboard. SATA was introduced in the year 2000.
- SATA-3 supports data transfer speeds of up to 600 megabytes per second.

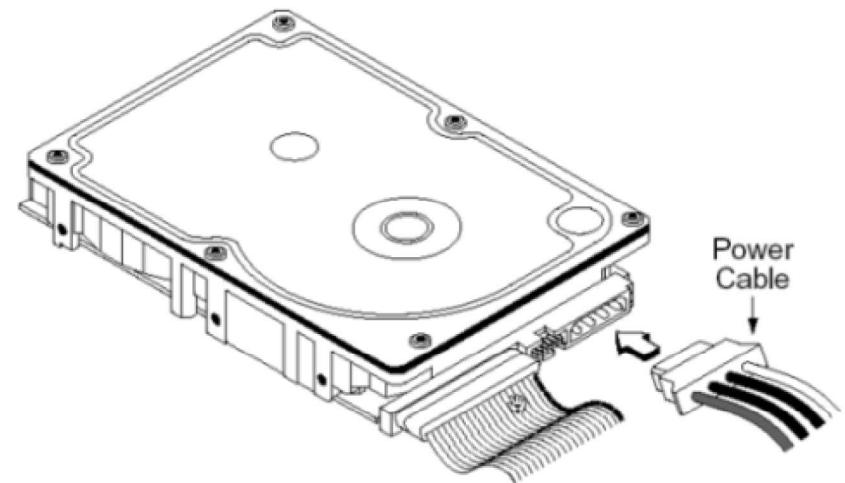


Storage Interfaces ... 2



Serial Attached SCSI (Small Computer System Interface) (SAS)

- For servers, SAS version 3 supports data transfer rates of 12 gigabits per second.



Storage Interfaces ... 3



- **Non-Volatile Memory Express (NVMe)** interface is a logical interface standard developed to better support SSDs and is typically used with the PCIe interface (**peripheral component interconnect express interface** provides high-speed data transfer internal to computer systems).
- Raw Data Transfer Rate for PCIe-6 is 64 GT/s, Bandwidth for x16 is 128 GB/s.
- Data Transfer for NVMe is 4GB/s.



www.shutterstock.com · 1468823192

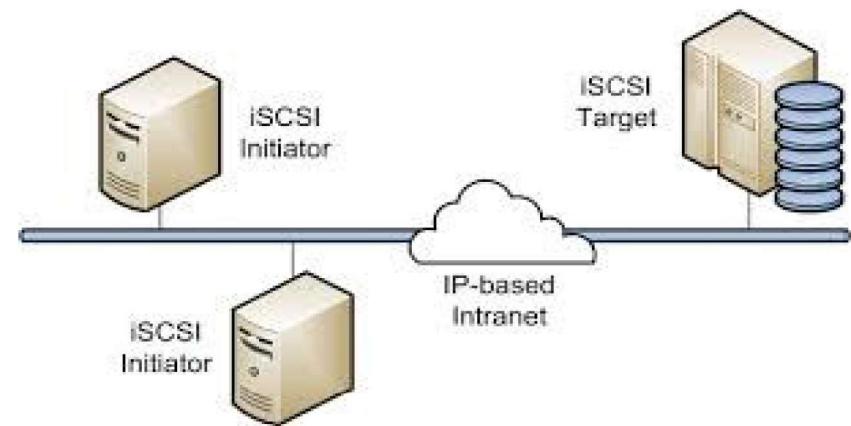


www.shutterstock.com · 1710261649

Storage Interfaces ... 4



- iSCSI (Internet Small Computer Systems Interface), is an Internet Protocol (IP) based **storage** networking standard for linking data **storage** facilities.
- An ***initiator*** functions as an iSCSI client.
- An iSCSI **target** is often a dedicated network-connected hard disk storage device, but may also be a general-purpose computer.
- In the storage area network (SAN) architecture, large numbers of disks are connected by a high-speed network to a number of server computers.
- The disks are usually organized locally using a storage organization technique called *redundant arrays of independent disks* (RAID).

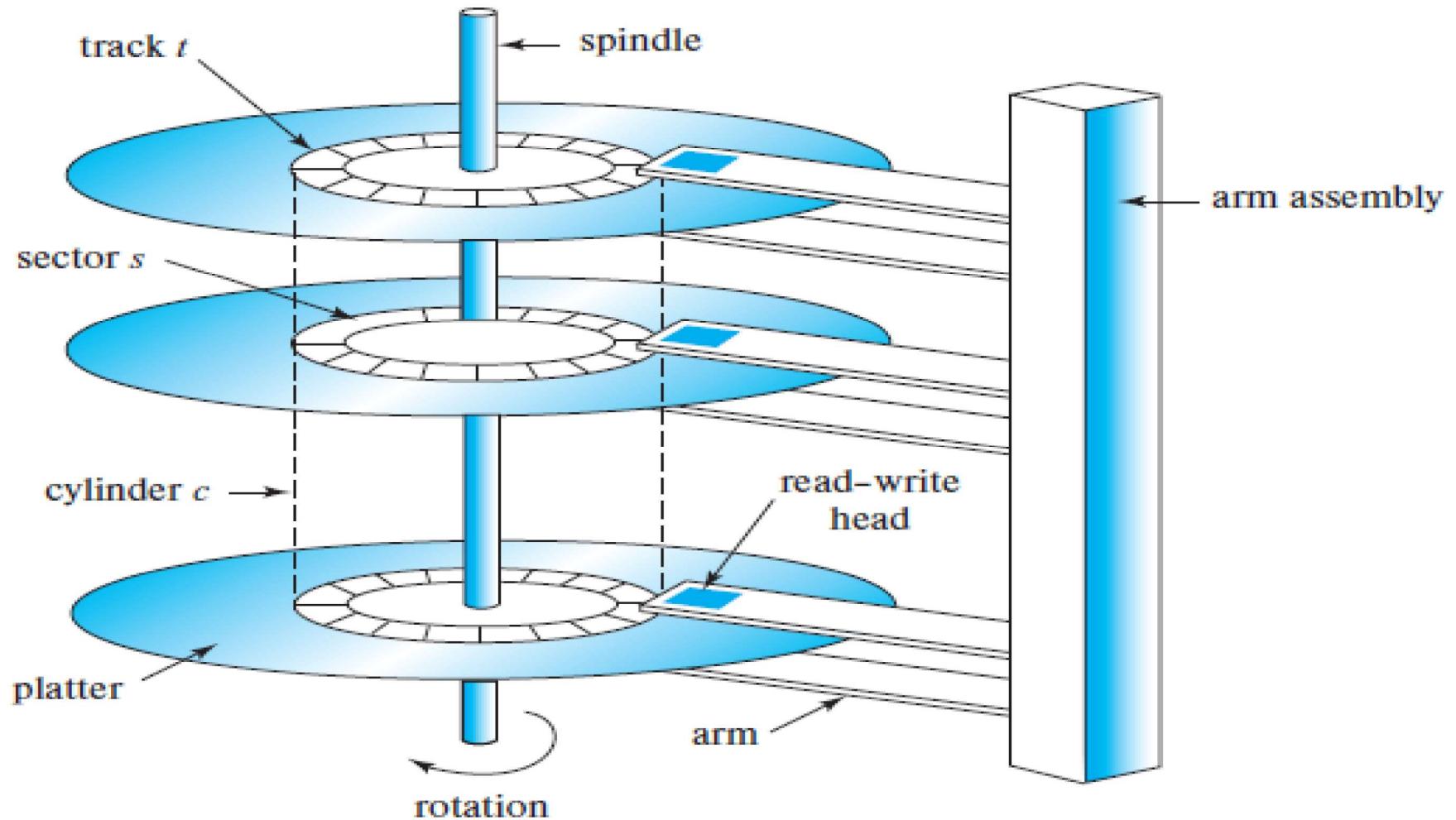




Storage Interfaces ... 5

- **Network attached storage (NAS)** is an alternative to SAN. NAS is much like SAN, except that instead of the networked storage appearing to be a large disk, it provides a file system interface using networked file system protocols.
- **Cloud storage**, where data are stored in the cloud and accessed via an API.
 - Cloud storage has a very high latency of tens to hundreds of milliseconds, if the data are not co-located with the database, and is thus not ideal as the underlying storage for databases.

Schematic diagram of a magnetic disk



Internals of an actual magnetic disk



Magnetic Disks: Features and Performance Measures



- Checksums
 - Disk controller interfaces
 - Remapping of bad sectors
 - Access time
 - Seek time
 - Average seek time
 - Rotational latency time
 - Average latency time
 - Data-transfer rate
 - Disk block
 - Page
 - Sequential access
 - Random access
 - I/O operations per second (IOPS)
 - Mean time to failure (MTTF)
-

Checksums



Secure Hash Algorithm 256

```
Windows PowerShell
Copyright (C) Microsoft Corporation. All rights reserved.

PS C:\Users\chris> Get-FileHash C:\Users\chris\File1.txt

Algorithm      Hash                                         Path
----          ----
SHA256        5A2EA7F56992A60B88206B5D8283A602EC8C9919947F2901230AFF8ECB637D16   C:\User

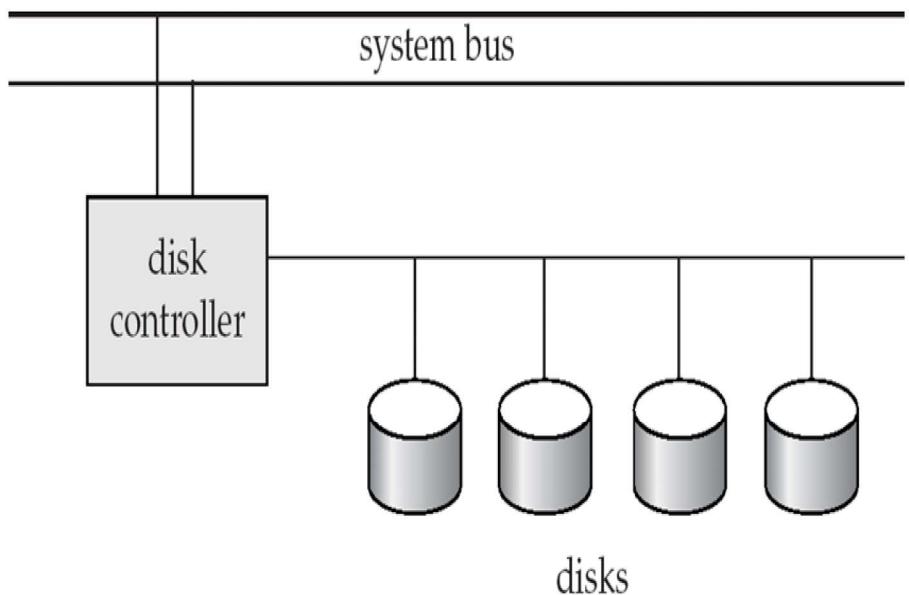
PS C:\Users\chris> Get-FileHash C:\Users\chris\File2.txt

Algorithm      Hash                                         Path
----          ----
SHA256        034A1CA76B1752D59496D86BE6FD781721E84B8CDC06BB5393828B2FFF211A2D   C:\User

The image shows two side-by-side Notepad windows. The left window is titled 'File1.txt - Notepad' and contains the text 'Words words words!'. The right window is titled 'File2.txt - Notepad' and contains the text 'Words words words.'. Both windows have standard Windows-style menus at the top.
```

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.)

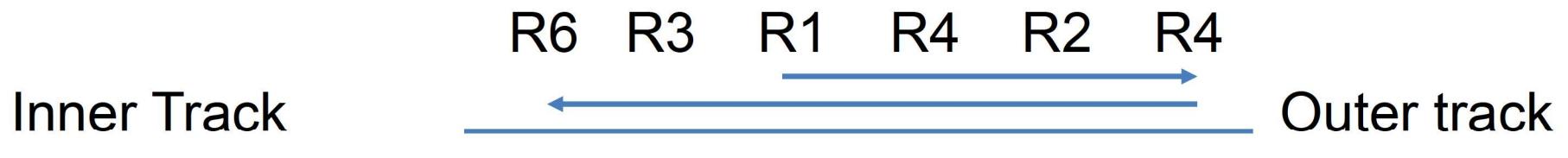


- **Mean time to failure (MTTF)** – the average time the disk is expected to run continuously without any failure.
 - Typically 3 to 5 years
 - Probability of failure of new disks is quite low, corresponding to a “theoretical MTTF” of 500,000 to 1,200,000 hours for a new disk
 - E.g., an MTTF of 1,200,000 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 500×10^6 hours (or 57,000 years) for a mirrored pair of disks (ignoring dependent failure modes)

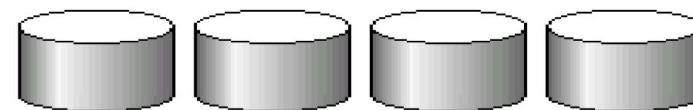
Improvement in Performance via Parallelism



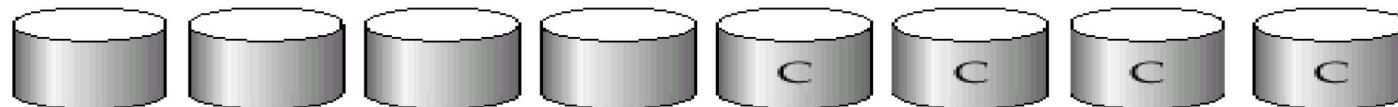
- Two main goals of parallelism in a disk system:
 1. Load balance multiple small accesses to increase throughput
 2. Parallelize large accesses to reduce response time.
- Improve transfer rate by striping data across multiple disks.
- **Bit-level striping** – split the bits of each byte across multiple 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

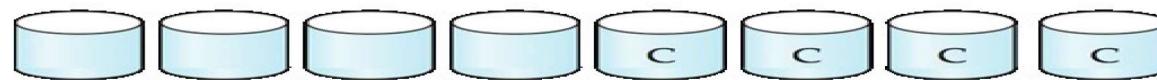


(b) RAID 1: mirrored disks

RAID Levels



(a) RAID 0: nonredundant striping



(b) RAID 1: mirrored disks



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



(d) RAID 3: bit-interleaved parity



(e) RAID 4: block-interleaved parity



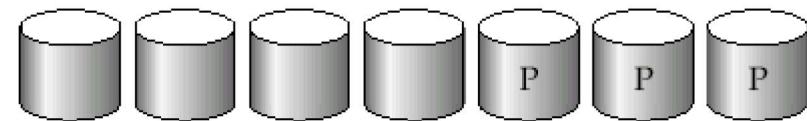
(f) RAID 5: block-interleaved distributed parity



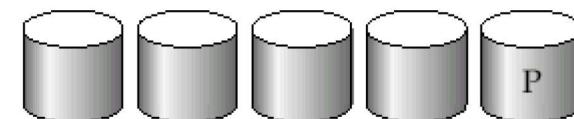
(g) RAID 6: P + Q redundancy

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)



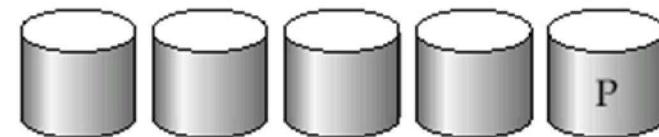
(c) RAID 2: memory-style error-correcting codes



(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



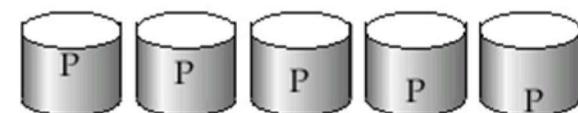
RAID Levels (Cont.)

- **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.

