

*W4111 – Introduction to Databases
Section 002, Fall 2023
Lecture 9: Module II, NoSQL, Part 2*



Contents

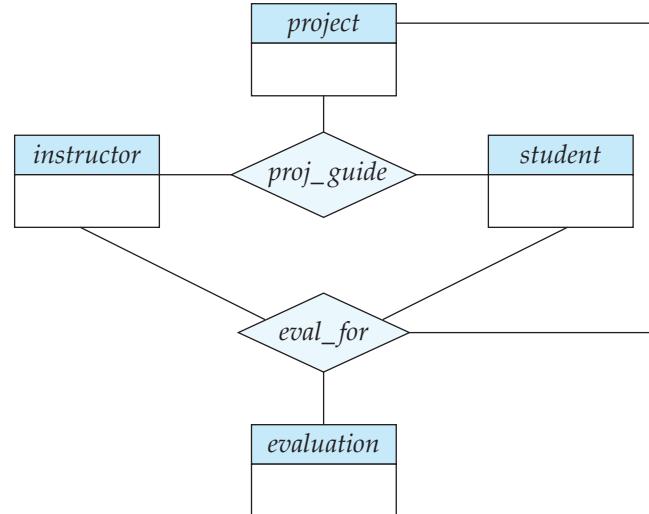
ER Modeling

Aggregation



Aggregation

- Consider the ternary relationship *proj_guide*, which we saw earlier
- Suppose we want to record evaluations of a student by a guide on a project





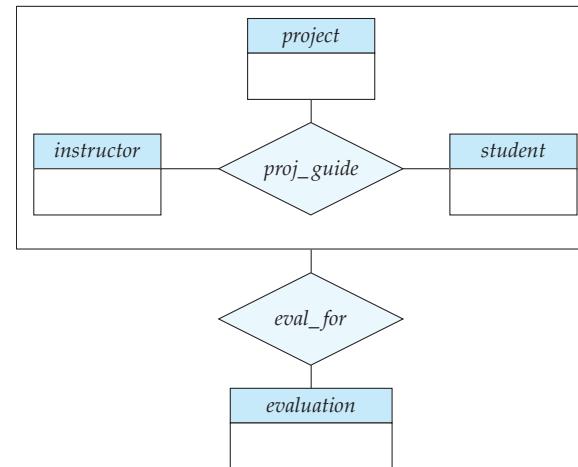
Aggregation (Cont.)

- Relationship sets *eval_for* and *proj_guide* represent overlapping information
 - Every *eval_for* relationship corresponds to a *proj_guide* relationship
 - However, some *proj_guide* relationships may not correspond to any *eval_for* relationships
 - So we can't discard the *proj_guide* relationship
- Eliminate this redundancy via *aggregation*
 - Treat relationship as an abstract entity
 - Allows relationships between relationships
 - Abstraction of relationship into new entity



Aggregation (Cont.)

- Eliminate this redundancy via *aggregation* without introducing redundancy, the following diagram represents:
 - A student is guided by a particular instructor on a particular project
 - A student, instructor, project combination may have an associated evaluation





Reduction to Relational Schemas

- To represent aggregation, create a schema containing
 - Primary key of the aggregated relationship,
 - The primary key of the associated entity set
 - Any descriptive attributes
- In our example:
 - The schema *eval_for* is:
$$\text{eval_for} (s_ID, project_id, i_ID, evaluation_id)$$
 - The schema *proj_guide* is redundant.

Module II and NoSQL

(Continued Part 2)

Reminder

Course Modules – Reminder

Course Overview

Each section of W4111 is slightly different based on student interest and professor's focus. There is a common, core syllabus. Professors cover topics in different orders and grouping based on teaching style.

This section of W4111 has four modules:

- **Foundational concepts (50% of semester):** This module covers concepts like data models, relational model, relational databases and applications, schema, normalization, ... The module focuses on the relational model and relational databases. The concepts are critical and foundational for all types of databases and data centric applications.
- **Database management system architecture and implementation (10%):** This module covers the software architecture, algorithms and implementation techniques that allow [databases management systems](#) to deliver functions. Topics include memory hierarchy, storage systems, caching/buffer pools, indexes, query processing, query optimization, transaction processing, isolation and concurrency control.
- **NoSQL – “Not Only SQL” databases (20%):** This module provides motivation for [“NoSQL”](#) data models and databases, and covers examples and use cases. The module also includes cloud databases and databases-as-a-service.
- **Data Enabled Decision Support (20%):** This module covers data warehouses, data import and cleanse, OLAP, Pivot Tables, Star Schema, reporting and visualization, and provides an overview of analysis techniques, e.g. clustering, classification, analysis, mining.

Module II – DBMS Architecture and

What is a database? In essence a database is nothing more than a collection of information that exists over a long period of time, often many years. In common parlance, the term *database* refers to a collection of data that is managed by a DBMS. The DBMS is expected to:

- 
1. Allow users to create new databases and specify their *schemas* (logical structure of the data), using a specialized *data-definition language*.

Covered for the relational model.

Database Systems: The Complete Book (2nd Edition)

by [Hector Garcia-Molina](#) (Author), [Jeffrey D. Ullman](#) (Author), [Jennifer Widom](#) (Author)

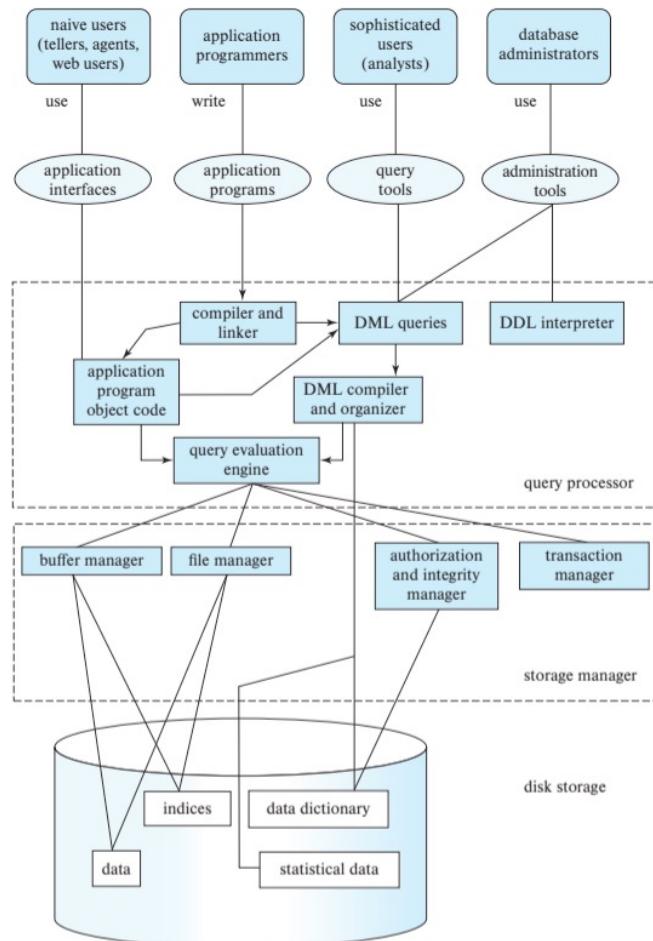
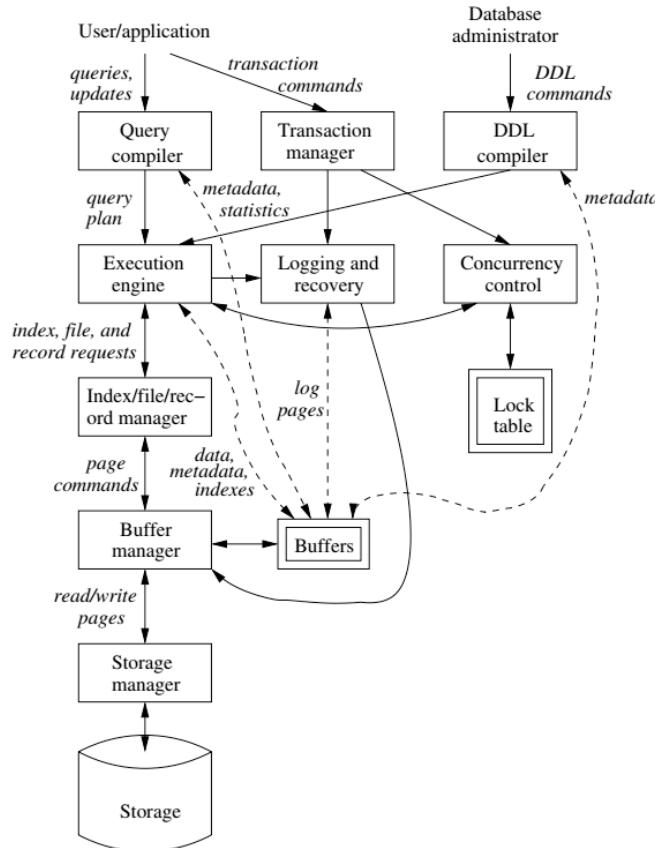
Module II – DBMS Architecture and

- 
- 
2. Give users the ability to *query* the data (a “query” is database lingo for a question about the data) and modify the data, using an appropriate language, often called a *query language* or *data-manipulation language*.
 3. Support the storage of very large amounts of data — many terabytes or more — over a long period of time, allowing efficient access to the data for queries and database modifications.
 4. Enable *durability*, the recovery of the database in the face of failures, errors of many kinds, or intentional misuse.
 5. Control access to data from many users at once, without allowing unexpected interactions among users (called *isolation*) and without actions on the data to be performed partially but not completely (called *atomicity*).

