

# **DATABASE MANAGEMENT SYSTEMS**

## **Module - IV**

### **Physical Design**

#### **Data Storage Structures**



# Outline: Data Storage Structures

- Database Storage Architecture
- Storage Organization (Records, Pages and Files)
  - File Organization
    - Fixed-Length Records
    - Variable-Length Records
    - Storing Large Objects
  - Organization of Records in Files
    - Heap File Organization
    - Sequential File Organization
    - Multitable Clustering File Organization
    - Partitioning
- Database Metadata (Data-Dictionary Storage)
- Database Buffers
  - Buffer Manager
  - Buffer-Replacement Strategies
  - Reordering of Writes and Recovery



# Database Storage Architecture

- **Persistent** (or **non-volatile**) data are stored on **non-volatile storage**, which is typically **magnetic disk** or **SSD (Solid-State Drive)**.
- **Magnetic disks** as well as **SSDs** are **block structured devices**, that is, data are read or written in units of a block.
- In contrast, **databases** deal with **records**, which are usually much smaller than a block (although in some cases **records** may have **attributes** that are very large).



# Database Storage Architecture

- Most databases use **operating system files** as an **intermediate layer** for **storing records**, which abstract away some details of the **underlying blocks**.

## *File Organization*

- However, to ensure **efficient access**, as well as to support **recovery from failures**, **databases** must continue to be **aware of blocks**.
- In ***File Organization*** – we study how individual records are stored in files, taking block structure into account.



# Database Storage Architecture

## *Organization of Records in Files*

- Given a **set of records**, the next decision lies in how to organize them in the **file structure**; for **example**, they may be stored in sorted order, in the order they are created, or in an arbitrary order.
- We study several alternative file organizations



# Database Storage Architecture

## *Data-Dictionary Storage*

- Describes how databases organize data about the relational schemas as well as storage organization, in the **data dictionary**.
- Information in the **data dictionary** is crucial for many tasks, for example, to locate and retrieve records of a relation when given the name of the relation.



# Database Storage Architecture

## *Database Buffer*

- For a **CPU** to access data, it must be in **main memory**, whereas **persistent data** must be resident on **non-volatile storage** such as **magnetic disks** or **SSDs**.
- For **databases** that are larger than **main memory**, which is the usual case, data must be fetched from **nonvolatile storage** and saved back if it is updated.
- **Database Buffer** describes how databases use a region of memory called the database buffer to store blocks that are fetched from **non-volatile storage**.



# Outline: Data Storage Structures

- Database Storage Architecture
- Storage Organization (Records, Pages and Files)
- File Organization
  - Fixed-Length Records
  - Variable-Length Records
  - Storing Large Objects
- Organization of Records in Files
  - Heap File Organization
  - Sequential File Organization
  - Multitable Clustering File Organization
  - Partitioning
- Database Metadata (Data-Dictionary Storage)
- Database Buffers
  - Buffer Manager
  - Buffer-Replacement Strategies
  - Reordering of Writes and Recovery





# 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

- We assume that records are smaller than a disk block

.



# 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

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



# Fixed-Length Records

- 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 3 deleted

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



# Fixed-Length Records

- 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 3 deleted and replaced by record 11

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



# Fixed-Length Records

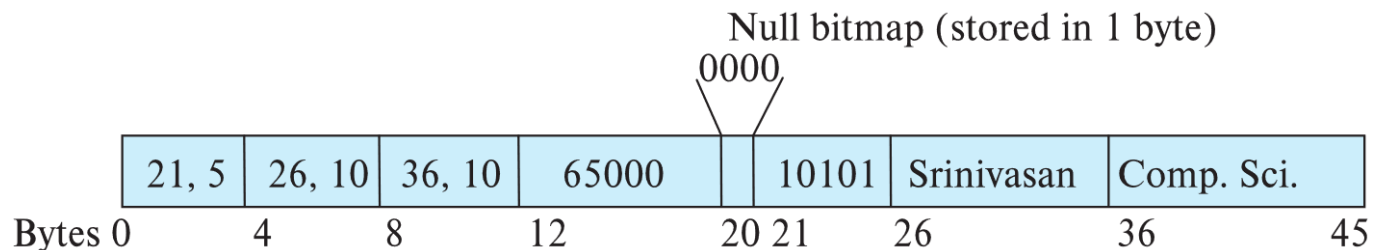
- Deletion of record  $i$ : alternatives:
  - move records  $i + 1, \dots, n$  to  $i, \dots, n - 1$
  - move record  $n$  to  $i$
  - **do not move records, but link all free records on a *free list***