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



(f) RAID 5: block-interleaved distributed parity

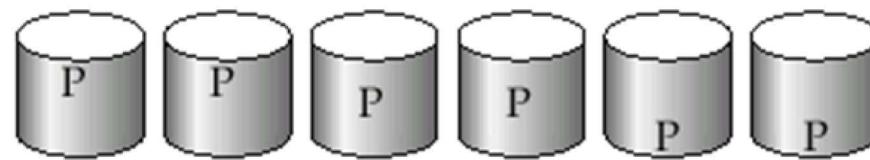
RAID Levels

Disk 1	Disk 2	Disk 3	Disk 4
B_1	B_2	B_3	B_4
P_1	B_5	B_6	B_7
B_8	P_2	B_9	B_{10}
:	:	:	:

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



**File Organization,
Record
Organization and
Storage
Access**



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 relationsThis 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

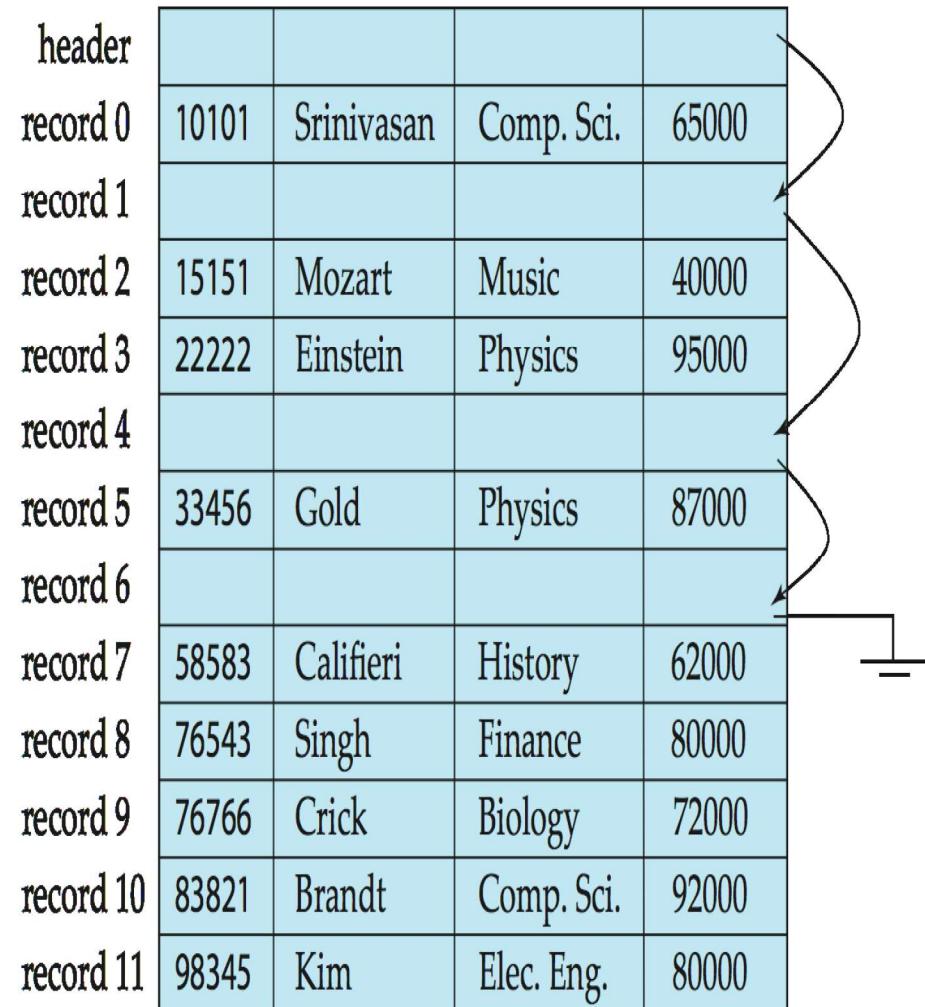
Deleting record 3 and moving last record

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

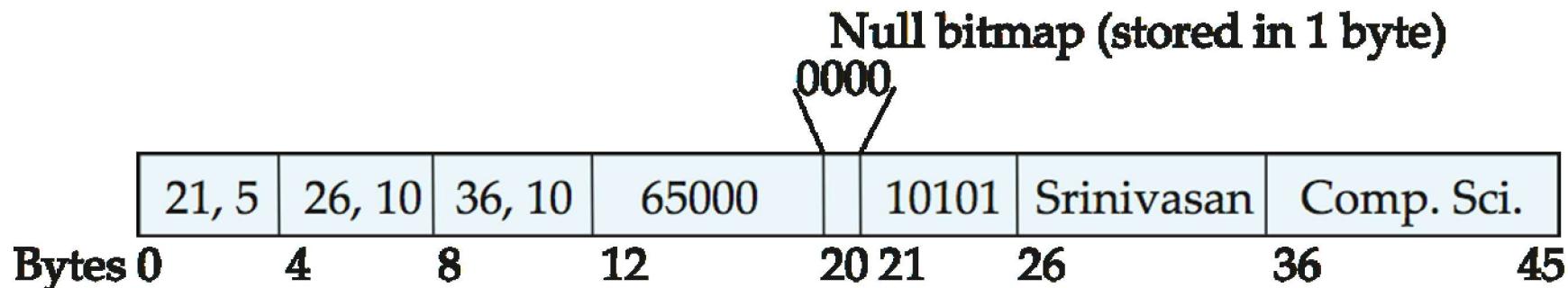
- Store the address of the first deleted record in the file header.
- Use this first record to store the address of the second deleted record, and so on
- Can think of these stored addresses as **pointers** since they “point” to the location of a record.
- More space efficient representation: reuse space 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



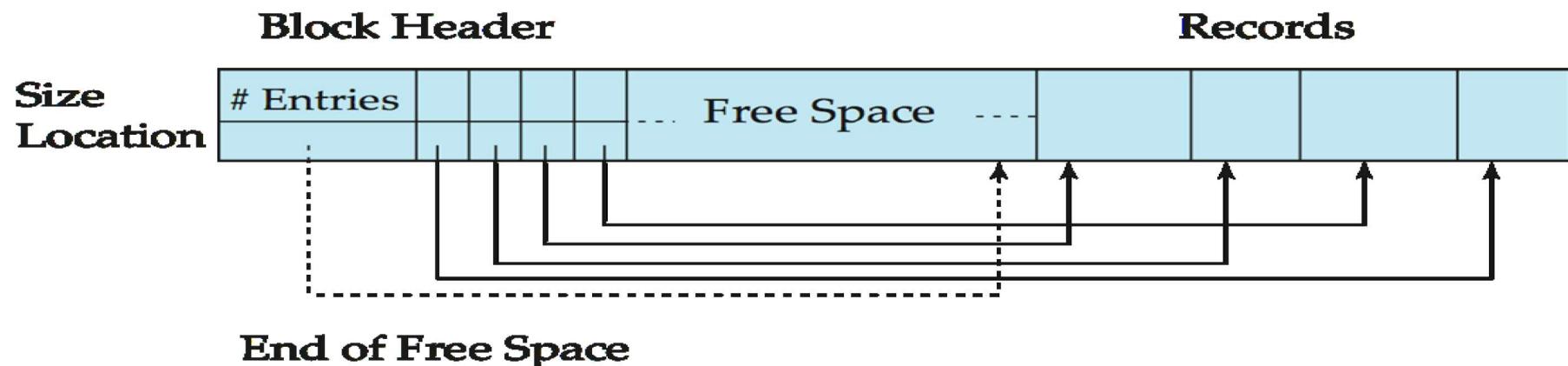
Variable-Length Records

- 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

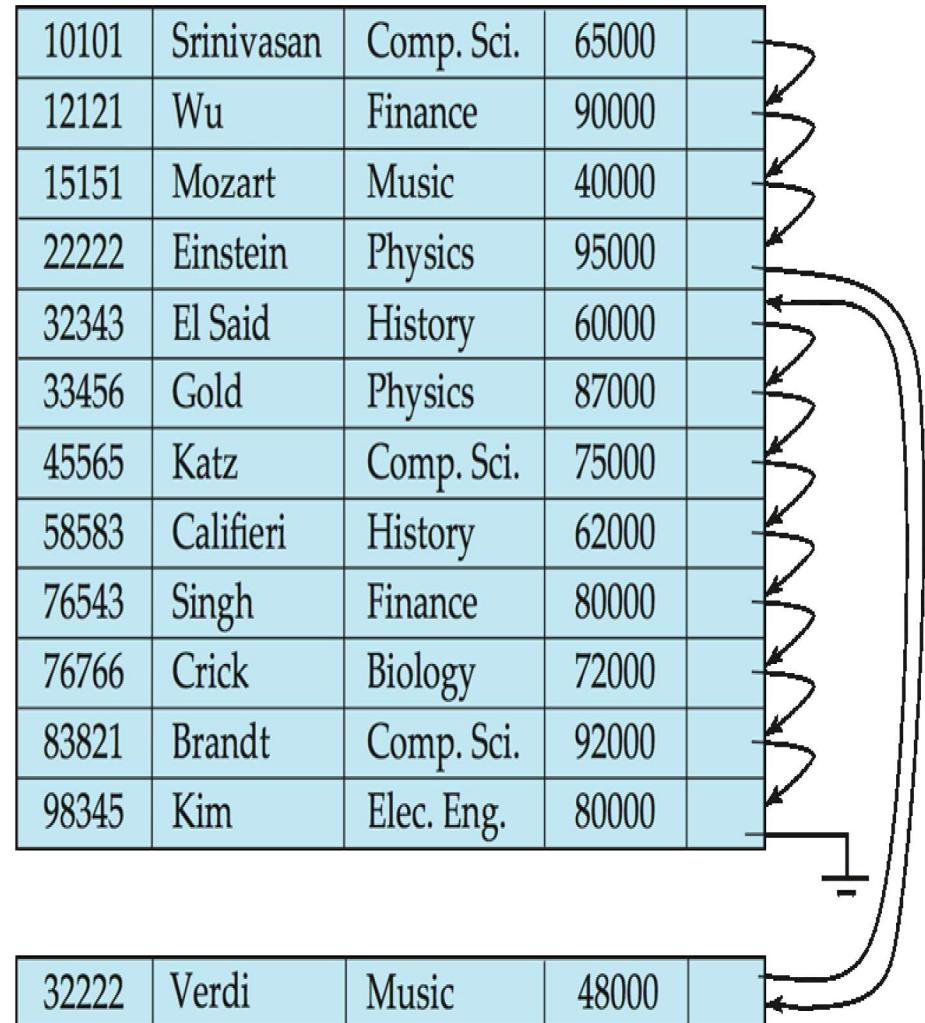
- Suitable for applications that require sequential processing of the entire file
- The records in the file are ordered by a **search-key**

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



Multi-table 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 instructors
- 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	





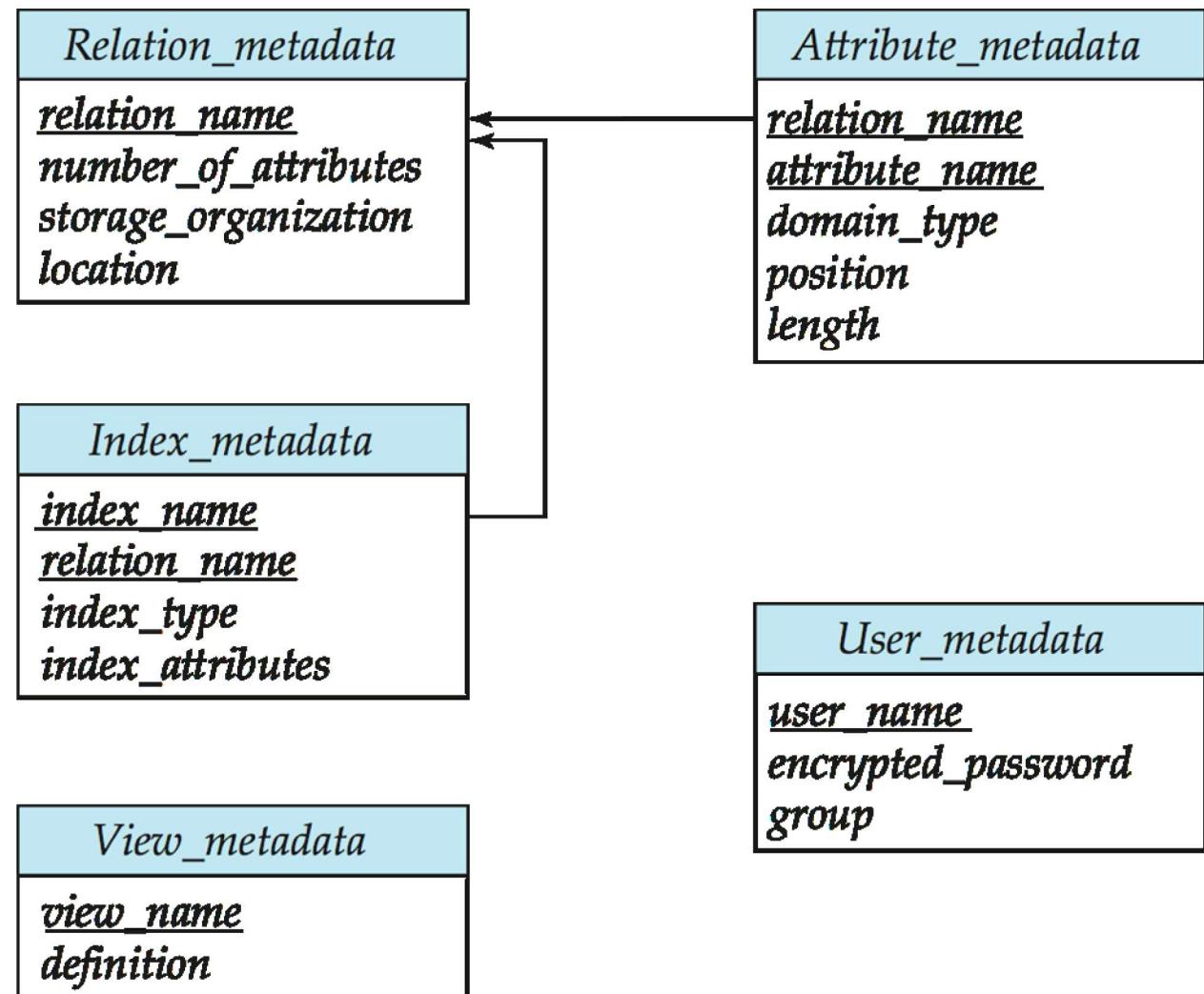
Data Dictionary Storage

- 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

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 details in subsequent sessions)



Indexing and Hashing

Indexing and Hashing

Topics to be covered:

- Basic Concepts
- Ordered Indices
- B⁺-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** – an attribute (or a 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



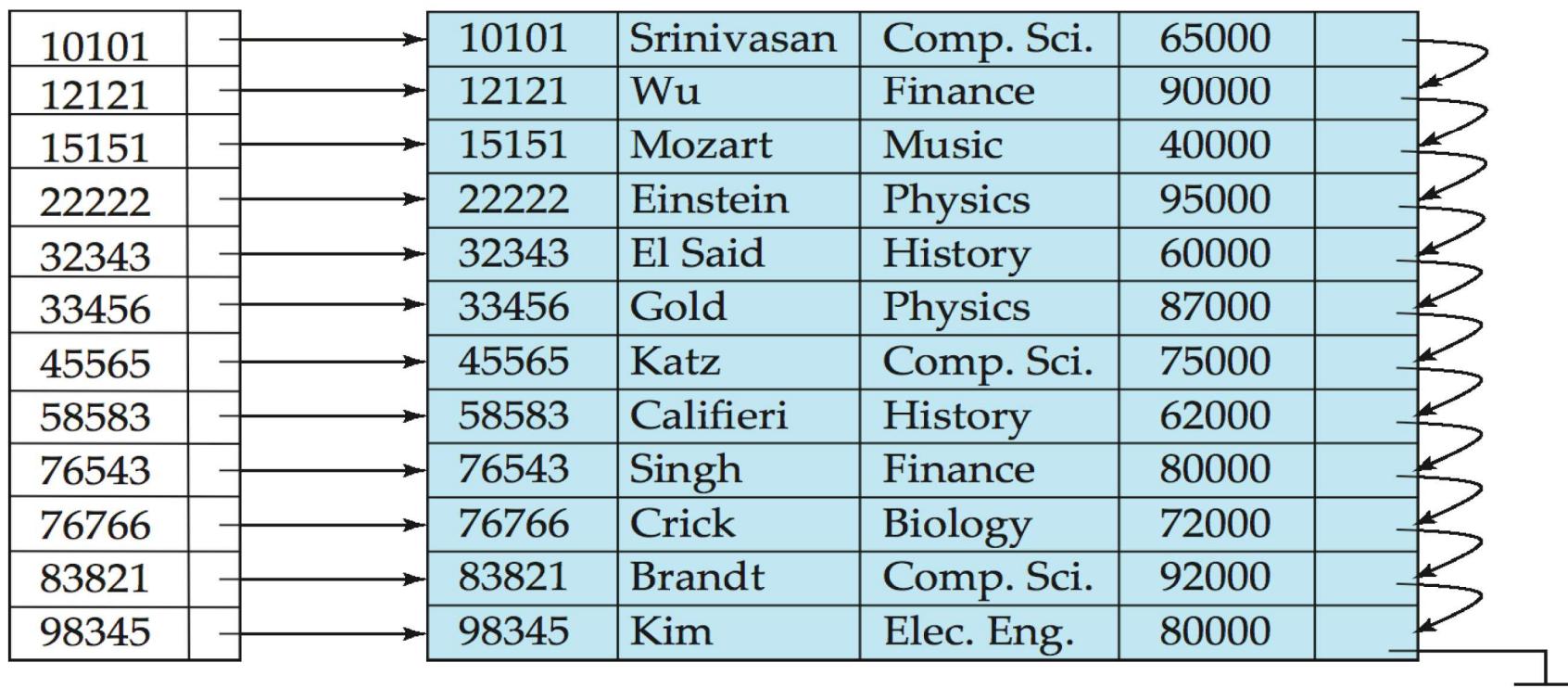
Ordered Indices

- 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

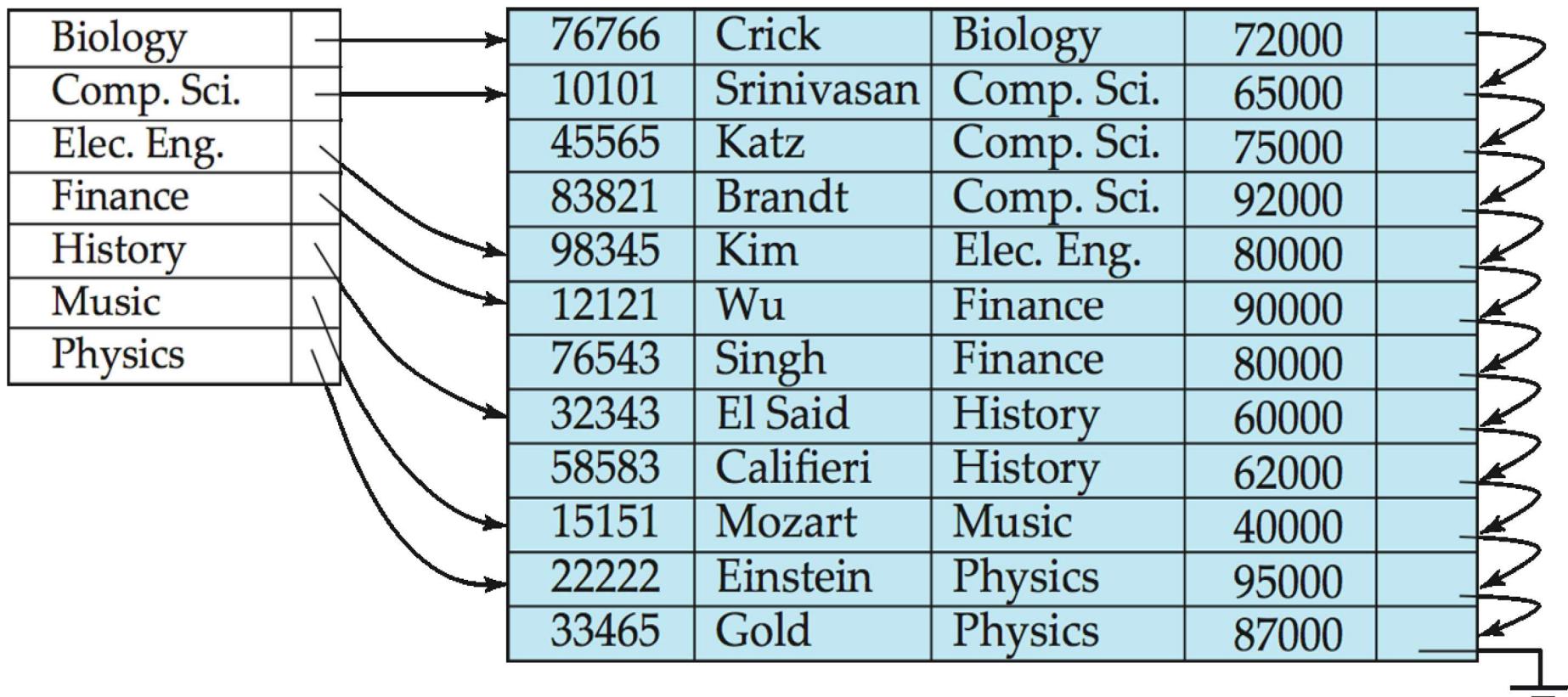
- **Dense index** — Index record appears for every search-key value in the file.