Database Systems: The Complete Book (2nd Edition)

by [Hector Garcia-Molina](#) (Author), [Jeffrey D. Ullman](#) (Author), [Jennifer Widom](#) (Author)

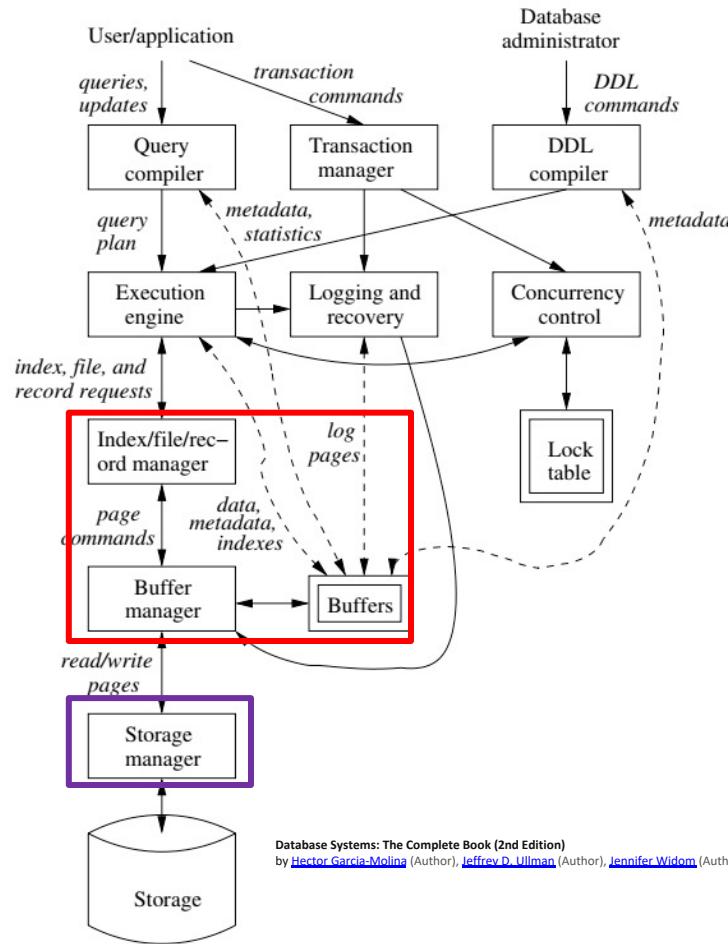
DBMS Arch.



Data Management

Today

- Load/save things quickly.
- Find things quickly.



Data Storage Structures

(Database Systems Concepts, V7, Ch. 13)

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

Terminology

- A tuple in a relation maps to a *record*. Records may be
 - *Fixed length*
 - *Variable length*
 - *Variable format* (which we will see in Neo4J, DynamoDB, etc).
- A *block*
 - Is the unit of transfer between disks and memory (buffer pools).
 - Contains multiple records, usually but not always from the same relation.
- The *database address space* contains
 - All of the blocks and records that the database manages
 - Including blocks/records containing data
 - And blocks/records containing free space.

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

From: Database System Concepts, 7th Ed.

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

From: Database System Concepts, 7th Ed.

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

From: Database System Concepts, 7th Ed.

Fixed-Length Records

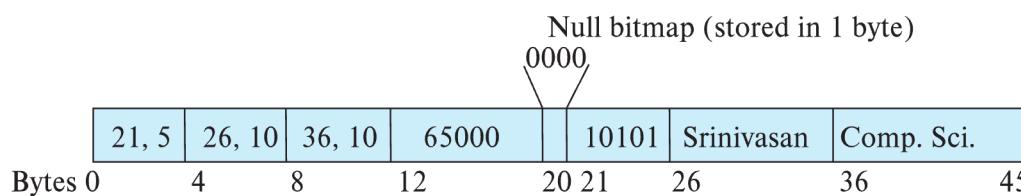
- Deletion of record i : alternatives:
 - move records $i + 1, \dots, n$ to $i, \dots, n - 1$
 - move record n to i
 - **do not move records, but link all free records on a *free list***

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

From: Database System Concepts, 7th Ed.

Variable-Length Records

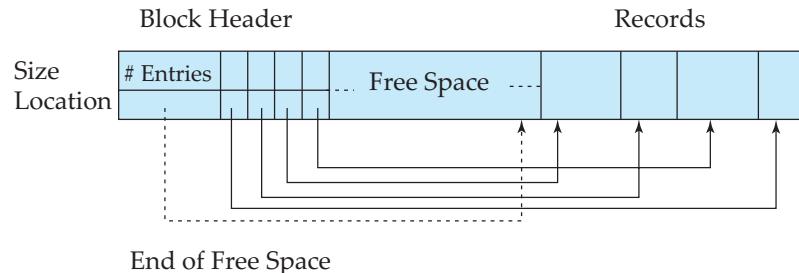
- Variable-length records arise in database systems in several ways:
 - Storage of multiple record types in a file.
 - Record types that allow variable lengths for one or more fields such as strings (`varchar`)
 - Record types that allow repeating fields (used in some older data models).
- Attributes 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



From: Database System Concepts, 7th Ed.

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.



From: Database System Concepts, 7th Ed.

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

From: Database System Concepts, 7th Ed.

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
 - More on this in Chapter 14
- **Hashing** – a hash function computed on search key; the result specifies in which block of the file the record should be placed
 - More on this in Chapter 14

From: Database System Concepts, 7th Ed.

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)

From: Database System Concepts, 7th Ed.

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

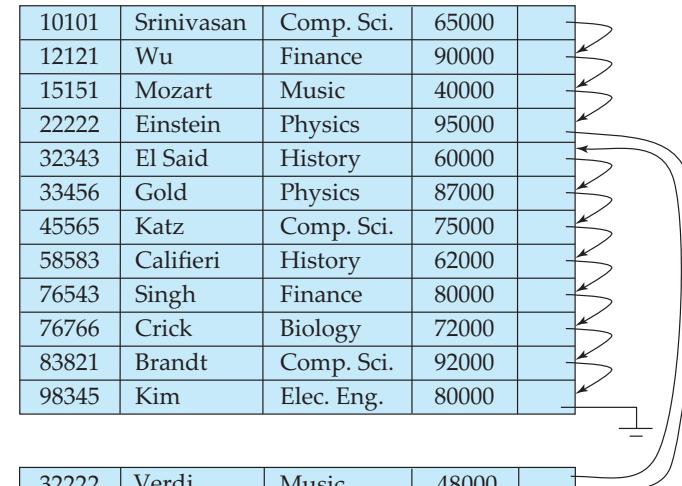
10101	Srinivasan	Comp. Sci.	65000	
12121	Wu	Finance	90000	
15151	Mozart	Music	40000	
22222	Einstein	Physics	95000	
32343	El Said	History	60000	
33456	Gold	Physics	87000	
45565	Katz	Comp. Sci.	75000	
58583	Califieri	History	62000	
76543	Singh	Finance	80000	
76766	Crick	Biology	72000	
83821	Brandt	Comp.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

From: Database System Concepts, 7th Ed.

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

From: Database System Concepts, 7th Ed.



From: Database System Concepts, 7th Ed.

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

From: Database System Concepts, 7th Ed.

