

Introduction to Data Curation:

Databases, Data Management, Data Cleaning,

Tim Mattson

FIELDS ARRANGED BY PURITY

→
MORE PURE

SOCIOLOGY IS
JUST APPLIED
PSYCHOLOGY

PSYCHOLOGY IS
JUST APPLIED
BIOLOGY.

BIOLOGY IS
JUST APPLIED
CHEMISTRY

WHICH IS JUST
APPLIED PHYSICS.
IT'S NICE TO
BE ON TOP.

OH, HEY, I DIDN'T
SEE YOU GUYS ALL
THE WAY OVER THERE.



SOCIOLOGISTS



PSYCHOLOGISTS



BIOLOGISTS



CHEMISTS

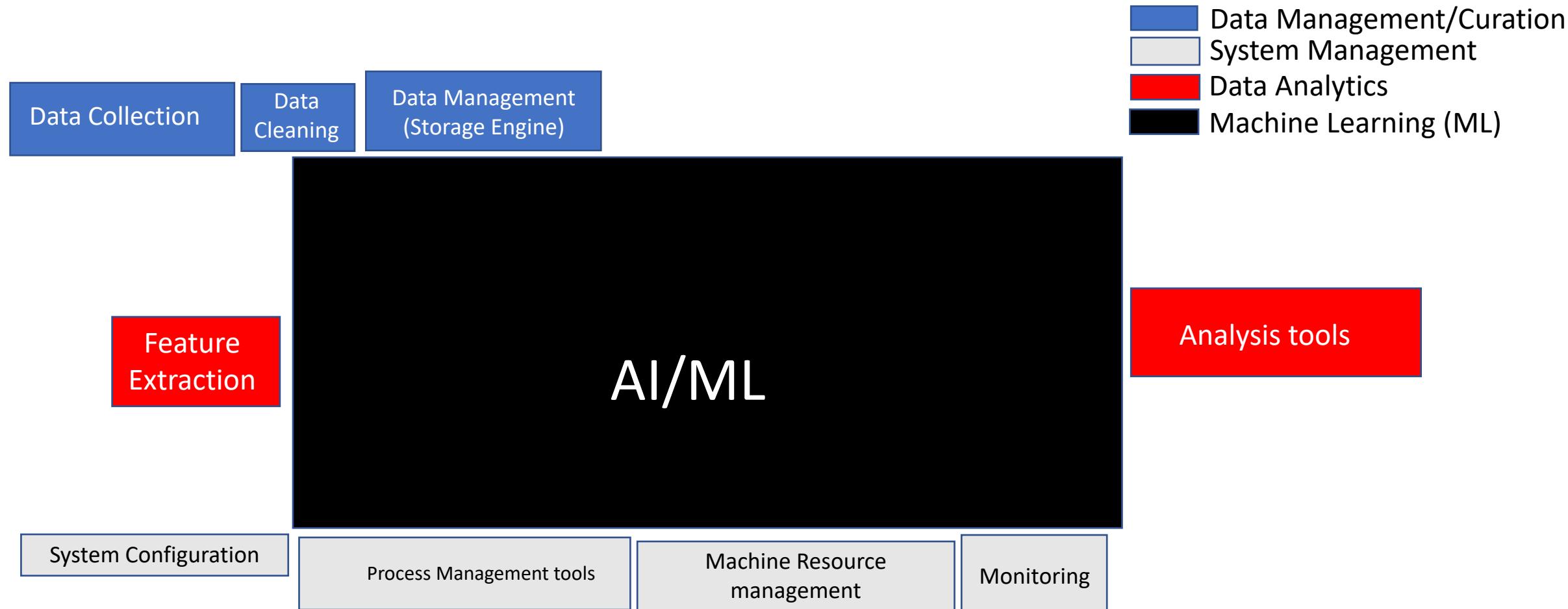


PHYSICISTS



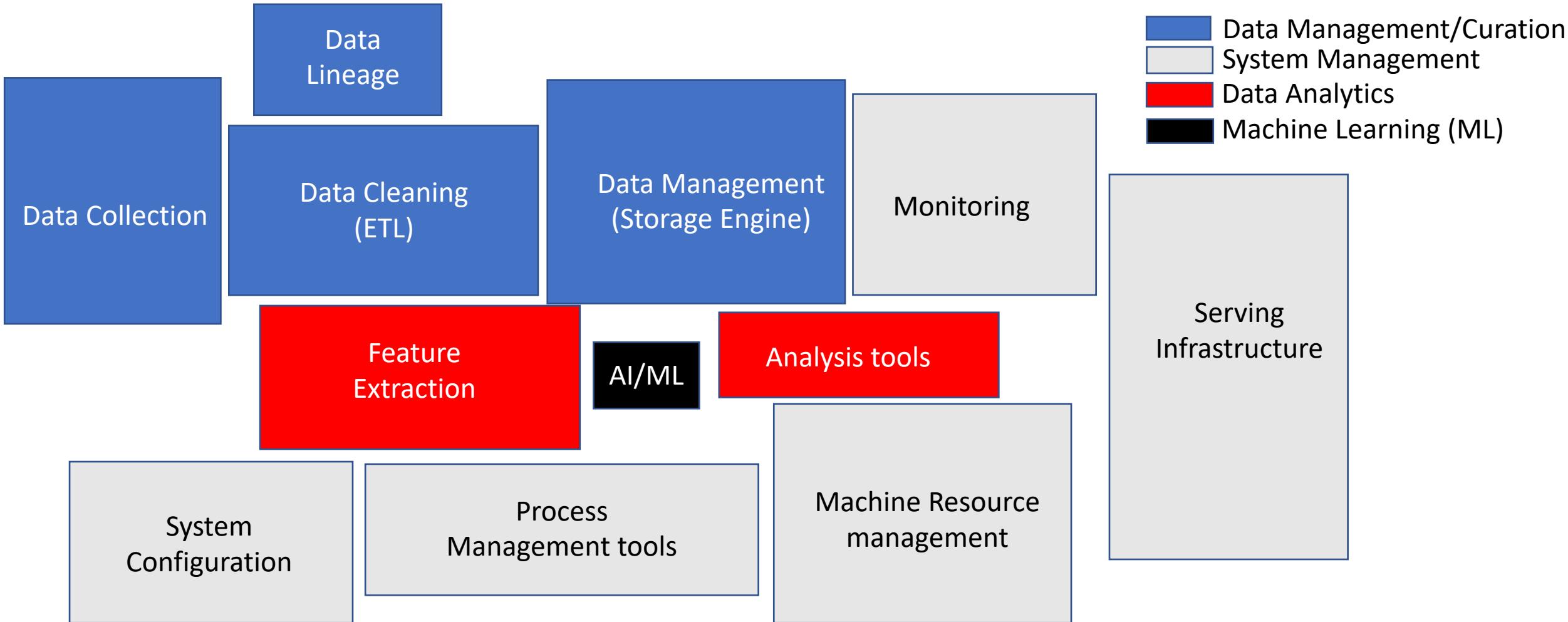
MATHEMATICIANS

The hype suggests a world looks like this ...



Area of rectangles approximates effort
spent working in the indicated domain
within a data sciences workflow

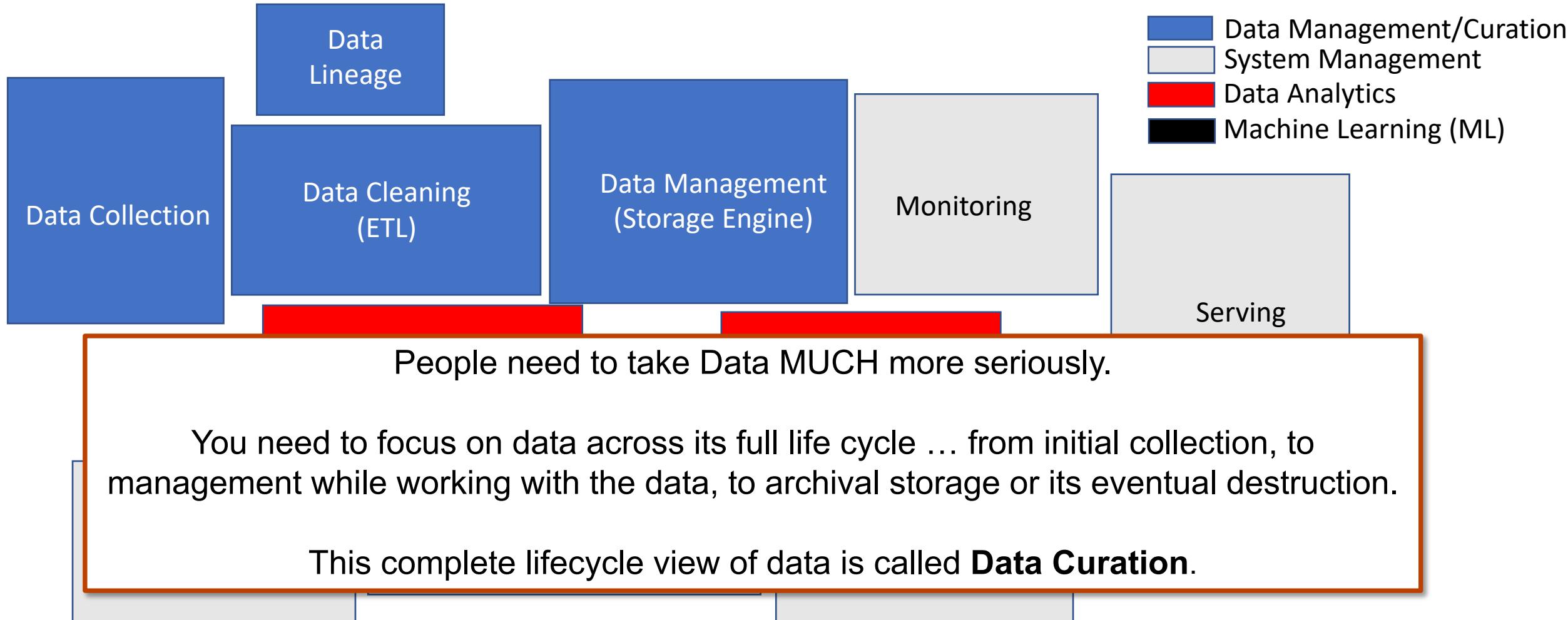
AI/ML developers see the world like this ...



AI/ML code is 2 to 5 percent of the code written by a data scientist*. Most of the code is “glue code” to manage data and system components

*based on figure 1 from “Hidden Technical Debt in Machine Learning Systems”, by D. Sculley et. al. from Google.

Users of AI/ML see the world like this ...



AI/ML code is 2 to 5 percent of the code written by a data scientist*. Most of the code is “glue code” to manage data and system components

*based on figure 1 from “Hidden Technical Debt in Machine Learning Systems”, by D. Sculley et. al. from Google.

Outline

- ➡ • Motivation: Why everyone needs a database management system?
- Database Technology: from ancient history to today
- Data Curation in the sciences
- My quest: One Algebra to rule them all

Raw Data from my camera

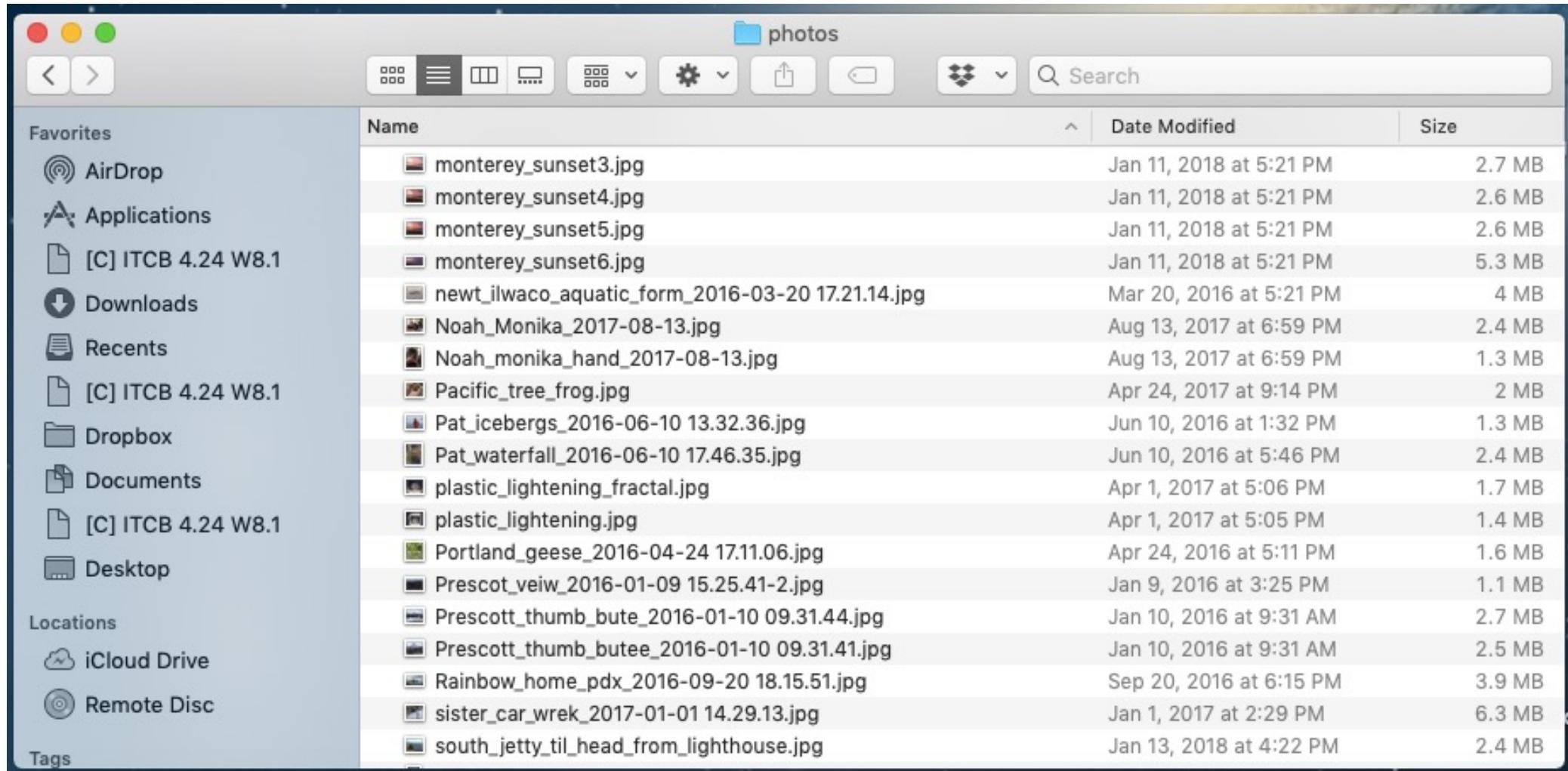
Unstructured, flat files

Data about the data is called metadata

Name	Date Modified	Size
img1.jpg	Aug 29, 2015 at 12:31 PM	3.6 MB
img2.jpg	Aug 29, 2015 at 12:40 PM	8.2 MB
img3.jpg	Aug 29, 2015 at 12:33 PM	9 MB
img4.jpg	Aug 29, 2015 at 12:16 PM	7.6 MB
img5.jpg	Aug 29, 2015 at 12:52 PM	9 MB
img6.jpg	Feb 19, 2016 at 12:42 PM	2.6 MB
img7.jpg	Aug 29, 2015 at 12:46 PM	3.2 MB
img8.jpg	Aug 29, 2015 at 5:40 PM	2.9 MB
img9.jpg	Aug 29, 2015 at 12:16 PM	7.5 MB
img10.jpg	Dec 19, 2017 at 4:19 PM	1.2 MB
img11.jpg	Aug 29, 2015 at 12:51 PM	5.4 MB
img12.jpg	Dec 25, 2017 at 4:29 PM	1.2 MB
img13.jpg	Dec 19, 2017 at 4:24 PM	1.4 MB
img14.jpg	Aug 29, 2015 at 12:33 PM	5.9 MB
img15.jpg	Aug 29, 2015 at 12:57 PM	10.2 MB
img16.jpg	Aug 29, 2015 at 1:07 PM	10.8 MB
img17.jpg	Aug 29, 2015 at 1:49 PM	7.2 MB
img19.jpg	Aug 29, 2015 at 12:31 PM	3.5 MB
img20.jpg	Aug 29, 2015 at 12:16 PM	6.6 MB

Raw Data from my camera

Unstructured, flat files. Attributes of the data captured in file names

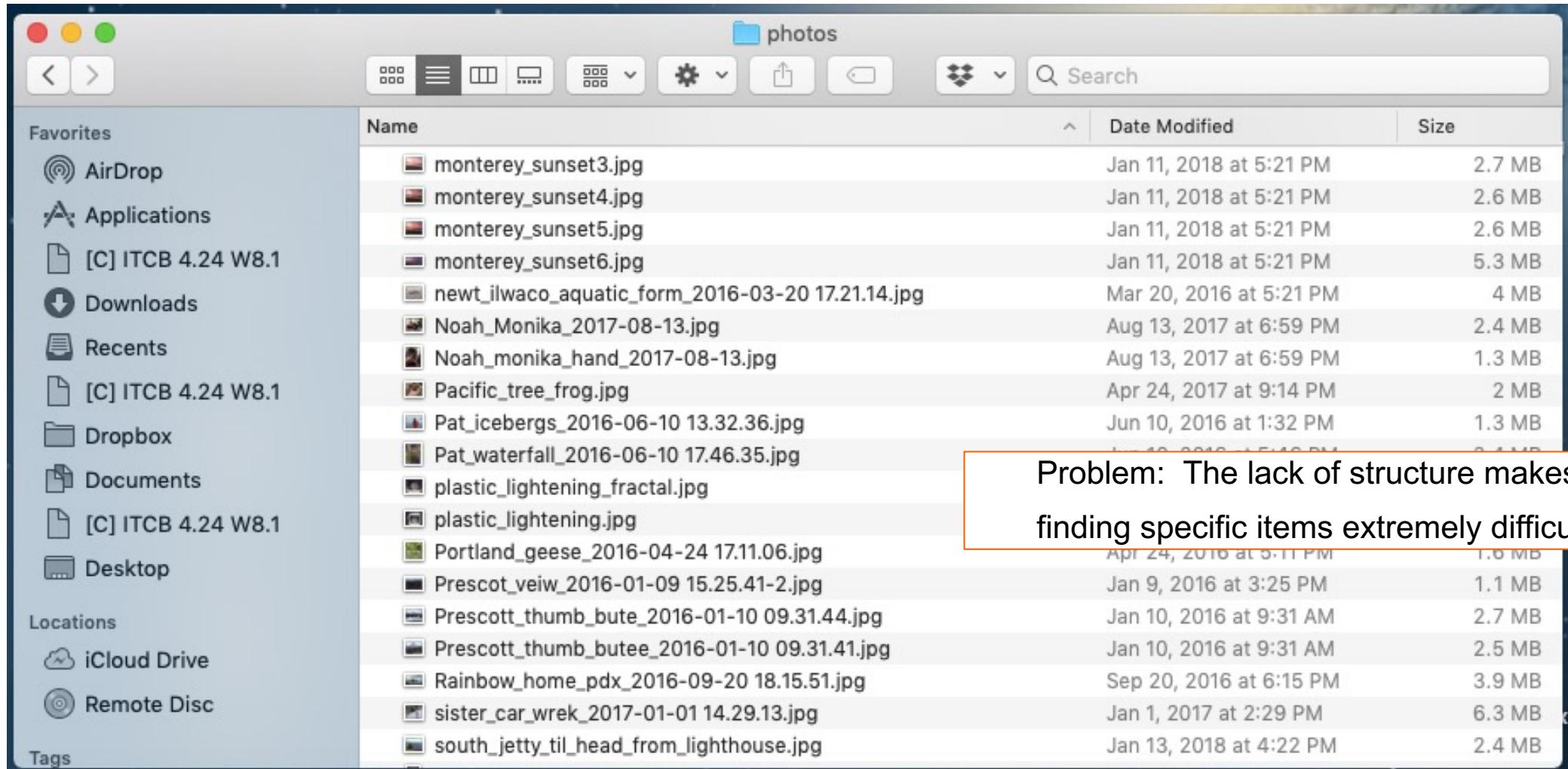


A screenshot of a Mac OS X Photos application window titled "photos". The window shows a list of image files in a table format. The columns are "Name", "Date Modified", and "Size". The "Name" column lists various image files with descriptive names like "monterey_sunset3.jpg", "Noah_Monika_2017-08-13.jpg", and "Rainbow_home_pdx_2016-09-20 18.15.51.jpg". The "Date Modified" column shows the last modification date for each file, and the "Size" column shows the file size in megabytes. The sidebar on the left contains "Favorites" and "Locations" sections, with items such as AirDrop, Applications, Downloads, Recents, and Desktop.