E.g. index on *ID* attribute of *instructor* relation



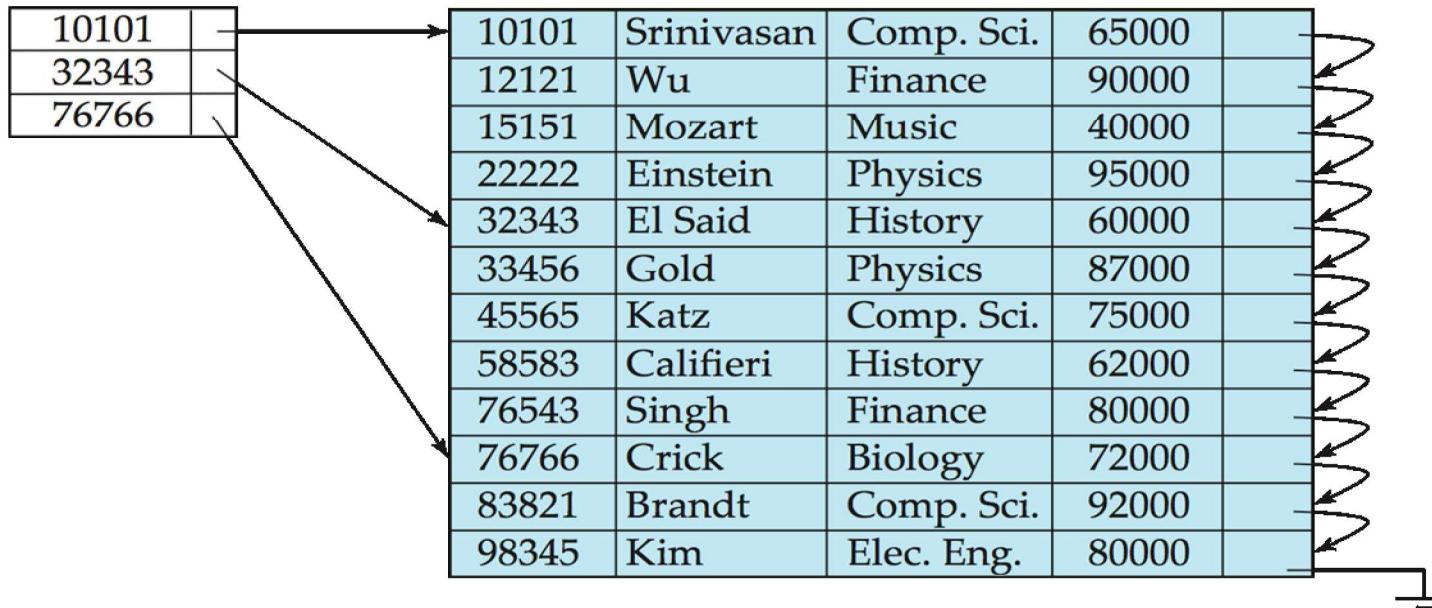
Dense Index Files (Cont.)

- Dense index on *dept_name*, with *instructor* file sorted on *dept_name*



Sparse Index Files

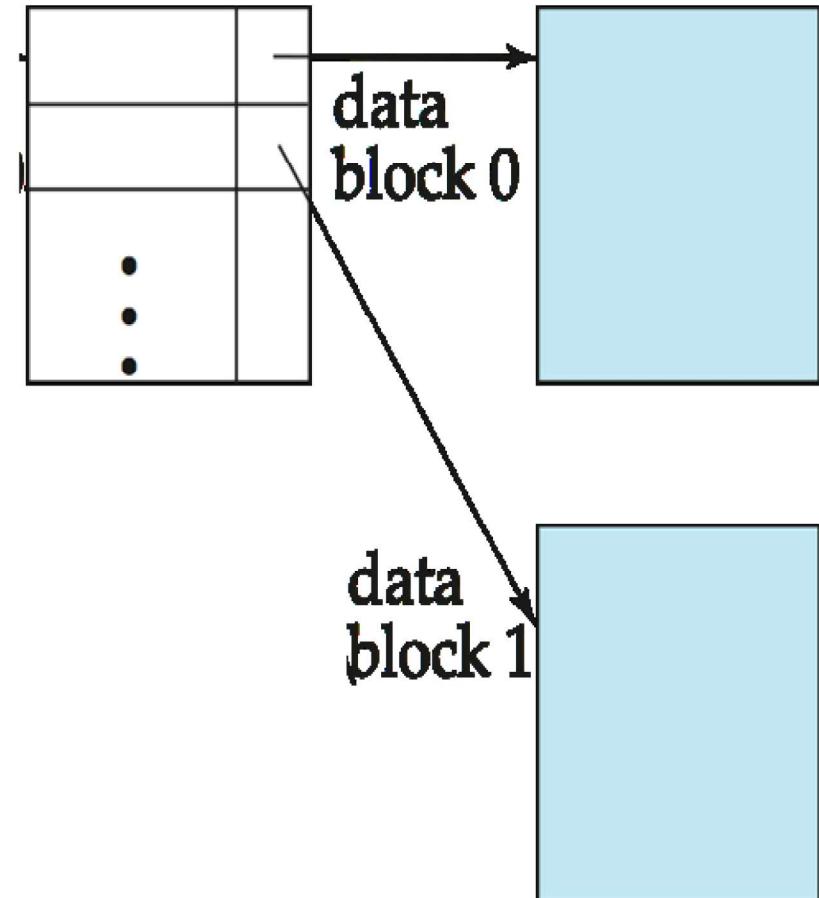
- **Sparse Index:** contains index records for only some search-key values.
 - Applicable when records are sequentially ordered on search-key
- To locate a record with search-key value K we:
 - Find index record with largest search-key value $< K$
 - Search file sequentially starting at the record to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



Sparse Index Files (Cont.)



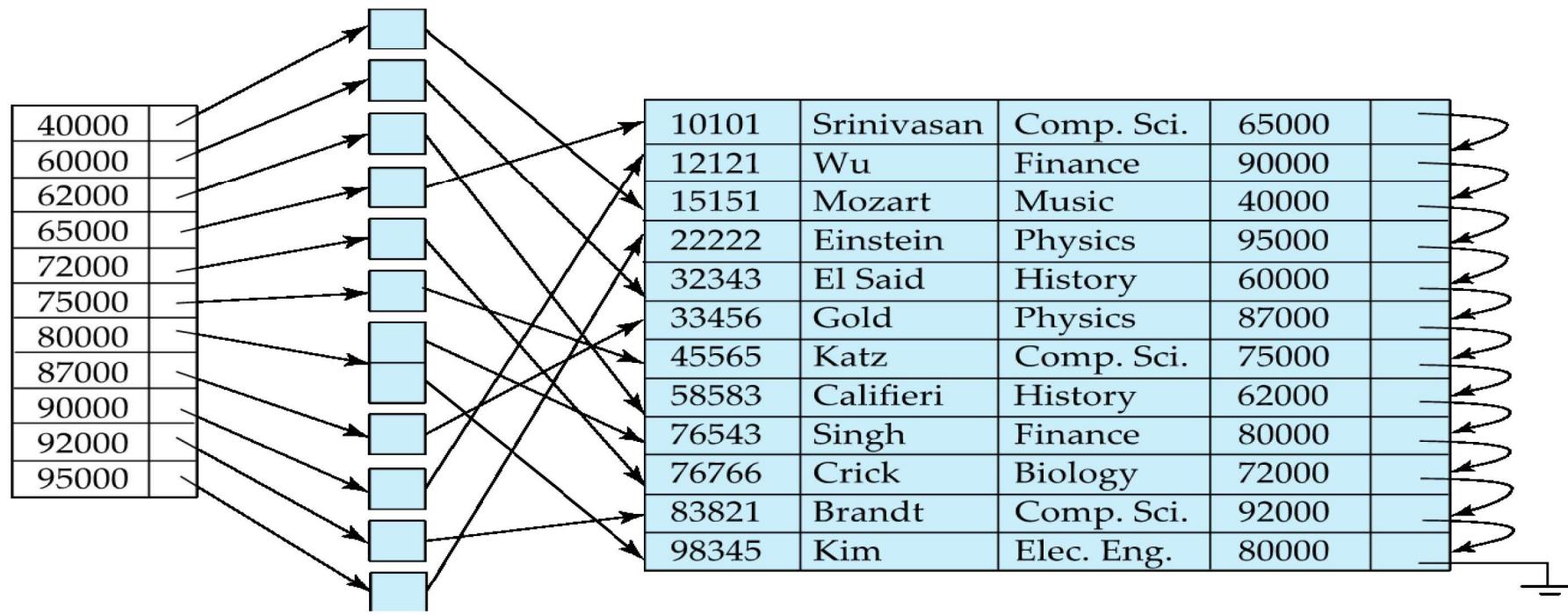
- Compared to dense indices:
 - Less space and less maintenance overhead for insertions and deletions.
 - Generally slower than dense index for locating records.
- **Good tradeoff:** sparse index with an index entry for every block in file, corresponding to least search-key value in the block.



Secondary Indices Example

- 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

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





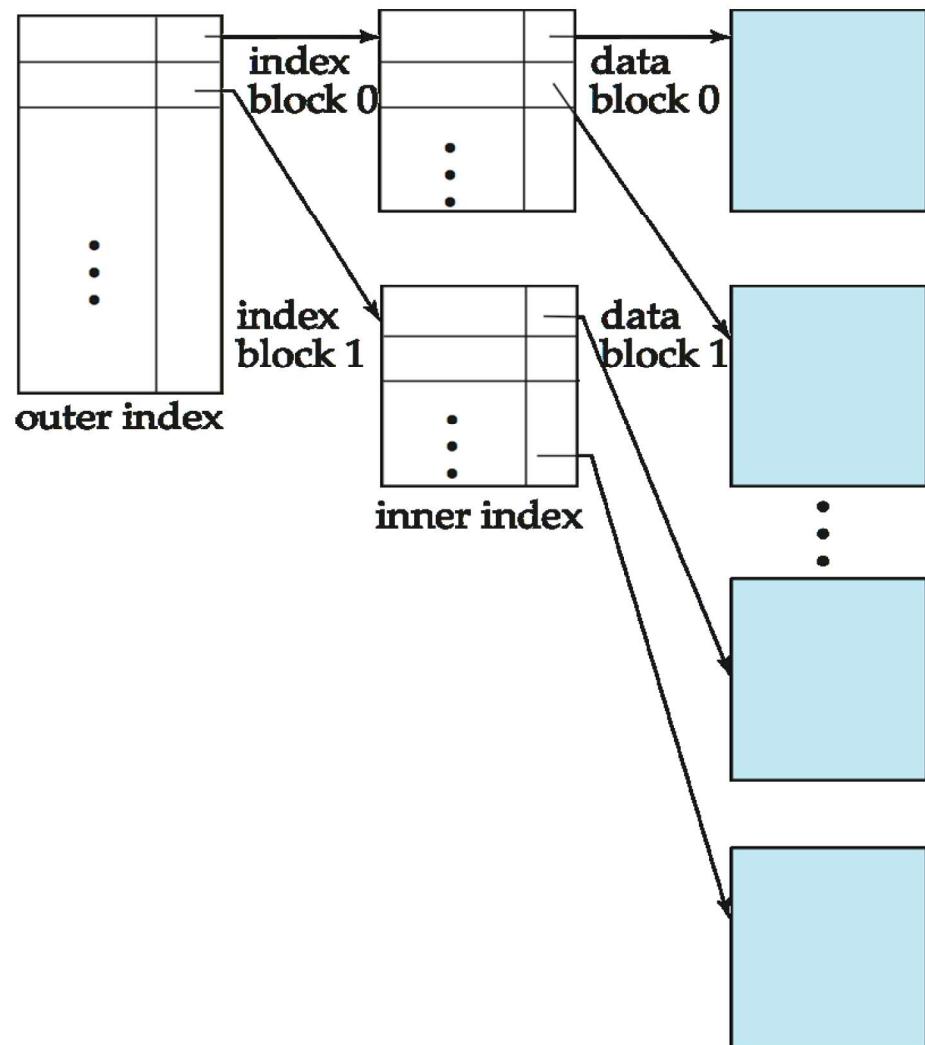
Primary and Secondary Indices

- 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

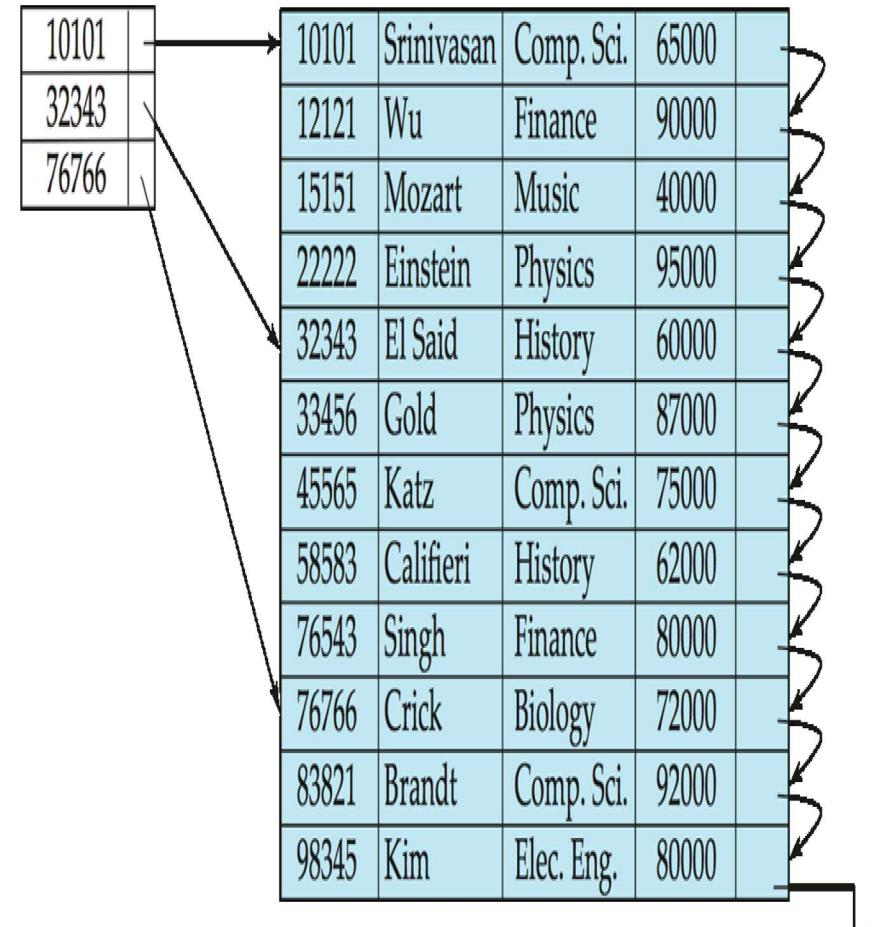


- If primary index does not fit in memory, access becomes expensive.
- Solution: treat primary index kept on disk as a sequential file and construct a sparse index on it.
 - outer index – a sparse index of primary index
 - inner index – the primary index file
- If even outer index is too large to fit 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.



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.