Column-Oriented Storage

- Also known as **columnar representation**
- Store each attribute of a relation separately
- Example

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

From: Database System Concepts, 7th Ed.

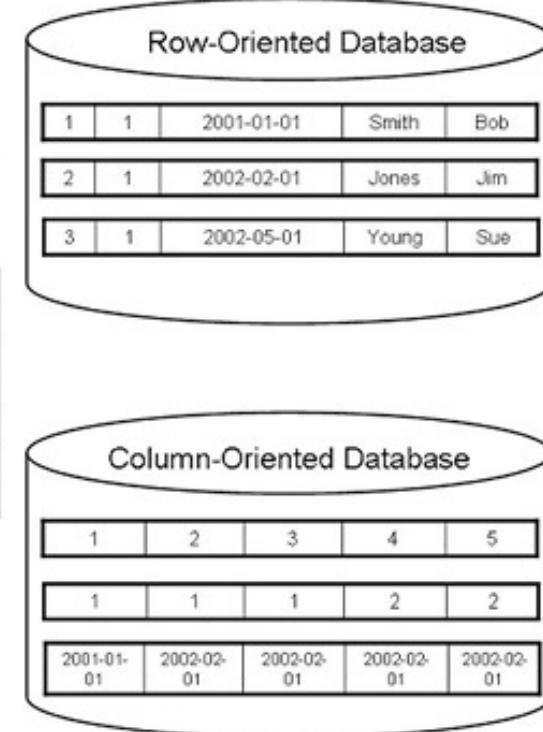
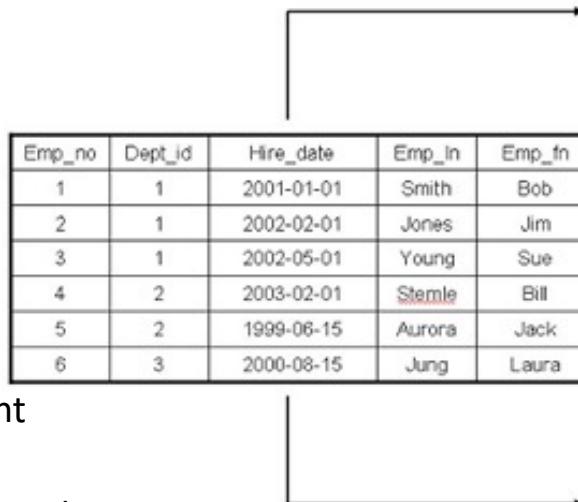
Columnar Representation

- Benefits:
 - Reduced IO if only some attributes are accessed
 - Improved CPU cache performance
 - Improved compression
 - **Vector processing** on modern CPU architectures
- Drawbacks
 - Cost of tuple reconstruction from columnar representation
 - Cost of tuple deletion and update
 - Cost of decompression
- Columnar representation found to be more efficient for decision support than row-oriented representation
- Traditional row-oriented representation preferable for transaction processing
- Some databases support both representations
 - Called **hybrid row/column stores**

From: Database System Concepts, 7th Ed.

Row vs Column

- Columnar and Row are both
 - Relational
 - Support SQL operations
- But differ in data storage
 - Row keeps row data together in blocks.
 - Columnar keeps column data together in blocks.
- This determines performance for different types of query, e.g.
 - Columnar is extremely powerful for BI scenarios
 - Aggregation ops, e.g. SUM, AVG
 - PROJECT (do not load all of the row) to get a few columns
 - Row is powerful for OLTP. Transaction typically create and retrieve
 - One row at a time
 - All the columns of a single row.

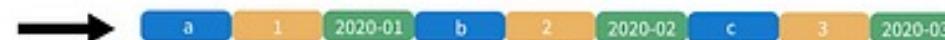


Columnar File Representation

Row-Based Storage Layout

String	Int	Date
a	1	2020-01
b	2	2020-02
c	3	2020-03

<https://towardsdatascience.com/demystifying-the-parquet-file-format-13adb0206705>



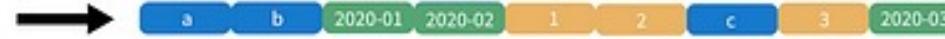
Column-Based Storage Layout

String	Int	Date
a	1	2020-01
b	2	2020-02
c	3	2020-03



Hybrid-Based Storage Layout (row group size = 2)

String	Int	Date
a	1	2020-01
b	2	2020-02
c	3	2020-03



"Final group only has 1 row"

Apache parquet is an open-source file format that provides efficient storage and fast read speed. It uses a hybrid storage format which sequentially stores chunks of columns, lending to high performance when selecting and filtering data. On top of strong compression algorithm support ([snappy](#), [gzip](#), [LZO](#)), it also provides some clever tricks for reducing file scans and encoding repeat variables.

Indexes



Basic Concepts

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

search-key	pointer
------------	---------
- 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”.



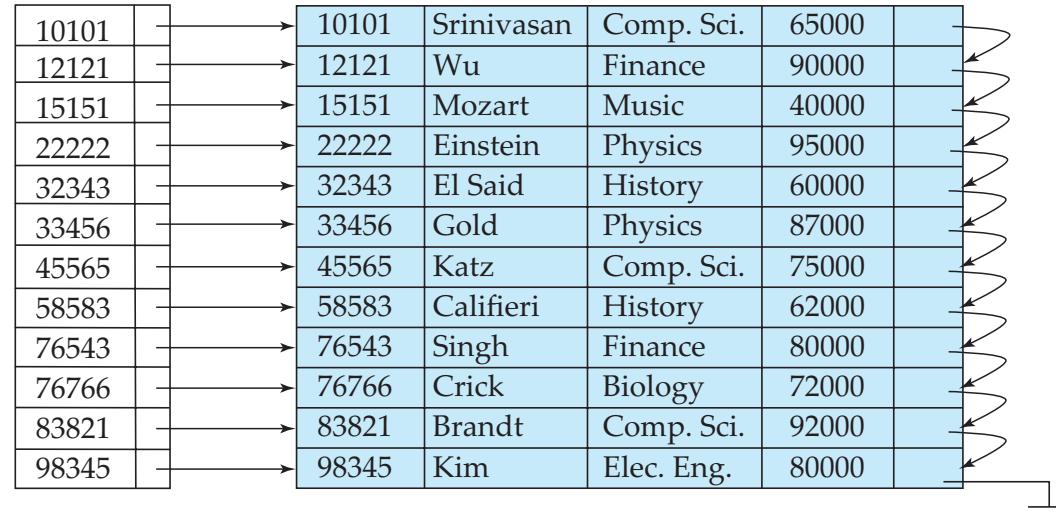
Ordered Indices

- In an **ordered index**, index entries are stored sorted on the search key value.
- **Clustering index:** in a sequentially ordered file, the index whose search key specifies the sequential order of the file.
 - Also called **primary 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 **nonclustering index**.
- **Index-sequential file:** sequential file ordered on a search key, with a clustering index on the search key.



Dense Index Files

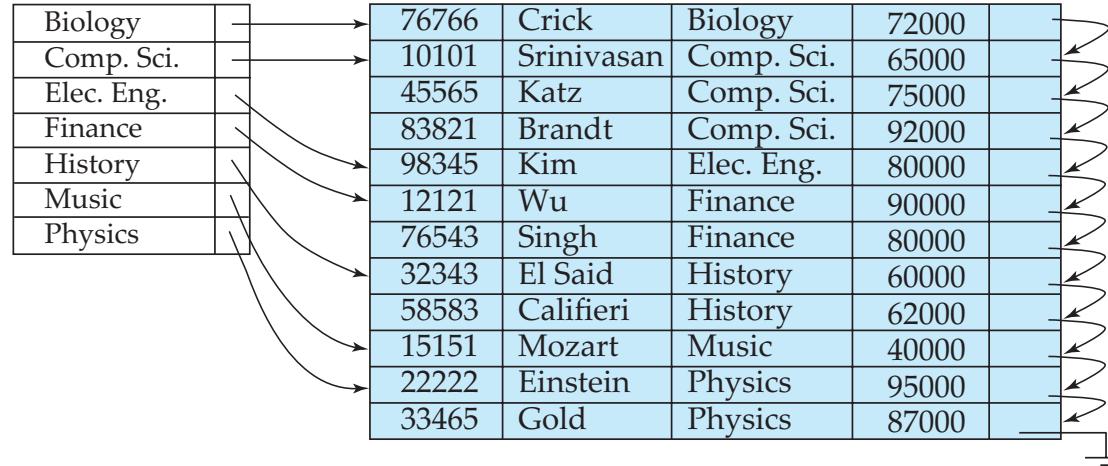
- **Dense index** — Index record appears for every search-key value in the file.
- E.g. index on *ID* attribute of *instructor* relation





Dense Index Files (Cont.)

