# Intro to Database Systems (15-445/645) Of Database Storage Part 2



### **ADMINISTRIVIA**

Project 0 due yesterday.

Homework 1 due this Friday, Feb 3<sup>rd</sup>.

Project 1 will be released later today.

### **DISK-ORIENTED ARCHITECTURE**

The DBMS assumes that the primary storage location of the database is on non-volatile disk.

The DBMS's components manage the movement of data between non-volatile and volatile storage.

## **PAGE-ORIENTED ARCHITECTURE**

#### Insert a new tuple:

- $\rightarrow$  Check page directory to find a page with a free slot.
- $\rightarrow$  Retrieve the page from disk (if not in memory).
- $\rightarrow$  Check slot array to find empty space in page that will fit.

#### Update an existing tuple using its record id:

- $\rightarrow$  Check page directory to find location of page.
- $\rightarrow$  Retrieve the page from disk (if not in memory).
- $\rightarrow$  Find offset in page using slot array.
- $\rightarrow$  Overwrite existing data (if new data fits).

SCMU·DB 15-445/645 (Spring 2023

## DISCUSSION

### Problems with the slotted page design

- $\rightarrow$  Fragmentation
- $\rightarrow$  Useless Disk I/O
- $\rightarrow$  Random Disk I/O (e.g., update 20 tuples on 20 pages)

What if the DBMS <u>cannot</u> overwrite data in pages and could only create new pages? → Examples: Some cloud storage, HDFS

SCMU·DB 15-445/645 (Spring 2023

### **TODAY'S AGENDA**

Log-Structured Storage Data Representation System Catalogs

## LOG-STRUCTURED STORAGE

DBMS stores log records that contain changes to tuples (PUT, DELETE).

- → Each log record must contain the tuple's unique identifier.
- $\rightarrow$  Put records contain the tuple contents.
- $\rightarrow$  Deletes marks the tuple as deleted.

As the application makes changes to the database, the DBMS appends log records to the end of the file without checking previous log records.



### LOG-STRUCTURED STORAGE

When the page gets full, the DBMS writes it out disk and starts filling up the next page with records.  $\rightarrow$  All disk writes are sequential.  $\rightarrow$  On-disk pages are immutable.

PUT #105 $\{y_2\}=c$	
PUT #102 {val=d <sub>1</sub> }	
DEL #101	
DEL #102	
PUT #105 {val=c <sub>3</sub> }	

## LOG-STRUCTURED STORAGE

To read a tuple with a given id, the DBMS finds the newest log record corresponding to that id.  $\rightarrow$  Scan log from newest to oldest.

Maintain an index that maps a tuple id to the newest log record.

- $\rightarrow$  If log record is in-memory, just read it.
- $\rightarrow$  If log record is on a disk page, retrieve it.
- $\rightarrow$  We will discuss indexes in two weeks.



**CMU·DB** 15-445/645 (Spring 2023



### LOG-STRUCTURED COMPACTION

After a page is compacted, the DBMS does <u>not</u> need to maintain temporal ordering of records within the page.
→ Each tuple id is guaranteed to appear at most once in the page.

The DBMS can instead sort the page based on id order to improve efficiency of future look-ups. → Called <u>Sorted String Tables</u> (SSTables)



11



## DISCUSSION

Log-structured storage managers are more common today. This is partly due to the proliferation of RocksDB.

#### What are some downsides of this approach?

- $\rightarrow$  Write-Amplification
- $\rightarrow$  Compaction is Expensive

RocksDB levelDB LevelDB LevelDB VugabyteDB Vggraph Dgraph CockroachDB

📿 cassandra

13

### **TUPLE STORAGE**

A tuple is essentially a sequence of bytes. It's the job of the DBMS to interpret those bytes into attribute types and values.

The DBMS's catalogs contain the schema information about tables that the system uses to figure out the tuple's layout.

ECMU-DB 15-445/645 (Spring 2023

### DATA REPRESENTATION

15

### INTEGER/BIGINT/SMALLINT/TINYINT

 $\rightarrow$  Same as in C/C++

FLOAT/REAL vs. NUMERIC/DECIMAL

→ IEEE-754 Standard / Fixed-point Decimals

### VARCHAR/VARBINARY/TEXT/BLOB

 $\rightarrow$  Header with length, followed by data bytes.

 $\rightarrow$  Need to worry about collations / sorting.

#### TIME/DATE/TIMESTAMP

 $\rightarrow$  32/64-bit integer of (micro)seconds since Unix epoch

### VARIABLE PRECISION NUMBERS

Inexact, variable-precision numeric type that uses the "native" C/C++ types.  $\rightarrow$  Examples: FLOAT, REAL/DOUBLE

Store directly as specified by **IEEE-754**.

Typically faster than arbitrary precision numbers but can have rounding errors...

ECMU·DB 15-445/645 (Spring 2023

### VARIABLE PRECISION NUMBERS



17

### **FIXED PRECISION NUMBERS**

Numeric data types with (potentially) arbitrary precision and scale. Used when rounding errors are unacceptable.

→ Example: NUMERIC, DECIMAL

Many different implementations.