10101	Srinivasan	Comp. Sci.	65000
32343	Wu	Finance	90000
76766	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



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

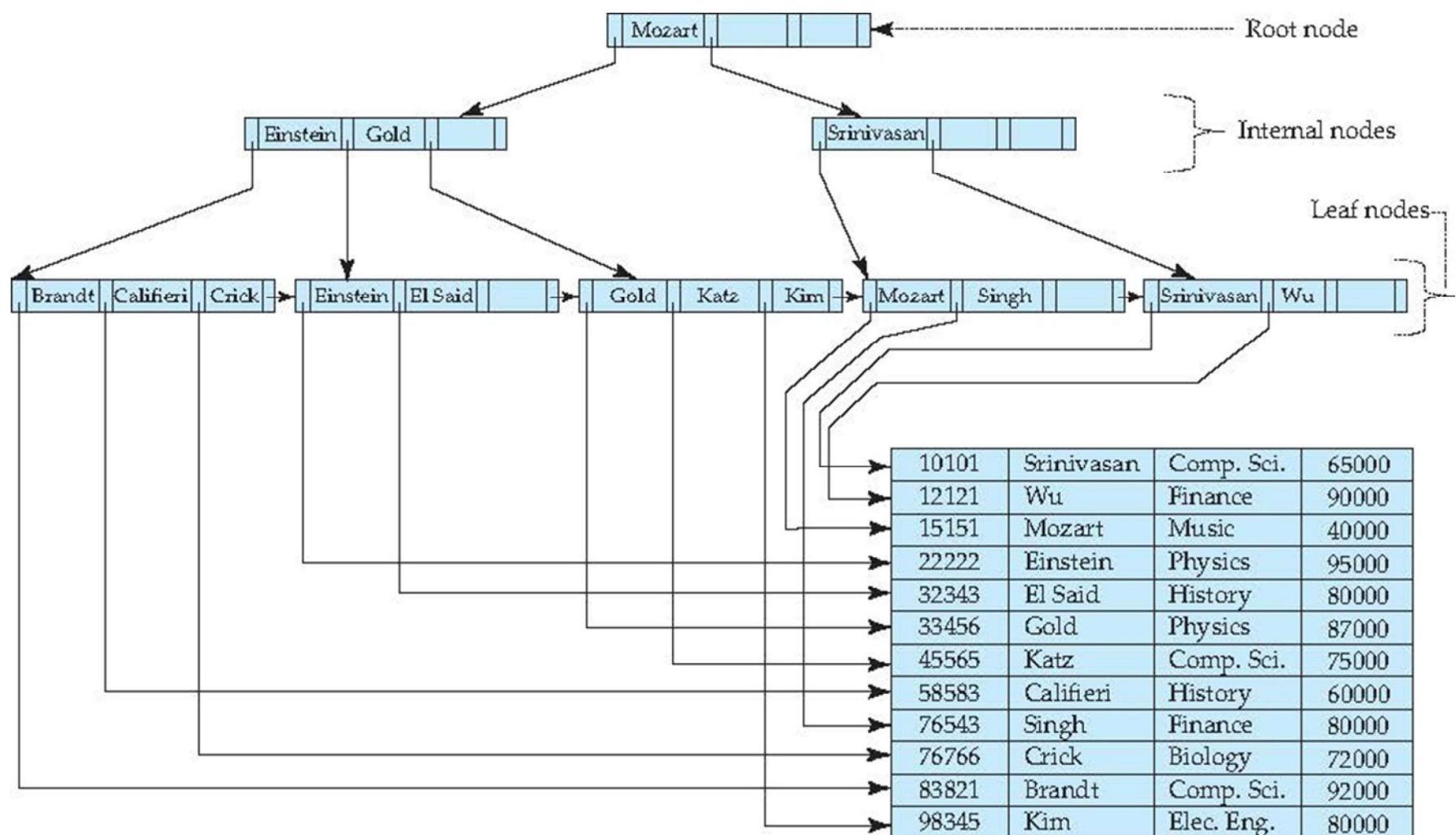
- Frequently, one wants to find all the records whose values in a certain field (which is not the search-key of the primary index) satisfy some condition.
 - Example 1: In the *instructor* relation stored sequentially by ID, we may want to find all instructors in a particular department
 - Example 2: 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⁺-Tree Index Files

- B⁺-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⁺-tree index files:
 - automatically reorganizes itself with small, local, changes, in the cases of insertions and deletions.
 - Reorganization of entire file is not required to maintain performance.
- (Minor) disadvantage of B⁺-trees:
 - extra insertion and deletion overhead, space overhead.
- Advantages of B⁺-trees outweigh disadvantages
 - B⁺-trees are used extensively

Example of B+ tree





B⁺-Tree Index Files (Cont.)

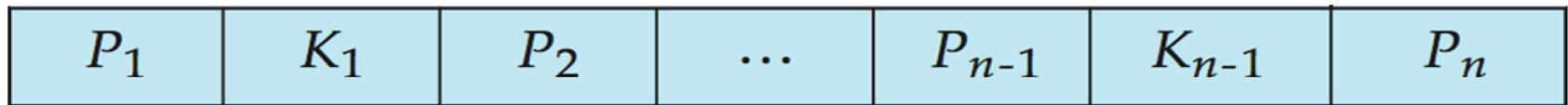
A B⁺-tree is a rooted tree, satisfying the following properties:

- 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



- 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

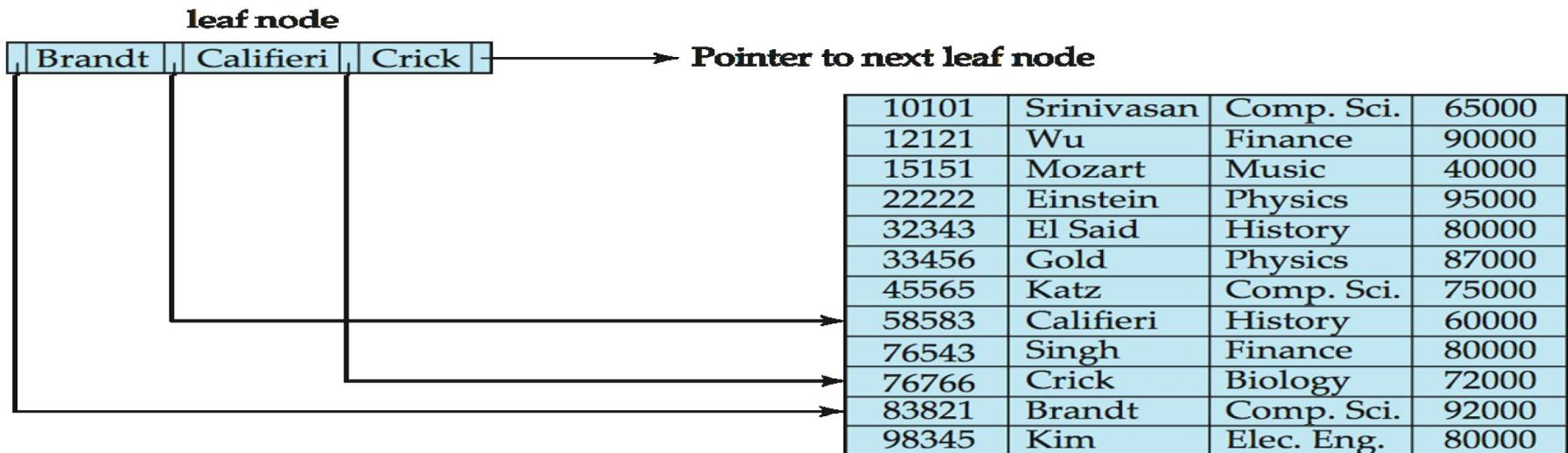
$$K_1 < K_2 < K_3 < \dots < K_{n-1}$$

(Initially assume no duplicate keys, discussions on duplicates later)

Leaf Nodes in B⁺-Trees

Properties of 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⁺-Trees

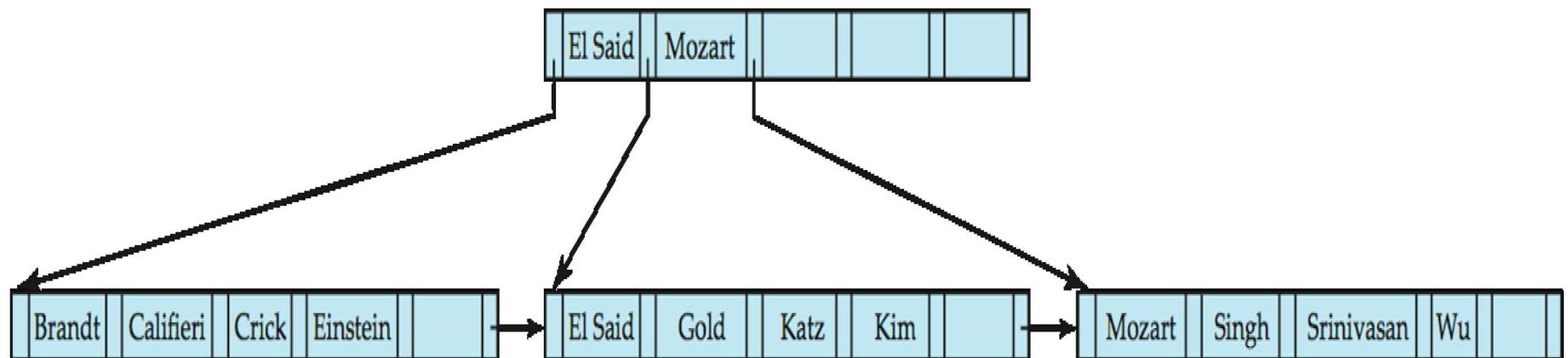
- Non leaf nodes form a multi-level sparse index on the leaf nodes. For a non-leaf node with m pointers:
 - All the search-keys in the subtree to which P_1 points are less than K_1
 - For $2 \leq i \leq n - 1$, all the search-keys in the subtree to which P_i points have values greater than or equal to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

B⁺-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.



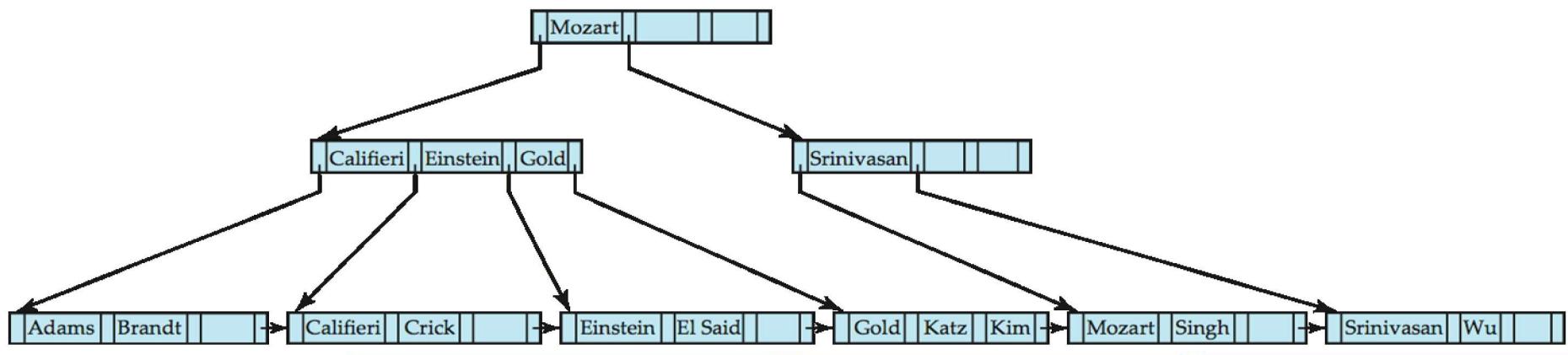
Observations about B⁺-trees



- 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. (i.e. Geometric progression)
 - 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
-

Queries on B⁺-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 $C = \text{right sibling of } C$ 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

Queries on B⁺-Trees (Cont.)



- If there are K search-key values in the file, the height of the tree is no more than $\lceil \log_{\lceil n/2 \rceil}(K) \rceil$.
- 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.
-

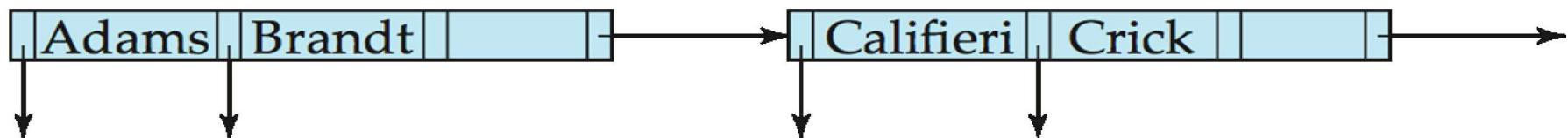
Updates on B⁺-Trees: Insertion (Cont.)

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.

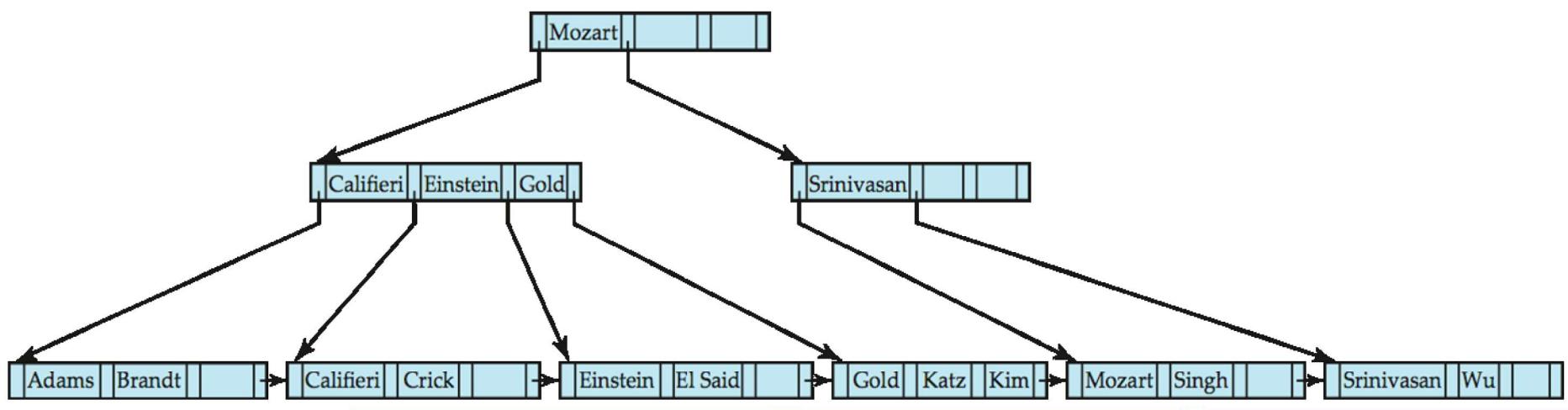
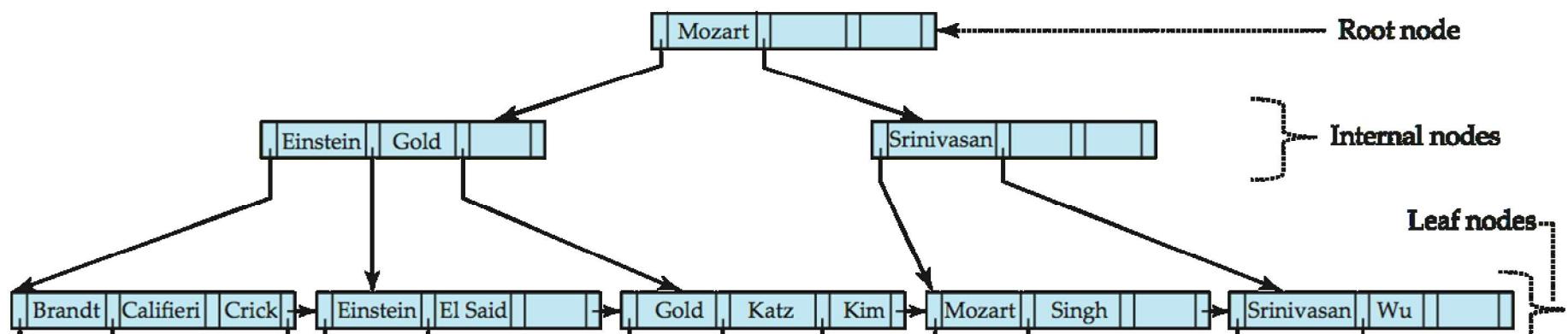
- In the worst case the root node may be split increasing the height of tree by 1.