Name	Date Modified	Size
monterey_sunset3.jpg	Jan 11, 2018 at 5:21 PM	2.7 MB
monterey_sunset4.jpg	Jan 11, 2018 at 5:21 PM	2.6 MB
monterey_sunset5.jpg	Jan 11, 2018 at 5:21 PM	2.6 MB
monterey_sunset6.jpg	Jan 11, 2018 at 5:21 PM	5.3 MB
newt_ilwaco_aquatic_form_2016-03-20 17.21.14.jpg	Mar 20, 2016 at 5:21 PM	4 MB
Noah_Monika_2017-08-13.jpg	Aug 13, 2017 at 6:59 PM	2.4 MB
Noah_monika_hand_2017-08-13.jpg	Aug 13, 2017 at 6:59 PM	1.3 MB
Pacific_tree_frog.jpg	Apr 24, 2017 at 9:14 PM	2 MB
Pat_icebergs_2016-06-10 13.32.36.jpg	Jun 10, 2016 at 1:32 PM	1.3 MB
Pat_waterfall_2016-06-10 17.46.35.jpg	Jun 10, 2016 at 5:46 PM	2.4 MB
plastic_lightening_fractal.jpg	Apr 1, 2017 at 5:06 PM	1.7 MB
plastic_lightening.jpg	Apr 1, 2017 at 5:05 PM	1.4 MB
Portland_geese_2016-04-24 17.11.06.jpg	Apr 24, 2016 at 5:11 PM	1.6 MB
Prescot_veiw_2016-01-09 15.25.41-2.jpg	Jan 9, 2016 at 3:25 PM	1.1 MB
Prescott_thumb_bute_2016-01-10 09.31.44.jpg	Jan 10, 2016 at 9:31 AM	2.7 MB
Prescott_thumb_butee_2016-01-10 09.31.41.jpg	Jan 10, 2016 at 9:31 AM	2.5 MB
Rainbow_home_pdx_2016-09-20 18.15.51.jpg	Sep 20, 2016 at 6:15 PM	3.9 MB
sister_car_wrek_2017-01-01 14.29.13.jpg	Jan 1, 2017 at 2:29 PM	6.3 MB
south_jetty_til_head_from_lighthouse.jpg	Jan 13, 2018 at 4:22 PM	2.4 MB

Raw Data from my camera

Unstructured, flat files. Attributes of the data captured in file names



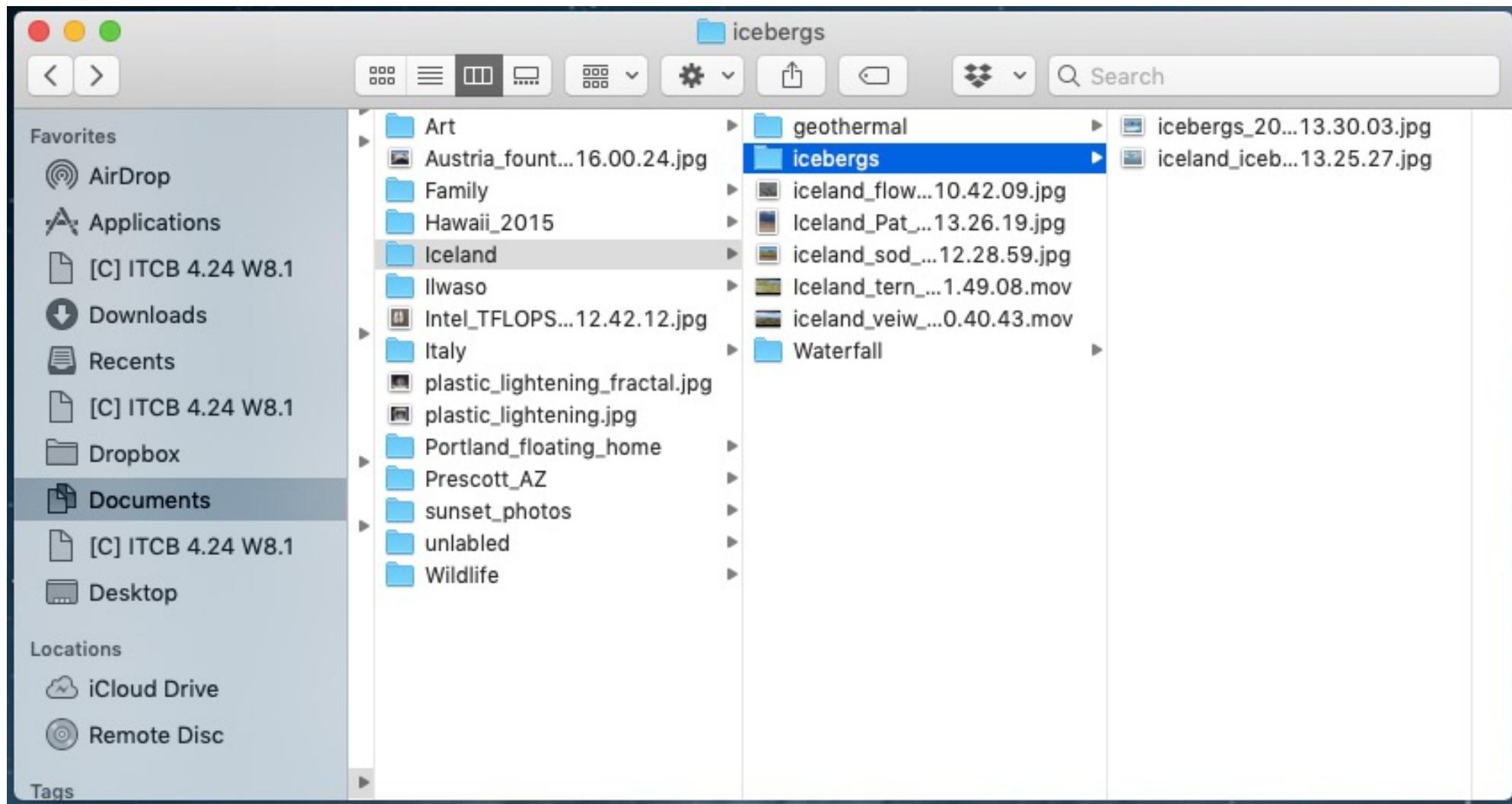
The screenshot shows a Mac OS X Photos application window titled "photos". The sidebar on the left contains "Favorites" (Airdrop, Applications, Downloads, Recents) and "Locations" (iCloud Drive, Remote Disc). The main pane displays a list of image files in a table format:

Name	Date Modified	Size
monterey_sunset3.jpg	Jan 11, 2018 at 5:21 PM	2.7 MB
monterey_sunset4.jpg	Jan 11, 2018 at 5:21 PM	2.6 MB
monterey_sunset5.jpg	Jan 11, 2018 at 5:21 PM	2.6 MB
monterey_sunset6.jpg	Jan 11, 2018 at 5:21 PM	5.3 MB
newt_ilwaco_aquatic_form_2016-03-20 17.21.14.jpg	Mar 20, 2016 at 5:21 PM	4 MB
Noah_Monika_2017-08-13.jpg	Aug 13, 2017 at 6:59 PM	2.4 MB
Noah_monika_hand_2017-08-13.jpg	Aug 13, 2017 at 6:59 PM	1.3 MB
Pacific_tree_frog.jpg	Apr 24, 2017 at 9:14 PM	2 MB
Pat_icebergs_2016-06-10 13.32.36.jpg	Jun 10, 2016 at 1:32 PM	1.3 MB
Pat_waterfall_2016-06-10 17.46.35.jpg		
plastic_lightening_fractal.jpg		
plastic_lightening.jpg		
Portland_geese_2016-04-24 17.11.06.jpg		
Prescot_veiw_2016-01-09 15.25.41-2.jpg	Jan 9, 2016 at 3:25 PM	1.1 MB
Prescott_thumb_bute_2016-01-10 09.31.44.jpg	Jan 10, 2016 at 9:31 AM	2.7 MB
Prescott_thumb_butee_2016-01-10 09.31.41.jpg	Jan 10, 2016 at 9:31 AM	2.5 MB
Rainbow_home_pdx_2016-09-20 18.15.51.jpg	Sep 20, 2016 at 6:15 PM	3.9 MB
sister_car_wrek_2017-01-01 14.29.13.jpg	Jan 1, 2017 at 2:29 PM	6.3 MB
south_jetty_til_head_from_lighthouse.jpg	Jan 13, 2018 at 4:22 PM	2.4 MB

Problem: The lack of structure makes finding specific items extremely difficult ...

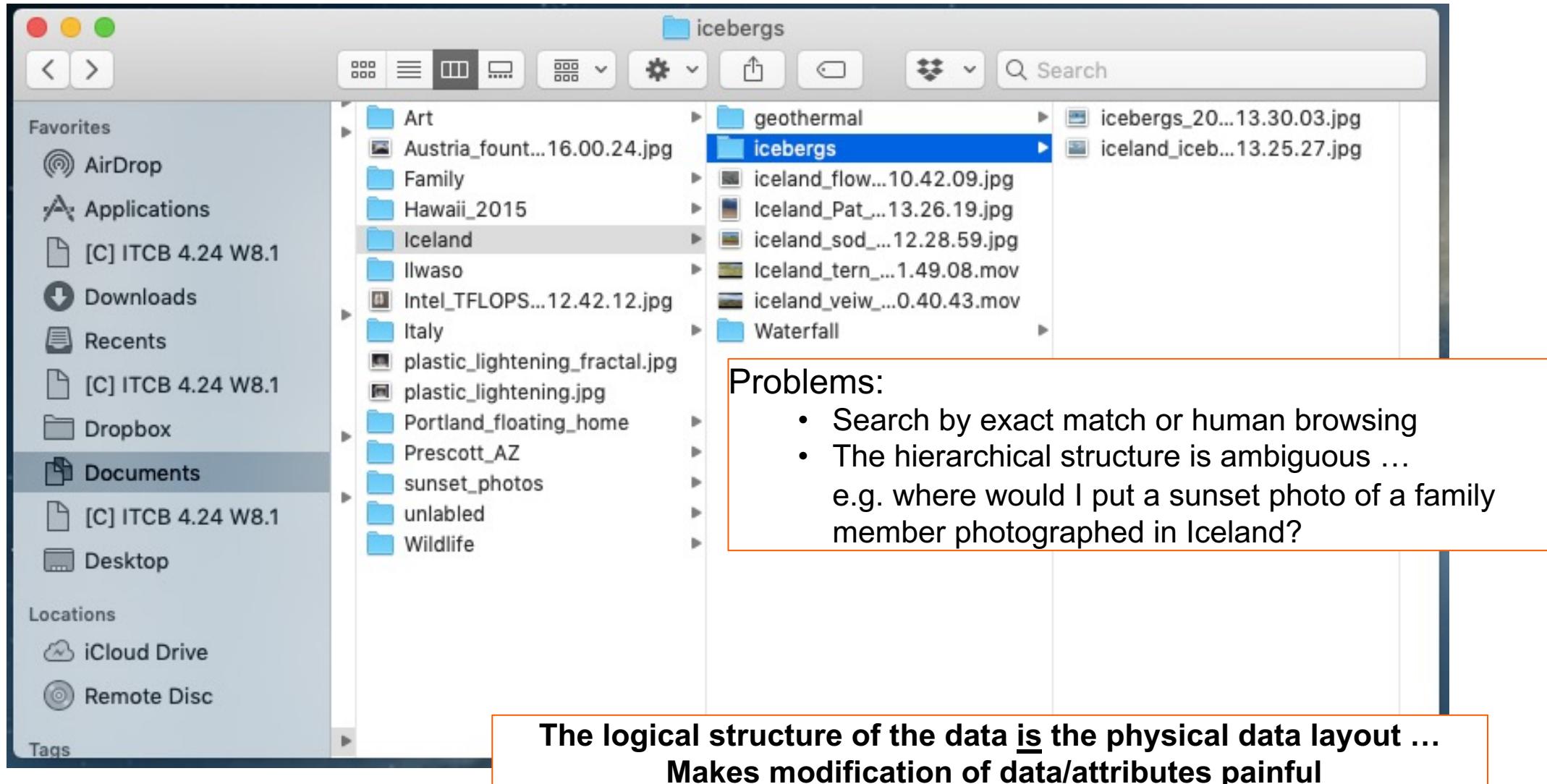
Raw Data from my camera

Hierarchical structure. Data Attributes in folder/file names



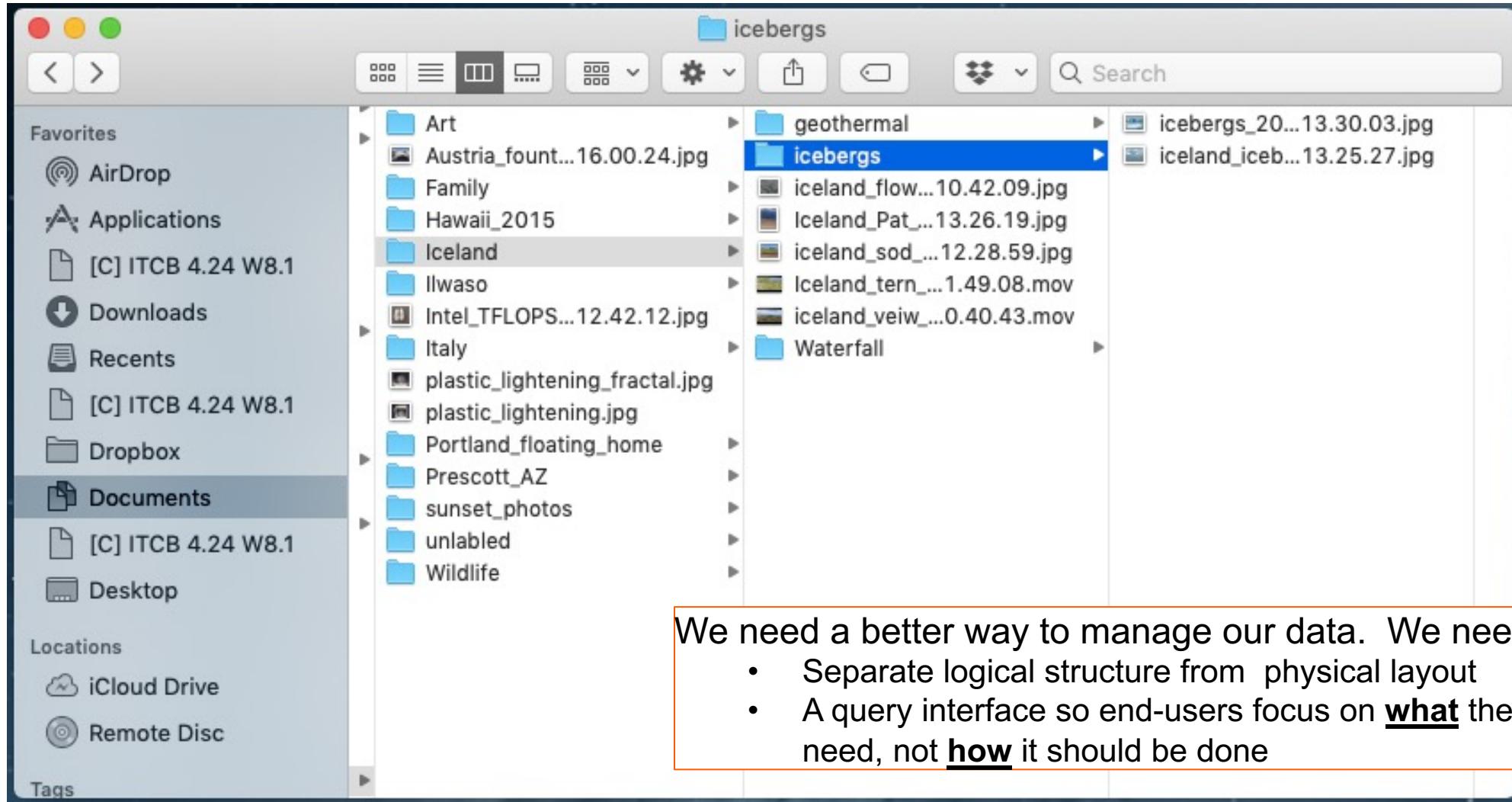
Raw Data from my camera

Hierarchical structure. Data Attributes in folder/file names

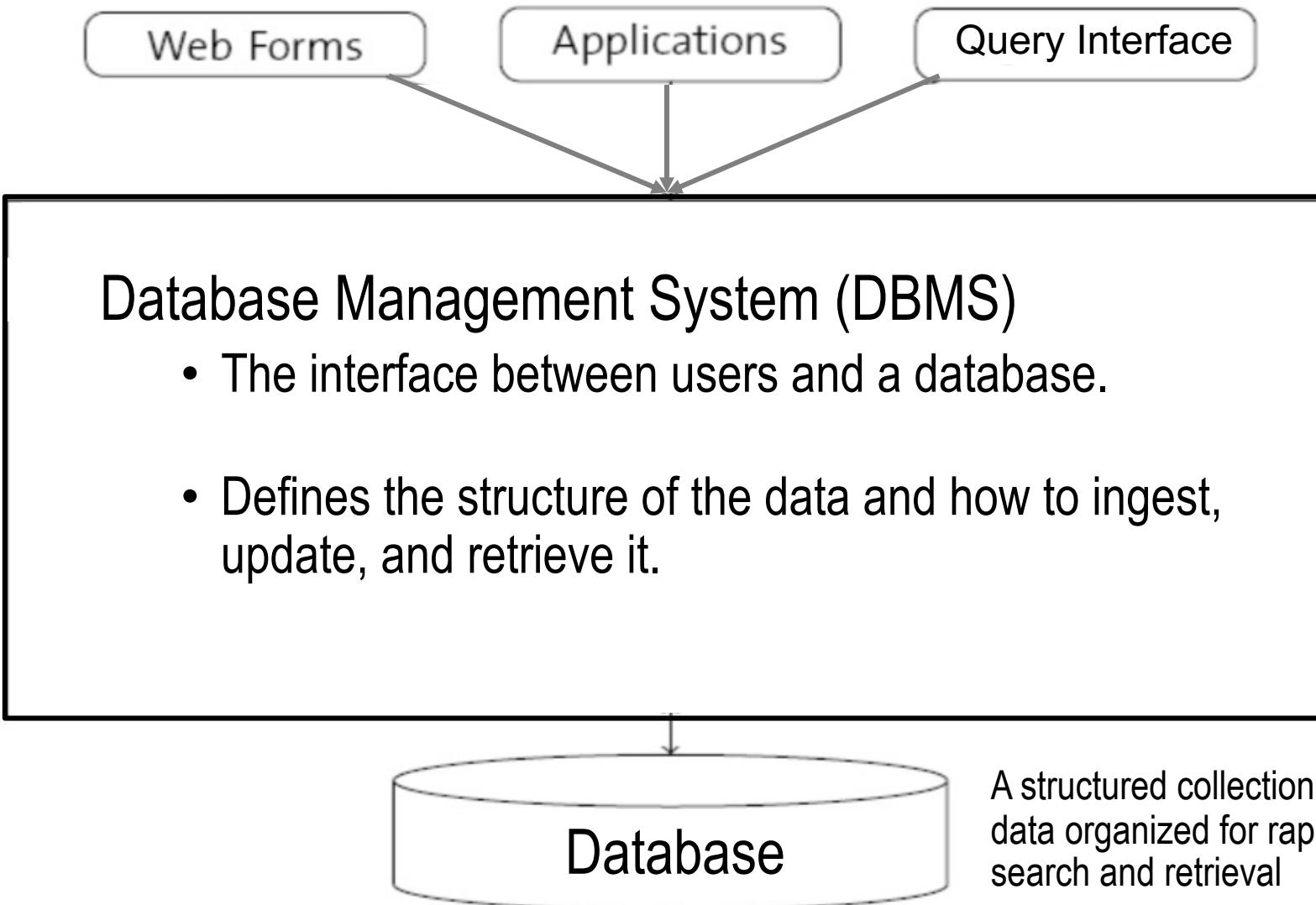


Raw Data from my camera

Hierarchical structure. Data Attributes in folder/file names



Databases and Database Management Systems



The DBMS separates how you work with Data from how the Data is stored.