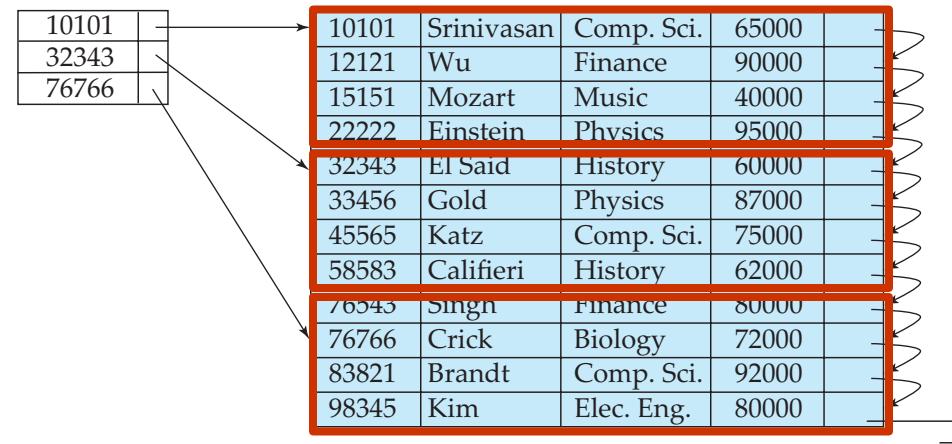
- Dense index on *dept_name*, with *instructor* file sorted on *dept_name*





Sparse Index Files

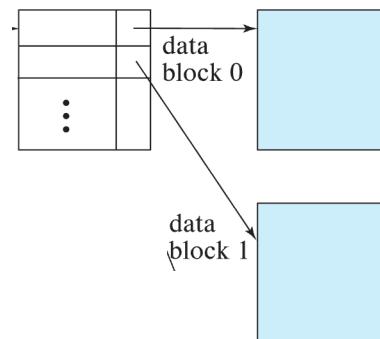
- **Sparse Index:** contains index records for only some search-key values.
 - Applicable when records are sequentially ordered on search-key
- To locate a record with search-key value K we:
 - Find index record with largest search-key value $< K$
 - Search file sequentially starting at the record to 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:**
 - for clustered index: sparse index with an index entry for every block in file, corresponding to least search-key value in the block.

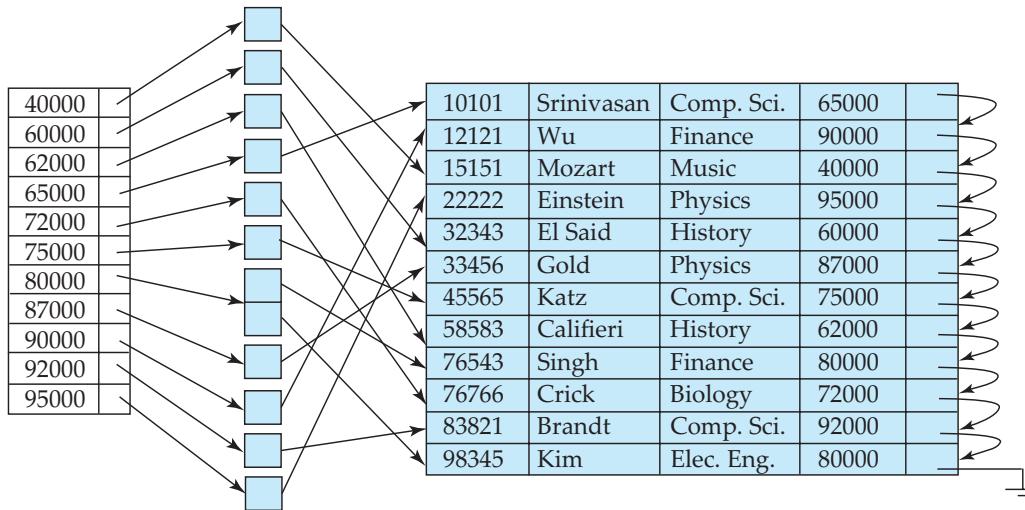


- For unclustered index: sparse index on top of dense index (multilevel index)



Secondary Indices Example

- Secondary index on salary field of instructor



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

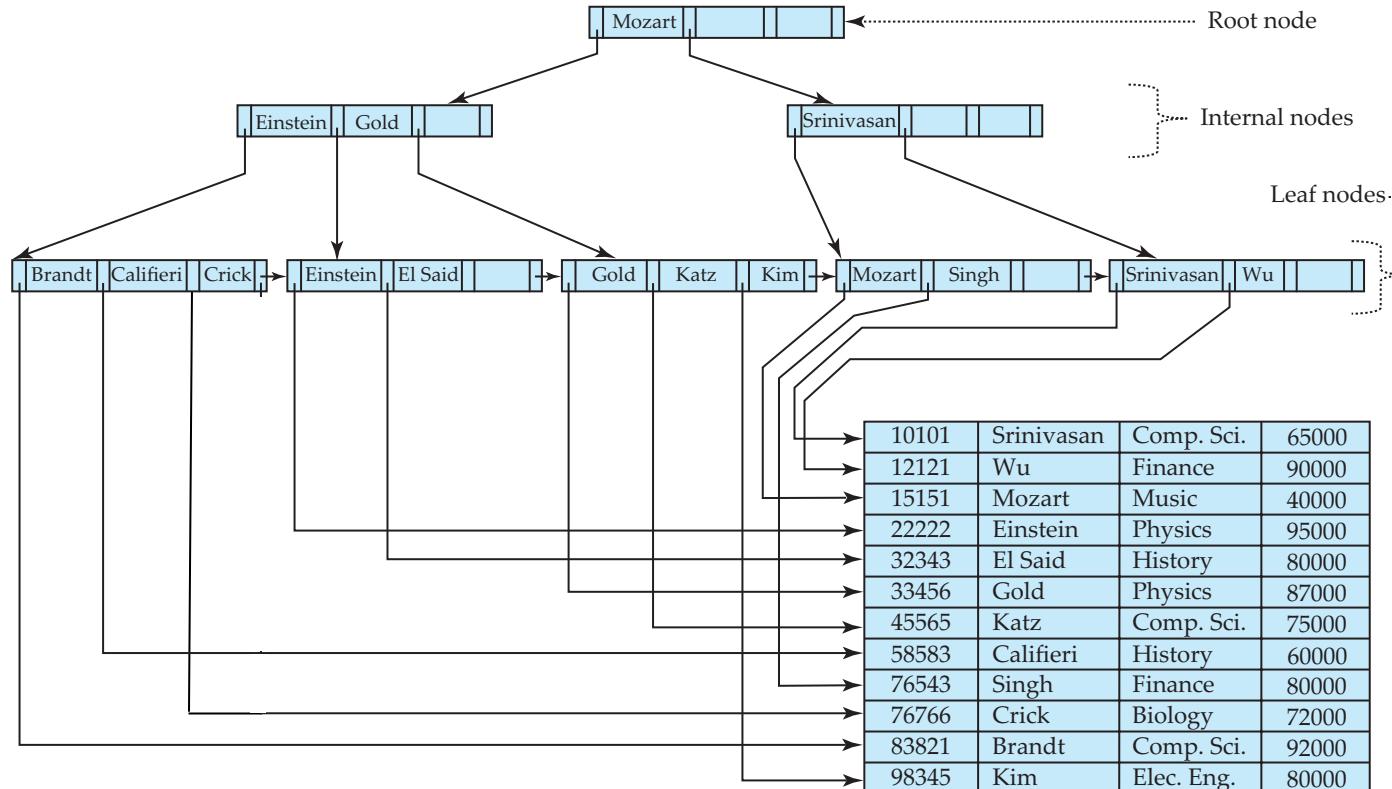


Indices on Multiple Keys

- **Composite search key**
 - E.g., index on *instructor* relation on attributes (*name*, *ID*)
 - Values are sorted lexicographically
 - E.g. (John, 12121) < (John, 13514) and (John, 13514) < (Peter, 11223)
 - Can query on just *name*, or on (*name*, *ID*)
- (nameLast, nameFirst, birthyear)
 - nameLast [nameLast = “Ferguson”] [nameLast like “Fer%”]
 - nameLast, nameFirst
 - nameLast, nameFirst, birthyear
- NOT and index on
 - nameFirst, nameLast
 - birthyear
 - nameLast like [%er%]