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



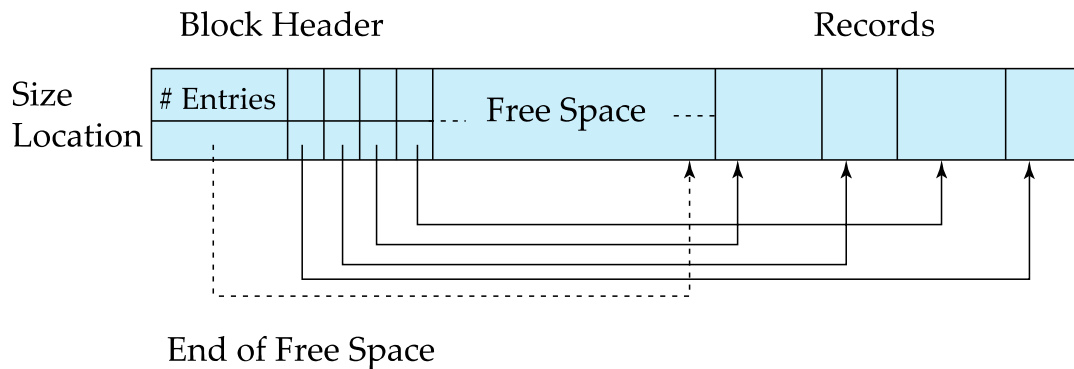
# Variable-Length Records

- 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.



# Storing Large Objects

- E.g., blob/clob types
- Records must be smaller than pages
- Alternatives:
  - Store as files in file systems
  - Store as files managed by database
  - Break into pieces and store in multiple tuples in separate relation
    - PostgreSQL TOAST





# Outline: Data Storage Structures

- Database Storage Architecture
- Storage Organization (Records, Pages and Files)
- File Organization
  - Fixed-Length Records
  - Variable-Length Records
  - Storing Large Objects
- Organization of Records in Files
  - Heap File Organization
  - Sequential File Organization
  - Multitable Clustering File Organization
  - Partitioning
- Database Metadata (Data-Dictionary Storage)
- Database Buffers
  - Buffer Manager
  - Buffer-Replacement Strategies
  - Reordering of Writes and Recovery



# Organization of Records in Files

- **Heap** – 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
- 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
- **B+-tree file organization**
  - Ordered storage even with inserts/deletes
- **Hashing** – a hash function computed on search key; the result specifies in which block of the file the record should be placed



# Heap File Organization

- Records can be placed anywhere in the file where there is free space
- Records usually do not move once allocated
- Important to be able to efficiently find free space within file
- **Free-space map**
  - Array with 1 entry per block. Each entry is a few bits to a byte, and records fraction of block that is free
  - In example below, 3 bits per block, value divided by 8 indicates fraction of block that is free

4	2	1	4	7	3	6	5	1	2	0	1	1	0	5	6
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

- Can have second-level free-space map
- In example below, each entry stores maximum from 4 entries of first-level free-space map

4	7	2	6
---	---	---	---

- Free space map written to disk periodically, OK to have wrong (old) values for some entries (will be detected and fixed)



