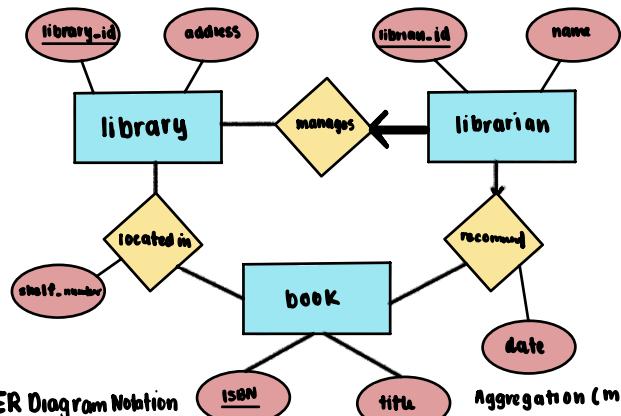


## Entity Relationship Diagram



can = "optional"

## Functional Dependency + Normalization

$X \rightarrow Y$  iff:  $t.X = u.X \rightarrow t.Y = u.Y$

$X \rightarrow Z$  iff:  $t.X = u.X + t.Y = u.Y \rightarrow t.Z = u.Z$

Normalization splitting/decomposing relations minimize redundancy

Functional dependencies constraints b/w two sets of attributes

More mongo

db.tabu.aggregate([

```

    { $group: { _id: "field", newField: { $avgFunc: "$field" } },
      { $match: { column: { $gte: 150000 } },
        { $sort: { field: -1 } },
        { $project: { field: 0 } },
        { $limit: 10 }
    }
  ]
)
```

Aggregation (ML): db.paintings.aggregate([

```

    { $group: { _id: "artistId", totalPrice: { $sum: "$purchasePrice" } },
      { $match: { "totalPrice": { $gte: 10000000 } },
        { $sort: { "totalPrice": -1 } },
        { $project: { _id: 0 } }
    }
  ]
)
```

Updates (ML): db.paintings.updateMany([

```

    { "purchasePrice": { $gt: 5000000 },
      { $inc: { "purchasePrice": 1000000 },
        { $set: { "status": "appreciated" } }
    }
  ]
)
```

Pros (+)

- remove redund.
- min update / delete anomalies

Cons (-)

- junks are costly

## ER Diagram Notation

Edge / Arrow	Bounds	Use Case
—	[ 0-many ]	"not home, up to nul"
—	[ 1-many ]	"at least one"
→	[ 0-1 ]	"at most one"
→	[ 1-1 ]	"exactly one"

## Olap in PostgreSQL

year	month	sales
2022	January	100
2022	February	150
2022	March	200
2023	January	180
2023	February	120

SELECT year, month, sum(sales)

FROM sales-data

GROUP BY ROLLUP(year, month);

SELECT year, month, sum(sales)

FROM sales-data

GROUP BY CUBE (year, month);

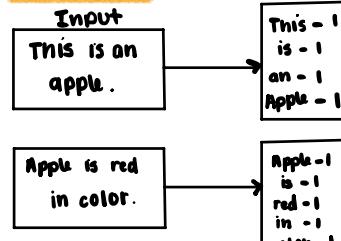
ROLLUP - generates a result that shows aggregates for a hierarchy of values

CUBE - generates a result set that shows all possible combos of aggs for the specified group by elements

+	NULL	January	200
	NULL	February	270
	NULL	March	200

year	month	sum
2022	January	100
2022	February	150
2022	March	200
2022	NULL	450
2023	January	180
2023	February	120
2023	NULL	300
NULL	NULL	750

## MapReduce



## shuffle

This - 1
is - 1
an - 1
Apple - 1
Apple - 1
red - 1
in - 1
color - 1

## Reducer

This - 1
is - 2
an - 1
Apple - 2
red - 1
in - 1
color - 1

## Output

This - 1
is - 2
an - 1
Apple - 2
red - 1
in - 1
color - 1

## MongoDB

Q: list of mems for movieId 192

db.movies.agg([{\$unwind: "\$cast"},

{ \$match: { movieId: { \$eq: 192 } } } ])