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

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

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

- The search-keys in a node are ordered

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

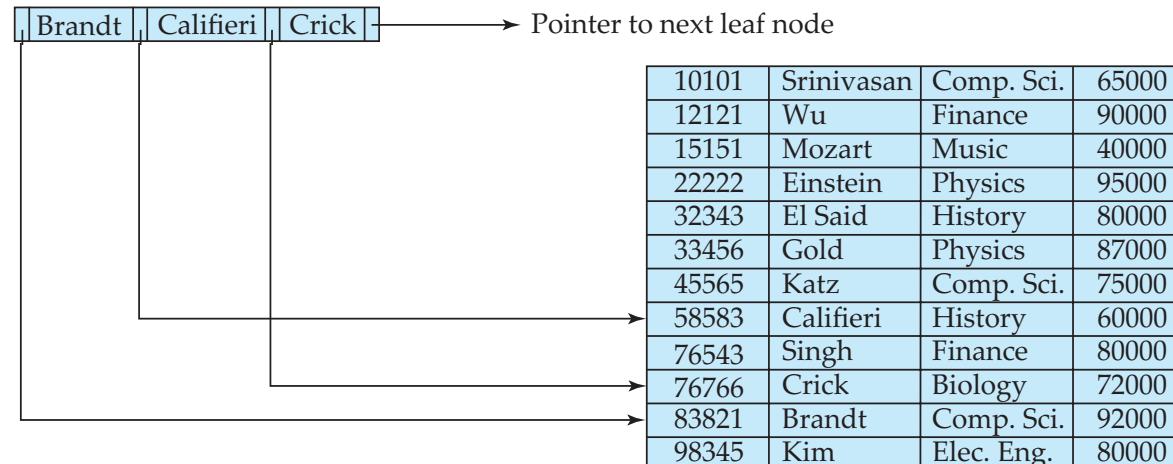
(Initially assume no duplicate keys, address duplicates later)



Leaf Nodes in B⁺-Trees

Properties of a leaf node:

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





Non-Leaf Nodes in B⁺-Trees

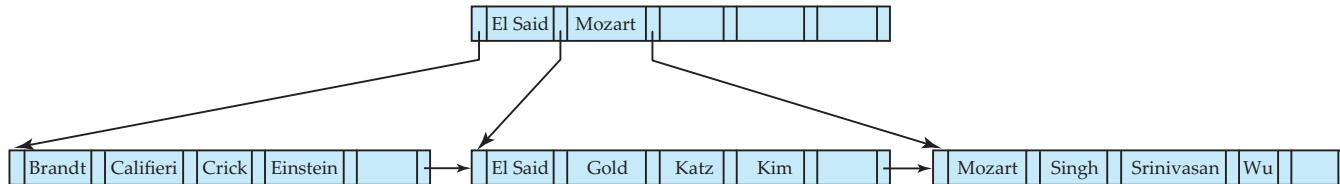
- Non leaf nodes form a multi-level sparse index on the leaf nodes. For a non-leaf node with m pointers:
 - All the search-keys in the subtree to which P_1 points are less than K_1
 - For $2 \leq i \leq n - 1$, all the search-keys in the subtree to which P_i points have values greater than or equal to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

P_1	K_1	P_2	\dots	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.
 - If there are K search-key values in the file, the tree height is no more than $\lceil \log_{\lceil n/2 \rceil}(K) \rceil$
 - thus searches can be conducted efficiently.
- Insertions and deletions to the main file can be handled efficiently, as the index can be restructured in logarithmic time (as we shall see).

Show the Simulator

<https://www.cs.usfca.edu/~galles/visualization/BPlusTree.html>



Hashing



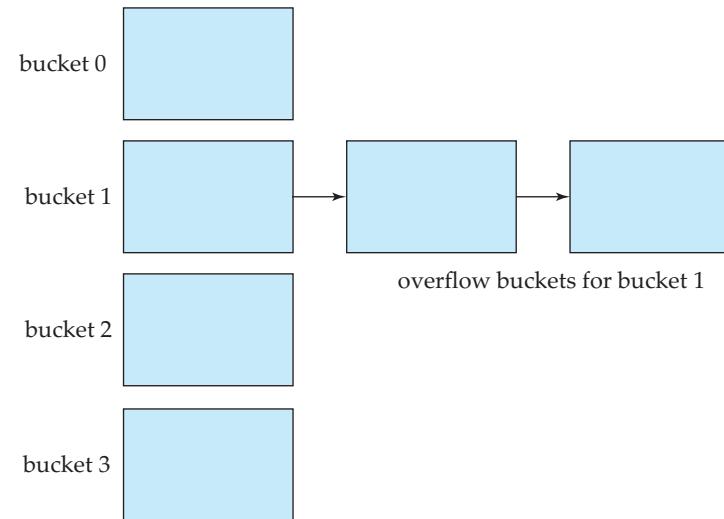
Static Hashing

- A **bucket** is a unit of storage containing one or more entries (a bucket is typically a disk block).
 - we obtain the bucket of an entry 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 entries for access, insertion as well as deletion.
- Entries with different search-key values may be mapped to the same bucket; thus entire bucket has to be searched sequentially to locate an entry.
- In a **hash index**, buckets store entries with pointers to records
- In a **hash file-organization** buckets store records



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 addressing** (also called **closed hashing** or **open hashing** depending on the book you use)
 - An alternative, called **open addressing** (also called **open hashing** or **closed hashing** depending on the book you use) which does not use overflow buckets, is not suitable for database applications.





Example of Hash File Organization

Hash file organization of *instructor* file, using *dept_name* as key.

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



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.

Show the Simulator

<http://iswsa.acm.org/mphf/openDSAPerfectHashAnimation/perfectHashAV.html>
<https://opendsa-server.cs.vt.edu/ODSA/AV/Development/hashAV.html>

NoSQL

Concepts

Overview (I) (<https://en.wikipedia.org/wiki/NoSQL>)

A **NoSQL** (originally referring to "non SQL" or "non relational")^[1] database provides a mechanism for storage and retrieval of data that is modeled in **means other than the tabular relations used in relational databases**. Such databases have existed since the late 1960s, but did not obtain the "NoSQL" moniker until a surge of popularity in the early twenty-first century,^[2] triggered by the needs of Web 2.0 companies such as Facebook, Google, and Amazon.com.^{[3][4][5]} NoSQL databases are increasingly used in big data and real-time web applications.^[6] NoSQL systems are also sometimes called "**Not only SQL**" to emphasize that they may support SQL-like query languages.^{[7][8]}

Motivations for this approach include: simplicity of design, simpler "horizontal scaling" to clusters of machines (which is a problem for relational databases),^[2] and finer control over availability. The data structures used by NoSQL databases (e.g. key-value, wide column, graph, or document) are different from those used by default in relational databases, making **some operations faster in NoSQL**. The particular suitability of a given NoSQL database depends on the problem it must solve. Sometimes the **data structures used by NoSQL databases are also viewed as "more flexible"** than relational database tables.^[9]

Overview (I) (<https://en.wikipedia.org/wiki/NoSQL>)

Many NoSQL stores compromise [consistency](#) (in the sense of the [CAP theorem](#)) in favor of availability, partition tolerance, and speed. Barriers to the greater adoption of NoSQL stores include the use of low-level query languages (instead of SQL, for instance the lack of ability to perform ad-hoc joins across tables), lack of standardized interfaces, and huge previous investments in existing relational databases.^[10] Most NoSQL stores lack true [ACID](#) transactions,

Instead, most NoSQL databases offer a concept of "eventual consistency" in which database changes are propagated to all nodes "eventually" (typically within milliseconds) so queries for data might not return updated data immediately or might result in reading data that is not accurate, a problem known as stale reads.^[11] Additionally, some NoSQL systems may exhibit lost writes and other forms of [data loss](#).^[12] Fortunately, some NoSQL systems provide concepts such as [write-ahead logging](#) to avoid data loss.^[13] For [distributed transaction processing](#) across multiple databases, data consistency is an even bigger challenge that is difficult for both NoSQL and relational databases. Even current relational databases "do not allow referential integrity constraints to span databases."^[14]

Simplistic Classification

(<https://medium.com/swlh/4-types-of-nosql-databases-d88ad21f7d3b>)

Relational is the foundational model.