# Sequential File Organization

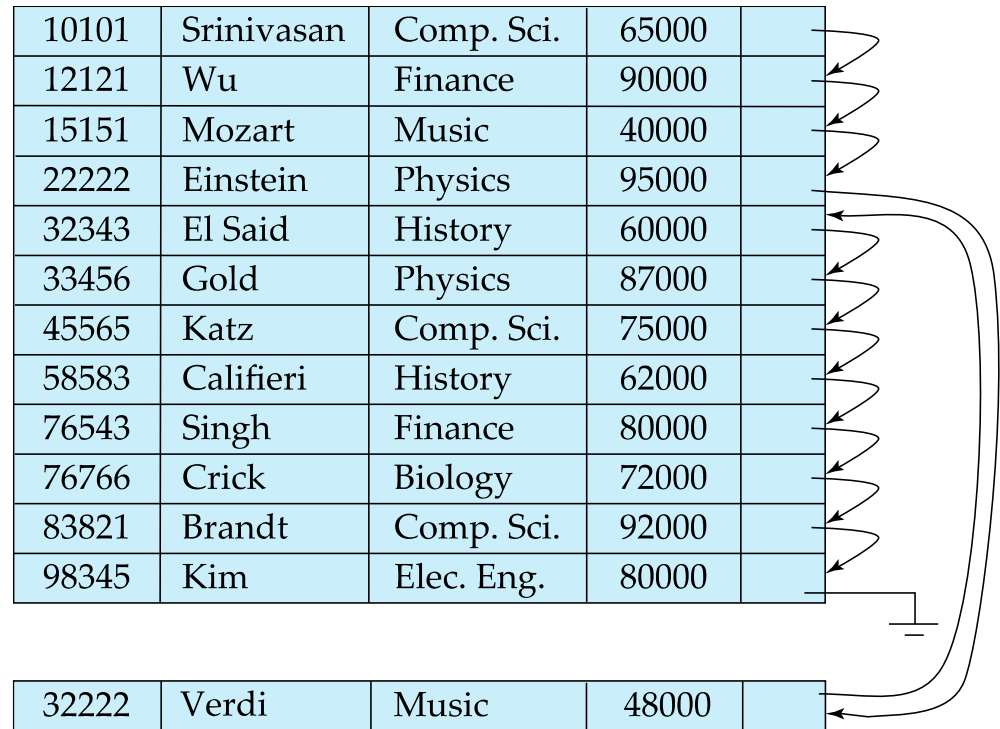
- Suitable for applications that require sequential processing of the entire file
- The records in the file are ordered by a search-key

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	
10101	Srinivasan	Comp. Sci.	65000
45565	Katz	Comp. Sci.	75000
83821	Brandt	Comp. Sci.	92000
Physics	Watson	70000	
33456	Gold	Physics	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



# Partitioning

- **Table partitioning:** Records in a relation can be partitioned into smaller relations that are stored separately
- E.g., *transaction* relation may be partitioned into *transaction\_2018*, *transaction\_2019*, etc.
- Queries written on *transaction* must access records in all partitions
  - Unless query has a selection such as *year=2019*, in which case only one partition is needed
- Partitioning
  - Reduces costs of some operations such as free space management
  - Allows different partitions to be stored on different storage devices
    - E.g., *transaction* partition for current year on SSD, for older years on magnetic disk





# Outline: Data Storage Structures

- Database Storage Architecture
- **Storage Organization (Records, Pages and Files)**
- File Organization
  - Fixed-Length Records
  - Variable-Length Records
  - Storing Large Objects
- Organization of Records in Files
  - Heap File Organization
  - Sequential File Organization
  - Multitable Clustering File Organization
  - Partitioning
- **Database Metadata (Data-Dictionary Storage)**
- Database Buffers
  - Buffer Manager
  - Buffer-Replacement Strategies
  - Reordering of Writes and Recovery



