# **15-721** DATABASE SYSTEMS

#### Lecture #09 – Storage Models & Data Layout

Andy Pavlo // Carnegie Mellon University // Spring 2016

#### TODAY'S AGENDA

In-Memory Data Layout Storage Models Project #2: Performance Profiling

#### DATA ORGANIZATION



# DATA ORGANIZATION

One can think of an in-memory database as just a large array of bytes.

 $\rightarrow$  The schema tells the DBMS how to convert the bytes into the appropriate type.

Each tuple is prefixed with a header that contains its meta-data.

Storing tuples with just their fixed-length data makes it easy to compute the starting point of any tuple.



#### **INTEGER/BIGINT/SMALLINT/TINYINT**

 $\rightarrow$  C/C++ Representation

#### NUMERIC

 $\rightarrow$  <u>IEEE-754</u> Standard

#### VARCHAR/VARBINARY/TEXT/BLOB

- $\rightarrow$  Pointer to other location if type is  $\geq$ 64-bits
- → Header with length and address to next location (if segmented), followed by data bytes.

#### TIME/DATE/TIMESTAMP

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



char[]



```
CREATE TABLE JoySux (
   id INT PRIMARY KEY,
   value BIGINT
);
```

char[]				
header	id	value		



char[]				
header	id	value		



<pre>char[]</pre>				
header	id	value		







#### Choice #1: Special Values

→ Designate a value to represent NULL for a particular data type (e.g., INT32\_MIN).

#### Choice #2: Null Column Bitmap Header

 $\rightarrow$  Store a bitmap in the tuple header that specifies what attributes are null.

#### **Choice #3: Per Attribute Null Flag**

- $\rightarrow$  Store a flag that marks that a value is null.
- $\rightarrow$  Have to use more space than just a single bit because this messes up with word alignment.

#### **Integer Numbers**

Data Type	Size	Size (Not Null)	Synonyms	Min Value	Max Value
BOOL	2 bytes	1 byte	BOOLEAN	0	1
BIT	9 bytes	8 bytes			
TINYINT	2 bytes	1 byte		-128	127
SMALLINT	4 bytes	2 bytes		-32768	32767
MEDIUMINT	4 bytes	3 bytes		-8388608	8388607
INT	8 bytes	4 bytes	INTEGER	-2147483648	2147483647
BIGINT	12 bytes	8 bytes		-2 ** 63	(2 ** 63) - 1

this messes up with word alignment.

#### **Integer Numbers**

0120	Size (Not Null)	Synonyms	Min Value	Max Value
2 bytes	1 byte	BOOLEAN	0	1
9 bytes	8 bytes			
2 bytes	1 byte		-128	127
4 bytes	2 bytes		-32768	32767
4 bytes	3 bytes		-8388608	8388607
8 bytes	4 bytes	INTEGER	-2147483648	2147483647
12 bytes	8 bytes		-2 ** 63	(2 ** 63) - 1
	2 bytes 9 bytes 2 bytes 4 bytes 4 bytes 8 bytes 12 bytes	2 bytes1 byte9 bytes8 bytes2 bytes1 byte4 bytes2 bytes4 bytes3 bytes8 bytes4 bytes12 bytes8 bytes	2 bytes1 byteBOOLEAN9 bytes8 bytes12 bytes1 byte14 bytes2 bytes14 bytes3 bytes18 bytes4 bytes112 bytes8 bytes1	2 bytes1 byteBOOLEAN09 bytes8 bytesII2 bytes1 byteIII4 bytes2 bytesIII4 bytes3 bytesIII8 bytes4 bytesIINTEGERI12 bytes8 bytesIII

this messes up with word alignment.

#### Choice #1: Special Values

→ Designate a value to represent NULL for a particular data type (e.g., INT32\_MIN).

#### Choice #2: Null Column Bitmap Header

 $\rightarrow$  Store a bitmap in the tuple header that specifies what attributes are null.

#### **Choice #3: Per Attribute Null Flag**

- $\rightarrow$  Store a flag that marks that a value is null.
- $\rightarrow$  Have to use more space than just a single bit because this messes up with word alignment.