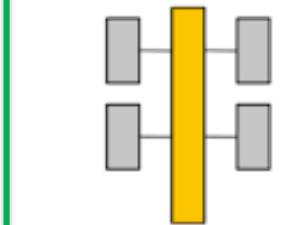
We covered graphs and examples.

SQL Database

Relational



Analytical (OLAP)

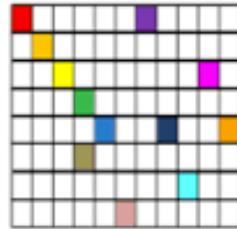


We will see OLAP in a future lecture.

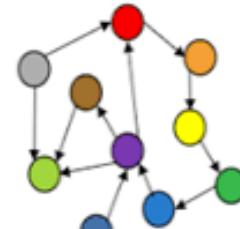
Subject of this lecture and part of HW4

NoSQL Database

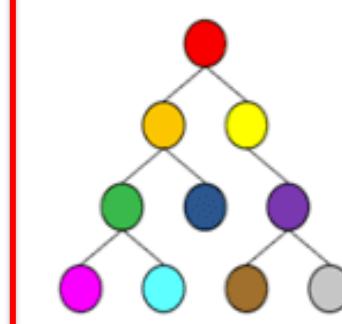
Column-Family



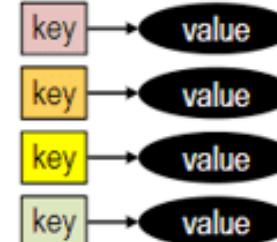
Graph



Document



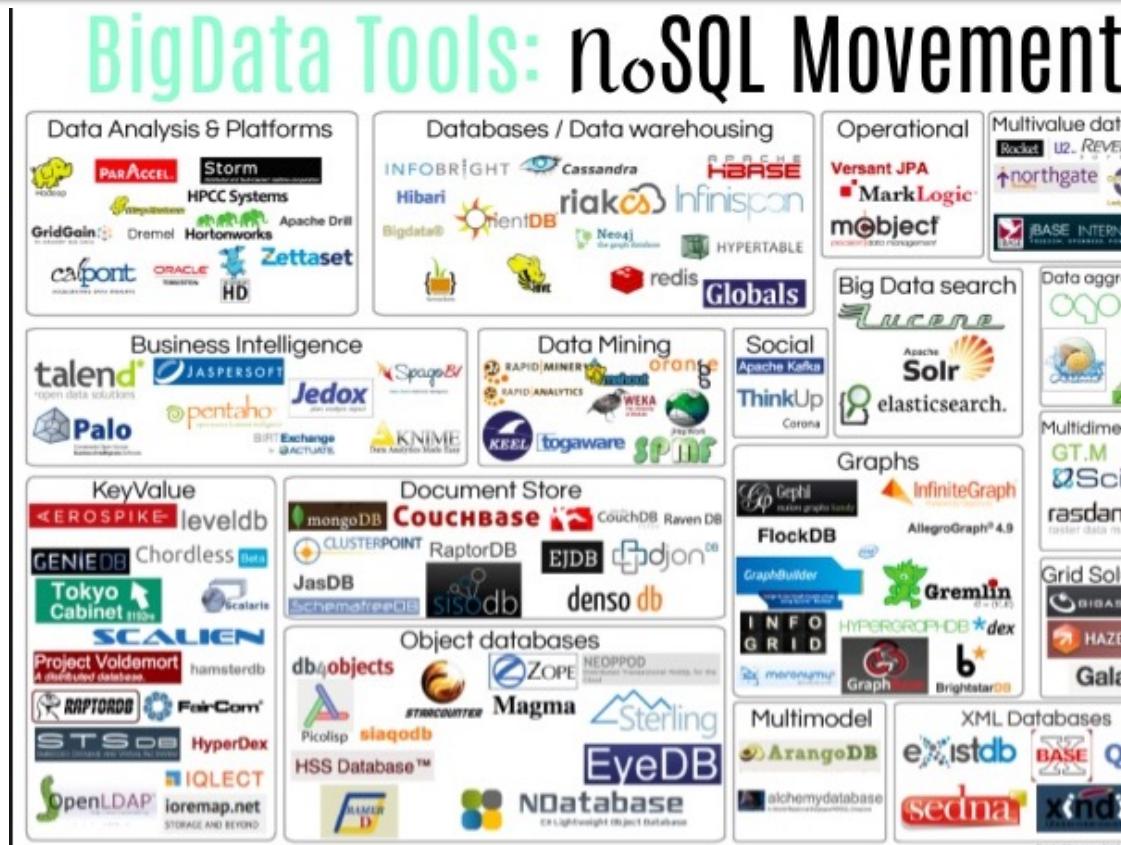
Key-Value



One Taxonomy

Document Database	Graph Databases
   	  The Distributed Graph Database
Wide Column Stores	Key-Value Databases
   	    

Another Taxonomy



Use Cases

Motivations

- Massive write performance.
- Fast key value look ups.
- Flexible schema and data types.
- No single point of failure.
- Fast prototyping and development.
- Out of the box scalability.
- Easy maintenance.

What is wrong with SQL/Relational?

- Nothing. One size fits all? Not really.
- Impedance mismatch. – Object Relational Mapping doesn't work quite well.
- Rigid schema design.
- Harder to scale.
- Replication.
- Joins across multiple nodes? Hard.
- How does RDMS handle data growth? Hard.
- Need for a DBA.
- Many programmers are already familiar with it.
- Transactions and ACID make development easy.
- Lots of tools to use.

Documents

Example Document – An Order

FOOD ORDER FORM TEMPLATE

Company Name
123 Main Street
Hamilton, OH 44416
(321) 456-7890
Email Address
Point of Contact
web address

YOUR LOGO

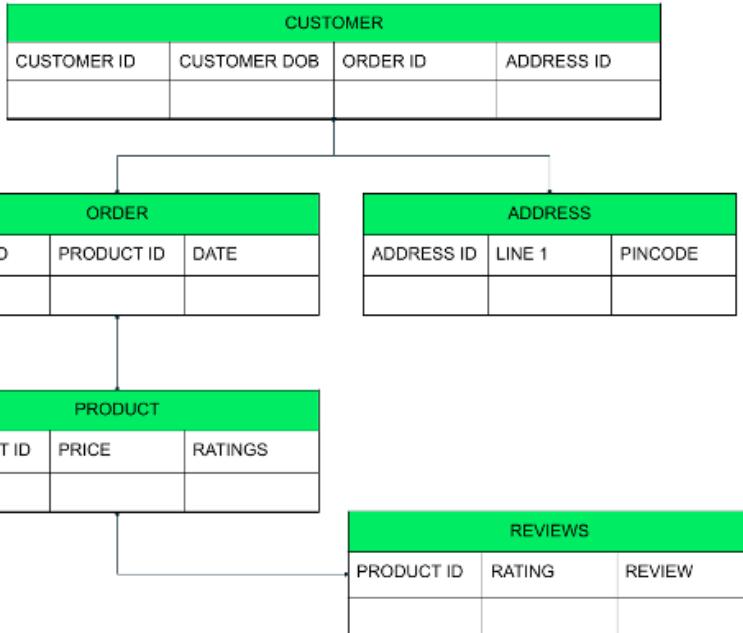
ORDER FORM

CUSTOMER		ORDER NO.	ORDER DATE
ATTN: Name / Dept		DATE NEEDED	TIME NEEDED
Company Name		ORDER RECEIVED BY	
123 Main Street			
Hamilton, OH 44416			
(321) 456-7890			
Email Address			
DESCRIPTION			
		UNIT PRICE	QUANTITY
		\$ -	
		\$ -	
		\$ -	
		\$ -	
		\$ -	
		\$ -	
		\$ -	
		\$ -	
		\$ -	
TOTALS			
enter percentage		TAX RATE	0.000%
enter initial pymt amount		TOTAL TAX	\$ -
		DELIVERY	\$ -
		GRAND TOTAL	\$ -
		LESS PAYMENT	\$ -
THANK YOU!			