Why a DBMS is so important: Big Data in the Real World

- Consider patient data in an Intensive Care Unit (e.g. MIMIC II data set*)

- EKG traces
- Blood oxygen
- Blood pressure
- EEG traces



- Demographic
- Caregiver notes
- Medical charts
- Lab test results
- Xray, MRI, etc.

Why a DBMS is so important: Big Data in the Real World

- Consider patient data in an Intensive Care Unit (e.g. MIMIC II data set*)



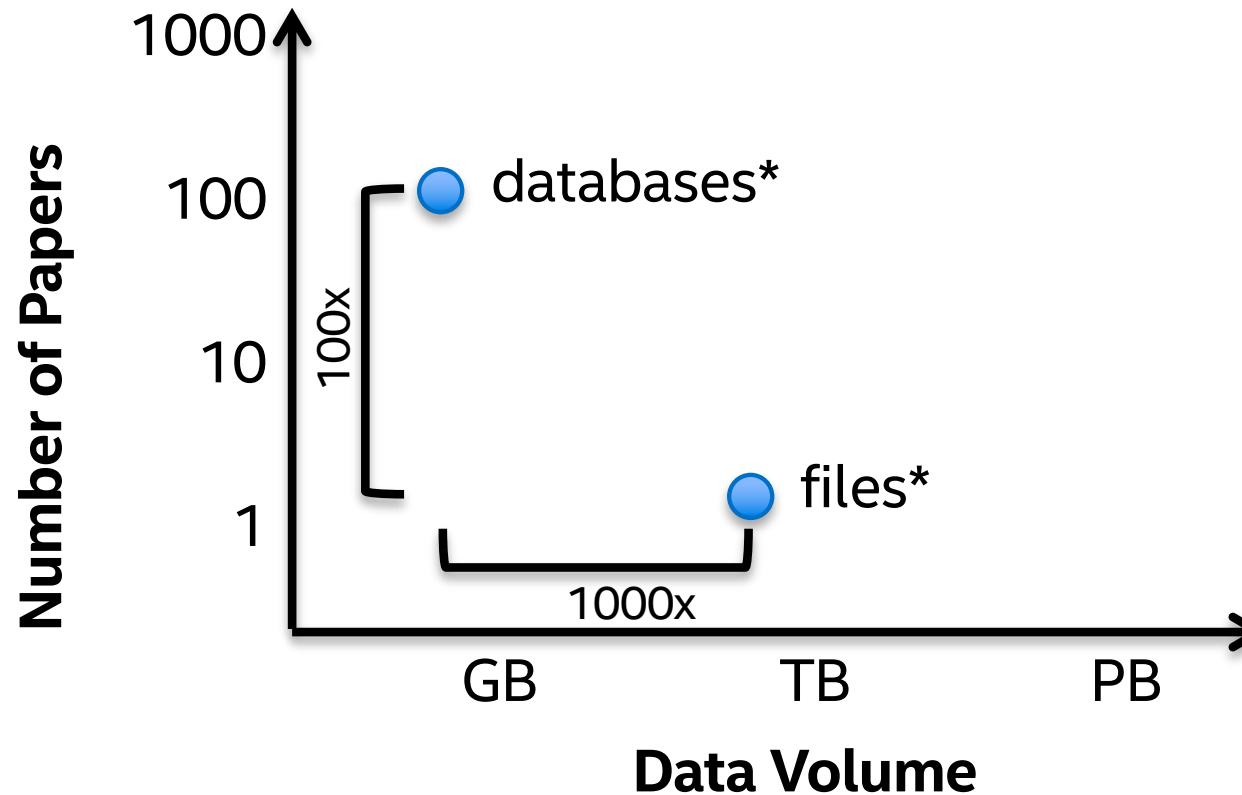
Arrays	• EKG traces	tables
Time Series	• Blood oxygen	documents
Time Series	• Blood pressure	tables
Arrays	• EEG traces	tables
		• Lab test results
		• Xray, MRI, etc.
		#images

Time series and tabular data: Structured data in a Database.

Other data? Flat files

Why a DBMS is so important: Big Data in the Real World

Analysis of published MIMICII papers, 2015



Storing data as flat files
is roughly equivalent to
deleting the data

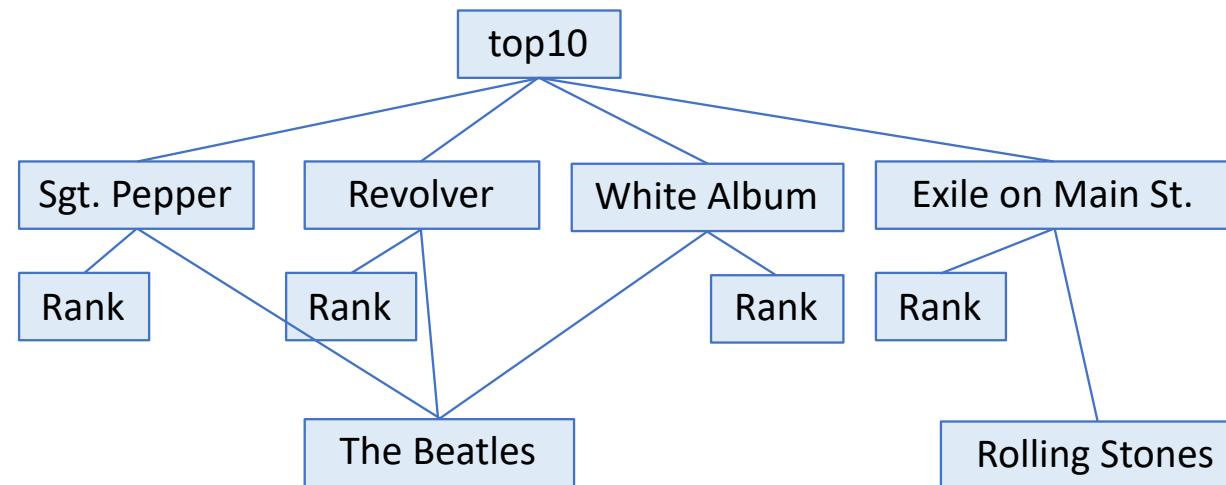
We must bring the power of databases to all data

Outline

- Motivation: Why everyone needs a database management system?
- • Database Technology: from ancient history to today
- Data Curation in the sciences
- My quest: One Algebra to rule them all

DBMS ancient history (1960s and 70s)

- In the 60's
 - Data in **Flat files** tied to a program (usually written in COBOL)
- Late 60s and throughout the 70s Hierarchical and Network models



- These approaches were expensive, difficult to adapt to changing data, and lacked a standard query interface for users.

The Relational Model of Databases

- In 1970 Edgar Codd (IBM) published one of the most important papers in the history of computer science.
- It defined a formal algebra* for building databases ... the **relational model**.
 - Object: A relation.
 - A set of tuples that share a set of attributes.
 - The set of attributes is defined by a schema
 - A relation is typically represented as a table.
 - A set of operators that act on relations. This set includes:
 - Select σ
 - Join \bowtie
 - Rename ρ
 - Project π

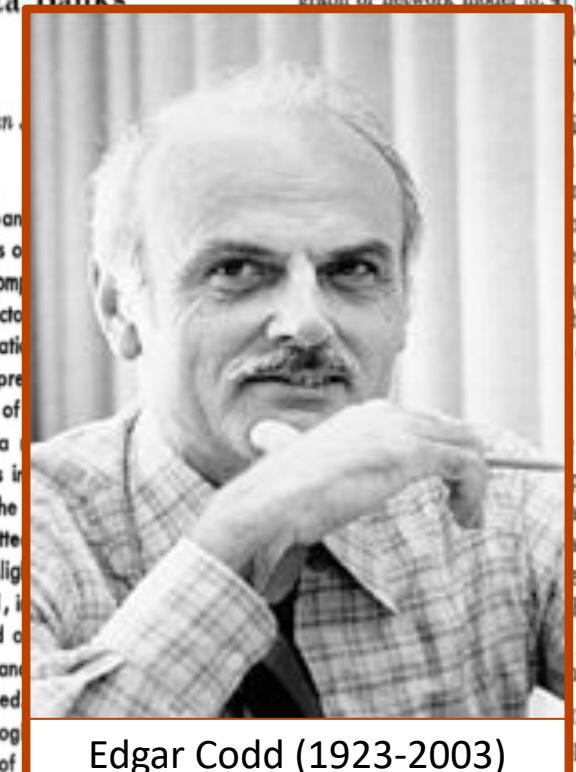
* Note: An “algebra” is a set of objects, operators that act on those objects, and rules for how those operators interact with each other

Information Retrieval

P. BAXENDALE, Editor

A Relational Model of Data for Large Shared Data Banks

E. F. CODD
IBM Research Laboratory, San



The relational view (or model) of data described in Section 1 appears to be superior in several respects to the graph or network model [3, 4] presently in vogue for non-a means of describing data—that is, without superim-for machine representationides a basis for a high level maximal independence be-and machine representa-on the other.

The relational view is that it g derivability, redundancy, nese are discussed in Section other hand, has spawned a least of which is mistaking for the derivation of rela-on the “connection trap”). permits a clearer evalua-tions of present formatted ative merits (from a logical esentations of data within a his clearer perspective are paper. Implementations of al model are not discussed.

IN PRESENT SYSTEMS option tables in recently de-represents a major advance

Edgar Codd (1923-2003)

KEY WORDS AND PHRASES: data bank, data base, data structure, data organization, hierarchies of data, networks of data, relations, derivability, redundancy, consistency, composition, join, retrieval language, predicate calculus, security, data integrity

CR CATEGORIES: 3.70, 3.73, 3.75, 4.90, 4.92, 4.93

Communications of the ACM, vol 13, no. 6 p. 337, 1970

Database Queries and the Relational Model

- Users interact with the relational database by issuing queries.
- Codd proposed elegant and mathematically rich procedural queries:

$$\pi_{e.name} \left(\sigma_{e.salary > m.salary} \left(\rho_e(\text{employee}) \bowtie_{e.manager = m.name} \rho_m(\text{employee}) \right) \right)$$

- When applied to this relation:

Name	Salary	Manager
Smith	45,000	Harker
Jones	40,000	Smith
Baker	50,000	Smith
Nelson	55,000	Baker

- The output is **Baker** and **Nelson** employees who earn more than their managers

Queries and the Structured Query Language

- Codd's notation was too obtuse for the general user

$$\pi_{e.name} \left(\sigma_{e.salary > m.salary} \left(\rho_e(\text{employee}) \bowtie_{e.manager = m.name} \rho_m(\text{employee}) \right) \right)$$

- In 1974, Codd's colleagues at IBM (Ray Boyce and Don Chamberlin) created a Query Language for Codd's relational Algebra called the Structured Query Language (SQL ... Pronounced as Sequel)
- The SQL Query equivalent to Codd's notation above reads:

```
select e.name  
from employee e, employee m  
where e.manager = m.name and e.salary > m.salary
```

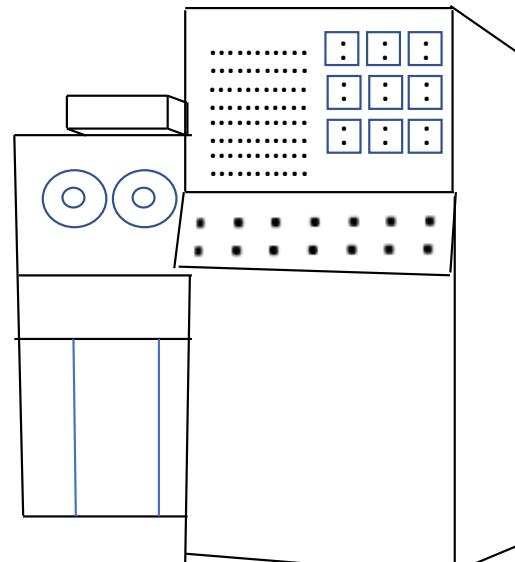
- Codd's notation is procedural. SQL is Declarative ... That difference is an important reason for why SQL is the most successful DSL of all time.

DBMS History: A platform perspective

60's to 70's

Flat-files to network models

Custom + emerging vendors



Mainframe

The early days of database technology ... centrally maintained, “big-iron” mainframe computers.

RDBMS: A deeper dive into Relations

- Example: the following relations come from the Rolling Stones top 500 albums database

Rolling Stone top 10 Albums (**top10**)

Number	Album
1	Sgt. Pepper's Lonely Hearts Club Band
2	Pet Sounds
3	Revolver
4	Highway 61 Revisited
5	Rubber Soul
6	What's Going On
7	Exile on Main St.
8	London Calling
9	Blonde on Blonde
10	The Beatles ("The White Album")

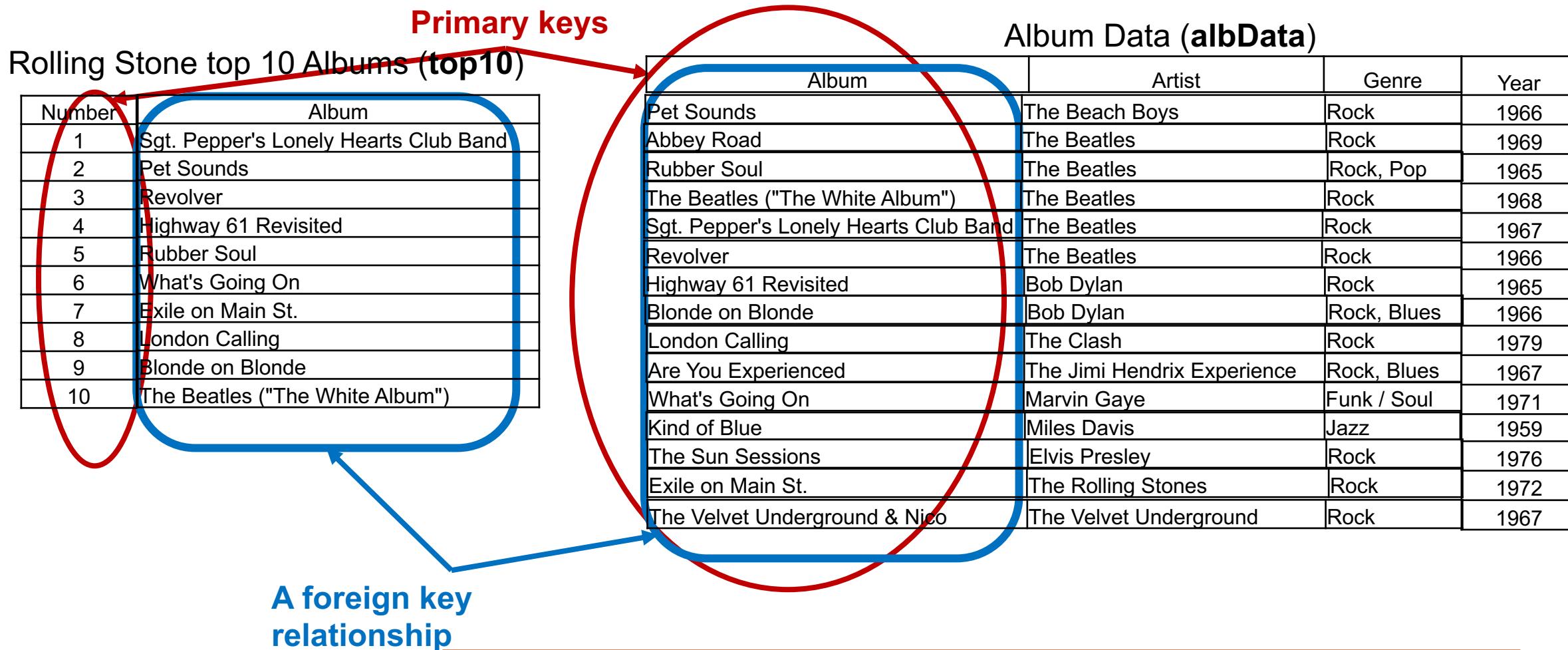
Each Row is a **record**
Each column is an **attribute**

Album Data (**albData**)

Album	Artist	Genre	Year
Pet Sounds	The Beach Boys	Rock	1966
Abbey Road	The Beatles	Rock	1969
Rubber Soul	The Beatles	Rock, Pop	1965
The Beatles ("The White Album")	The Beatles	Rock	1968
Sgt. Pepper's Lonely Hearts Club Band	The Beatles	Rock	1967
Revolver	The Beatles	Rock	1966
Highway 61 Revisited	Bob Dylan	Rock	1965
Blonde on Blonde	Bob Dylan	Rock, Blues	1966
London Calling	The Clash	Rock	1979
Are You Experienced	The Jimi Hendrix Experience	Rock, Blues	1967
What's Going On	Marvin Gaye	Funk / Soul	1971
Kind of Blue	Miles Davis	Jazz	1959
The Sun Sessions	Elvis Presley	Rock	1976
Exile on Main St.	The Rolling Stones	Rock	1972
The Velvet Underground & Nico	The Velvet Underground	Rock	1967