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

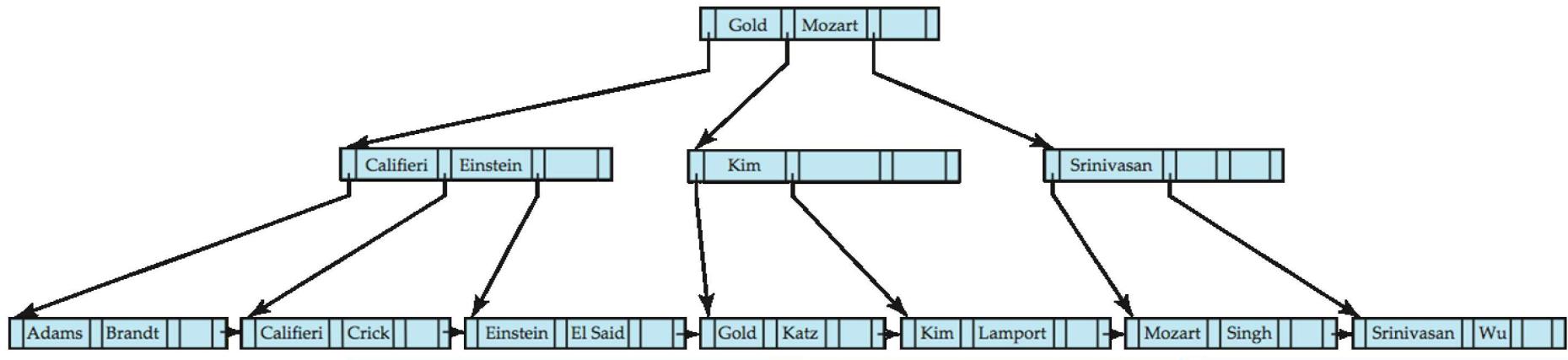
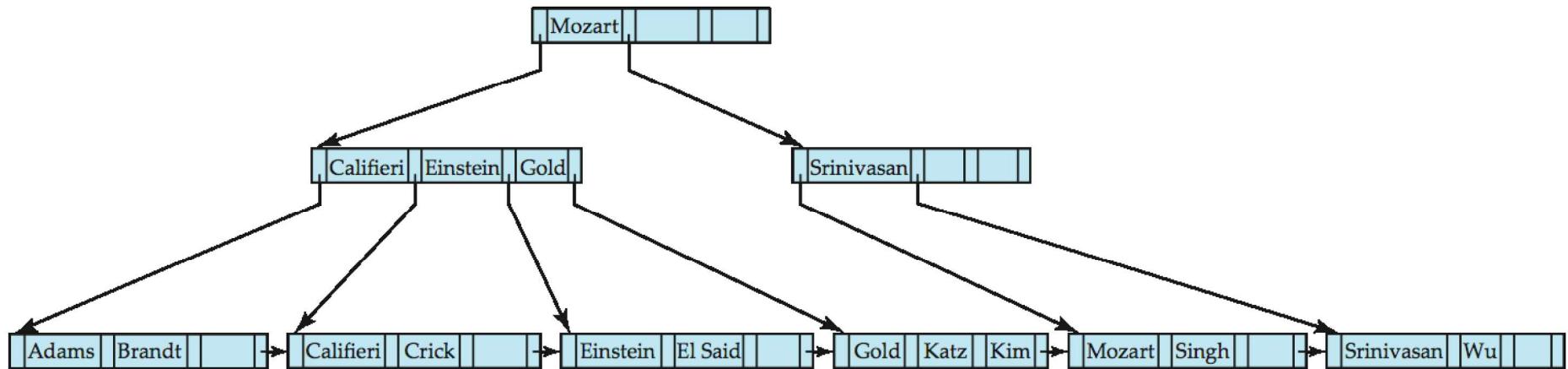
B⁺-Tree Insertion

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



B⁺-Tree Insertion

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

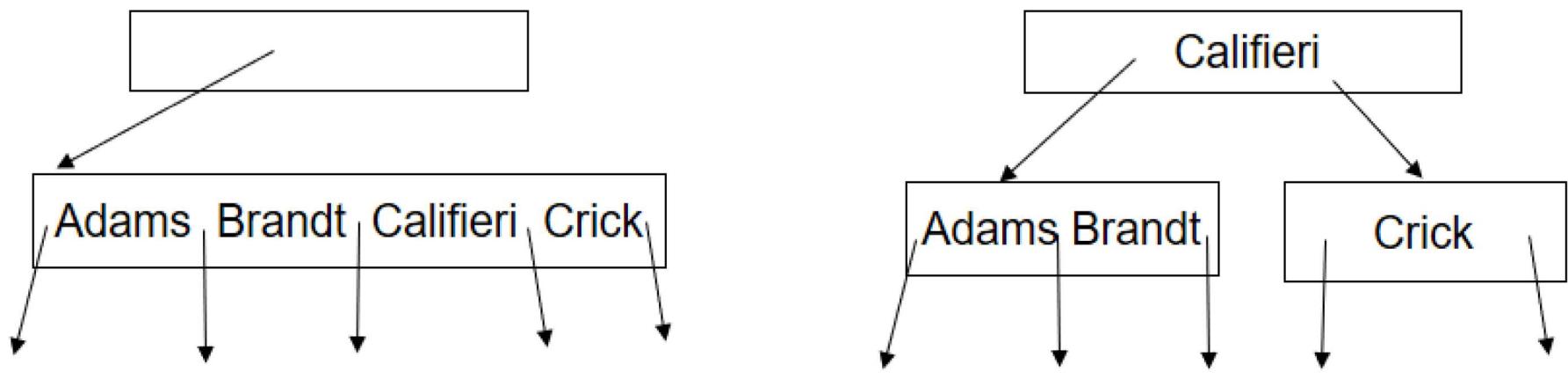


Insertion in B⁺-Trees (Cont.)



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

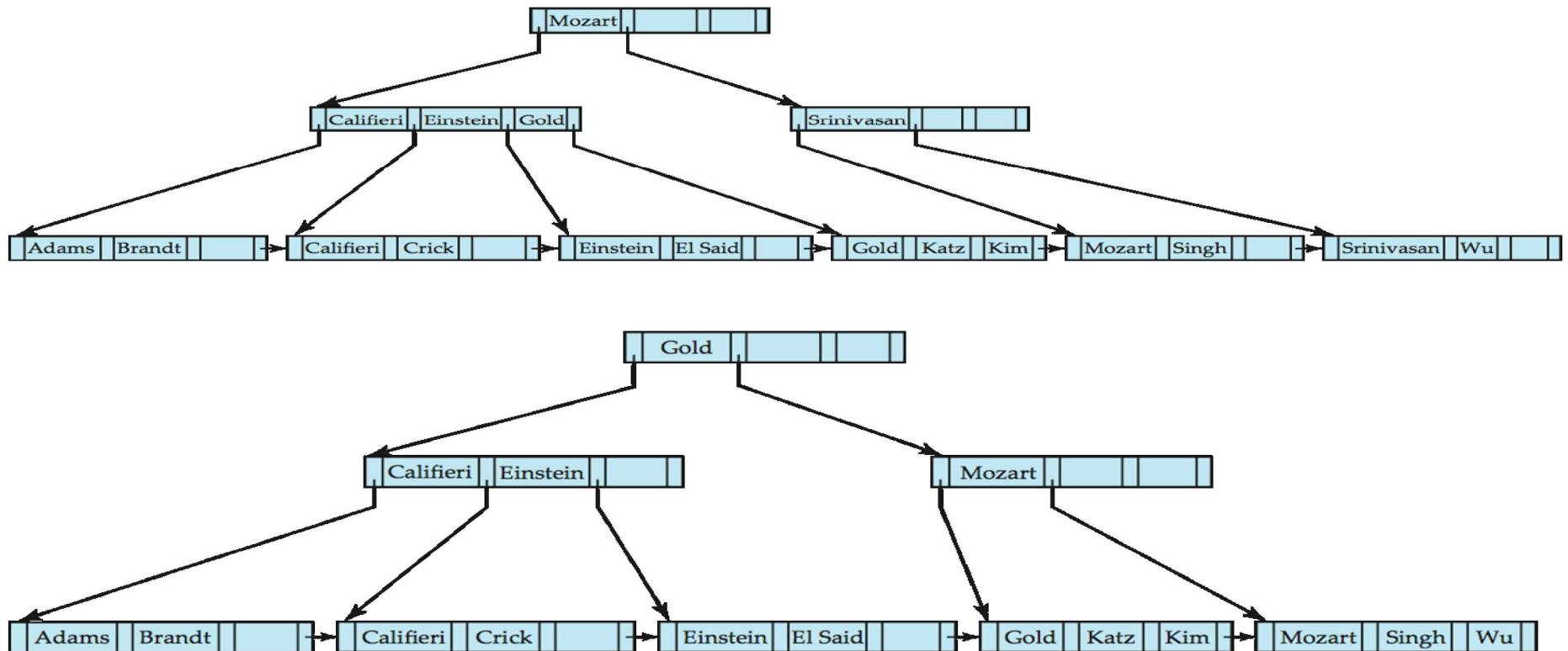


Examples of B⁺-Tree Deletion



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

Before and after deleting “Srinivasan”

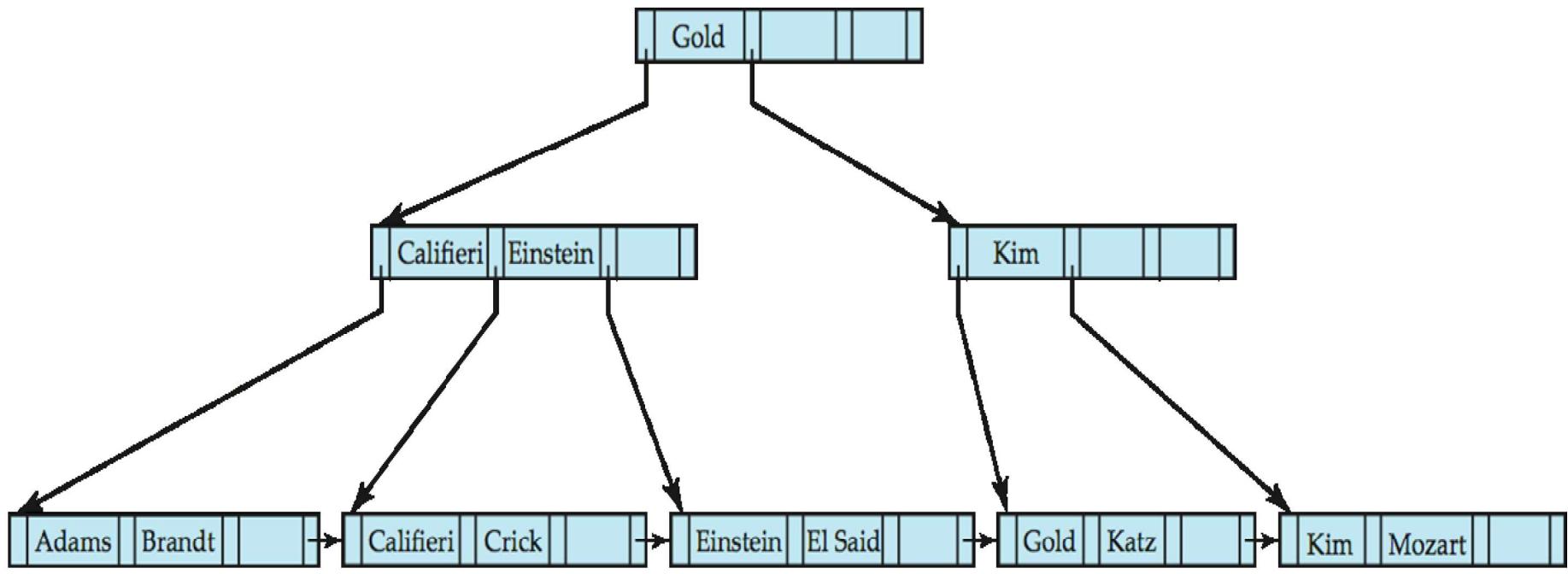


Examples of B⁺-Tree Deletion (Cont.)



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

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

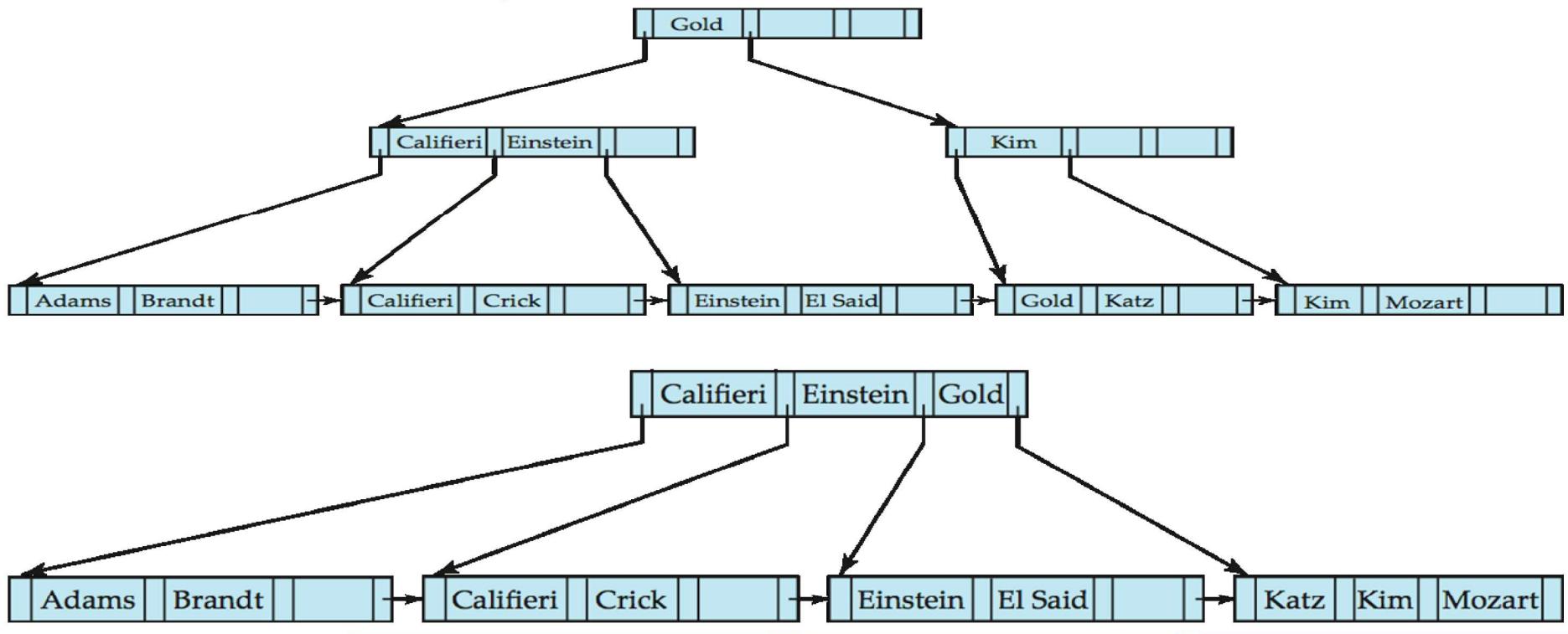


Example of B⁺-tree Deletion (Cont.)



Before and after deletion of “Gold” from earlier example

- Node with Gold and Katz became under-full, and was merged with its sibling
- Parent node becomes under-full, 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

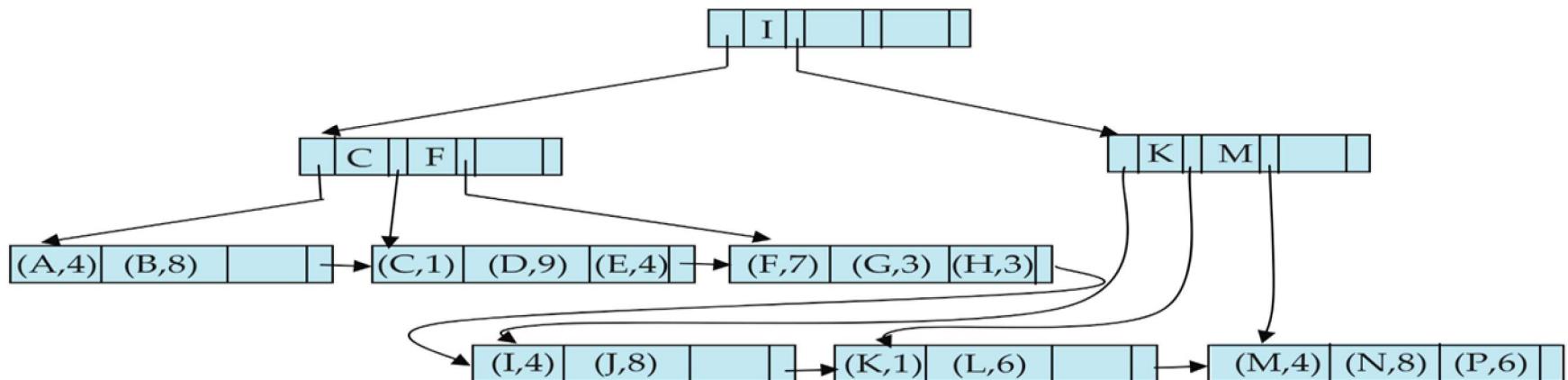
- Index file degradation problem is solved by using B⁺-Tree indices.
- Data file degradation problem is solved by using B⁺-Tree File Organization.
- The leaf nodes in a B⁺-tree file organization store records, instead of pointers.
- Leaf nodes are still required to be half full
 - Since records are larger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

B⁺-Tree File Organization (Cont.)



- Good space utilization must be done.
- 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 $\lfloor 2n/3 \rfloor$ entries.

Example of B⁺-tree File Organization





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⁺-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

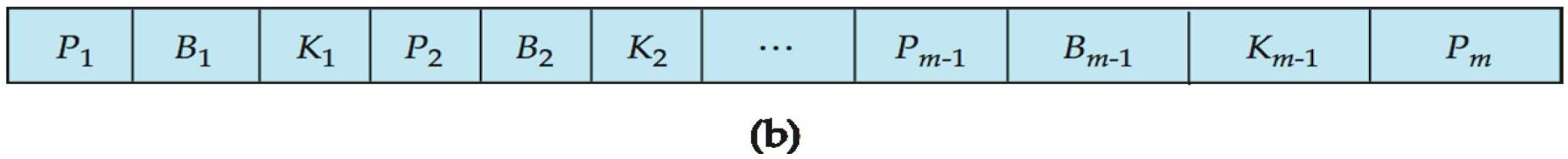
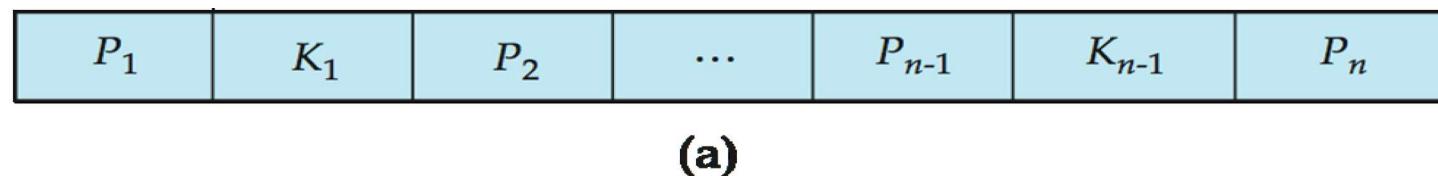
Bulk Loading and Bottom-Up Build