# Data Dictionary (Metadata) Storage

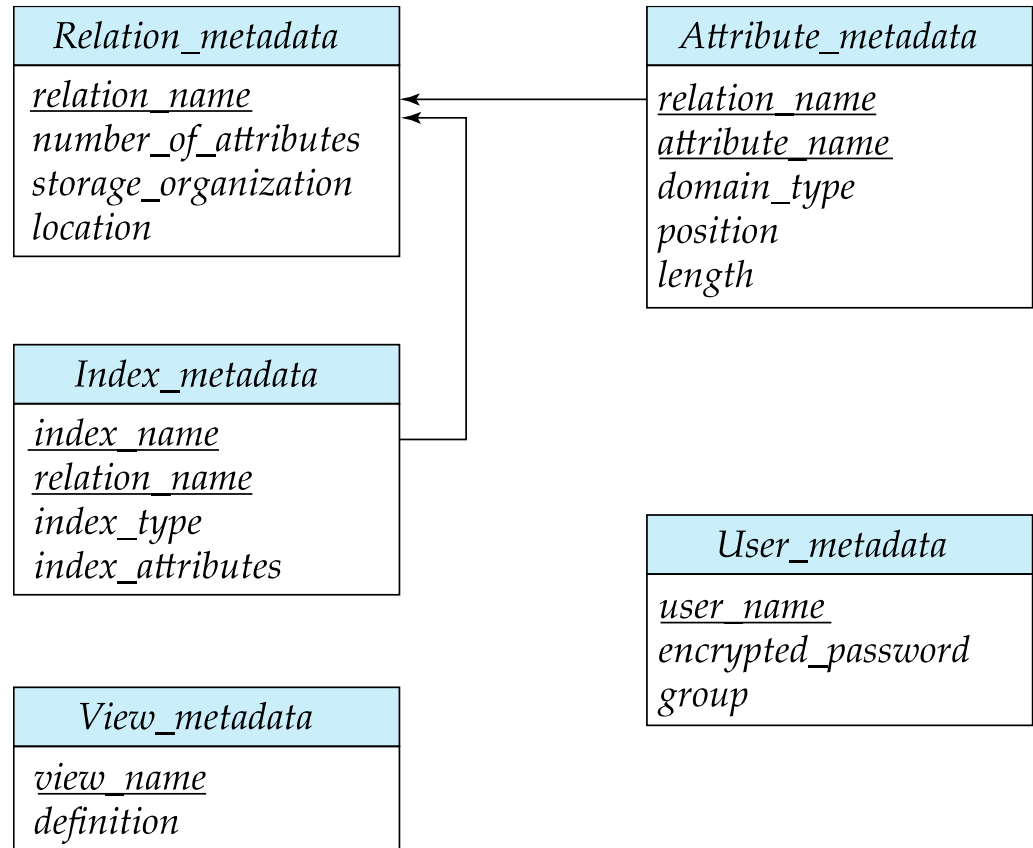
The **Data dictionary** (also called **system catalog**) stores **metadata**; that is, data about data, such as

- 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



# Relational Representation of System Metadata

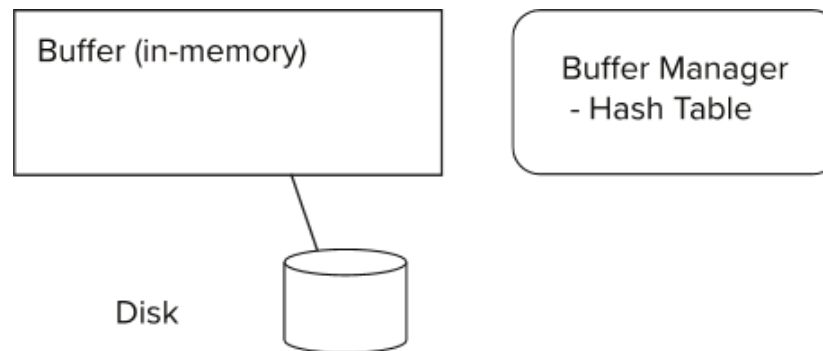
- Relational representation on disk
- Specialized data structures designed for efficient access, in memory





# Storage Access

- 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.