RDBMS: The concept of a Schema

- Schema: defines how the data is organized ... (1) the column labels/types, (2) the primary key (uniquely identifies a record), and (3) the foreign keys (columns that map onto primary keys of other relations)



A complex database might need hundreds of tables.
Defining consistent schema for all those tables is extremely difficult

Relational Database example with SQL

Input two tables

Number	Album
1	Sgt. Pepper's Lonely Hearts Club Band
2	Pet Sounds
3	Revolver
4	Highway 61 Revisited
5	Rubber Soul
6	The Beatles ("The White Album")
7	Sgt. Pepper's Lonely Hearts Club Band
8	Revolver
9	Highway 61 Revisited
10	Blonde on Blonde
	London Calling
	Exile on Main St.
	Are You Experienced
	Kind of Blue
	The Sun Sessions
	Blonde on Blonde
	London Calling
	Exile on Main St.
	The Velvet Underground & Nico
	The Velvet Underground



```
SELECT top10.Number, top10.Album, albData.Artist, albData.year  
From top10  
INNER JOIN albData  
ON top10.Album = albData.Album;
```

(top10 and albData)

Rolling Stone top 10 Albums (**top10**)

Number	Album
1	Sgt. Pepper's Lonely Hearts Club Band
2	Pet Sounds
3	Revolver
4	Highway 61 Revisited
5	Rubber Soul
6	What's Going On
7	Exile on Main St.
8	London Calling
9	Blonde on Blonde
10	The Beatles ("The White Album")

Album Data (**albData**)

Album	Artist	Genre	Year
Pet Sounds	The Beach Boys	Rock	1966
Abbey Road	The Beatles	Rock	1969
Rubber Soul	The Beatles	Rock, Pop	1965
The Beatles ("The White Album")	The Beatles	Rock	1968
Sgt. Pepper's Lonely Hearts Club Band	The Beatles	Rock	1967
Revolver	The Beatles	Rock	1966
Highway 61 Revisited	Bob Dylan	Rock	1965
Blonde on Blonde	Bob Dylan	Rock, Blues	1966
London Calling	The Clash	Rock	1979
Are You Experienced	The Jimi Hendrix Experience	Rock, Blues	1967
What's Going On	Marvin Gaye	Funk / Soul	1971
Kind of Blue	Miles Davis	Jazz	1959
The Sun Sessions	Elvis Presley	Rock	1976
Exile on Main St.	The Rolling Stones	Rock	1972
The Velvet Underground & Nico	The Velvet Underground	Rock	1967

Relational Database example with SQL

Input two tables

Album			
Number	Album	Artist	Year
1	Sgt. Pepper's Lonely Hearts Club Band	The Beatles	1967
2	Pet Sounds	The Beach Boys	1966
3	Revolver	The Beatles	1966
4	Highway 61 Revisited	Bob Dylan	1965
5	Rubber Soul	The Beatles	1965
6	What's Going On	Marvin Gaye	1971
7	Exile on Main St.	The Rolling Stones	1972
8	London Calling	The Clash	1979
9	Blonde on Blonde	Bob Dylan	1966
10	The Beatles ("The White Album")	The Beatles	1968

```
SELECT top10.Number, top10.Album, albData.Artist, albData.year
From top10
INNER JOIN albData
ON top10.Album = albData.Album;
```

(top10 and albData)

Output
a new
table

Number	Album	Artist	Year
1	Sgt. Pepper's Lonely Hearts Club Band	The Beatles	1967
2	Pet Sounds	The Beach Boys	1966
3	Revolver	The Beatles	1966
4	Highway 61 Revisited	Bob Dylan	1965
5	Rubber Soul	The Beatles	1965
6	What's Going On	Marvin Gaye	1971
7	Exile on Main St.	The Rolling Stones	1972
8	London Calling	The Clash	1979
9	Blonde on Blonde	Bob Dylan	1966
10	The Beatles ("The White Album")	The Beatles	1968

Rolling Stone top 10 Albums (top10)

Number	Album
1	Sgt. Pepper's Lonely Hearts Club Band
2	Pet Sounds
3	Revolver
4	Highway 61 Revisited
5	Rubber Soul
6	What's Going On
7	Exile on Main St.
8	London Calling
9	Blonde on Blonde
10	The Beatles ("The White Album")

Album Data (albData)

Album	Artist	Genre	Year
Pet Sounds	The Beach Boys	Rock	1966
Abbey Road	The Beatles	Rock	1969
Rubber Soul	The Beatles	Rock, Pop	1965
The Beatles ("The White Album")	The Beatles	Rock	1968
		Rock	1967
		Rock	1966
		Rock	1965
		Rock, Blues	1966
		Rock	1979
		Rock, Blues	1967
		Funk / Soul	1971
		Jazz	1959
		Rock	1976
		Rock	1972
		Rock	1967

Number	Album	Artist	Year
1	Sgt. Pepper's Lonely Hearts Club Band	The Beatles	1967
2	Pet Sounds	The Beach Boys	1966
3	Revolver	The Beatles	1966
4	Highway 61 Revisited	Bob Dylan	1965
5	Rubber Soul	The Beatles	1965
6	What's Going On	Marvin Gaye	1971
7	Exile on Main St.	The Rolling Stones	1972
8	London Calling	The Clash	1979
9	Blonde on Blonde	Bob Dylan	1966
10	The Beatles ("The White Album")	The Beatles	1968

Online Transaction Processing (OLTP)

- When you use an ATM (bancomat), you expect that your transaction will be:

Atomic: the transaction either completes or it doesn't happen.

Consistent: The data is always in a valid state.



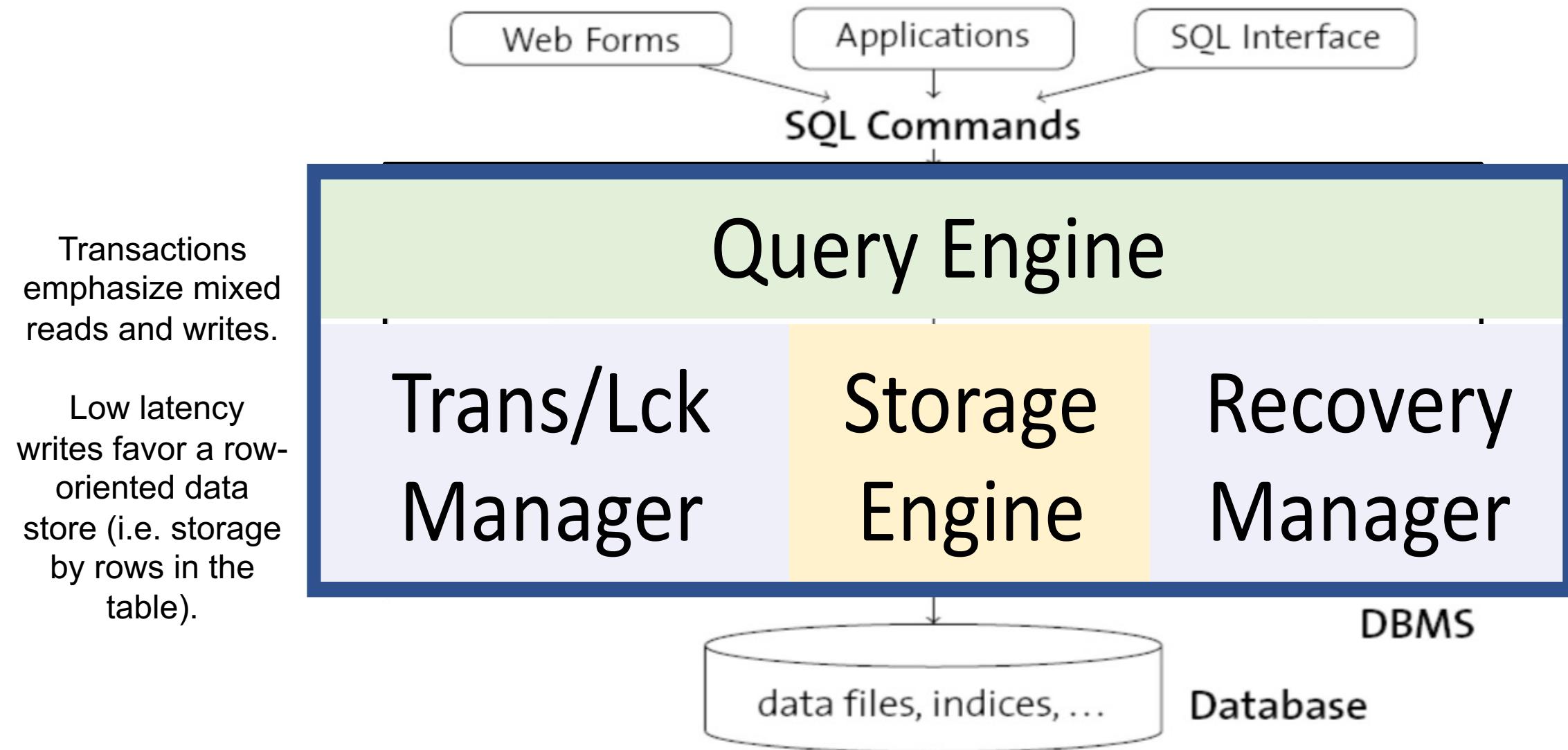
Isolated: Transactions cannot interfere with each other.

Durable: Once committed, the effect of a transaction is always there ... regardless of system crashes.

A Database Management Systems that supports these conditions is said to be **ACID compliant**

ACID compliance is essential for a multiuser, online transaction Processing System (e.g. an ATM)

Structure of a Modern Relational DBMS

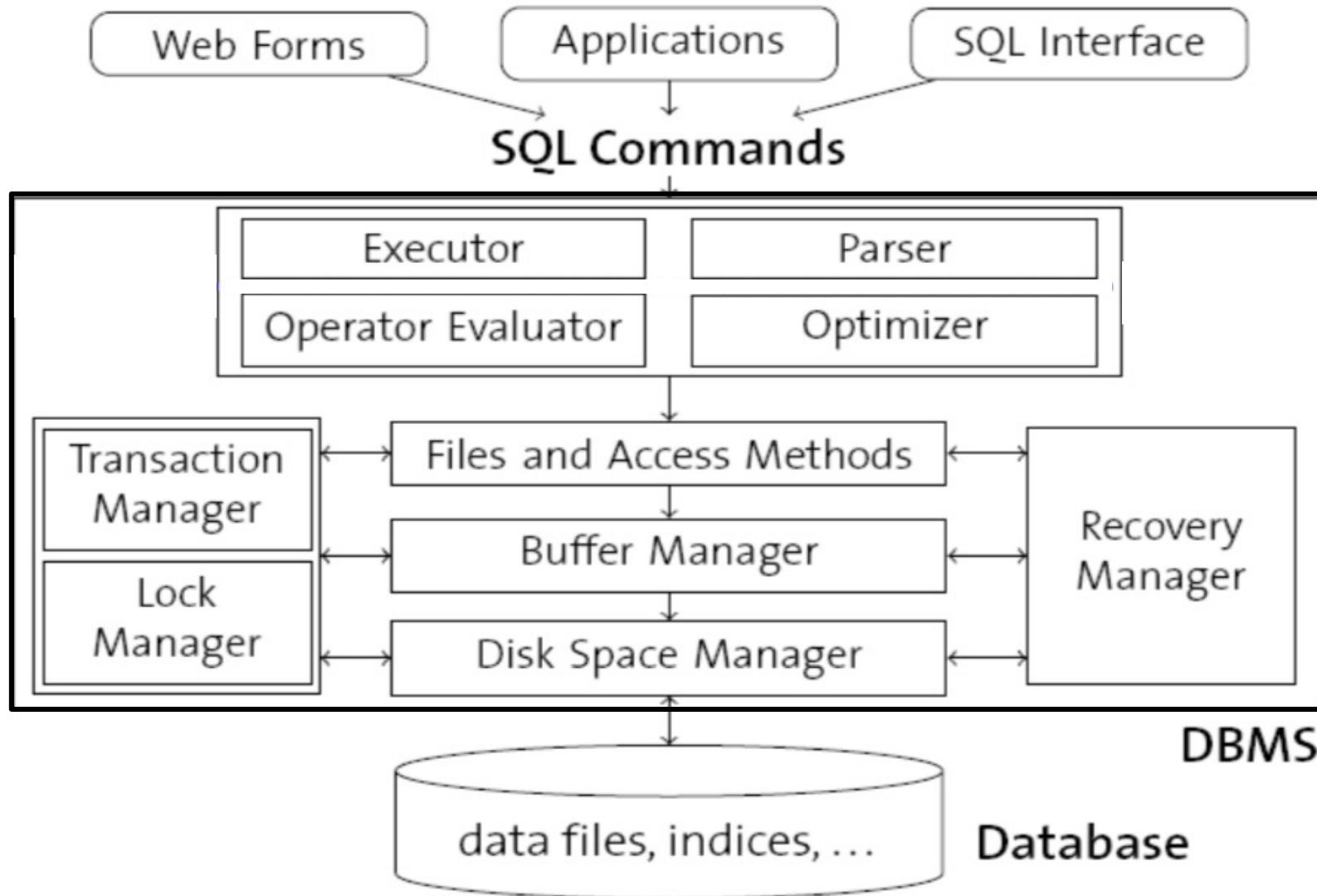


Source: Nesime Tatbul, information systems, ETH, 2012
And Ramakrishnam/Gehrke: "Database managementSystems", McGraw-Hill 2003

Structure of a Modern Relational DBMS

Transactions emphasize mixed reads and writes.

Low latency writes favor a row-oriented data store (i.e. storage by rows in the table).



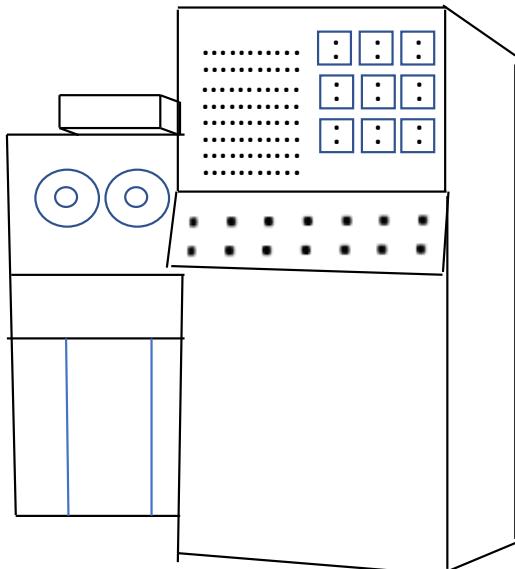
Trans: Transaction

Lck: Lock

Source: Nesime Tatbul, information systems, ETH, 2012
And Ramakrishnam/Gehrke: "Database management Systems", McGraw-Hill 2012

DBMS History: A platform perspective

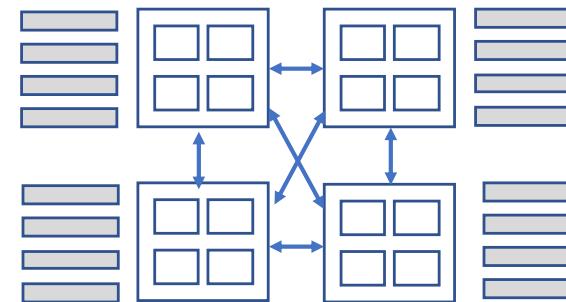
60's to 70's
Flat-files to network models
Custom + emerging vendors



Mainframe

80's to 90's

Relational models
RDBMS vendors (Oracle and friends)



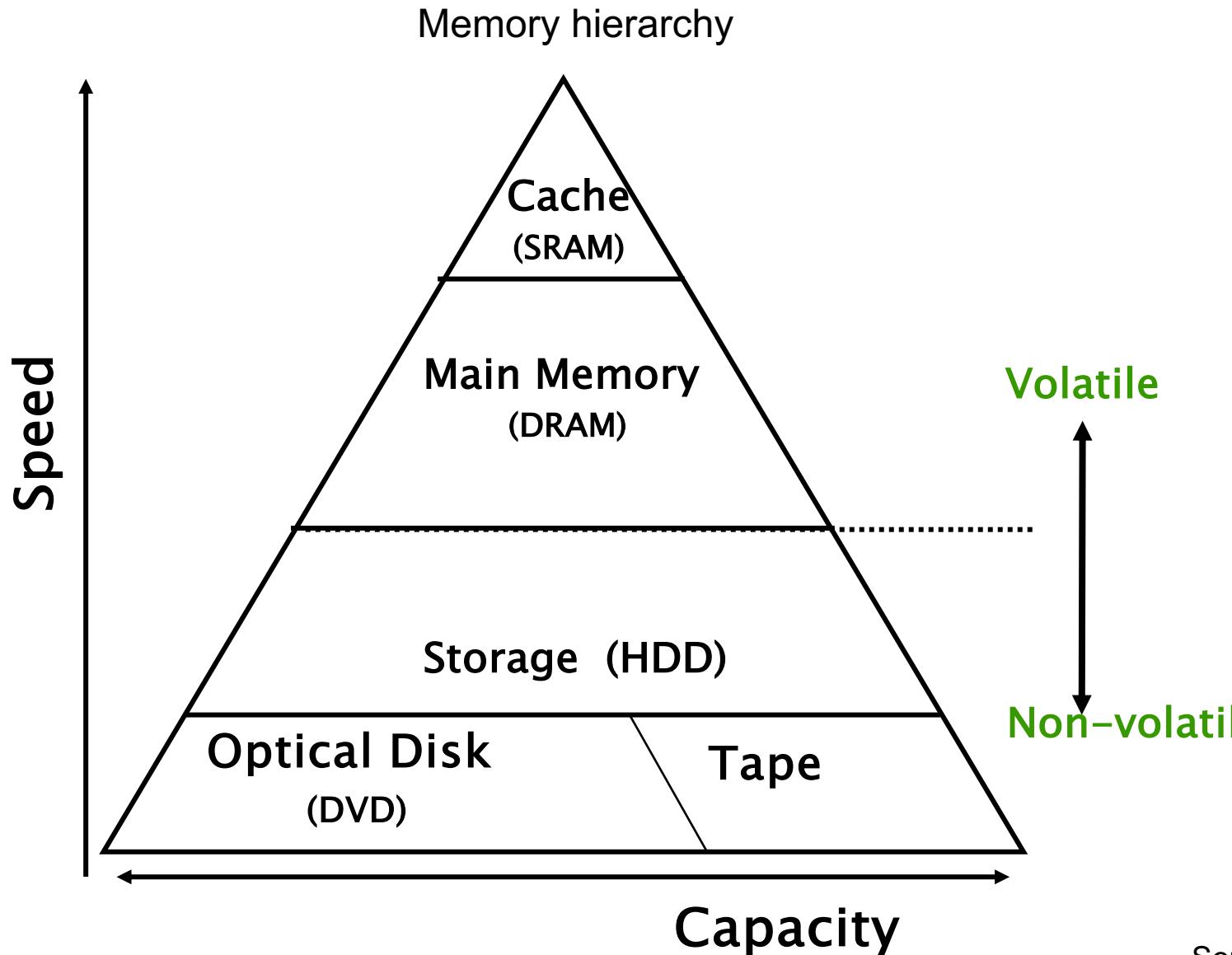
Shared Address-space Multiprocessor Server

Databases to support ACID compliance while handling larger problems depended on multithreading and large address spaces.