Q: movies w/ rating higher than 7

db.ratings.find({ "ratingInRating": { \$gt: 7 } })

Q: Leo DiCaprio movies

db.movies.find({ cast: "Leo DiCaprio" })

## Transactions

Time	1	2	3	4	5	6	7	w(c)
T1	R(A)							
T2		W(B)						
T3			W(A)	W(B)	W(C)			

conflicting Actions

R<sub>1</sub>(A), W<sub>2</sub>(B)  
W<sub>3</sub>(A), W<sub>4</sub>(A)  
R<sub>1</sub>(A), W<sub>3</sub>(A)  
W<sub>2</sub>(B), W<sub>3</sub>(B)  
W<sub>3</sub>(C), W<sub>1</sub>(C)

not conflict serializable

T<sub>1</sub> → T<sub>2</sub> → T<sub>3</sub>

## Vocab

Shared (S) Lock: allow mul trans to read; prevents writing

Exclusive (X) Lock: exec. rights to read + write

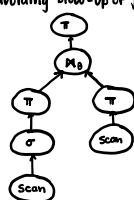
Strict 2-Phase Locking: growing (acquire) + shrinking (release)

Deadlock: when two transactions wait for held resources

### Query Optimization

useful for reducing the size of relational instances and avoiding blow-up of joins  
**predicate pushdown:** pre-select rows before join  
**projection pushdown:** pre-select columns join

ex.)  
`SELECT c.name, s.name, s.population  
FROM candidate AS c  
INNER JOIN states AS s  
ON c.homestateID != s.sid  
WHERE s.population > 800000 AND`



3 components of window function:

1. **PARTITION BY:** groups the rows so that functions will be calculated within these groups instead of the entire set of rows

2. **ORDER BY:** Sort rows in window frame if the order of rows is important (e.g. running totals)

3. **RANGE:** Specify the window frame's relation to current row

`RANK() OVER(PARTITION BY ORDER BY DESC) AS rank`

ex.)  
`SELECT student_name, order_id, SUM(price)  
OVER (PARTITION BY student_id) AS total_spent  
FROM orders  
ORDER BY order_id;`

### Summary

#### Views

↳ virtual: output not stored, computed on demand

CTEs: used within larger queries

materialized: stored on disk, periodically refreshed

### ACID Principles: properties of database transactions

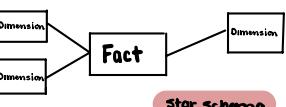
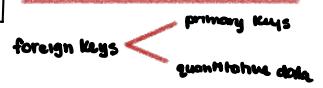
**Atomicity:** "all or nothing" nature of transactions

**Consistency:** ensures valid states (i.e. integrity constraints like uniqueness, PKs)

**Isolation:** ensures that the concurrent execution of diff. transactions leaves database in same state as if transactions executed sequentially

**Durability:** ensures that once transaction is committed, effects → permanent even in the case of system failure

### Physical models of Data Cubes



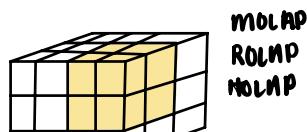
- still one fact table
- many dimension tables
- dimension table connected to others

### Data Cube OLAP Queries

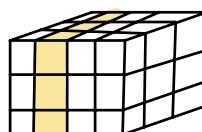
OL (Analytical) P: reads / sums large volumes of multidimensional data

OL (Transactional) P: read, insert, update, delete data; smaller storage

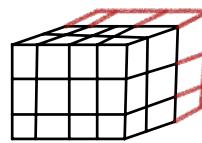
↳ each element in the cube map the unique intersection of all combo of values along dim.



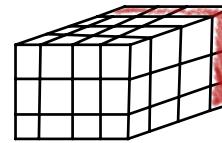
**Dice:** get range partition on one or more dimensions (subcube)



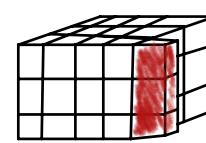
**Slice:** to get equality condition on a dimension



**roll-up:** decrease granularity to sum at higher level of dim hierarchy  
**Summarize**



**drill-down:** increase granularity to sum at low levels of a dim hierarchy  
**lower level**



**Cross-tab:** pivot table, dist of 1 var relative to another

OLAP systems fundamentally provide comprehensive, summarized views of data in data warehouses. Sum will happen across multiple dimensions of data...