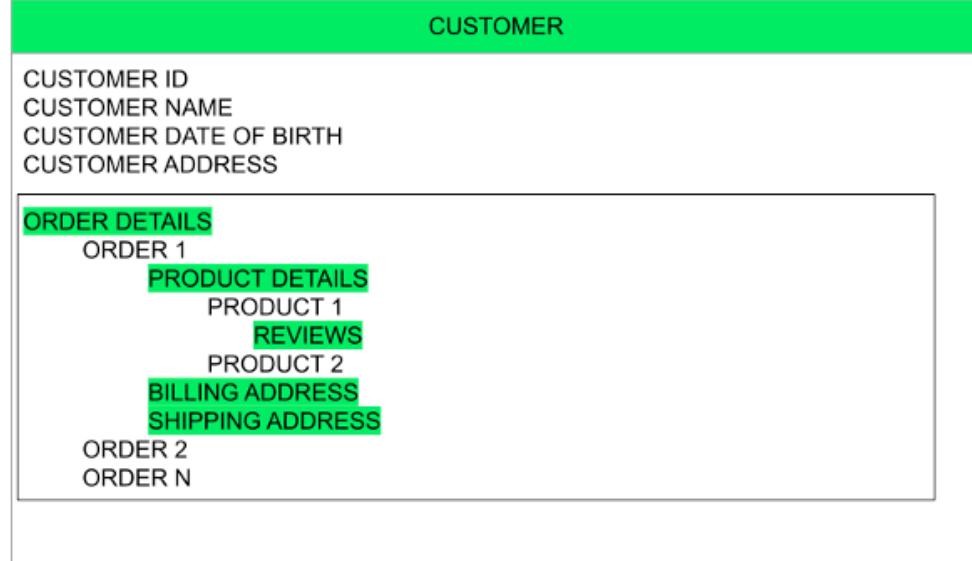
- There are 5 entity types on the form:
 - A “copy of” customer information.
 - Links (via productCode) to products
 - Order
 - OrderDetails
 - Comments
- But OrderDetails and Comments are somehow different from the others.
 - These are arrays of objects
 - That are sort of “inside” the order.
- These are weak entities, but relational does not always handle these well.

Relational vs Document

Relational



Relational



The key difference?

- Relational DB attributes are atomic.
- Document DB attributes may be *multi-valued* and *complex*.

UI Interface and Logical Data Model

{

Columbia University

CUSTOMER NUMBER

10001

CONTACT NAME

Lord Voldemort

CONTACT NUMBER

+1 212-555-6666

Customer Address (THIS WAS INSANELY PAINFUL TO LAYOUT)

STREET 1

520 W 120th St.

STREET 2

4th Floor

CITY

New York

STATE

NY

COUNTRY

US

POSTAL CODE

10027

I nearly went "postal" doing the previous layout and almost failed all of you!

NUMBER	STATUS	ORDER DATE	REQUIRED DATE	SHIPPED DATE	DETAILS
21	Shipped	2022-01-01	2022-02-01	2022-01-28	<button>Expand</button>

LINE NUMBER PRODUCT CODE QUANTITY PRICE EACH

0	P0	0	0
1	P1	2	3
2	P2	4	6

NUMBER	STATUS	ORDER DATE	REQUIRED DATE	SHIPPED DATE	DETAILS
22	Pending	2022-03-01	2022-04-01	0000000	<button>Expand</button>

customerName: "Columbia University",

contact: {

name: "Lord Voldemort",

phone: "+1 212-555-6666"

},

address: {

line1: "520 W 120th St."

line2: "Floor 4",

... ...

}

orders: [

{

number: ...,

... ...

details: [

... ...

]

}]

Switch to Notebook

MongoDB

MongoDB Concepts

RDBMS	MongoDB
Database	Database
Table	Collection
Tuple/Row	Document
column	Field
Table Join	Embedded Documents
Primary Key	Primary Key (Default key <code>_id</code> provided by MongoDB itself)
Database Server, Client, Tools, Packages	
mysqld/Oracle	<code>mongod</code>
mysql/sqlplus	<code>mongo</code>
DataGrip	Compass
pymysql	<code>pymongo</code>

Core Operations

Basic Operations:

- Create database
- Create collection
- Create-Retrieve-Update-Delete (CRUD):
 - Create: insert()
 - Retrieve:
 - find()
 - find_one()
 - Update: update()
 - Delete: remove()

More Advanced Concepts:

- Limit
- Sort
- Aggregation Pipelines
 - Merge
 - Union
 - Lookup
 - Match
 - Merge
 - Sample
 -

We will just cover the basics for now and may cover more things in HW or other lectures.

find()

- Note:
 - MongoDB uses a more `pymysql` approach, e.g. an API, than pure declarative languages like SQL.
 - The parameters for `find()` are where the declarative language appears.
- The basic forms of `find()` and `find_one()` have two parameters:
 - *filter expression*
 - *Project expression*
- You can use the Compass tool and screen captures for some HW and exam answers.
- What if I want the answer in a Jupyter Notebook?

The screenshot shows the MongoDB Compass interface with the following details:

- Collection:** GOT.seasons
- Documents:** 8 (TOTAL SIZE 1.0MB, AVG. SIZE 130.2KB)
- Indexes:** 1 (TOTAL SIZE 20.0KB, AVG. SIZE 20.0KB)
- Filter:** `{"episodes.scenes.location": "The Dothraki Sea"}`
- Project:** `{ season: 1, "episodes.episodeNum":1, "episodes.episodeLink": 1, "episodes.episodeTitle": 1 }`
- Sort:** `{ field: -1 }`
- Collation:** `{ locale: 'simple' }`
- Options:** MAX TIME MS 60000, SKIP 0, LIMIT 0
- Results:** Displaying documents 1 - 3 of 3
- Document 1:** _id: ObjectId("60577e50c68b67110968b6d1"), season: "1", episodes: Array (9 elements), each element containing episodeNum: 1, episodeTitle: "Winter Is Coming", episodeLink: "/title/tt1480055/", followed by 8 more objects.
- Document 2:** _id: ObjectId("60577e50c68b67110968b6d5"), season: "5", episodes: Array (1 element).
- Document 3:** _id: ObjectId("60577e50c68b67110968b6d6"), season: "5", episodes: Array (1 element).

Generate Code

The screenshot shows the MongoDB Compass interface. On the left, the sidebar lists databases (Local, HOST localhost:27017, CLUSTER Standalone, EDITION MongoDB 4.2.6 Community) and collections (GOT, characters, seasons, admin, config, db, fantasy_baseball, local). The 'seasons' collection is selected. In the main area, the 'GOT.seasons' document list shows two documents with IDs 60577e50c68b67110968b6d5 and 60577e50c68b67110968b6d6. A modal dialog titled 'Export Query To Language' is open, showing a query builder interface. The 'My Query:' section contains the following MongoDB query:1 \n2 { \n3 "episodes.scenes.location": "The Dothraki Sea" \n4 }

```
The 'Export Query To' dropdown is set to 'PYTHON 3'. The generated Python code is:
```

1 # Requires the PyMongo package.\n2 # https://api.mongodb.com/python/current\n3\n4 client = MongoClient('mongodb://localhost:27017/?\n5 \n6 \n7 \n8 \n9 \n10 \n11 \n12 \n13 \n14 \n15 \n16 result = client['GOT']['seasons'].find(\n filter=filter,\n projection=\n { 'season': 1,\n 'episodes.episodeNum': 1,\n 'episodes.episodeLink': 1,\n 'episodes.episodeTitle': 1\n }\n)\n\nfor season in result:\n print(season)\n\n for episode in season['episodes']:\n print(episode)

```
Below the code, there are two checkboxes: 'Include Import Statements' (unchecked) and 'Include Driver Syntax' (checked). A red box highlights the '...' button in the top right corner of the modal.
```

- Choose Export to Language.
- Copy into the notebook.
- The export has an option to include all the connection setup, choosing DB,
- Switch to Notebook

Result is not Quite You Expect

GOT.seasons

DOCUMENTS 8 TOTAL SIZE 1.0MB AVG. SIZE 130.2KB INDEXES 1 TOTAL SIZE 20.0KB AVG. SIZE 20.0KB

Documents Aggregations Schema Explain Plan Indexes Validation

FILTER: {"episodes.scenes.location": "The Dothraki Sea"}
PROJECT: { season: 1, "episodes.episodeNum":1, "episodes.episodeLink": 1, "episodes.episodeTitle": 1}
SORT: { field: -1 }
COLLATION: { locale: 'simple' }

FIND RESET ...