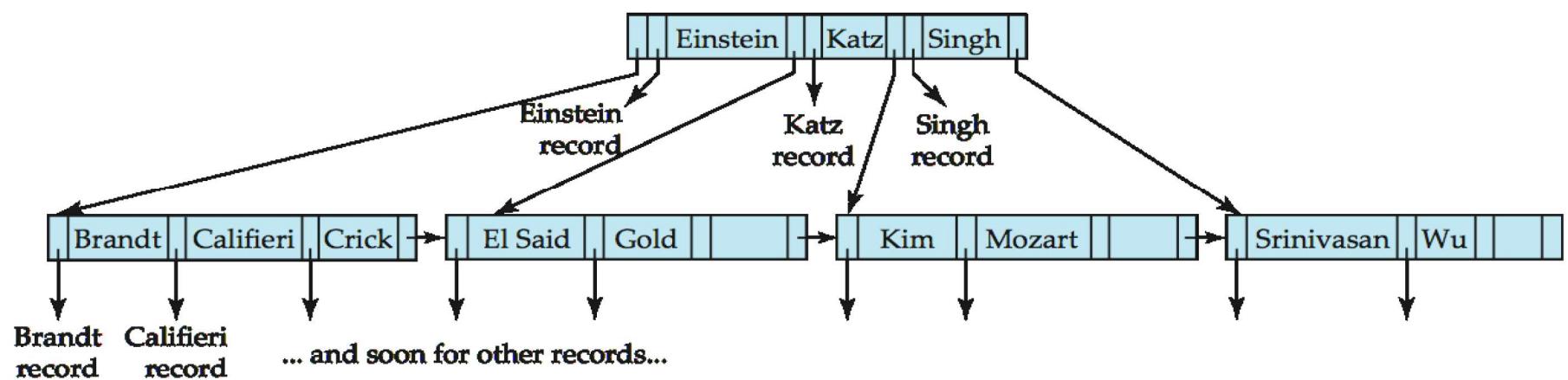
- Inserting entries one-at-a-time into a B⁺-tree requires ≥ 1 IO per entry
 - assuming leaf level does not fit in memory
 - can be very inefficient for loading a large number of entries at a time (**bulk loading**)
 - Efficient alternative 1:
 - sort entries first (using efficient external-memory sort algorithms to be discussed later)
 - 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⁺-tree construction**
 - As before sort entries
 - And then create tree layer-by-layer, starting with leaf level
 - Implemented as part of bulk-load utility by most database systems
-

B-Tree Index Files

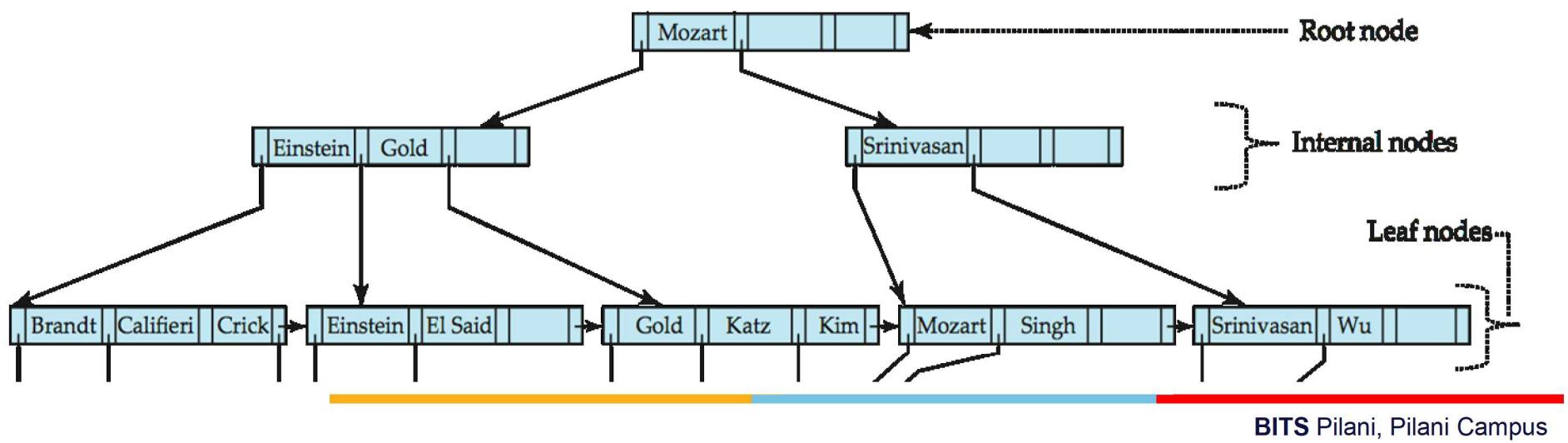
- Similar to B+-tree, but B-tree allows search-key values to appear only once; eliminates redundant storage of search keys.
- Search keys in non-leaf nodes appear nowhere else in the B-tree; an additional pointer field for each search key in a non-leaf node must be included.
- Generalized B-tree leaf node [see (a) below]
- Non-leaf 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⁺-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⁺-Tree
 - Insertion and deletion more complicated than in B⁺-Trees
 - Implementation is harder than B⁺-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 = “Finance” 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 = “Finance” and salary < 80000
 - But cannot efficiently handle
where dept_name < “Finance” and balance = 80000
 - May fetch many records that satisfy the first but not the second condition
-



Other Features

□ Covering indices

- All the requested columns are available with in the index, the index is covering the query
- 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



Hashing



Introduction

- Hashing is a widely used technique for building indices in main memory;
- Such indices may be transiently created to process a join operation or may be a permanent structure in a main memory database.
- Hashing has also been used as a way of organizing records in a file, although hash file organizations are not very widely used.
- We consider:
 1. In-memory hash indices, and
 2. Disk-based hashing



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

- 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. $h(\text{Music}) = 1$ $h(\text{History}) = 2$
 $h(\text{Physics}) = 3$ $h(\text{Elec. Eng.}) = 3$

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 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.



Polynomial rolling hash function

- $\text{hash} = (s[0]*P^0 + s[1]*P^1 + \dots + s[m]*P^m) \bmod M$
- where P and M are some positive numbers.
- The $s[0], s[1], s[2] \dots s[n-1]$ are the values assigned to each character in English alphabet ($a \rightarrow 1, b \rightarrow 2, \dots z \rightarrow 26$).
- P: The value of P can be any prime number roughly equal to the number of different characters used. For example: if the input string contains only lowercase letters of the English alphabet, then $P = 31$ is the appropriate value of P.

If the input string contains both uppercase and lowercase letters, then $P = 53$ is an appropriate option.

M: the probability of two random strings colliding is inversely proportional to m, Hence m should be a large prime number.

$M = 10^9 + 9$ is a good choice.



Applying hash function to a character string K

```
Temp = 1
/* K = string key */
For i=1 to 20 do
    temp = temp*code(K[i]) mod M;
    /* code(x) returns code e.g. ANSI/ASCII of x */
hash_address = temp mode M
```



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***.



Overflow chaining

- **Overflow chaining** = closed addressing = closed hashing
- If the bucket does not have enough space, a **bucket overflow** is said to occur. We handle bucket overflow by using overflow buckets (many chained overflow buckets possible for a usual bucket).
- So that the probability of bucket overflow is reduced, the number of buckets is chosen to be $(n_r f_r) * (1 + d)$, where n_r denotes the number of records, f_r denotes the number of records per bucket, d is a fudge factor, typically around 0.2. With a fudge factor of 0.2, about 20 percent of the space in the buckets will be empty.
- The hash index can be rebuilt with an increased number of buckets. This may cause disruption of normal processing when there are large number of records.
- So, the number of buckets to be increased in a more incremental fashion. Such schemes are called **dynamic hashing** (e.g. linear hashing and extendable hashing) techniques.

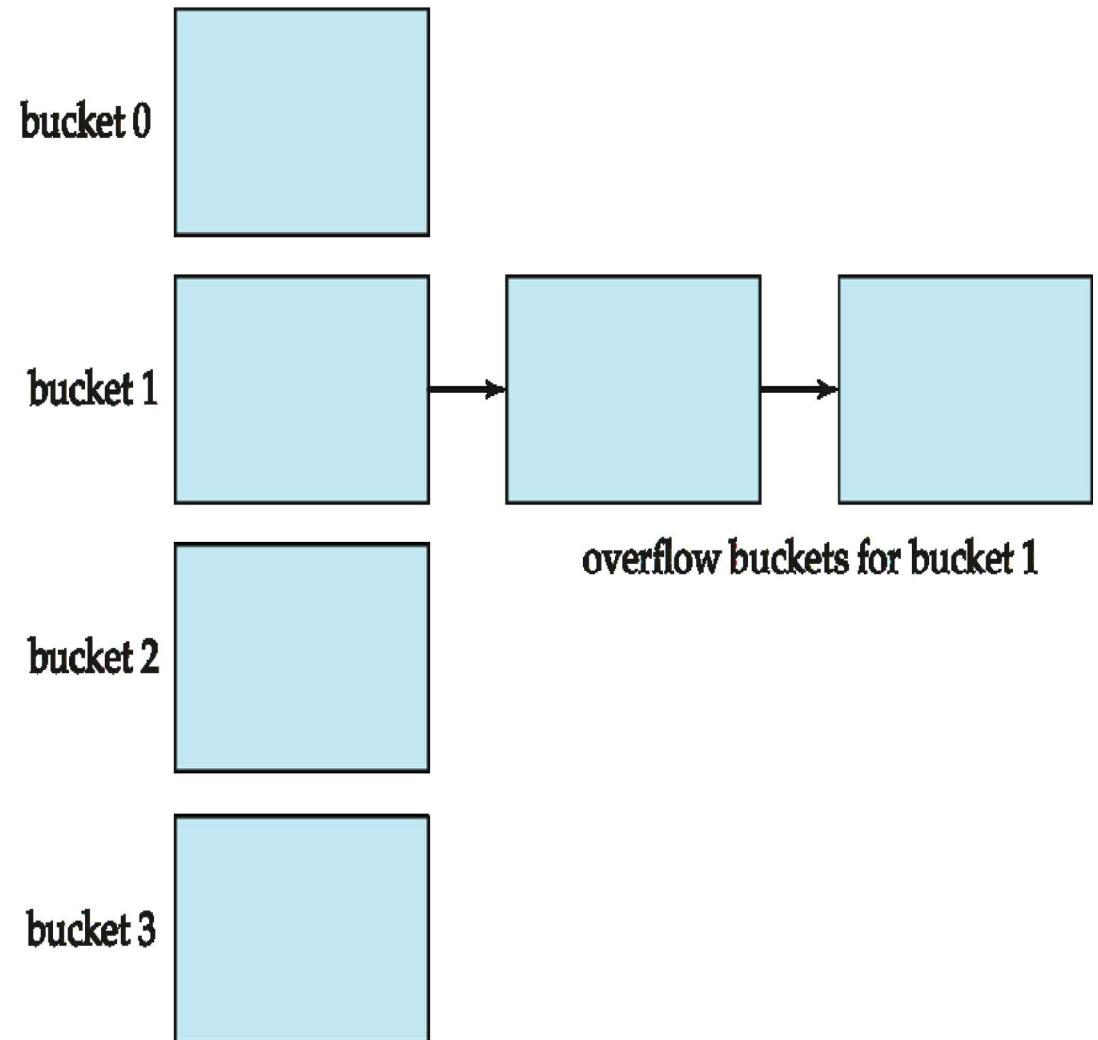
Handling of Bucket Overflows (Cont.)

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.

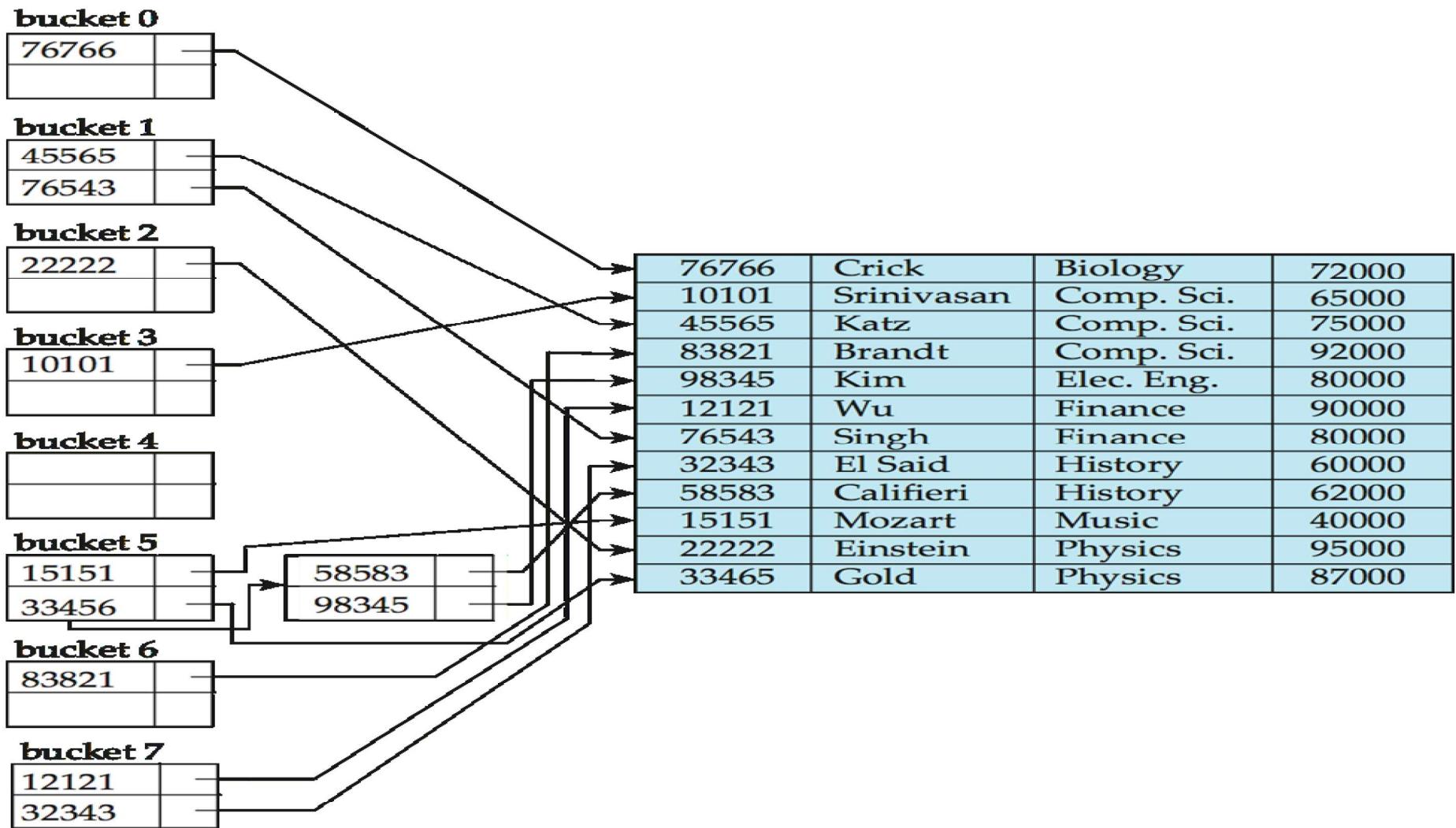




Hashing for Integers

- The most commonly used method for hashing integers is called *modular hashing*.
- We choose the number of buckets, M, to be prime, and, for any positive integer key k, compute the remainder when dividing k by M.
- This function is very easy to compute ($k \% M$, and is effective in dispersing the keys evenly between 0 and M-1).
- **The multiplication method:** In multiplication method, we multiply the key k by a constant real number c in the range $0 < c < 1$ and extract the *fractional part of $k * c$* .
- Then we multiply this value by table_size m and take the floor of the result.

Example of Hash Index: hash index on *instructor*, on attribute *ID*



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.



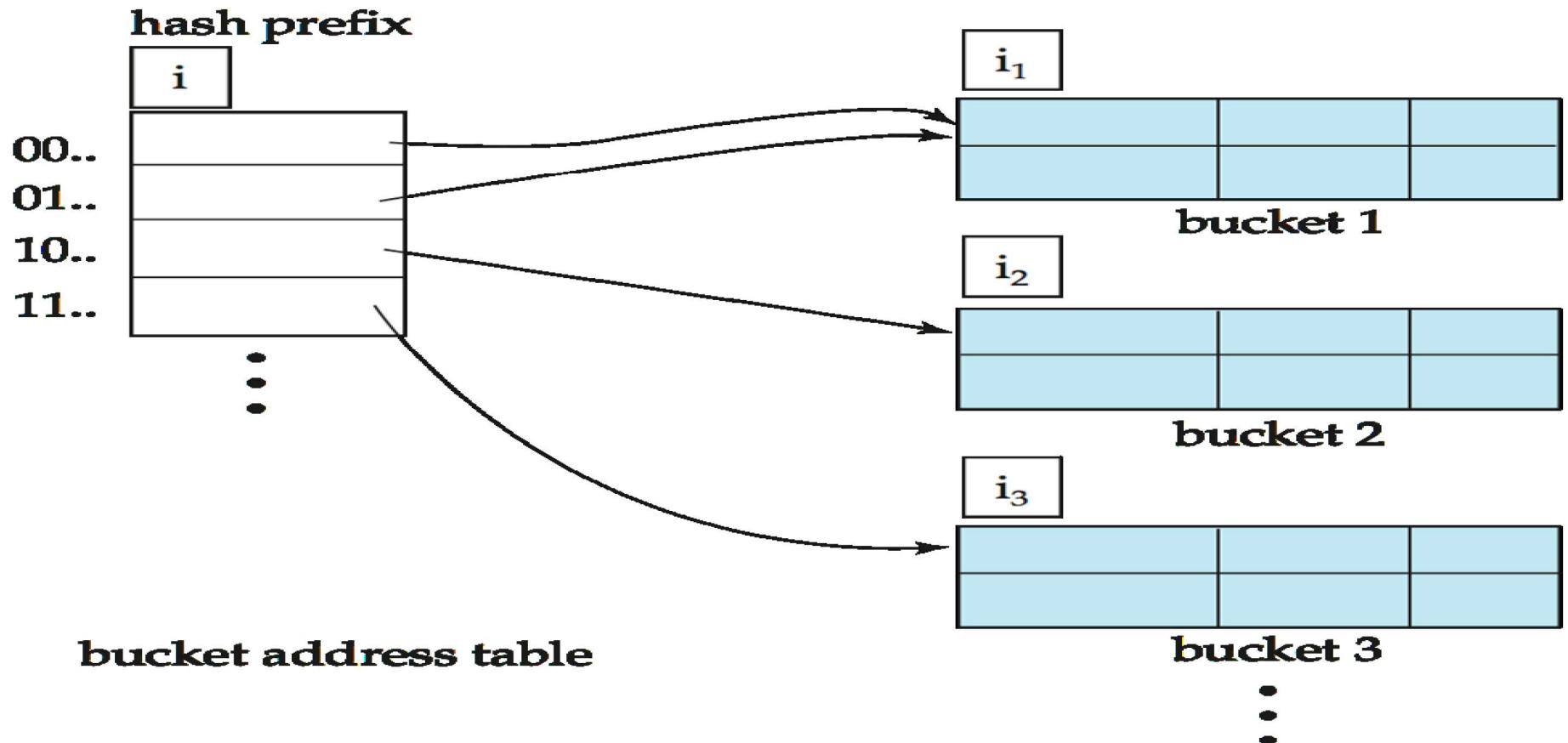
Dynamic Hashing

- 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 (i.e. merging in some way) and splitting of buckets.