Scale-up: increasing processors in a single physical address space

RDBS: Storing the data

- Traditional databases live “on disk”.
- Consider memory latencies:



CPU: <nanosecond
Memory: ~80 ns
3D Xpoint: ~300ns
Optane SSD: ~10 usec
NAND SSD: ~80 usec
Disk: ~6 milliseconds

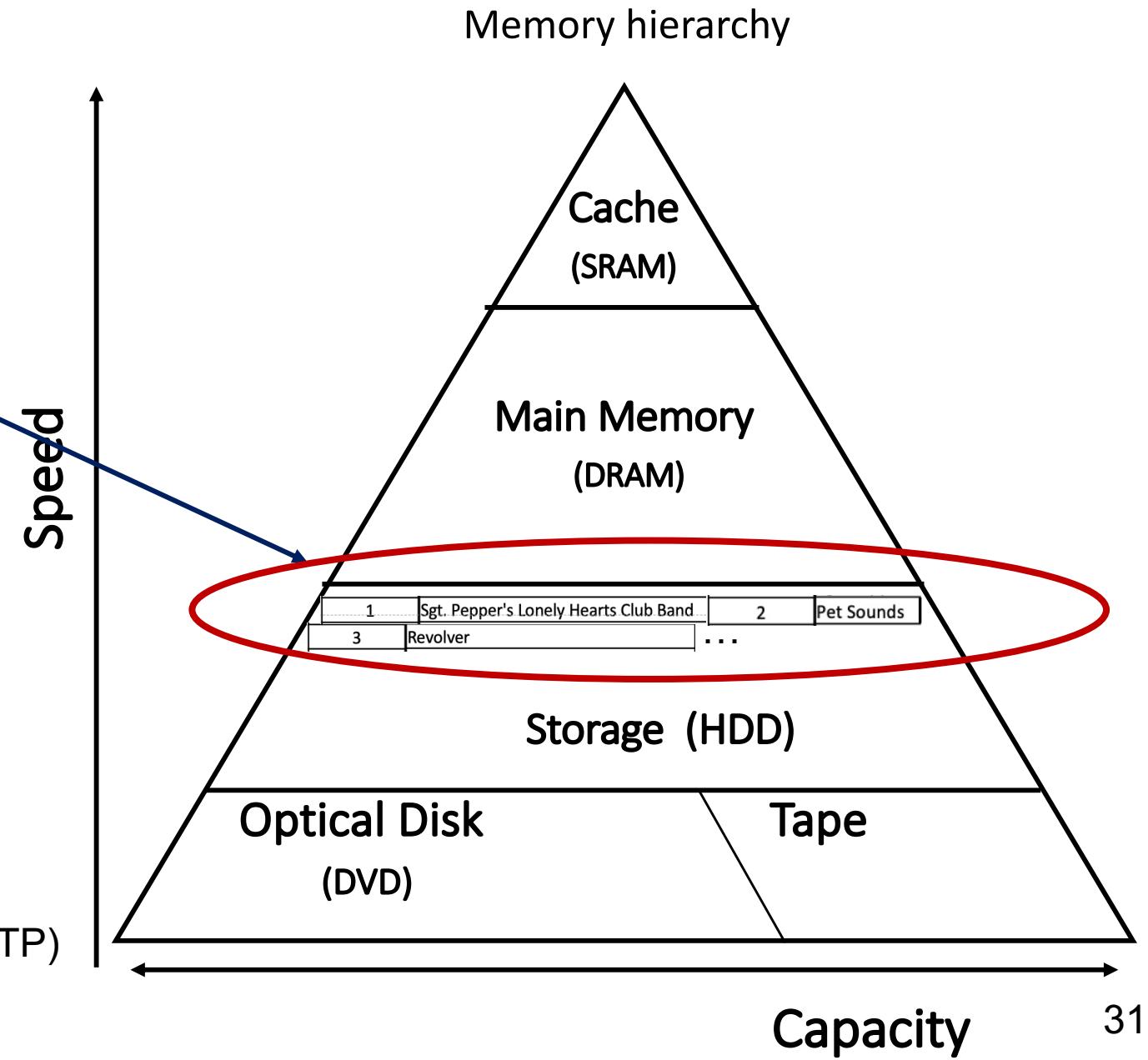
- Disk access is slow compared to CPU speeds.
- Disk access is the primary performance bottleneck in traditional RDBMS
- How should we layout data on disk?

RDBMS: Storing Relations ... Row Store

- Organize by row ... one record after another

Table: Rolling Stone top 10 Albums

Number	Album
1	Sgt. Pepper's Lonely Hearts Club Band
2	Pet Sounds
3	Revolver
4	Highway 61 Revisited
5	Rubber Soul
6	What's Going On
7	Exile on Main St.
8	London Calling
9	Blonde on Blonde
10	The Beatles ("The White Album")



A **row store** is optimized for writing records into a database.

Preferred for **Online Transaction Processing (OLTP)**

RDBMS: Storing Relations ... Column Store

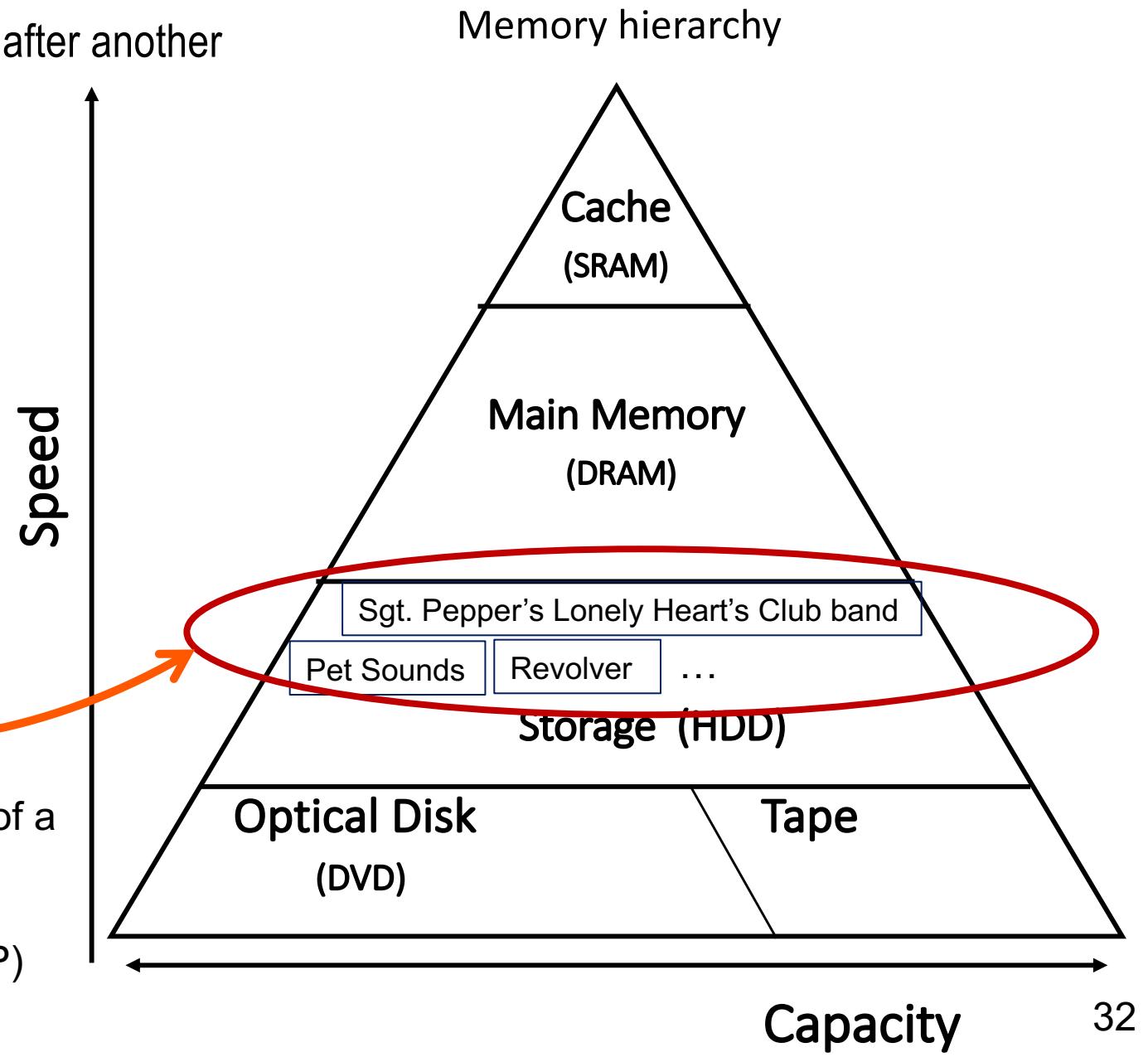
- Organize by columns (attributes) ... one attribute after another

Table: Rolling Stone top 10 Albums

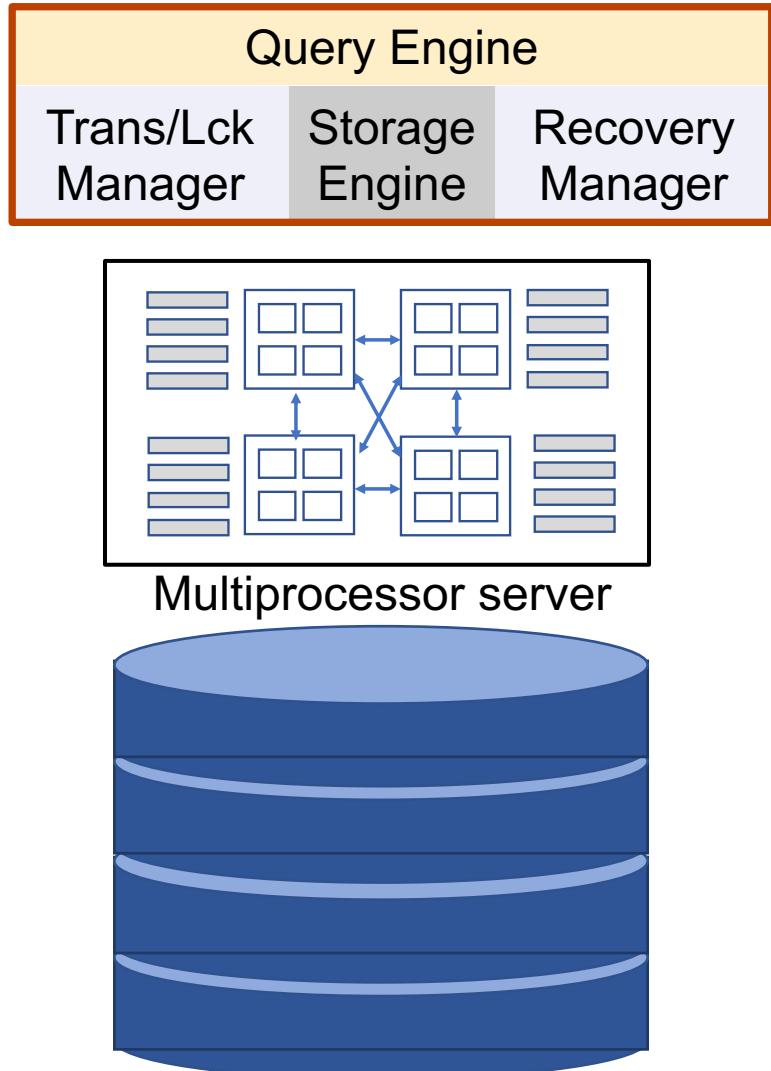
Number	Album
1	Sgt. Pepper's Lonely Hearts Club Band
2	Pet Sounds
3	Revolver
4	Highway 61 Revisited
5	Rubber Soul
6	What's Going On
7	Exile on Main St.
8	London Calling
9	Blonde on Blonde
10	The Beatles ("The White Album")

A *column store* is optimized for reading a set of attributes as is often done in computing properties of a set of records.

Preferred for **Online Analytical Processing (OLAP)**



Databases in the Internet Age

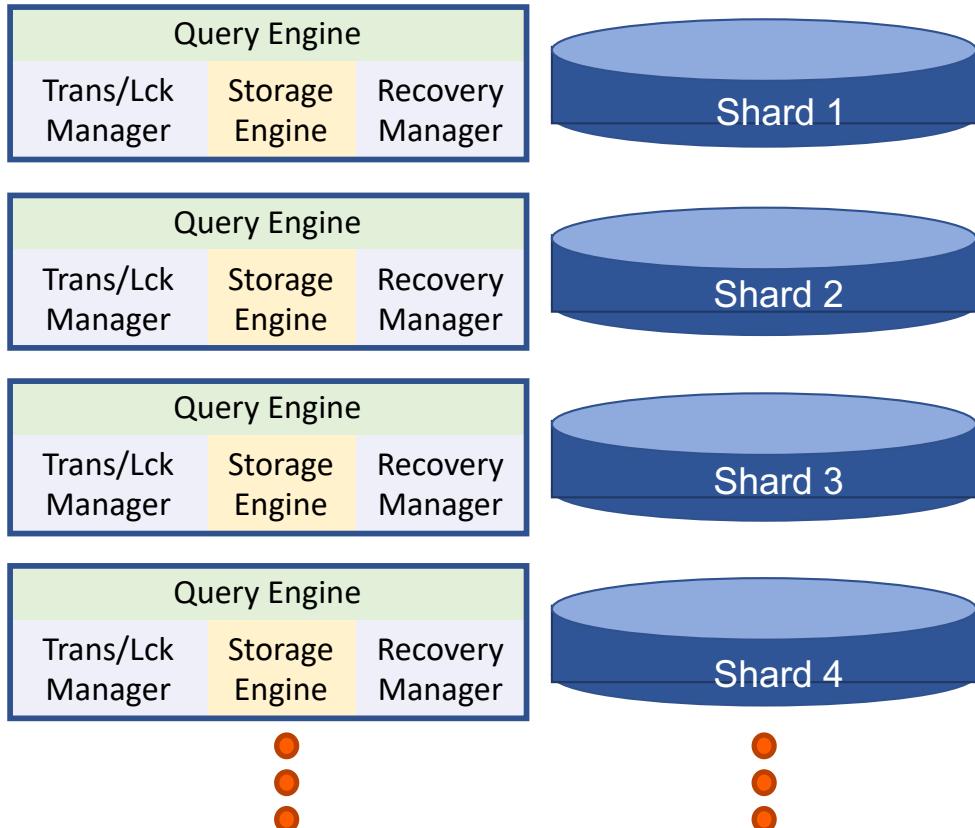


RDBMS Challenges

- Must define Schemas upfront
- Scale-up → expensive multiprocessor servers
- Scalability limited: You can only connect so-many processors to a single shared memory.
- ACID compliance is REALLY HARD to scale.

Internet-scale problems → Extreme scalability and unstructured data → moving beyond RDMS/Scale-up

Databases in the Internet Age: the birth of NoSQL

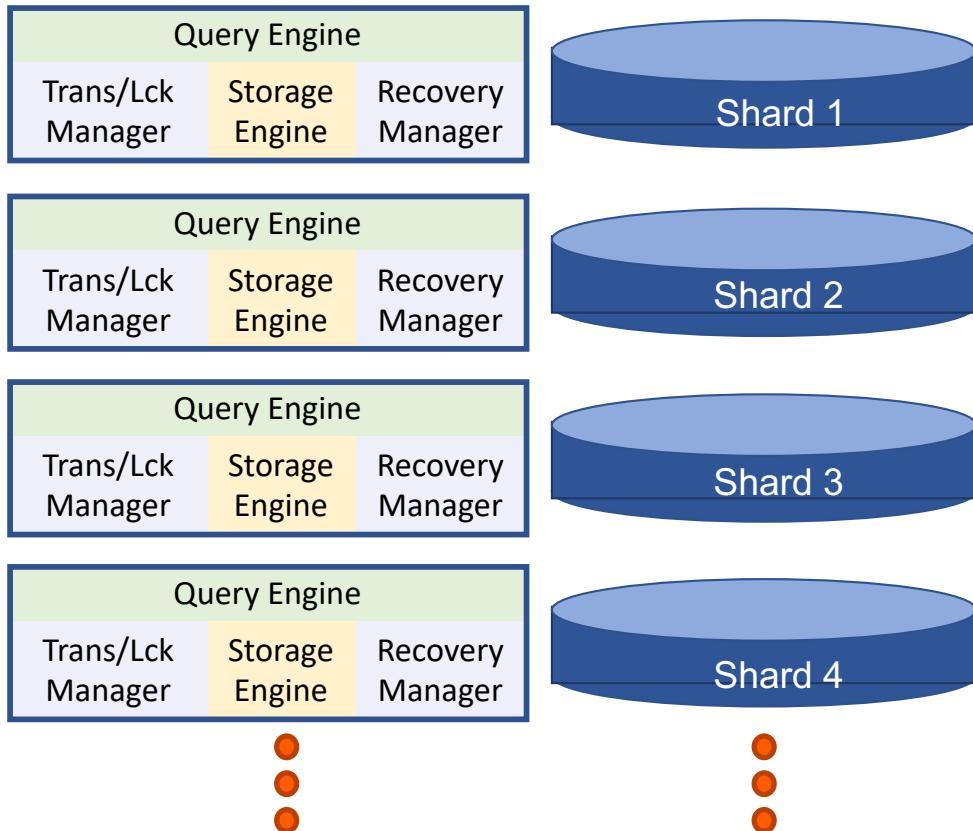


NoSQL

- Break up the data into chunks called Shards
- Distribute shards among multiple DBMS instances
- Scale-out solutions: use “shared nothing” clusters
- Eventual consistency ... violates ACID, but more scalable
- Uses a simple schema free-algebra .. Key-value store

Scale-out: Lots of servers connected by a network (no shared address space). Scalability limited by \$\$\$ and electric bill.