- → Example: Store in an exact, variable-length binary representation with additional meta-data.
- $\rightarrow$  Can be less expensive if you give up arbitrary precision.

ECMU-DB 15-445/645 (Spring 2023





### LARGE VALUES

Most DBMSs don't allow a tuple to exceed the size of a single page.

To store values that are larger than a page, the DBMS uses separate <u>overflow</u> storage pages. → Postgres: TOAST (>2KB) → MySQL: Overflow (>1/2 size of page) → SQL Server: Overflow (>size of page)



21

### **EXTERNAL VALUE**

Some systems allow you to store a really large value in an external file. Treated as a **BLOB** type. → Oracle: **BFILE** data type → Microsoft: **FILESTREAM** data type

The DBMS <u>cannot</u> manipulate the contents of an external file.
→ No durability protections.
→ No transaction protections.

#### To BLOB or Not To BLOB: Large Object Storage in a Database or a Filesystem?

Russell Sears<sup>2</sup>, Catharine van Ingen<sup>1</sup>, Jim Gray<sup>1</sup> 1: Microsoft Research, 2: University of California at Berkeley sears@cs.berkeley.edu, vanIngen@microsoft.com, gray@microsoft.com MSR-TR-2006-45 April 2006 Revised June 2006

#### Abstract

Application designers must decide whether to store large objects (BLOBs) in a filesystem or in a database. Generally, this decision is based on factors such as application simplicity or manageability. Often, system performance affects these factors.

Following an encounter network, and a status of the problem of the

Of course, this depends on the particular filesystem, database system, and workload in question. This study shows that when comparing the NTFS file system and SQL Server 2005 database system on a create, (read, replace)\* delete workload, BLOBs smaller than 256KB are more efficiently handled by SQL Server, while NTFS is more efficient BLOBS larger than IMB. Of course, this break-even point will vary among different database systems, filesystems, and workloads. By measuring the performance of a storage server

b) incasting the performance of a storage server workload typical of web applications which use gev/put photocols such as WebDAV [WebDAV], we found that be the ak-even point depends on many factors. However, our experiments suggest that storage age, the ains of bytes in deleted or replaced objects to bytes in like objects, is dominant. As storage age increases, fragmentation tends to increases. The filesystem we established benefit from incorporating ideas from filesystem world benefit from incorporating ideas from filesystem architecture oversely, filesystem performance may be impresed by using database techniques to handle small index for the file of the storage storage storage.

Surprisingly, for these studies, when average object size is held constant, the distribution of object sizes did not significantly affect performance. We also found that, in addition to low percentage free space, a low ratio of free space to average object size leads to fragmentation and performance degradation.

#### 1. Introduction

Application data objects are getting larger as digital media becomes ubiquitous. Furthermore, the increasing popularity of web services and other network applications means that systems that once managed static archives or fushed" objects now manage frequently modified versions of application data as it is being created and updated. Rather than updating these objects, the archive either stores multiple versions of the objects (the V of WebDAV stands for "versioning"), or simply does wholesale replacement (as in SharePoint Team Services (SharePoint).

Application designers have the choice of storing large objects as files in the filesystem, as BLOBs (binary large objects) in a database, or as a combination of both. Only folklore is available regarding the tradeoffs - often the design decision is based on which technology the designer knows best. Most designers will tell you that a database is probably best for small binary objects and that that files are best for large objects. But, what is the break-even point? What are the tradeoffs?

What are the traacotist' This article characterizes the performance of an abstracted write-intensive web application that deals with relatively large objects. Two versions of the system are compared; objects as relational database to store large objects, while the other version stores the objects as files in the filesystem. We measure how performance changes over time as the storage becomes fragmented. The article concludes by describing and quantifying the factors that a designer should consider then picking a storage system. It also suggests filesystem and database improvements for large object support.

One surprising (to us at least) conclusion of our work is that storage fragmentation is the main determinant of the break-even point in the tradeoff. Therefore, much of our work and much of this article focuses on storage fragmentation issues. In essence, filesystems seem to have better fragmentation handling than databases and this drives the break-even point down from about IMB to about 266KB.

2

### SYSTEM CATALOGS

23

# A DBMS stores meta-data about databases in its internal catalogs.

- $\rightarrow$  Tables, columns, indexes, views
- $\rightarrow$  Users, permissions
- $\rightarrow$  Internal statistics

Almost every DBMS stores the database's catalog inside itself (i.e., as tables).

- $\rightarrow$  Wrap object abstraction around tuples.
- $\rightarrow$  Specialized code for "bootstrapping" catalog tables.

### SYSTEM CATALOGS

24

You can query the DBMS's internal INFORMATION\_SCHEMA catalog to get info about the database.

→ ANSI standard set of read-only views that provide info about all the tables, views, columns, and procedures in a database

DBMSs also have non-standard shortcuts to retrieve this information.

ECMU-DB 15-445/645 (Spring 2023

### **ACCESSING TABLE SCHEMA**

25

#### List all the tables in the current database:



### **ACCESSING TABLE SCHEMA**

#### List all the tables in the student table:



### CONCLUSION

Log-structured storage is an alternative approach to the page-oriented architecture we discussed last class.

The storage manager is not entirely independent from the rest of the DBMS.