# NOTICE

The truth is that you only need to worry about word-alignment for cache lines (e.g., 64 bytes).

I'm going to show you the basic idea using 64bit words since it's easier to see...

```
CREATE TABLE JoySux (
   id INT PRIMARY KEY,
   cdate TIMESTAMP,
   color CHAR(2),
   zipcode INT
);
```



```
CREATE TABLE JoySux (
    id INT PRIMARY KEY,
    cdate TIMESTAMP,
    color CHAR(2),
    zipcode INT
);
```



```
CREATE TABLE JoySux (
32-bits id INT PRIMARY KEY,
cdate TIMESTAMP,
color CHAR(2),
zipcode INT
);
```

All attributes in a tuple must be word aligned to enable the CPU to access it without any unexpected behavior or additional work.

char[]

```
CREATE TABLE JoySux (
32-bits id INT PRIMARY KEY,
64-bits cdate TIMESTAMP,
                                 id
                                        cdate
      color CHAR(2),
      zipcode INT
                                64-bit Word 64-bit Word 64-bit Word 64-bit Word
    );
```

```
CREATE TABLE JoySux (
32-bits id INT PRIMARY KEY,
64-bits cdate TIMESTAMP,
16-bits color CHAR(2),
zipcode INT
);
```



```
CREATE TABLE JoySux (
32-bits id INT PRIMARY KEY,
64-bits cdate TIMESTAMP,
16-bits color CHAR(2),
32-bits zipcode INT
);
```



```
CREATE TABLE JoySux (
32-bits id INT PRIMARY KEY,
64-bits cdate TIMESTAMP,
16-bits color CHAR(2),
32-bits zipcode INT
);
```



If the CPU fetches a 64-bit value that is not word-aligned, it has four choices:

- → Execute two reads to load the appropriate parts of the data word and reassemble them.
- → Read some unexpected combination of bytes assembled into a 64-bit word.
- $\rightarrow$  Throw an exception

```
CREATE TABLE JoySux (
32-bits id INT PRIMARY KEY,
64-bits cdate TIMESTAMP,
16-bits color CHAR(2),
32-bits zipcode INT
);
```



```
CREATE TABLE JoySux (
32-bits id INT PRIMARY KEY,
64-bits cdate TIMESTAMP,
16-bits color CHAR(2),
32-bits zipcode INT
);
```



#### STORAGE MODELS

N-ary Storage Model (NSM) Decomposition Storage Model (DSM) Hybrid Storage Model

# N-ARY STORAGE MODEL (NSM)

The DBMS stores all of the attributes for a single tuple contiguously.

Ideal for OLTP workloads where txns tend to operate only on an individual entity and insert-heavy workloads.

Use the tuple-at-a-time iterator model.



#### NSM PHYSICAL STORAGE

#### **Choice #1: Heap-Organized Tables**

- $\rightarrow$  Tuples are stored in blocks called a heap.
- $\rightarrow$  The heap does not necessarily define an order.

#### **Choice #2: Index-Organized Tables**

- $\rightarrow$  Tuples are stored in the index itself.
- $\rightarrow$  Not quite the same as a clustered index.

# CLUSTERED INDEXES

The table is stored in the sort order specified by the primary key.

 $\rightarrow$  Can be either heap- or index-organized storage.

Some DBMSs always use a clustered index.  $\rightarrow$  If a table doesn't include a pkey, the DBMS will automatically make a hidden row id pkey.

Other DBMSs cannot use them at all.  $\rightarrow$  A clustered index is non-practical in a MVCC DBMS using the **Insert Method**.

# N-ARY STORAGE MODEL (NSM)

#### Advantages

- $\rightarrow$  Fast inserts, updates, and deletes.
- $\rightarrow$  Good for queries that need the entire tuple.
- $\rightarrow$  Can use index-oriented physical storage.

#### Disadvantages

 $\rightarrow$  Not good for scanning large portions of the table and/or a subset of the attributes.

# DECOMPOSITION STORAGE MODEL (DSM)

The DBMS stores a single attribute for all tuples contiguously in a block of data.  $\rightarrow$  Sometimes also called <u>vertical partitioning</u>.