### Relational Algebra

Union : R ∪ S

Natural Join : R ⋈ D S

Intersection : R ∩ S

Projection (π) ex. π<sub>id, citation (stops)</sub> Push ↓

Selection (σ) ex. σ<sub>id, citation (age < 50 (stops))</sub>

Renaming (ρ) ex. ρ<sub>stops → data(id → person.id) (stops)</sub>

Cartesian/Cross Product (×) ex. σ<sub>stops.location = zip.location</sub> (stops × zip)

### DML (Data Manipulation Language)

#### Insertion

`INSERT INTO Relation VALUES (<list>)`

↳ ex. `INSERT INTO stops VALUES (5000, 'Asian', 24, 'West Oakland')`

#### Data Models

data frame - row / col labels

matrix - no labels

relations - column labels

tensors - generalization of matrix

- Rank 0, scalar

- Rank 1, vector

- Rank 2, matrix

- Rank 3, tensor

#### Updating

`UPDATE Relation`

`SET < list of new column value assignments >`

`WHERE < condition >`

↳ ex. `UPDATE stops`

`SET age = 10`

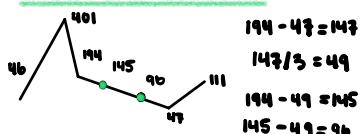
`WHERE age IS NULL ;`

#### Deleting

`DELETE FROM < relationname >`

`DELETE FROM < columnname > WHERE < condition >`

### General Interpolation



### DDL (Data Definition Language)

`CREATE`: defines schema

`DROP`: deletes both schema and instance

`ALTER`: redefines schema

↳ ex.) `CREATE TABLE cast_info`  
`(person_id INT FOREIGN KEY REFERENCES actor(actor_id))`

↳ ex.) `location VARCHAR(20) NOT NULL,`  
`arrest BOOLEAN DEFAULT False,`  
`PRIMARY KEY (stopID)`  
`UNIQUE (person ID, stopTime)`

### Multidimensional OLAP

cube stored as MD array

fast responses

expensive ETL time

pre-comp mat. slices

### Relational OLAP

cube stored as star schema

Supports gen SQL tools

no pre-computation

implemented in SQL

### Common Table Expression EX

WITH relevant\_authors AS (

SELECT \*  
FROM authors

WHERE debut\_year > 2010)

SELECT DISTINCT b.genre

FROM relevant\_authors AS a

INNER JOIN Books AS b ON a.author\_id;

### Materialized View EX

`CREATE MATERIALIZED VIEW _ AS`

### Subqueries EX

`SELECT name  
FROM authors  
WHERE id NOT IN (SELECT distinct id FROM Books);`

### Spreadsheets Conceptual Model

A spreadsheet workspace comprises many sheets

• each sheet has cells

• a spreadsheet is structured around cells

• cells contain one or more of:

↳ value

↳ formula (arithmetic / special functions)

### Access Control

`GRANT [privileges] ON object TO users [WITH GRANT OPTION]`

`REVOKE privileges ON object TO users [CASCADE]`

K  $10^3$  bytes clust. index ↑ sort time, order, min/max ↓ writing  
 M  $10^6$  views - memory - every reference  
 G  $10^9$  mat. views - disk - query must make (low recom)  
 T  $10^{12}$  (te - memory - same as query (visual))  
 P  $10^{15}$  table - disk - permanent data  
 E  $10^{18}$   
 page = 8kb

hierarchy - granularity lvl (t:me = mins/hrs, days)  
 numer. gran. - sig. digits, roll up → coarser (up in hier)  
 ↳ less specific  
 ✓ delete from child  
 X delete from parent (unless cascade)