MAX TIME MS: 60000
SKIP: 0 LIMIT: 0

VIEW

Displaying documents 1 - 3 of 3 < > C REFRESH

```
_id:ObjectId("60577e50c68b67110968b6d1")
season:"5"
episodes:Array
  ▾ 0:Object
    episodeNum:1
    episodeTitle:"Winter Is Coming"
    episodeLink:"/title/tt1480055"
  ▾ 1:Object
  ▾ 2:Object
  ▾ 3:Object
    episodeNum:4
    episodeTitle:"Cripples, Bastards, and Broken Things"
    episodeLink:"/title/tt1829963"
  ▾ 4:Object
  ▾ 5:Object
  ▾ 6:Object
  ▾ 7:Object
  ▾ 8:Object
  ▾ 9:Object

_id:ObjectId("60577e50c68b67110968b6d5")
season:"5"
episodes:Array

_id:ObjectId("60577e50c68b67110968b6d6")
season:"6"
episodes:Array
```

- The query returns documents that match.
 - The document is “Large” and has seasons and episodes and seasons.
 - If you do a \$project requesting episodes/episode content,
 - You get all episodes in the documents that match.
 - Not just the episodes with the scene/location.
 - Projecting array elements from arrays whose elements are arrays is complex and baffling.
- You also get back something (a cursor) that is iterable.

Result is not Quite You Expect

The screenshot shows the MongoDB Compass interface with a query results grid and the MongoDB shell output below it.

Query Results Grid:

_id	season	episodes
60577e50c68b6711096b6d1	"5"	[{"episodeNum": 1, "episodeTitle": "Winter Is Coming", "episodeLink": "/title/tt1480055"}, {"episodeNum": 2, "episodeTitle": "The Night King", "episodeLink": "/title/tt1480056"}, {"episodeNum": 3, "episodeTitle": "Cripples, Bastards, and Broken Men", "episodeLink": "/title/tt1829963"}, {"episodeNum": 4, "episodeTitle": "The Children", "episodeLink": "/title/tt1480057"}, {"episodeNum": 5, "episodeTitle": "The Mountain and the Viper", "episodeLink": "/title/tt1480058"}, {"episodeNum": 6, "episodeTitle": "The Red and the Black", "episodeLink": "/title/tt1480059"}, {"episodeNum": 7, "episodeTitle": "The Kingsguard Must Die", "episodeLink": "/title/tt1480060"}, {"episodeNum": 8, "episodeTitle": "The Rains of Castamere", "episodeLink": "/title/tt1480061"}, {"episodeNum": 9, "episodeTitle": "The Bells", "episodeLink": "/title/tt1480062"}]

MongoDB Shell Output:

```
_id:ObjectId("60577e50c68b6711096b6d1")
season:"5"
> episodes:Array
  ▾ 0:Object
    episodeNum:1
    episodeTitle:"Winter Is Coming"
    episodeLink:"/title/tt1480055/"
  ▾ 1:Object
  ▾ 2:Object
  ▾ 3:Object
    episodeNum:4
    episodeTitle:"Cripples, Bastards, and
    episodeLink:"/title/tt1829963/"
  ▾ 4:Object
  ▾ 5:Object
  ▾ 6:Object
  ▾ 7:Object
  ▾ 8:Object
  ▾ 9:Object

_id:ObjectId("60577e50c68b6711096b6d5")
season:"5"
> episodes:Array

_id:ObjectId("60577e50c68b6711096b6d6")
season:"6"
> episodes:Array
```

Net for HWs and exams:

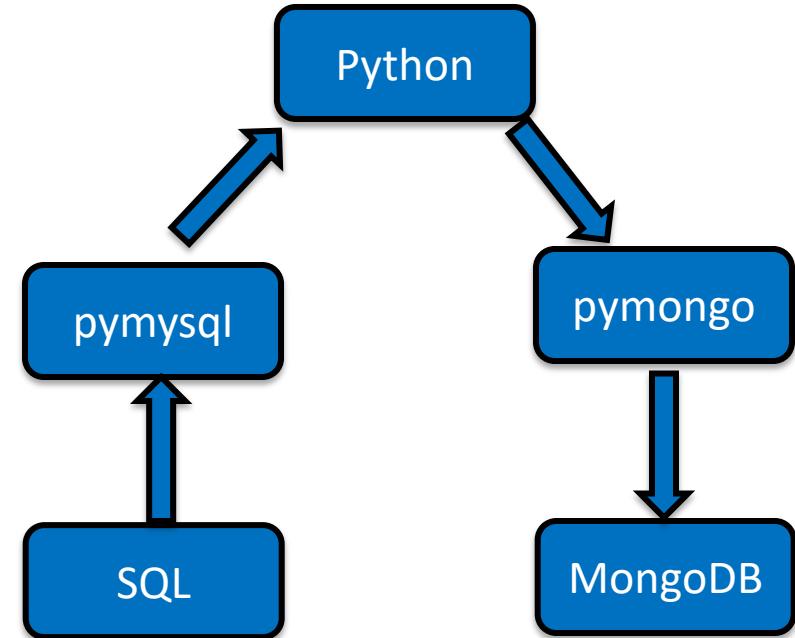
- We will keep the queries simple.
- The language is as complex as SQL, and we spent several weeks on the language.
- Let's take a look in the Jupyter Notebook on some create and insert functions.
- But, first, the datamodel impedance mismatch concept.

- The query returns documents that match.
 - The document is “Large” and has many episodes and seasons.
 - You do a \$project requesting episodes/episode content,
 - You get all episodes in the documents that match.
 - Not just the episodes with the scene/location.
 - Projecting array elements from arrays whose elements are arrays is complex and baffling.
 - You also get back something (a cursor) that is iterable.

More Fun – Data Types

MongoDB supports many datatypes. Some of them are –

- **String** – This is the most commonly used datatype to store the data. String in MongoDB must be UTF-8 valid.
- **Integer** – This type is used to store a numerical value. Integer can be 32 bit or 64 bit depending upon your server.
- **Boolean** – This type is used to store a boolean (true/ false) value.
- **Double** – This type is used to store floating point values.
- **Min/ Max keys** – This type is used to compare a value against the lowest and highest BSON elements.
- **Arrays** – This type is used to store arrays or list or multiple values into one key.
- **Timestamp** – timestamp. This can be handy for recording when a document has been modified or added.
- **Object** – This datatype is used for embedded documents.
- **Null** – This type is used to store a Null value.
- **Symbol** – This datatype is used identically to a string; however, it's generally reserved for languages that use a specific symbol type.
- **Date** – This datatype is used to store the current date or time in UNIX time format. You can specify your own date time by creating object of Date and passing day, month, year into it.
- **Object ID** – This datatype is used to store the document's ID.
- **Binary data** – This datatype is used to store binary data.
- **Code** – This datatype is used to store JavaScript code into the document.
- **Regular expression** – This datatype is used to store regular expression.



(Some) MongoDB CRUD Operations

- Create:
 - db.collection.insertOne()
 - db.collection.insertMany()
- Retrieve:
 - db.collection.find()
 - db.collection.findOne()
 - db.collection.findOneAndUpdate()
 -
- Update:
 - db.collection.updateOne()
 - db.collection.updateMany()
 - db.collection.replaceOne()
- Delete:
 - db.collection.deleteOne()
 - db.collection.deleteMany()

pymongo maps the camel case to _, e.g.

- findOne()
- find_one()

There are good online tutorials:

- https://www.tutorialspoint.com/python_data_access
- <https://www.tutorialspoint.com/mongodb/index.htm>

(Some) MongoDB Pipeline Operators

<https://www.slideshare.net/mongodb/s01-e04-analytics>

Aggregation operators

- Pipeline and Expression operators

Pipeline	Expression	Arithmetic	Conditional
\$match	\$addToSet	\$add	\$cond
\$sort	\$first	\$divide	\$ifNull
\$limit	\$last	\$mod	
\$skip	\$max	\$multiply	
\$project	\$min	\$subtract	
\$unwind	\$avg		Variables
\$group	\$push		
\$geoNear	\$sum		
\$text			\$let
\$search			\$map

Tip: Other operators for date, time, boolean and string manipulation

MongoDB Checkpoint

- You can see that MongoDB has powerful, sophisticated
 - Operators
 - Expressions
 - Pipelines
- We have only skimmed the surface. There is a lot more:
 - Indexes
 - Replication, Sharding
 - Embedded Map-Reduce support
 -
- We will explore a little more in subsequent lectures, homework,
- You will have to install MongoDB and Compass for HW4 and final exam.