Ideal for OLAP workloads where read-only queries perform large scans over a subset of the table's attributes.

Use the vector-at-a-time iterator model.

# DECOMPOSITION STORAGE MODEL (DSM)

- 1970s: Cantor DBMS
- 1980s: DSM Proposal
- **1990s:** SybaseIQ (in-memory only)
- 2000s: Vertica, Vectorwise, MonetDB
- 2010s: "The Big Three" Cloudera Impala, Amazon Redshift, SAP HANA, MemSQL



# CLUSTERED INDEXES

Some columnar DBMSs store data in sorted order to maximize compression.  $\rightarrow$  Bitmap indexes with RLE from last class

Vertica does not even use indexes because all columns are sorted.

# TUPLE IDENTIFICATION

# Choice #1: Fixed-length Offsets → Each value is the same length for an attribute. Choice #2: Embedded Tuple Ids → Each value is stored with its tuple id in a column.

Offsets



Embedded Ids



# DECOMPOSITION STORAGE MODEL (DSM)

#### Advantages

- $\rightarrow$  Reduces the amount wasted work because the DBMS only reads the data that it needs.
- $\rightarrow$  Better compression (last lecture).

#### Disadvantages

 $\rightarrow$  Slow for point queries, inserts, updates, and deletes because of tuple splitting/stitching.
### OBSERVATION

Data is "hot" when first entered into database  $\rightarrow$  A newly inserted tuple is more likely to be updated again the near future.

As a tuple ages, it is updated less frequently.  $\rightarrow$  At some point, a tuple is only accessed in read-only queries along with other tuples.

What if we want to use this data to make decisions that affect new txns?









CARNEGIE MELLON DATABASE GROUP

CMU 15-721 (Spring 2016)



**OLTP Data Silos** 

CARNEGIE MELLON

CMU 15-721 (Spring 2016)



CMU 15-721 (Spring 2016)





# HYBRID STORAGE MODEL

Single logical database instance that uses different storage models for hot and cold data.

Store new data in NSM for fast OLTP Migrate data to DSM for more efficient OLAP

# HYBRID STORAGE MODEL

#### **Choice #1: Separate Execution Engines**

 $\rightarrow$  Use separate execution engines that are optimized for either NSM or DSM databases.

### Choice #2: Single, Flexible Architecture

 $\rightarrow$  Use single execution engine that is able to efficiently operate on both NSM and DSM databases.

# SEPARATE EXECUTION ENGINES

Run separate "internal" DBMSs that each only operate on DSM or NSM data.

- $\rightarrow$  Need to combine query results from both engines to appear as a single logical database to the application.
- $\rightarrow$  Have to use a synchronization method (e.g., 2PC) if a txn spans execution engines.

Two approaches to do this:

→ **Fractured Mirrors** (Oracle, IBM)

 $\rightarrow$  **<u>Delta Store</u>** (SAP HANA)















### DELTA STORE

Stage updates to the database in an NSM table. A background thread migrates updates from delta store and applies them to DSM data.



### DELTA STORE

Stage updates to the database in an NSM table. A background thread migrates updates from delta store and applies them to DSM data.



### DELTA STORE

Stage updates to the database in an NSM table. A background thread migrates updates from delta store and applies them to DSM data.



# SINGLE, FLEXIBLE ARCHITECTURE

Use a single execution engine architecture that is able to operate on both NSM and DSM data.  $\rightarrow$  Don't need to store two copies of the database.  $\rightarrow$  Don't need to sync multiple database segments.

Note that a DBMS can use the delta-store for NSM data with a single architecture.

Examine the access patterns of queries and then dynamically reconfigure the database to optimize decomposition and layout.

Copies columns into a new layout that is optimized for each query.

- $\rightarrow$  Think of it like a mini fractured mirror.
- $\rightarrow$  Use query compilation to speed up operations.



UPDATE JoyStillSux
 SET B = 1234
 WHERE C = "xxx"

SELECT AVG(B)
FROM JoyStillSux
WHERE C = "yyy"

SELECT SUM(A)
FROM JoyStillSux

Α	В	С	D	



А	В	С	D





А	В	С	D





# $H_2O$ ADAPTIVE STORAGE

This approach is unable to handle updates to the database. It also unable to store tuples in the same table in a different layout.