# Outline: Data Storage Structures

- Database Storage Architecture
- Storage Organization (Records, Pages and Files)
- File Organization
  - Fixed-Length Records
  - Variable-Length Records
  - Storing Large Objects
- Organization of Records in Files
  - Heap File Organization
  - Sequential File Organization
  - Multitable Clustering File Organization
  - Partitioning
- Database Metadata (Data-Dictionary Storage)
- Database Buffers
  - Buffer Manager
  - Buffer-Replacement Strategies
  - Reordering of Writes and Recovery



# Database Buffer

- The size of **main memory** on **servers** has increased greatly over the years, and many **medium-sized databases** can fit in **memory**.
- However, a **server** has many demands on its **memory**, and the amount of memory it can give to a database may be much smaller than the database size even for **medium-sized databases**.
- And many **large databases** are much larger than the available memory on servers.
- Thus, even today, **database data** reside primarily on **disk** in most databases, and they must be brought into memory to be **read** or **updated**; **updated data blocks** must be written back to disk subsequently.



# Database Buffer

- Since **data access** from **disk** is much **slower** than **in-memory data access**, a major goal of the database system is to **minimize the number of block transfers** between the **disk and memory**.
- One way to **reduce the number of disk accesses** is to keep **as many blocks as possible** in main memory.
- The goal is to **maximize the chance** that, when a block is accessed, it is already in main memory, and, thus, no disk access is required.
- Since it is not possible to keep all blocks in **main memory**, we need to **manage the allocation of the space available** in main memory for the storage of blocks.



# Database Buffer

- The **buffer** is that part of main memory available for storage of copies of disk blocks.
- There is always a copy kept on disk of every block, but the copy on disk may be a version of the block older than the version in the **buffer**.
- The subsystem responsible for the allocation of buffer space is called the **buffer manager**.





# Buffer Manager

- Programs call on the buffer manager when they need a block from disk.
  - If the block is already in the buffer, buffer manager returns the address of the block in main memory
  - If the block is not in the buffer, the buffer manager
    - Allocates space in the buffer for the block
      - Replacing (throwing out) some other block, if required, to make space for the new block.
      - 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.
    - Reads the block from the disk to the buffer, and returns the address of the block in main memory to requester.



# Buffer Manager

- **Buffer replacement strategy** (details coming up!)
- **Pinned block:** memory block that is not allowed to be written back to disk
  - **Pin** done before reading/writing data from a block
  - **Unpin** done when read /write is complete
  - Multiple concurrent pin/unpin operations possible
    - Keep a pin count, buffer block can be evicted only if pin count = 0
- **Shared and exclusive locks on buffer**
  - Needed to prevent concurrent operations from reading page contents as they are moved/reorganized, and to ensure only one move/reorganize at a time
  - Readers get shared lock, updates to a block require exclusive lock
  - **Locking rules:**
    - Only one process can get exclusive lock at a time
    - Shared lock cannot be concurrently with exclusive lock
    - Multiple processes may be given shared lock concurrently



# 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
  - LRU can be bad for some queries
- 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
- Mixed strategy with hints on replacement strategy provided by the query optimizer is preferable
- Example of bad access pattern for LRU: 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 ...
```



# Buffer-Replacement Policies (Cont.)

- **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
- Operating system or buffer manager may reorder writes
  - Can lead to corruption of data structures on disk
    - E.g., linked list of blocks with missing block on disk
    - File systems perform consistency check to detect such situations
  - Careful ordering of writes can avoid many such problems



# Optimization of Disk Block Access (Cont.)

- Buffer managers support **forced output** of blocks for the purpose of recovery
- **Nonvolatile write buffers** speed up disk writes by writing blocks to a non-volatile RAM or flash buffer immediately
  - *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
- **Journaling file systems** write data in-order to NV-RAM or log disk
  - Reordering without journaling: risk of corruption of file system data