Database = records + record - size  
 Pages needed =  $\frac{\text{usable space per page}}{\text{bytes/records}}$  =  $\lambda \rightarrow \frac{\text{records}}{\lambda}$   
 df.melt = unpivot df.pivot (x, y, val)

Sparse = many NaNs

Epoch - 4 byte int  
 1 int + 16 = 2 bytes  
 32 = 4 bytes  
 64 = 8 bytes  
 String - timestamp - 20 bytes double 54 = 8 bytes  
 (1 byte for each spot) page = 8 kb  
 SQL timestamp = 8 bytes  
 Extract (Year from timestamp) = 1 byte  
 TZ - timestamp (epoch) at Time Zone 'UTC' as timestamp\_utc  
 String :: int → convert to int  
 int 32 =  $2^{32} = 4 \cdot 3 \times 10^9$  bytes = 28 bytes  
 16 = 65K 64 =  $10^{19}$  8 bits

Nested - small table, ind on join col X big X with index ( $n^2$ )  
 NOT RECOMMEND (!=) small table  
 Sort-merge - medium+ table, data sorted/index both on join  
 ↳ join needs to sort  
 Hash - large tables, no index, join join X range <, >

order by random() expensive (full table)  
 tablesample Bernoulli(p) factor (p for each row)  
 tablesample System(p) factor (selects whole page)

DF - rows, cols have labels  
 Mat = no labels  
 relational col. labels  
 sum c - scalar (univ)  
 1 - vector  
 2 - matrix X  
 3 - tensor (its all)

MAD = median of deviations  
 Hampel =  $1.4826 \cdot \text{MAD}$   
 2 Hampel = 2.0  
 MDL - takes default if possible  
 percen - inc (0.5) within group (order by devicetime)

Mongo  
 .find({ "key": "value" }, ... )  
 \$match: {  
 "lat": { "\$gte": 37, "\$lte": 21 } }  
 \$group ... arg\_star: { "\$avg": "\$stars" }  
 .updateMany() → '\$text'  
 '\$search': "a b c"  
 '\$set': { "field": "val" }  
 \$all, \$in, \$exists, \$unwind, \$lookup  
 \$insertMany \$updateMany w/ \$set join  
 \$lookup('value', where, return)

\$group: {  
 \_id: '\$device',  
 count: { \$sum: 1 } }  
 Conflict = order is not valid (T<sub>2</sub> between T<sub>1</sub>, functions)  
 T<sub>1</sub> [RJ WP] S T<sub>1</sub> [ ] [RJ WP] S T<sub>2</sub> [RP WP]  
 T<sub>2</sub> [ ] [RP WP] S if conflict  
 cycles = unserializable

Serial  
 T<sub>1</sub> [ ] [RJ WP] S T<sub>1</sub> [ ] [RJ WP] S Z  
 T<sub>2</sub> [RP WP] T<sub>2</sub> [RP WP]

T<sub>1</sub> [ ] [RJ WP] T<sub>1</sub> [ ] [RJ WP] X  
 T<sub>2</sub> [RP WP] T<sub>2</sub> [RD RJS WJS WP]

but serializable

Conflict = diff trans, same object, one is write  
 i.e.: T<sub>1</sub> → R1(x) T<sub>2</sub> → W2(x) T<sub>1</sub> → W1(x)

DBT from { \$ref: { 'name': ' } }

roll up = more general day → month  
 \$ = anchor (lexical)

mult X by 8 → X = x < < 3  
 2<sup>3</sup> splits Shuffled  
 cont divide

rollup(yr, mon, val) = 2023 Jun 100  
 n+1 2023 Feb 200  
 2023 null 300

cube(yr, mon, val) = 2023 Jun 100  
 2023 Feb 200  
 2023 null 300  
 2024 Jan 1000  
 null Jan 1200

cube math =  $(\# \text{unique} + 1) \cdot (\# \text{unique} + 1)$

serial = ordered = can do all T<sub>1</sub> x y z  
 and T<sub>2</sub> x y z

same type =  $M, T$

col. label =  $R, df$

row reference by address:  $DF, M, \Gamma, S$

Res. sample = SRS n big O, n rows  
total not known (streams)

Ex. res =