Databases in the Internet Age: the birth of NoSQL



Example of a NoSQL key-value store

Data as a set of tuples: <key,value>

<“Pet Sounds”, Rank=2>

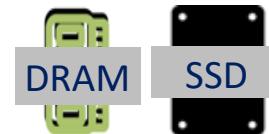
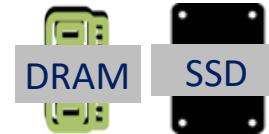
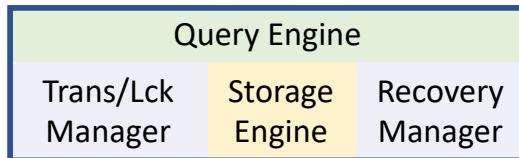
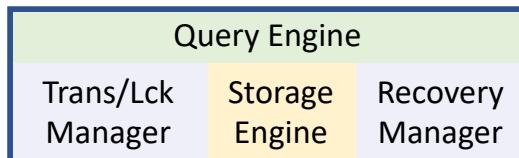
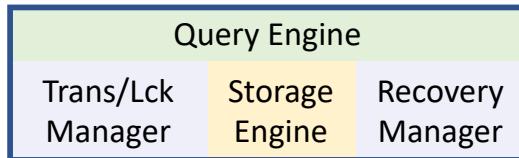
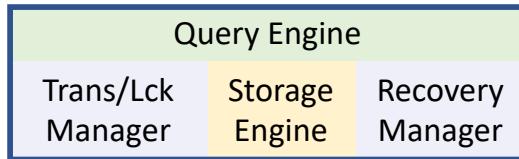
<“Rubber Soul”, Artist=“The Beatles”>

<“Rubber Soul”, Rank=5>

<“Blonde on Blonde”, Artist=“Bob Dylan”>

Scale-out: Lots of servers connected by a network (no shared address space).
Scalability limited by \$\$\$ and electric bill.

Databases in the Internet Age: the birth of NoSQL



- NoSQL originally used disk-based solutions such as Hadoop.
- Later moved to **in-memory storage engines** such as Spark with SSDs restricted to backing up persistent data.

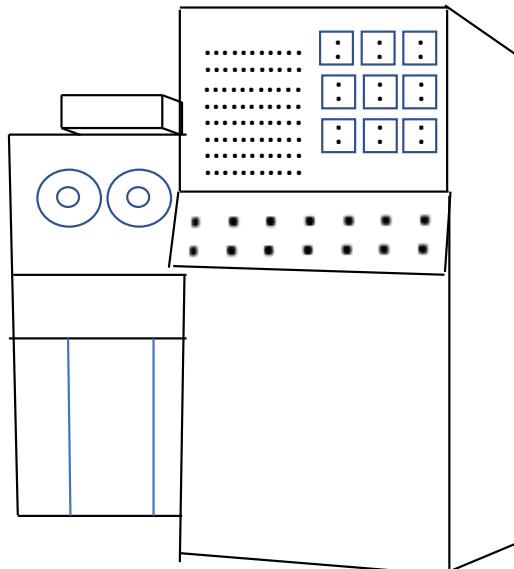
NoSQL - arbitrary scalability, reliability by replication, easy ingestion of new data, flexibility as data changes

DBMS History: A platform perspective

60's to 70's

Flat-files to network models

Custom + emerging vendors

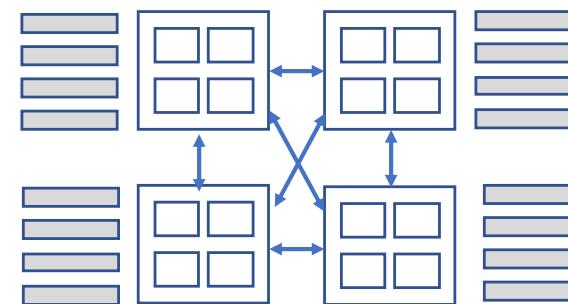


Mainframe

80's to 90's

Relational models

RDBMS vendors (Oracle and friends)

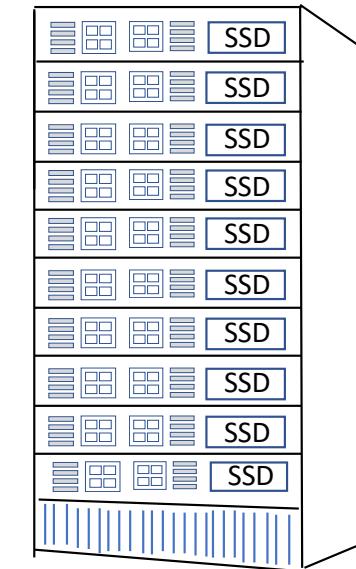


SMP Server

00's to 10's

NoSQL

Legacy RDBMS + a swarm of NoSQL vendors



Hyperconverged Infrastructure (HCI)
dual-processor Servers

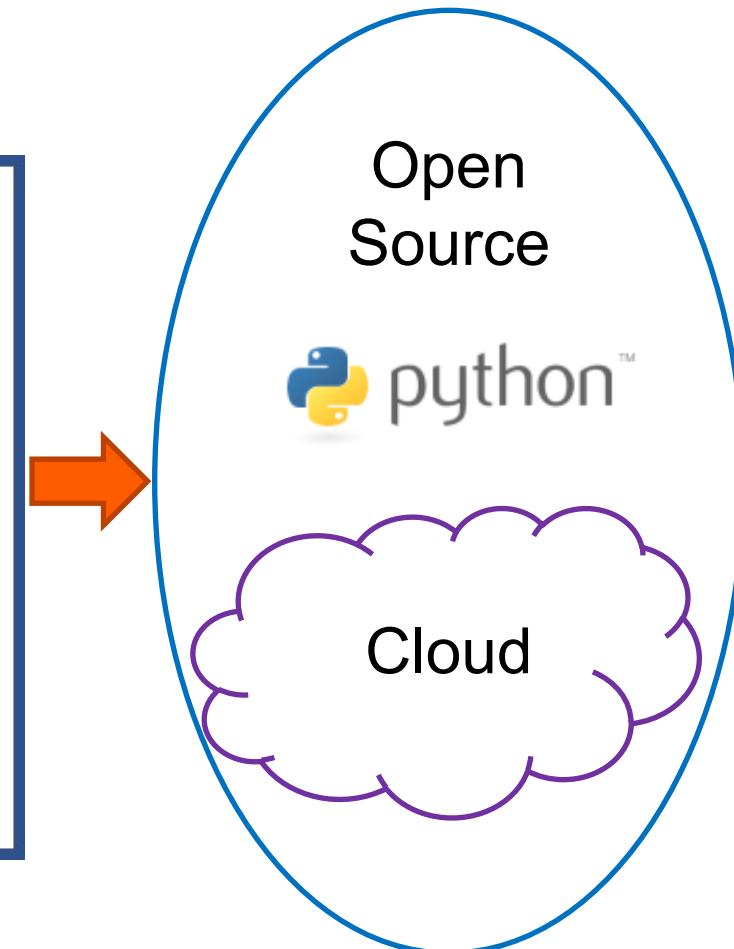
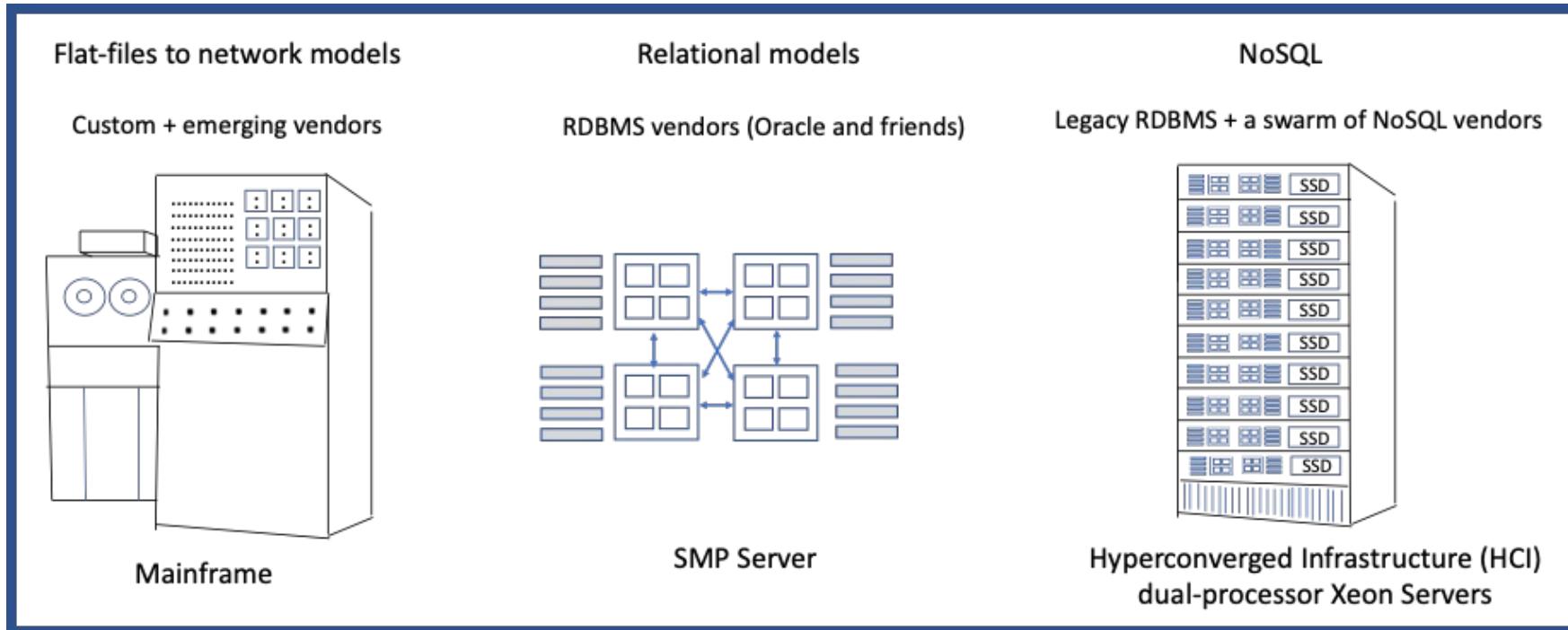
Four classes of systems in use today

DBMS Classes	When they emerged	Feature-set/notes
RDBMS/SQL	1980s	Online Transaction Processing (OLTP) with full ACID guarantees for transactions. Universal declarative query language (SQL)
OLAP data warehouses	2000s	Extract Transform Load (ETL) data from multiple sources and gather into a “single” system for Online Analytic Processing (OLAP). The birth of column stores data for high performance analytics
NoSQL (later ... not only SQL)	Mid-2000s	DBMS for the internet age. Pioneered (1) eventual consistency and relaxed ACID, (2) shard data for distributed systems, (3) in memory storage, and (4) replication
NewSQL	2010s	All the benefits of NoSQL but ACID is brought back with lock-free consistency models for OLTP applications. Pioneered Hybrid Transaction-Analytics Processing (HTAP)

Source: What's Really New with NewSQL, Pavlo and Aslett, SIGMOD Record, June 2016

ACID: **A**tomicity, **C**onsistency, **I**ndependence, **D**ependability

DBMS Platforms: Where are we going?



DBMS technology is moving into the cloud using FaaS (Function as a Service) to make the physical computers “invisible”.

Outline

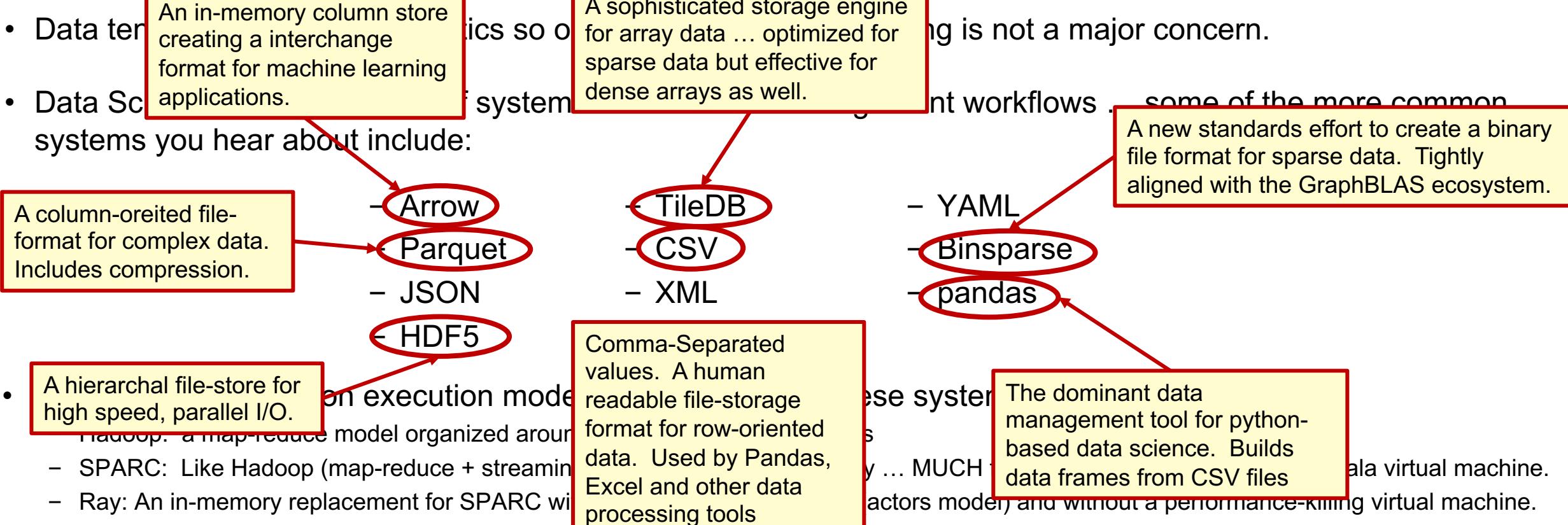
- Motivation: Why everyone needs a database management system?
- Database Technology: from ancient history to today
- • Data Curation in the sciences
- My quest: One Algebra to rule them all

How Data Scientists deal with data today ...

- A DBMS separates how data is stored from how it is processed. If data is arbitrarily complex, heterogeneous, and needs protections such as ACID, you must use a DataBase Management System such as Postgress.
- In the sciences, our data is potentially huge but the structure is much simpler than the more general data dealt with in the database world.
- Data tends to be used for analytics so online-transaction-processing is not a major concern.
- Data Scientists use a number of systems for their data management workflows ... some of the more common systems you hear about include:
 - Arrow
 - Parquet
 - JSON
 - HDF5
 - TileDB
 - CSV
 - XML
 - YAML
 - Binsparse
 - pandas
- There are also common execution models for working with these systems:
 - Hadoop: a map-reduce model organized around operating on data on disks
 - SPARC: Like Hadoop (map-reduce + streaming) but data stored in memory ... MUCH faster than Hadoop. Based on Scala virtual machine.
 - Ray: An in-memory replacement for SPARC with more flexible processing (actors model) and without a performance-killing virtual machine.

How Data Scientists deal with data today ...

- A DBMS separates how data is stored from how it is processed. If data is arbitrarily complex, heterogeneous, and needs protections such as ACID, you must use a DataBase Management System such as Postgress.
- In the sciences, our data is potentially huge but the structure is much simpler than the more general data dealt with in the database world.



Map-Reduce is a variation on SPMD and the Bulk Synchronous Pattern restricted to reductions for the communication phase

Managing Scientific Data

How Data Scientists deal with data today ...

- A DBMS separates how data is stored from how it is processed. If data is arbitrarily complex, heterogeneous, and needs protections such as ACID, you must use a DataBase Management System such as Postgress.
- In the sciences, our data is potentially huge but the structure is much simpler than the more general data dealt with in the database world.
- Data tends to be used for analytics so online-transaction-processing is not a major concern.
- Data Scientists use a number of systems for their data management workflows ... some of the more common systems you hear about include:

It would require an entire course to survey all of these

- 
- Arrow
 - Parquet
 - JSON
 - HDF5
 - TileDB
 - CSV
 - XML
 - YAML
 - Binsparse
 - pandas

- There are also common execution models for working with these systems:
 - Hadoop: a map-reduce model organized around operating on data on disks
 - SPARC: Like Hadoop (map-reduce + streaming) but data stored in memory ... MUCH faster than Hadoop. Based on Scala virtual machine.
 - Ray: An in-memory replacement for SPARC with more flexible processing (actors model) and without a performance-killing virtual machine.

Map-Reduce is a variation on SPMD and the Bulk Synchronous Pattern restricted to reductions for the communication phase

Managing Scientific Data

How Data Scientists deal with data today ...

- A DBMS separates how data is stored from how it is processed. If data is arbitrarily complex, heterogeneous, and needs protections such as ACID, you must use a DataBase Management System such as Postgress .
- In the sciences, our data is potentially huge but the structure is much simpler than the more general data dealt with in the database world.
- Data tends to be used for analytics so online-transaction-processing is not a major concern.

It would require an entire course to survey all of these ... so we'll just pick one to discuss to give you a feel for why it is so important to use a data storage system suited to your problem.

Data Scientists use a number of systems for their data management workflows ... some of the more common systems you hear about include:

- Arrow
- Parquet
- JSON
- HDF5
- TileDB
- CSV
- XML
- YAML
- Binsparse
- pandas

There are also common execution models for working with these systems:

Hadoop: a map-reduce model organized around operating on data on disks

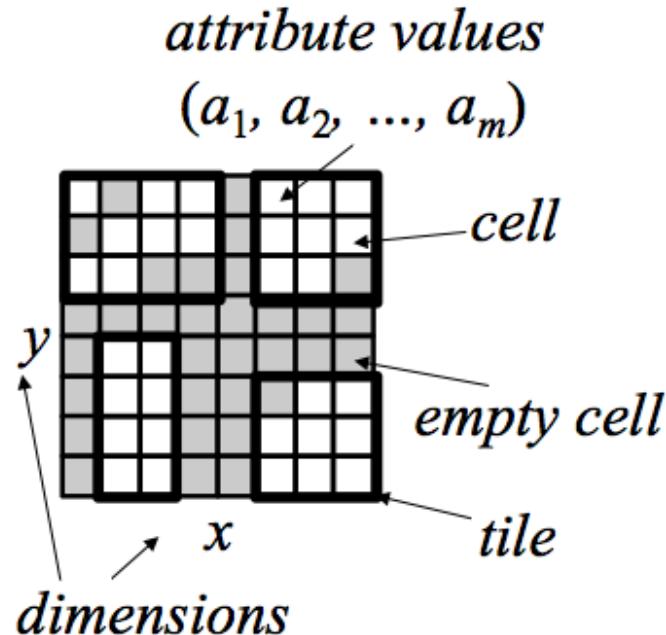
SPARC: Like Hadoop (map-reduce + streaming) but data stored in memory ... MUCH faster than Hadoop. Based on Scala virtual machine.

- Ray: An in-memory replacement for SPARC with more flexible processing (actors model) and without a performance-killing virtual machine.