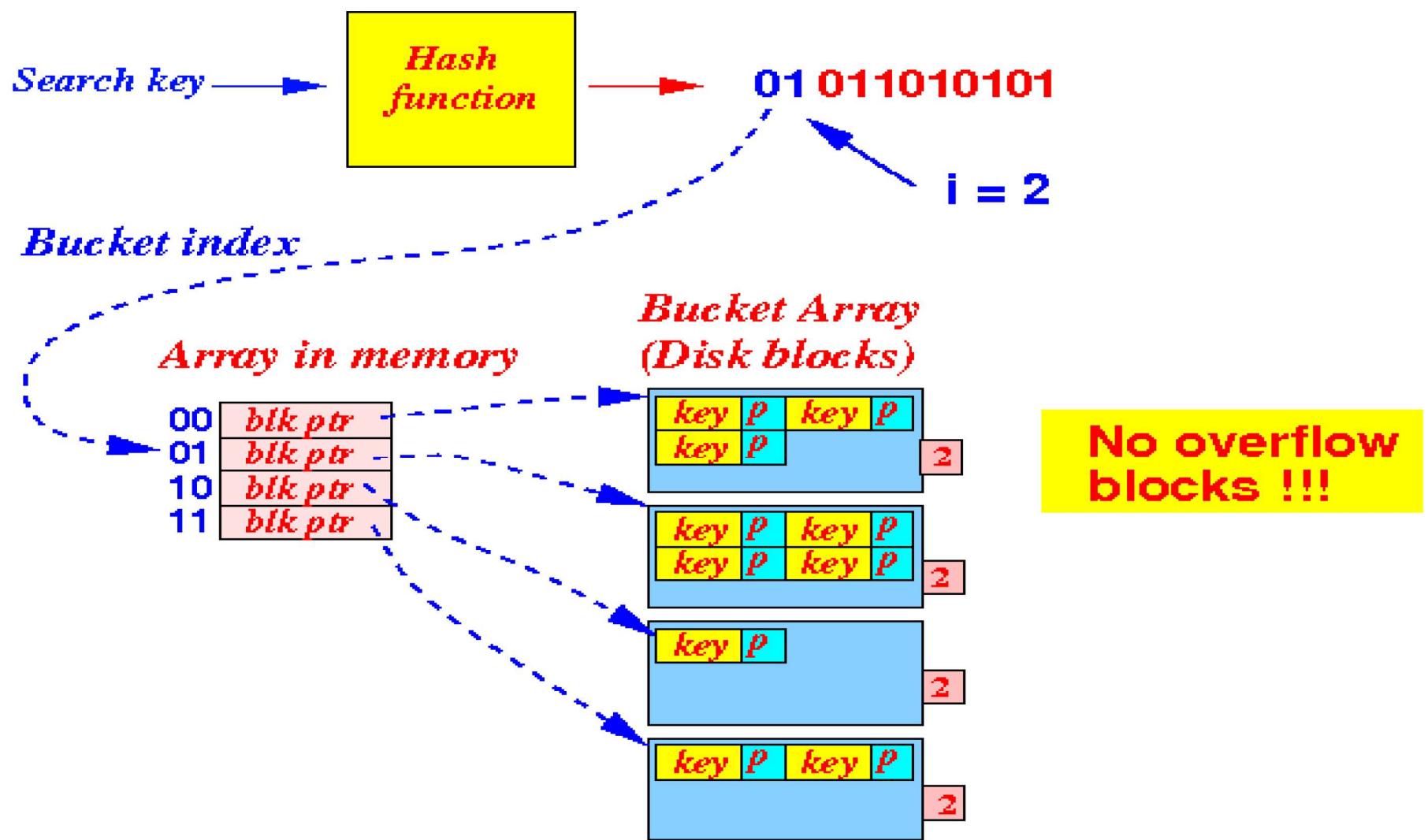
General Extendable Hash Structure



In this structure, $i_2 = i_3 = i$, whereas $i_1 = i - 1$.



Extendable Hash Indexing



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 $i_z = i_j + 1$
 - 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

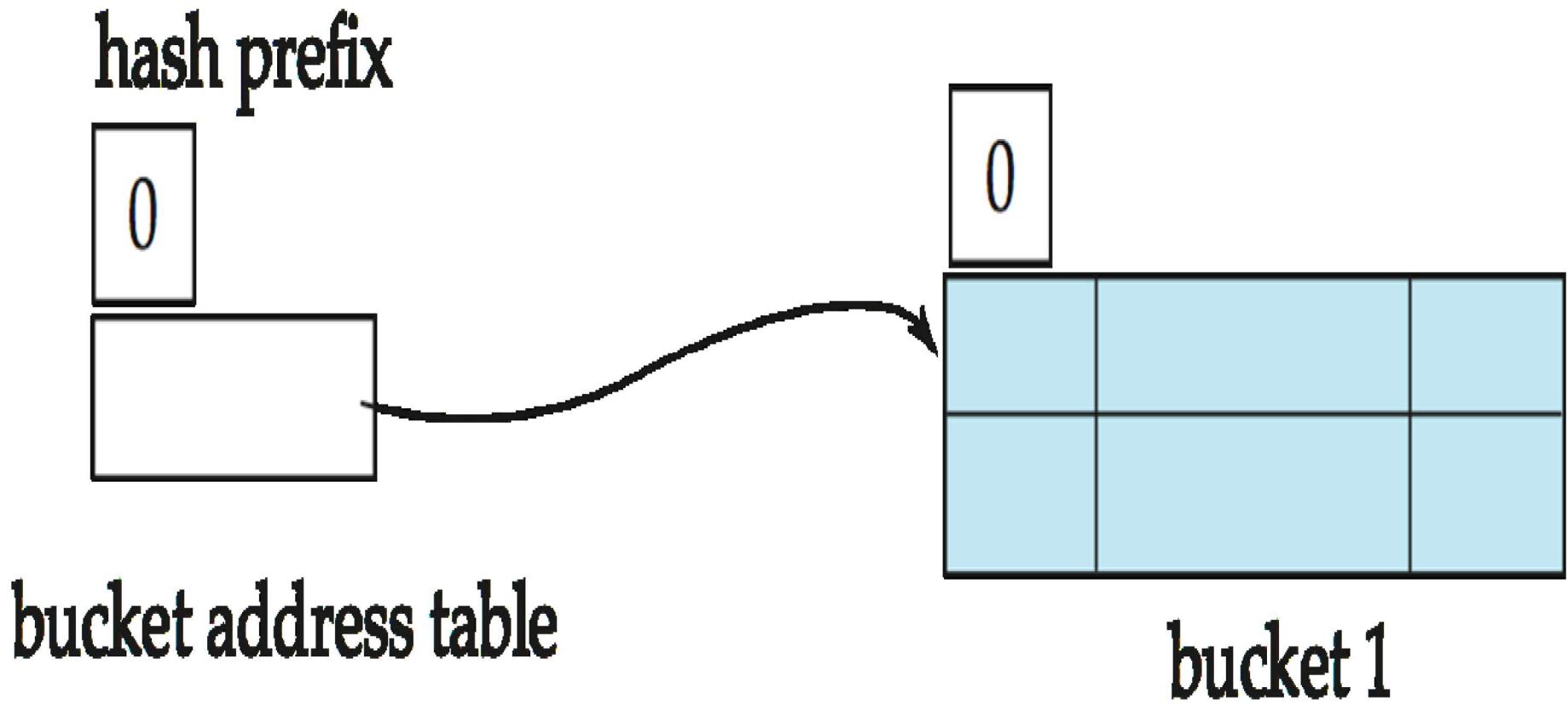


Use of Extendable Hash Structure: Example

<i>dept_name</i>	$h(dept_name)$
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

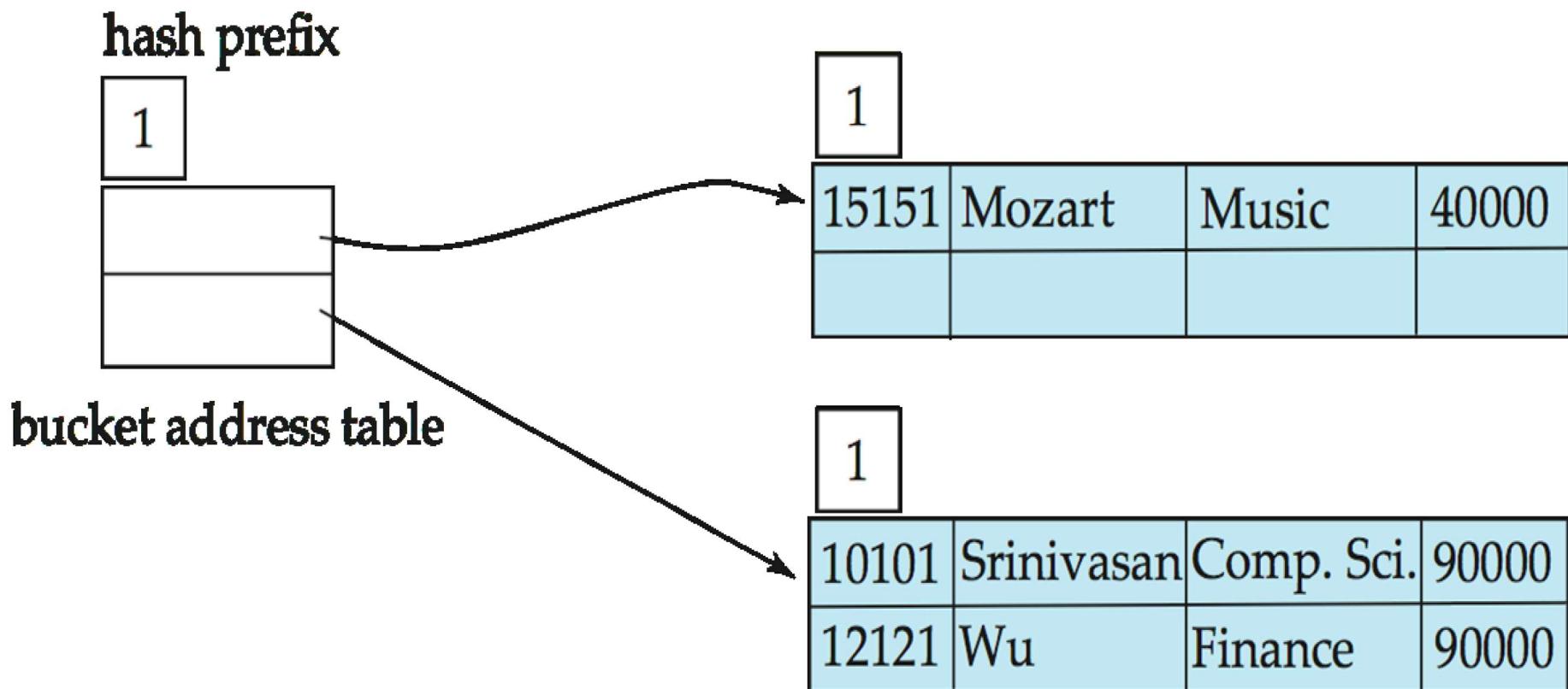
Example (cont.)

Initial Hash structure; bucket size = 2



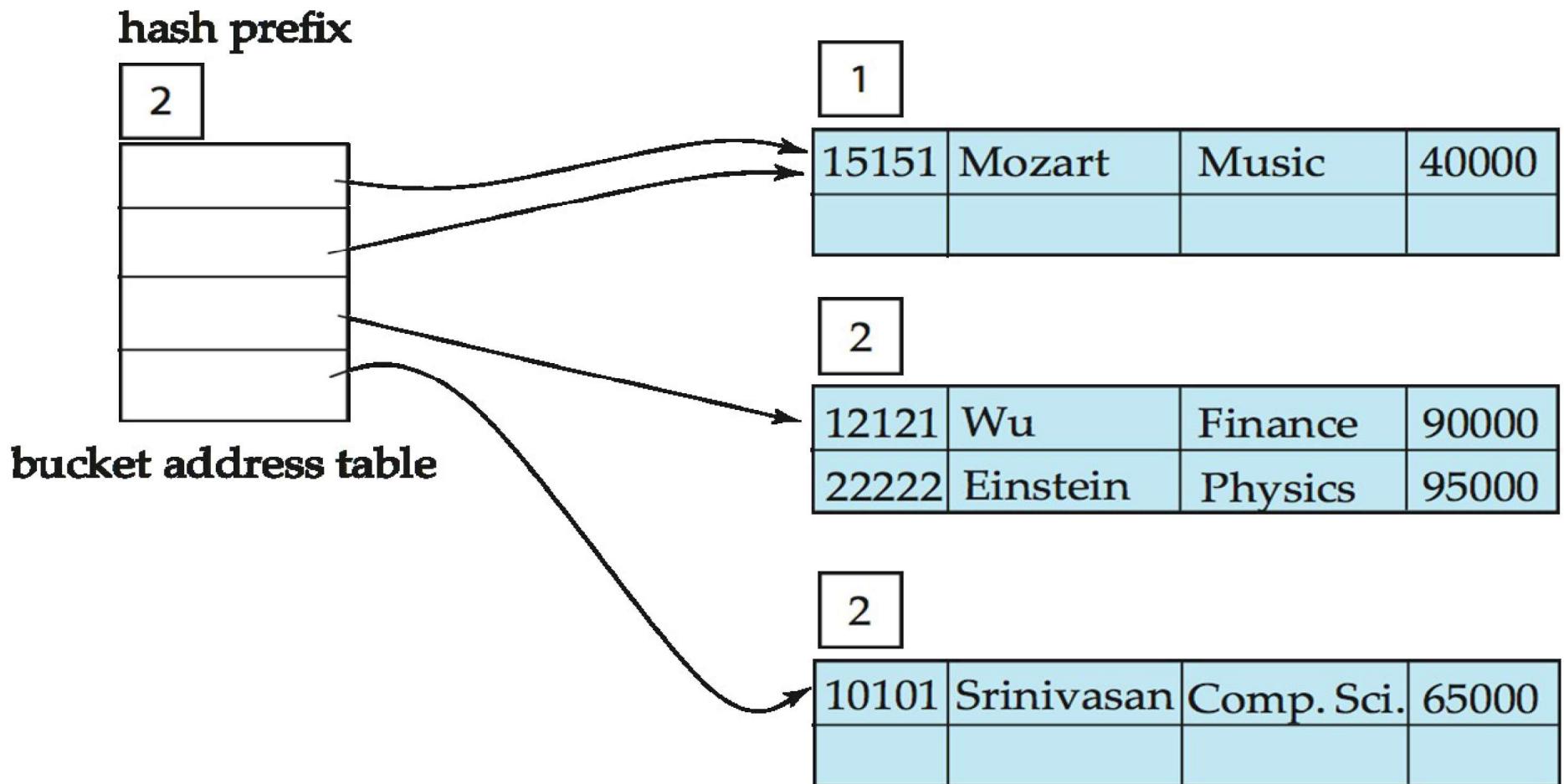
Example (cont.)

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



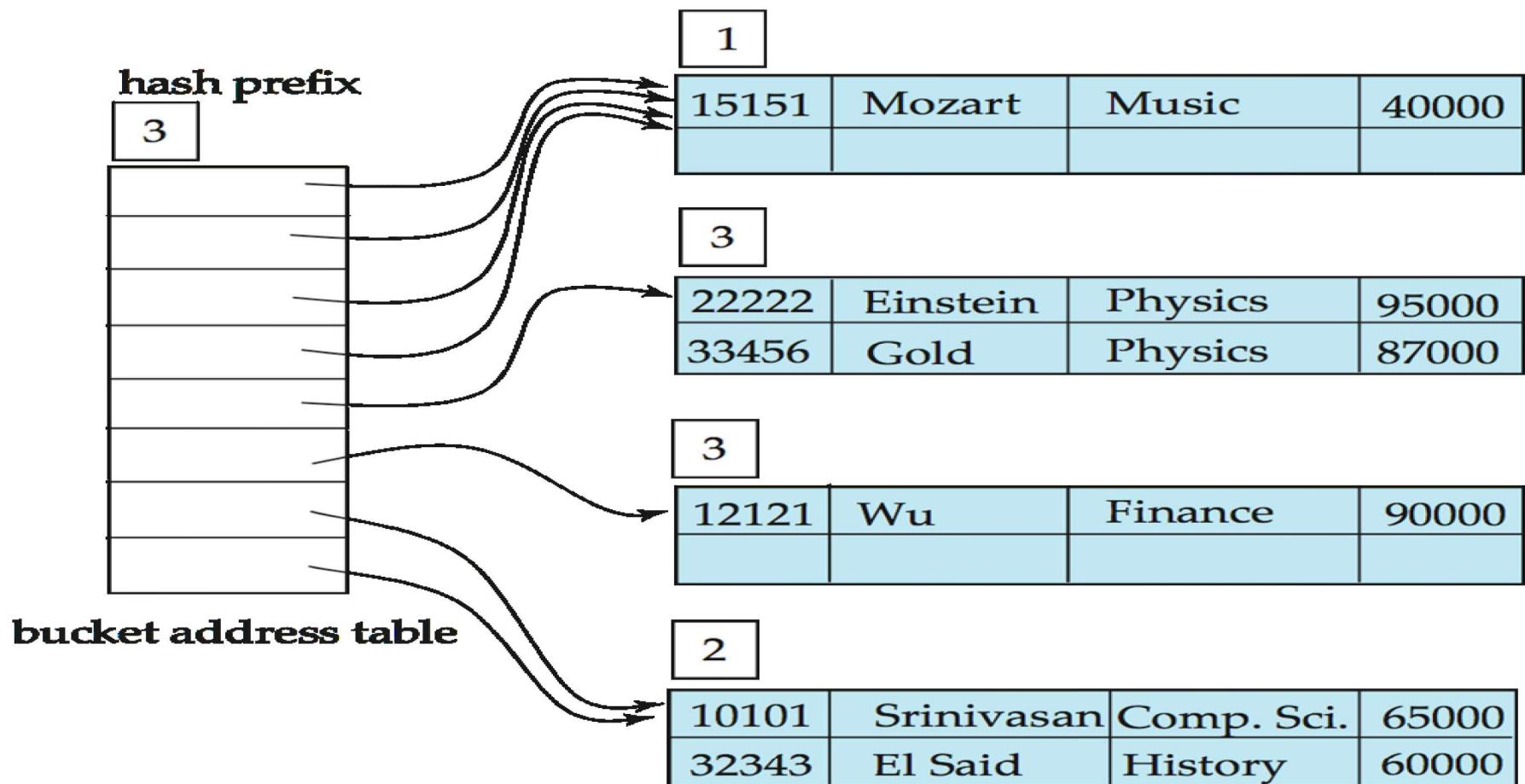
Example (cont.)