This is because they are missing the ability to categorize whether data is hot or cold...

UPDATE JoyStillSux
SET B = 1234
WHERE C = "xxx"

SELECT AVG(B)
FROM JoyStillSux
WHERE C = "yyy"

SELECT SUM(A)
FROM JoyStillSux

Α	В	С	D	





SELECT SUM(A)
FROM JoyStillSux

А	В	С	D	









# CATEGORIZING DATA

### Choice #1: Manual Approach

 $\rightarrow$  DBA specifies what tables should be stored as DSM.

### Choice #2: Off-line Approach

 $\rightarrow$  DBMS monitors access logs offline and then makes decision about what data to move to DSM.

### **Choice #3: On-line Approach**

 $\rightarrow$  DBMS tracks access patterns at runtime and then makes decision about what data to move to DSM.



### PARTING THOUGHTS

A flexible architecture that supports a hybrid storage model is the next major trend in DBMSs

This will enable relational DBMSs to support all known database workloads except for matrices in machine learning.

35





# MOTIVATION

Consider a hot program **Z** with two functions **foo** and **bar**.

How can we speed up Z with only a debugger ?  $\rightarrow$  Randomly pause it during execution  $\rightarrow$  Collect the function call stack

### RANDOM PAUSE METHOD

#### Consider this scenario

- $\rightarrow$  Collected 10 call stack samples
- $\rightarrow$  Say 6 out of the 10 samples were in **foo**

What percentage of time was spent in **foo**?  $\rightarrow$  Roughly 60% of the time was spent in **foo** 

 $\rightarrow$  Accuracy increases with # of samples

### AMDAHL'S LAW

Say we optimized **foo** to run 2 times faster What's the expected overall speedup ?

→ p = percentage of time spent in optimized task → s = speed up for the optimized task → Overall speedup = \_\_\_\_ = 1.4 times faster
# AMDAHL'S LAW

Say we optimized **foo** to run 2 times faster What's the expected overall speedup ?  $\rightarrow$  60% of time spent in **foo** drops in half  $\rightarrow$  40% of time spent in **bar** unaffected

→ p = percentage of time spent in optimized task
 → s = speed up for the optimized task
 → Overall speedup = \_\_\_\_ = 1.4 times faster

## AMDAHL'S LAW

1 0.6 2 +0.4 1 1 0.6 2 +0.4 0.6 2 0.6 0.6 2 2 0.6 2 +0.4 1 0.6 2 +0.4 = 1.4 times faster

1 0.6 2 +0.4 1 1 0.6 2 +0.4 0.6 2 0.6 0.6 2 2 0.6 2 +0.4 1 0.6 2 +0.4 = 1.4 times faster

Say we optimized **foo** to run 2 times faster What's the expected overall speedup ?  $\rightarrow$  60% of time spent in **foo** drops in half  $\rightarrow$  40% of time spent in **bar** unaffected

# PROFILING TOOLS FOR REAL

### Choice #1: Valgrind

- $\rightarrow$  Heavyweight instrumentation framework with a lot of tools
- $\rightarrow$  Sophisticated visualization tools

### Choice #2: Perf

- $\rightarrow$  Lightweight tool that can record different kinds of events
- $\rightarrow$  Console-oriented visualization tools

# CHOICE #1: VALGRIND

Instrumentation framework for building dynamic analysis tools → memcheck: a memory error detector

→ **callgrind**: a call-graph generating profiler

# CHOICE #1: VALGRIND

Instrumentation framework for building dynamic analysis tools → memcheck: a memory error detector

 $\rightarrow$  callgrind: a call-graph generating profiler

Using **callgrind** to profile the index test and Peloton in general:

\$ valgrind --tool=callgrind --trace-children=yes
./tests/index\_test

\$ valgrind --tool=callgrind --trace-children=yes
./build/src/peloton -D data &> /dev/null&

# KCACHEGRIND

### Profile data visualization tool

\$ kcachegrind callgrind.out.12345







CARP DAT/ callgrind.out [1] - Total Instruction Fetch Cost: 9 318 940 42



CARP DAT/ callgrind.out [1] - Total Instruction Fetch Cost: 9 318 940

42

# CHOICE #2: PERF

Tool for using the performance counters subsystem in Linux.