Map-Reduce is a variation on SPMD and the Bulk Synchronous Pattern restricted to reductions for the communication phase

TileDB: an array data storage manager

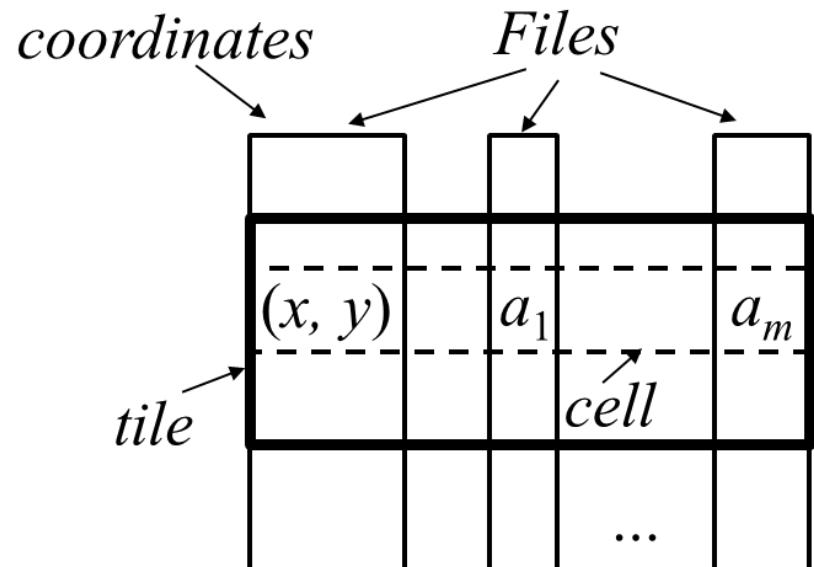
Logical representation



TileDB stores only the non-empty elements of sparse arrays

Sorts cells and packs them into groups of fixed capacity, called tiles

Physical representation



- The tile is the atomic unit of compression
- Their fixed capacity leads to balanced computations

TileDB: Updates at high speed with fragments

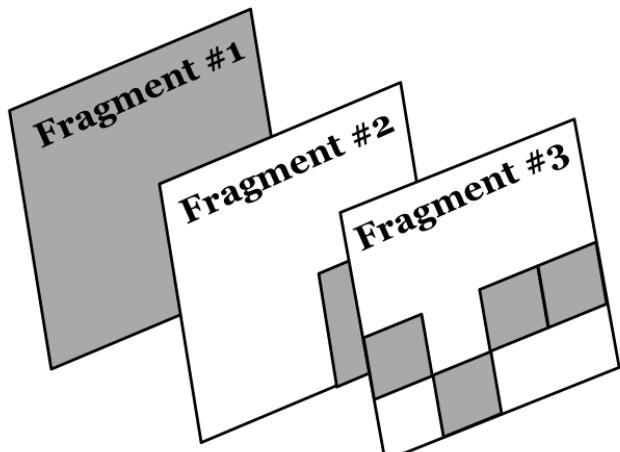
Writes are buffered in fragments

Fragment #1 (dense)				
	1	2	3	4
1	o a	1 bb	4 e	5 ff
2	2 ccc	3 dddd	6 ggg	7 hhhh
3	8 i	9 jj	12 m	13 nn
4	10 kkk	11 llll	14 ooo	15 pppp

Fragment #2 (dense)				
	1	2	3	4
1				
2				
3			112 M	113 NN
4			114 OOO	115 PPPP

Fragment #3 (sparse)				
	1	2	3	4
1				
2				
3	208 u		212 x	213 yy
4		211 wwww		

Batches up fragments for later consolidation in the background



Collective logical array view

	1	2	3	4
1	o a	1 bb	4 e	5 ff
2	2 ccc	3 dddd	6 ggg	7 hhhh
3	208 u	9 jj	212 x	213 yy
4	10 kkk	211 wwww	114 OOO	115 PPPP

Provides a consistent view for reads

TileDB Performance

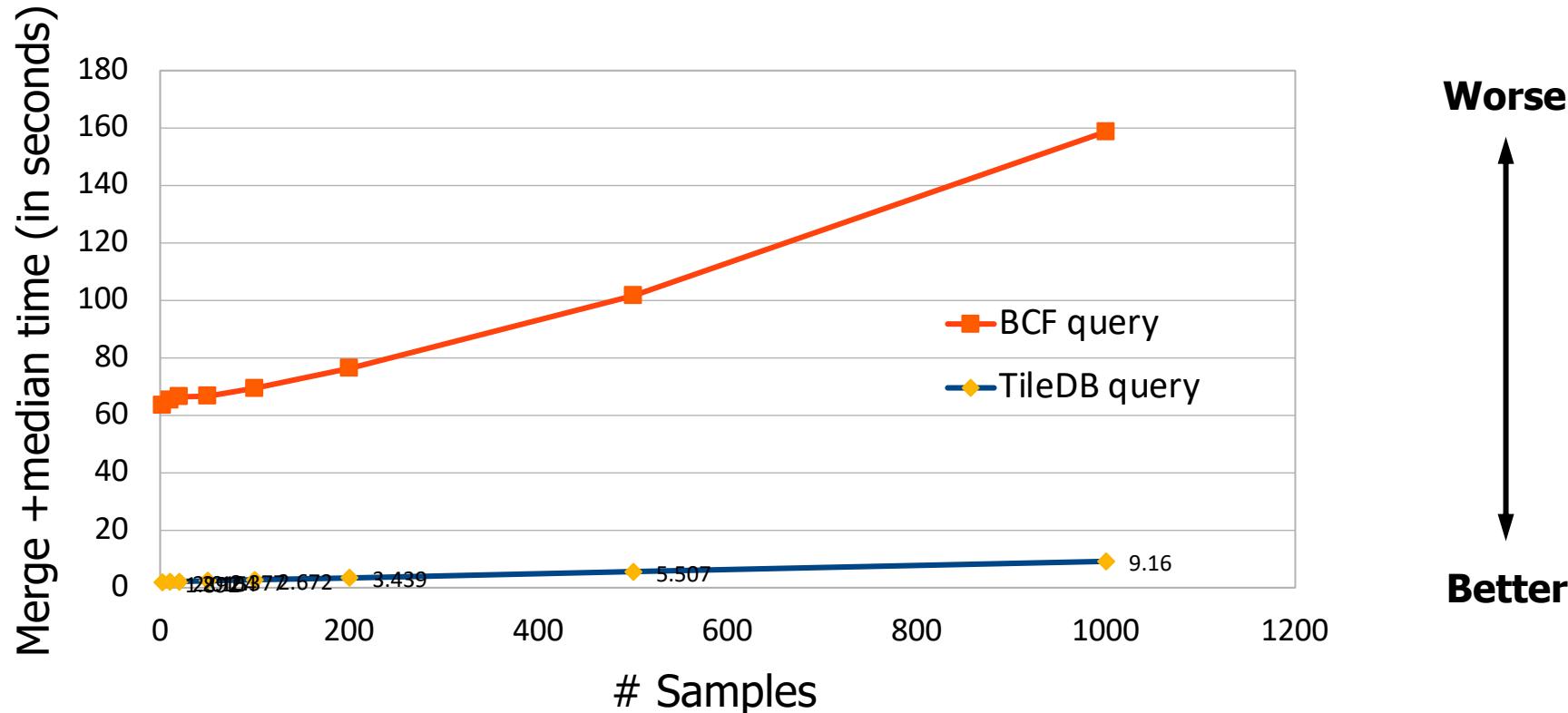
- Extensively benchmarked versus HDF5, SciDB and Vertica[%]
- Benchmarks heavily optimized for each storage engine with help from SciDB and HDF5 teams.
- **Takeaways:**
 - 2x-4x faster than HDF5 on dense reads and sequential writes
 - Orders of magnitude better than HDF5 on random writes
 - Orders of magnitude better than SciDB in all settings
 - Up to 40x faster than Vertica on dense arrays, 2x faster on sparse

Intel® Xeon™ platform with a 2.3 GHz 36-core CPU and 128 GB of RAM, running CentOS6. We utilized a 4 TB, 7200 rpm Western Digital HDD. SciDB v15.12, Vertica v7.02.0202, and HDF5 v.10.0.

[%]Source: Papadopoulos, Datta, Madden, Mattson, VLDB 2017

GenomicsDB: A Data Store optimized for Genomics built on top of TileDB

GenomicsDB combine gVCF operation + median (5K random positions)



BCF refers to the Broad processing pipeline highly optimized by Intel.

gVCF is the genomic Variant Call Format used in the Broad GATK platform for genomics

Intel® Xeon® E5 2697 v2 CPU, 12 cores, dual socket, 128 GB RAM, CentOS6.6, Western Digital 4 TB WD4000F9YZ-0 as a ZFS RAID0 pool.

Open-source and commercial versions of TileDB

The Universal Data Engine



↔ ↔ Pluggable Compute: Efficient APIs & Tool Integrations ↔ ↔

Cloud hosted for
data science
customers

Provides a cloud native
object-store for array data
structures

TileDB Cloud

- ❑ Access control and logging
- ❑ Serverless SQL, UDFs, task graphs
- ❑ Jupyter notebooks and dashboards

Unified data management
and easy serverless compute
at global scale

TileDB Embedded

- ❑ Data versioning & time traveling
- ❑ Columnar, cloud-optimized
- ❑ Parallel IO, rapid reads & writes

Open-source interoperable
storage with a universal
open-spec array format

Full featured
open-source
version for
researchers



Outline

- Motivation: Why everyone needs a database management system?
 - Database Technology: from ancient history to today
 - Data Curation in the sciences
- ➡ • My quest: One Algebra to rule them all

The importance of Algebras

- Remember how Codd's relational algebra revolutionized database management systems?

The Relational Model of Databases

- In 1970 Edgar Codd (IBM) published one of the most important papers in the history of computer science.
- It defined a formal algebra* for building databases ... the **relational model**.
 - Object: A relation.
 - A set of tuples that share a set of attributes.
 - The set of attributes is defined by a schema
 - A relation is typically represented as a table.
 - A set of operators that act on relations. This set includes:
 - Select σ
 - Join \bowtie
 - Rename ρ
 - Project π

* Note: An "algebra" is a set of objects, operators that act on those objects, and rules for how those operators interact with each other

Information Retrieval

P. BAXENDALE, Editor

A Relational Model of Data for Large Shared Data Banks

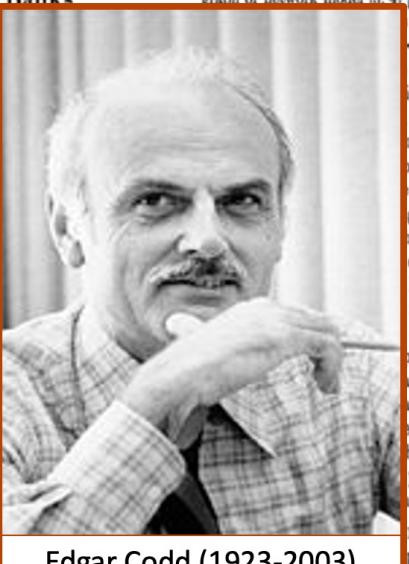
E. F. CODD
IBM Research Laboratory, San

Future users of large data banks will have to know how the data is organized (internal representation). A promising way to do this is to use a relational model. Such information is not a satisfactory answer to the question of what data is available at terminals and most applications. It is also not a satisfactory answer to the question of what data is available when the internal representation changes. It is not even a satisfactory answer to the question of what data is available when some aspects of the internal representation change. Changes in data are often needed as a result of changes in the environment, such as traffic and natural growth in the system.

Existing noninferential, formattable languages, such as COBOL, are not well suited for this purpose. They are designed for use with tree-structured files or similar hierarchical models of the data. In Section 1, it is shown that a relational model can be used to represent data in a more natural and convenient way. A model based on the relational model is introduced. A formal language for data base relations, called SQL, is introduced. A sublanguage of SQL, called the data sublanguage, is introduced. The data sublanguage is used to express queries on relations (other than logical relations) and to apply them to the problems of data representation in the user's model.

KEY WORDS AND PHRASES: data bank, data base, data structure, data organization, hierarchies of data, networks of data, relations, derivability, redundancy, consistency, composition, join, retrieval language, predicate calculus, security, data integrity

CROSS CATEGORIES: 3.70, 3.73, 3.75, 4.20, 4.22, 4.29



Edgar Codd (1923-2003)

The relational view (or model) of data described in Section 1 appears to be superior in several respects to the graph or network model [3, 4] presently in vogue for non-a means of describing data

— that is, without superimposing a means of describing data for machine representation — provides a basis for a high level maximal independence between data and machine representation on the other.

The relational view is that it provides a means of deriving derivability, redundancy, consistency, and normal forms. These are discussed in Section 1. The relational model, on the other hand, has spawned a number of variants, at least of which is misleading for the derivation of relational normal forms (the "connection trap").

The relational model permits a clearer evaluation

of the merits of present formatted

representations of data within a

more general perspective

than that of the relational model. Implementations of the relational model are not discussed.

IN PRESENT SYSTEMS

option tables in recently de-

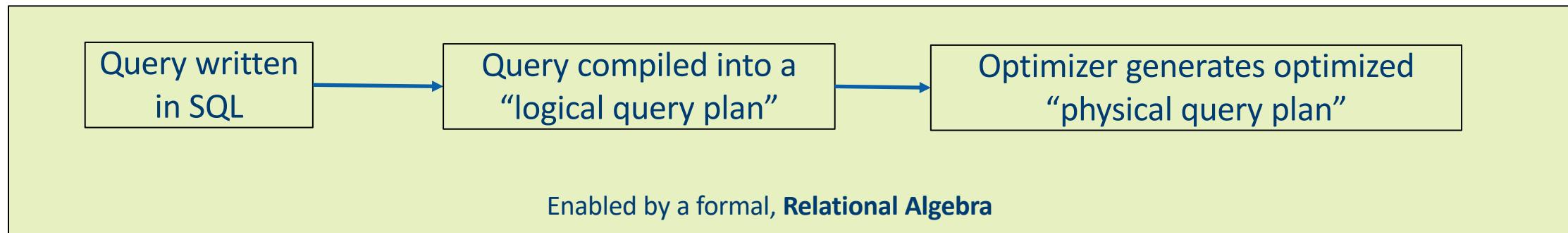
veloped systems represent a major advance

toward the goal of data independence [5, 6, 7]. Such tables facilitate changing certain characteristics of the data representation stored in a data bank. However, the variety of data representation characteristics which can be changed without logically impairing some application programs is

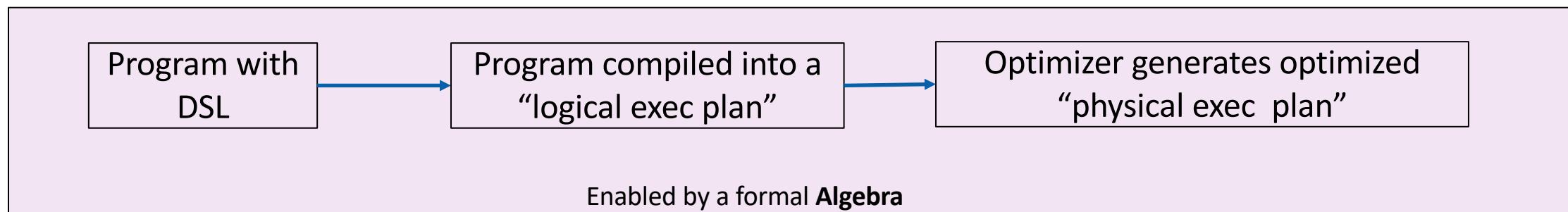
Communications of the ACM, vol 13, no. 6 p. 337, 1970

Productivity, Performance and Portability in one programming framework

- By the 1980s, database researchers at IBM and UC Berkeley exploited the declarative nature of SQL to build systems that delivered on the “3 Ps” ...



- Can we replicate this strategy for programming heterogeneous systems?



The lesson from Edgar Codd so long ago was the power of an algebra to unify disparate approaches to a problem.

Relational algebras are great at data management, but they suck at computation. It would be stupid to build a PDE solver around a relational algebra.

So if we want "one algebra to rule them all", what should be our algebra?

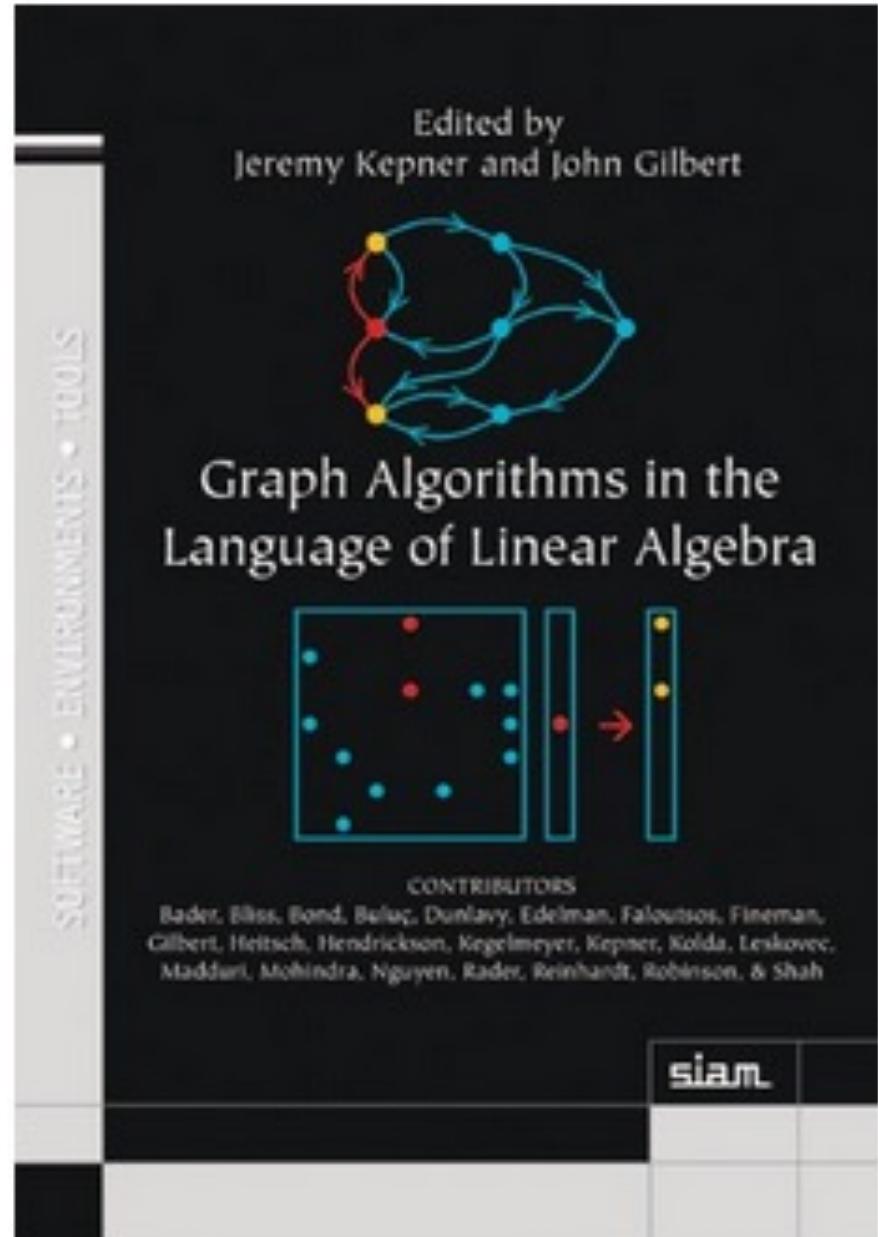
Linear Algebra: One Algebra to rule them all

- Computational physics is basically applied linear algebra
 - We create differential equations from the physics, discretize domains to replace derivatives with differences, and solve resulting algebraic equations.
 - Since the differential operators are replaced by modest sized stencils, the arrays in physics problems are sparse (with a small number of exceptions such as in ab initio quantum chemistry).
- Graphs are linear algebra, databases map onto linear algebra, science and engineering is linear algebra ... if you go deep enough, in almost any field, you end up doing linear algebra.
- All we need is a good library for Sparse Linear Algebra.

Sparse Linear Algebra

If it can do graph algorithms, it can do anything!