Hash structure after insertion of Einstein record



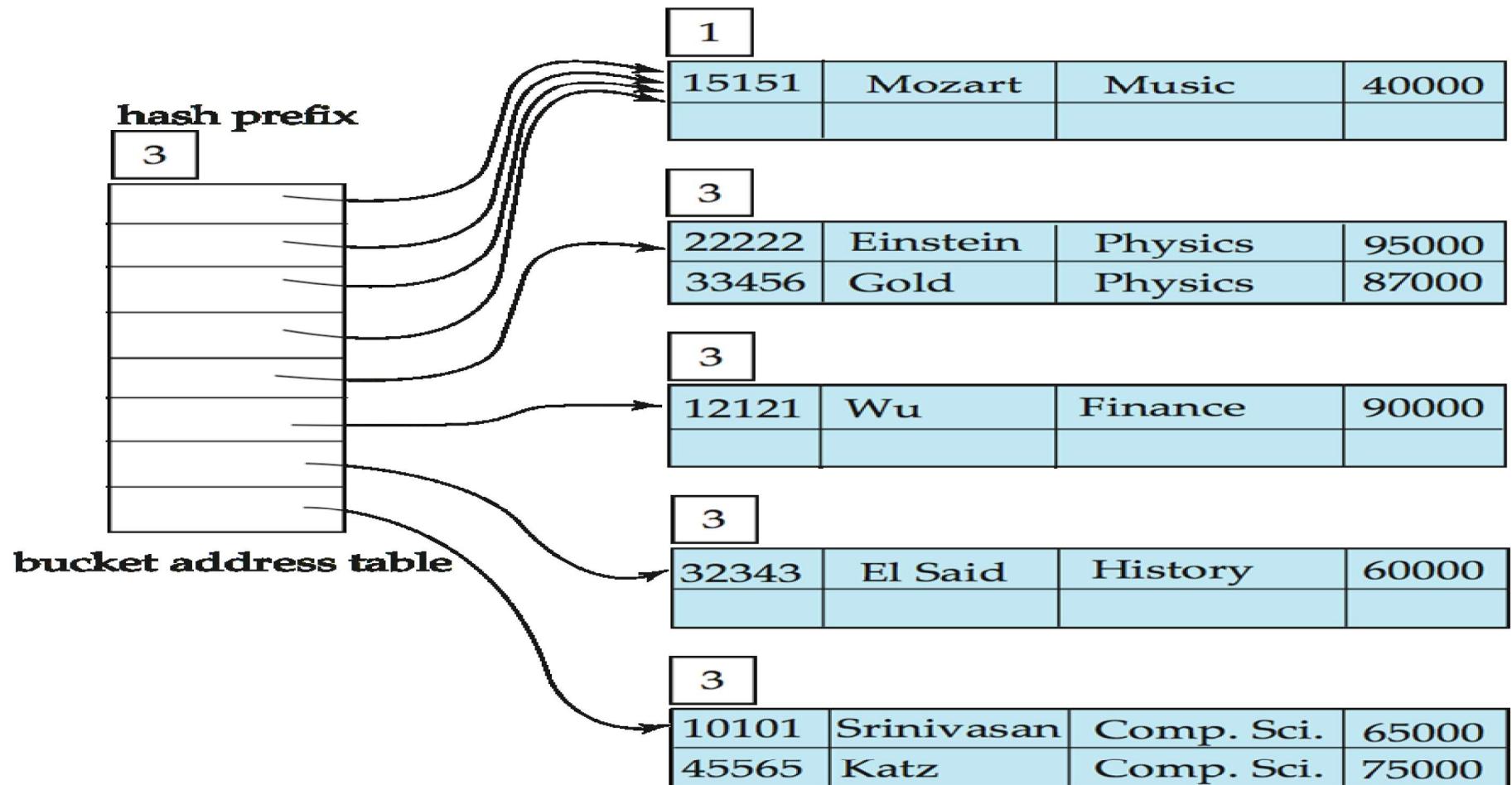
Example (cont.)

Hash structure after insertion of Gold and El Said records



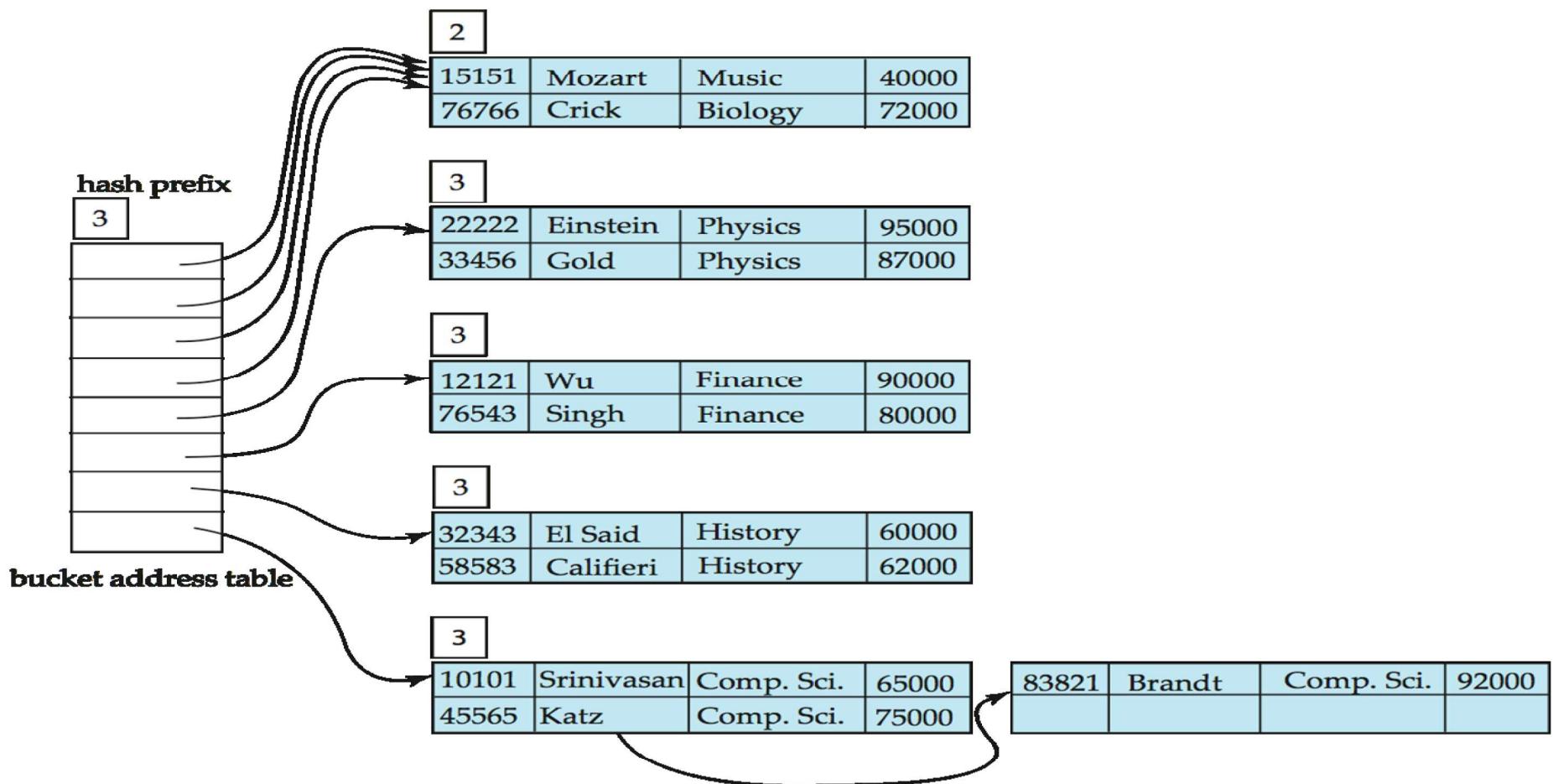
Example (cont.)

Hash structure after insertion of Katz record



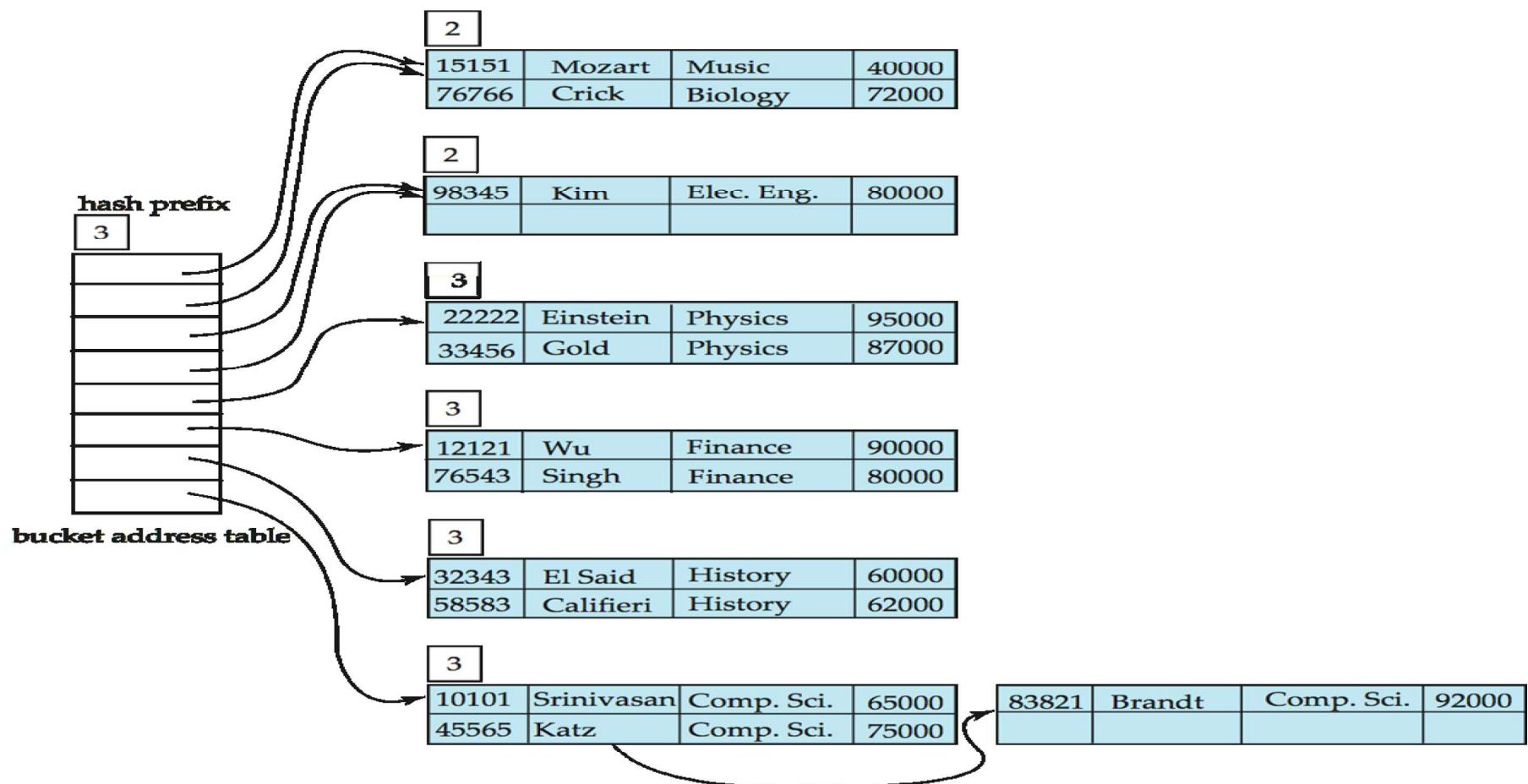
Example (cont.)

And after insertion of eleven records



Example (cont.)

And after insertion of Kim record in previous hash structure





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⁺-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 and **SQLServer** supports only B⁺-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 <i>gender</i>		Bitmaps for <i>income_level</i>	
m	10010	L1	10100
f	01101	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⁺-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⁺-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.
 - 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.



Partitioned Hashing

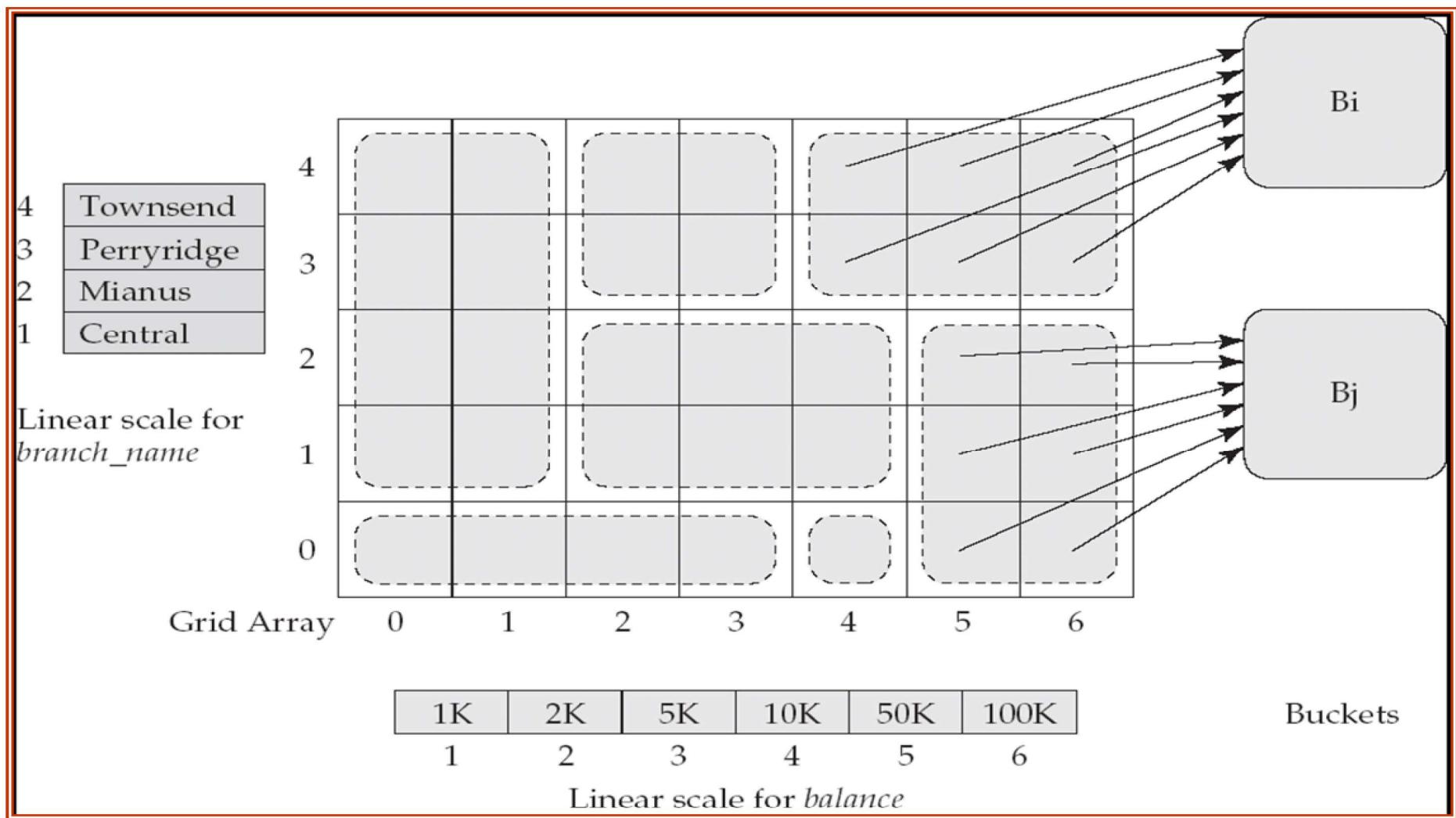
- Hash values are split into segments that depend on each attribute of the search-key.
- (A_1, A_2, \dots, A_n) for n attribute search-key
- Example: $n = 2$, for *customer*, search-key being (*customer-street*, *customer-city*)
- | search-key value | hash value |
|--------------------|------------|
| (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 are an alternative.



Thanks

Any questions please

Thanking you