→ -e = sample the event cycles at the user level only
 → -c = collect a sample every 2000 occurrences of event

\$ perf record -e cycles:u -c 2000
./tests/index\_test

Uses counters for tracking events  $\rightarrow$  On counter overflow, the kernel records a sample  $\rightarrow$  Sample contains info about program execution



## PERF VISUALIZATION

We can also use **perf** to visualize the generated profile for our application.

\$ perf report



### DEDE VICLALIZATION

#### × -Edit View Search Terminal Help File Samples: 56 of event 'cpu-clock:u', Event count (approx.): 56 index test ld-2.19.so [.] do lookup x index test ld-2.19.so dl lookup symbol x [.] index test ld-2.19.so dl relocate object check match.9458 index test ld-2.19.so [.] index test libstdc++.so.6.0.21 operator delete(void\*) [.] index test libstdc++.so.6.0.21

dynamic cast [.] index test libstdc++.so.6.0.21 operator new(unsigned long) index test libpelotonpg.so.0.0.0 [.] Json::Value::~Value() index test libpeloton.so.0.0.0 peloton::Value::CompareWithoutNull(peloton::Value) const index test libc-2.19.so int free [.] index test libc-2.19.so memcpy\_sse2\_unaligned [.] dl addr index test libc-2.19.so [.] index test libc-2.19.so libc dl error tsd index test ld-2.19.so strcmp [.] index test index test testing::TestEventListeners::TestEventListeners()

CARNEGIE MELLON DATABASE GROUP perf report

### DEDE MICHALIZATION

### x – 🗆

perf report

File Edit View Search Terminal Help

Samples:	56 of event	t 'cpu-clock:u', Event	count (approx.): 56
	index_test	ld-2.19.so	[.] do_lookup_x
	index_test	ld-2.19.so	[.] _dl_lookup_symbol_x
	index_test	ld-2.19.so	[.] _dl_relocate_object
	index_test	ld-2.19.so	[.] check_match.9458
	index_test	libstdc++.so.6.0.21	<pre>[.] operator delete(void*)</pre>
1.79%	index_test	libstdc++.so.6.0.21	[.]dynamic_cast
1.79%	index test	libstdc++.so.6.0.21	[.] operator new(unsigned long)
1.79%	index_test	libpelotonpg.so.0.0.0	[.] Json::Value::~Value()
1.79%	index_test	libpeloton.so.0.0.0	[.] peloton::Value::CompareWithoutNull(peloton::Value) const
1.79%	index test	libc-2.19.so	[.] _int_free
1.79%	index_test	libc-2.19.so	[.]memcpy_sse2_unaligned
1.79%	index_test	libc-2.19.so	[.] _dl_addr
1.79%	index test	libc-2.19.so	[.]libc_dl_error_tsd
1.79%	index test	ld-2.19.so	[.] strcmp
1.79%	index_test	index_test	<pre>[.] testing::TestEventListeners::TestEventListeners()</pre>

Cumulative Time Distribution

# PERF EVENTS

### Supports several other events like: $\rightarrow$ L1-dcache-load-misses $\rightarrow$ branch-misses

### To see a list of events:

\$ perf list

### Another usage example:

\$ perf record -e cycles,LLC-load-misses -c 2000
./tests/index\_test



# REFERENCES

### Valgrind

- $\rightarrow$  The Valgrind Quick Start Guide
- $\rightarrow$  <u>Callgrind</u>
- $\rightarrow$  <u>Kcachegrind</u>
- $\rightarrow$  <u>Tips for the Profiling/Optimization process</u>

### Perf

- $\rightarrow$  <u>Perf Tutorial</u>
- $\rightarrow$  <u>Perf Examples</u>
- $\rightarrow$  <u>Perf Analysis Tools</u>