- Graph algorithms can be represented in terms of Linear Algebra.
- This is important for Graphs, but it is also used for a wide range of applications ... from engineering codes to databases.
- We need the data structures and a fundamental set of building blocks from which we can construct algorithms ...
We need the GraphBLAS



GraphBLAS is a specification (graphblas.org)

Mathematical Foundations of the GraphBLAS

Jeremy Kepner (MIT Lincoln Laboratory Supercomputing Center), Peter Aaltonen (Indiana University),
David Bader (Georgia Institute of Technology), Aydin Buluç (Lawrence Berkeley National Laboratory),
Franz Franchetti (Carnegie Mellon University), John Gilbert (University of California, Santa Barbara),

Dylan Hutchison (University of Washington), Manoj Kumar (IBM),

Andrew Lumsdaine (Indiana University), Henning Meyerhenke (Karlsruhe Institute of Technology),

Scott McMillan (CMU Software Engineering Institute), Jose Moreira (IBM),

John D. Owens (University of California, Davis), Carl Yang (University of California, Davis),

Marcin Zalewski (Indiana University), Timothy Mattson (Intel)

IEEE HPEC 2016

Design of the GraphBLAS API for C

Aydin Buluç[†], Tim Mattson[‡], Scott McMillan[§], José Moreira[¶], Carl Yang^{*,†}

[†]*Computational Research Division, Lawrence Berkeley National Laboratory*

[‡]*Intel Corporation*

[§]*Software Engineering Institute, Carnegie Mellon University*

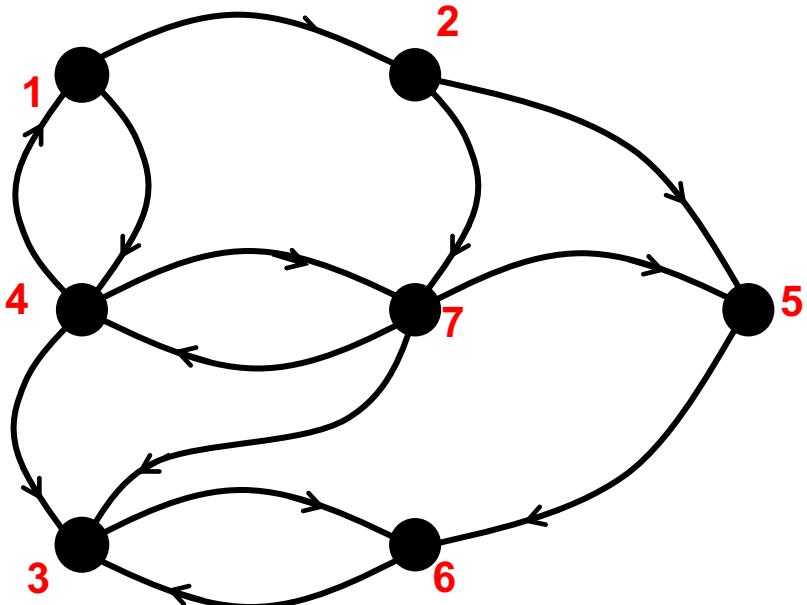
[¶]*IBM Corporation*

^{*}*Electrical and Computer Engineering Department, University of California, Davis, USA*

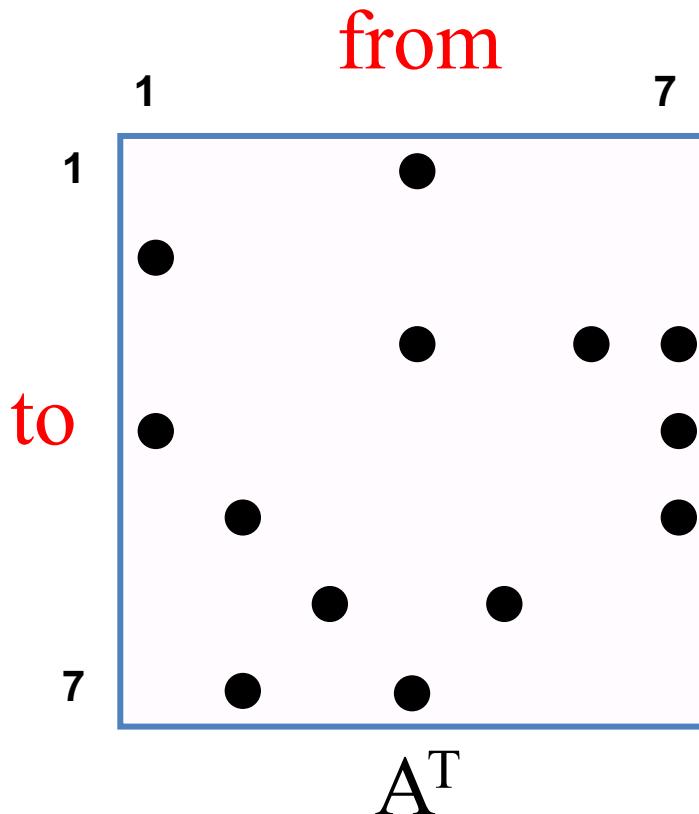
IEEE HPEC 2017

The official GraphBLAS C spec can be found at: www.graphblas.org

Graphs in the Language of Linear Algebra

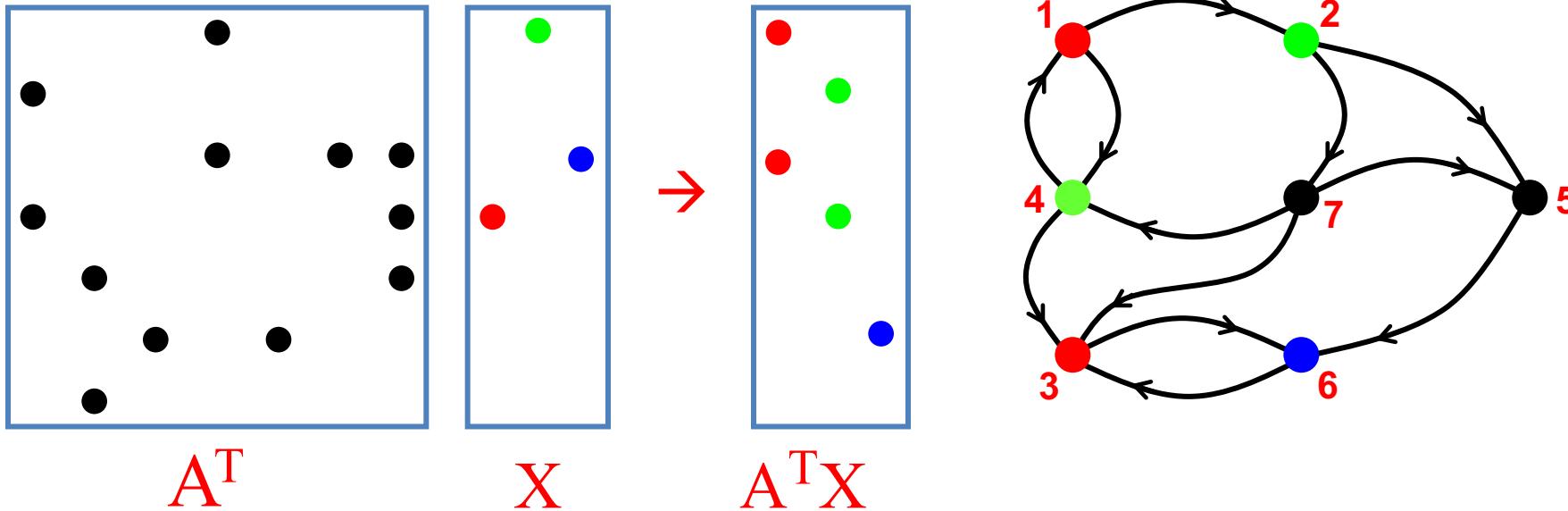


These two diagrams are equivalent representations of a graph.



$A =$ the adjacency matrix ... Elements denote edges between vertices

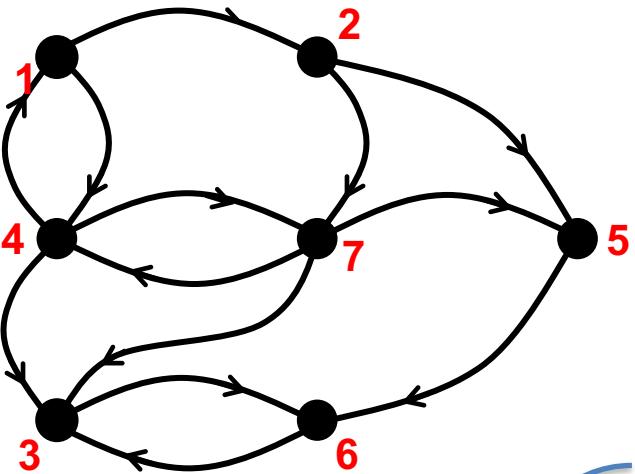
Multiple-source breadth-first search



- Sparse array representation => space efficient
- Sparse matrix-matrix multiplication => work efficient
- Three possible levels of parallelism: searches, vertices, edges

Multiplication of sparse matrices captures breadth first search and serves as the foundation of all algorithms based on BFS

Working with paths

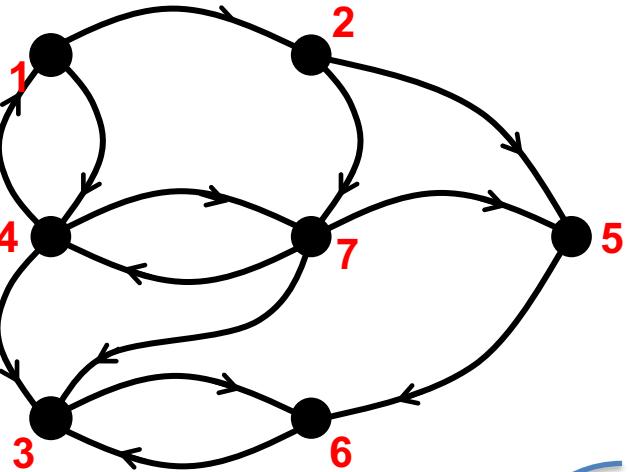


Consider the adjacency matrix with unit cost for “one hop” paths between vertices.

$$A =$$

0	1	0	1	0	0	0
0	0	0	0	1	0	1
0	0	0	0	0	1	0
1	0	1	0	0	0	1
0	0	0	0	0	1	0
0	0	1	0	0	0	0
0	0	1	1	1	0	0

Working with paths



A^2 finds all the “two hop” paths in the graph.

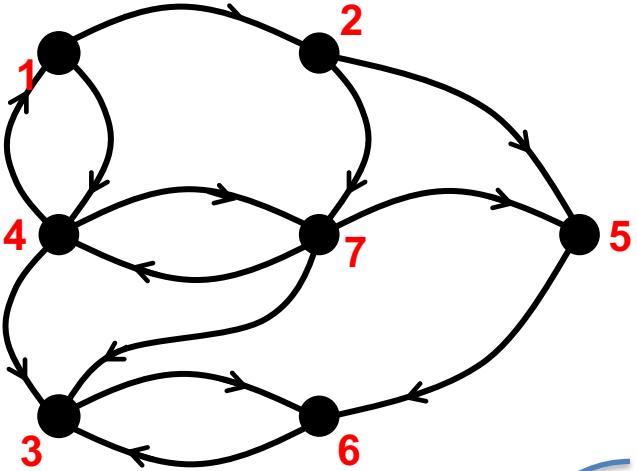
$$A \min.+ A = A^2 =$$

2	0	2	0	2	0	2
0	0	2	2	2	2	0
0	0	2	0	0	0	0
0	2	2	2	2	2	0
0	0	2	0	0	0	0
0	0	0	0	0	2	0
2	0	2	0	0	2	2

Same pattern through the matrices as familiar matrix multiply but:

- replace $+/*$ with $\min/+$
- Replace “zero” with identity of $\min(\infty)$

Working with paths

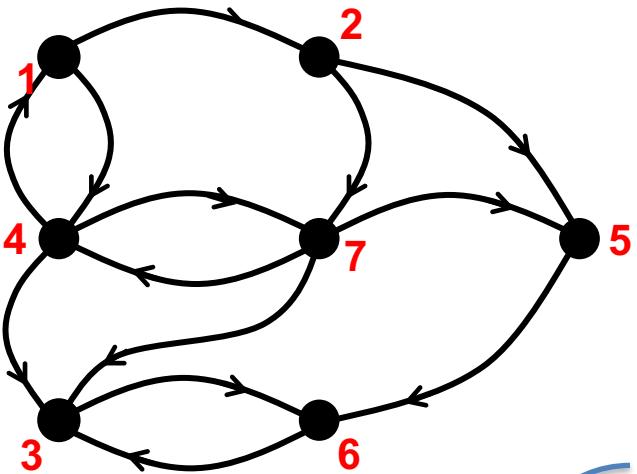


A^3 finds all the “three hop” paths in the graph.

$$A \text{ min.} + A^2 = A^3 =$$

0	3	3	3	3	3	0
3	0	3	0	0	3	3
0	0	0	0	0	3	0
3	0	3	0	3	3	3
0	0	0	0	0	3	0
0	0	3	0	0	0	0
0	3	3	3	3	3	0

Working with paths



Continue until the shortest path matrix no longer changes

In this case, Beyond A^4 the shortest paths don't change. We are done.

Shortest paths =

2	1	2	1	2	3	2
3	4	2	2	1	2	1
0	0	2	0	0	1	0
1	2	1	2	2	2	1
0	0	2	0	0	1	0
0	0	1	0	0	2	0
2	3	1	1	1	2	2

Generalizing Linear Algebra with Algebraic Semirings

- A semiring generalizes the operations of traditional linear algebra by replacing $(+, \cdot)$ with binary operations $(\text{Op1}, \text{Op2})$
 - Op1 and Op2 have identity elements sometimes called 0 and 1
 - Op1 and Op2 are associative.
 - Op1 is commutative, Op2 distributes over op1 from both left and right
 - The Op1 identify is an Op2 annihilator.

Generalizing Linear Algebra with Algebraic Semirings

- A semiring generalizes the operations of traditional linear algebra by replacing $(+, *)$ with binary operations $(\text{Op1}, \text{Op2})$
 - Op1 and Op2 have identity elements sometimes called 0 and 1
 - Op1 and Op2 are associative.
 - Op1 is commutative, Op2 distributes over op1 from both left and right
 - The Op1 identity is an Op2 annihilator.

$(R, +, *, 0, 1)$ Real Field	Standard operations in linear algebra
---------------------------------	---------------------------------------

Notation: $(R, +, *, 0, 1)$

Scalar type Op1 Op2 Identity Op1 Identity Op2

Generalizing Linear Algebra with Algebraic Semirings

- A semiring generalizes the operations of traditional linear algebra by replacing $(+, *)$ with binary operations $(\text{Op1}, \text{Op2})$
 - Op1 and Op2 have identity elements sometimes called 0 and 1
 - Op1 and Op2 are associative.
 - Op1 is commutative, Op2 distributes over op1 from both left and right
 - The Op1 identity is an Op2 annihilator.

$(R, +, *, 0, 1)$ Real Field	Standard operations in linear algebra
$(\{0,1\}, , \&, 0, 1)$ Boolean Semiring	Graph traversal algorithms
$(R \cup \{\infty\}, \min, +, \infty, 0)$ Tropical semiring	Shortest path algorithms
$(R \cup \{\infty\}, \min, *, \infty, 1)$	Selecting a subgraph or contracting nodes to form a quotient graph.

The GraphBLAS Operations

Operation Name	Mathematical Notation		
mxm	$C \langle M, z \rangle$	$=$	$C \odot A \oplus . \otimes B$
mxv	$w \langle m, z \rangle$	$=$	$w \odot A \oplus . \otimes u$
vxm	$w^T \langle m^T, z \rangle$	$=$	$w^T \odot u^T \oplus . \otimes A$
eWiseMult	$C \langle M, z \rangle$	$=$	$C \odot A \otimes B$
	$w \langle m, z \rangle$	$=$	$w \odot u \otimes v$
eWiseAdd	$C \langle M, z \rangle$	$=$	$C \odot A \oplus B$
	$w \langle m, z \rangle$	$=$	$w \odot u \oplus v$
reduce (row)	$w \langle m, z \rangle$	$=$	$w \odot [\oplus_j A(:, j)]$
reduce (scalar)	s	$=$	$s \odot [\oplus_{i,j} A(i, j)]$
	s	$=$	$s \odot [\oplus_i u(i)]$
apply	$C \langle M, z \rangle$	$=$	$C \odot f_u(A)$
	$w \langle m, z \rangle$	$=$	$w \odot f_u(u)$
transpose	$C \langle M, z \rangle$	$=$	$C \odot A^T$
extract	$C \langle M, z \rangle$	$=$	$C \odot A(i, j)$
	$w \langle m, z \rangle$	$=$	$w \odot u(i)$
assign	$C \langle M, z \rangle(i, j)$	$=$	$C(i, j) \odot A$
	$w \langle m, z \rangle(i)$	$=$	$w(i) \odot u$

$\langle M, m \rangle$ are write masks (Matrix/vector). $\langle z \rangle$ selects replace or combine for elements outside the mask.
 \odot is an accumulation operator.

Sparse arrays do science simulations (HPC people have been using them for years).

Sparse arrays do Graphs.

Sparse arrays do ML

... but they can also be used in databases.



A Minimalist Kernel for Linear and Relational Algebra

Shana Hutchison, Bill Howe, Dan Suciu
BeyondMR @SIGMOD, 19 May 2017





Objects: *Associative Tables*

Total functions from keys to
values with finite support

		Attributes		
		Keys	Values	
		k_1	k_2	
		[0]	[“”]	
		v_1	v_2	
a	37	7	'dan'	
a	20	0	" "	
b	25	0	'dylan'	
b	20	2	'bill'	

Annotations:

- A bracket labeled "Default Values" spans the last two columns of the table.
- Two arrows labeled "Support" point to the values 'dan' and 'bill' in the fourth column.

Operators:

UDFs: \otimes , \oplus , f
Think "Semiring"

Join



"horizontal
concat"

Union



"vertical
concat"

Extension

ext_f

"flatmap"

Join and Union adapted from:
M. Spight and V. Tropashko.
First steps in relational lattice. 2006.

Ext is a restricted form
of monadic bind

One algebra to rule them all

- This is very much “work in progress”.
- We know how to do engineering/scientific computing with sparse arrays.
- We have a sophisticated storage engine for sparse arrays.
- We know we can build a full featured database with sparse arrays
- We know in principle that we can indeed create “one algebra to rule them all”. There’s just a bit of engineering work needed to pull everything together.

Conclusion

- Long ago, in scientific computing we selected problems that did not involve much data ... the input/output behavior of our supercomputers was so awful we avoided I/O as much as possible.
- That is no longer the case ... currently, much of scientific computing involves data.
- Hence, a computational scientist is also a data scientist.
- In this lecture, we covered the core concepts to get you started in your journey into the depths of data science. We covered.
 - What is a database and why we need to make our data useful through database technology.
 - The importance of using data-storage engines instead of “flat files”.
 - Key trends in database technology
 - A very brief survey of key data science tools in use today