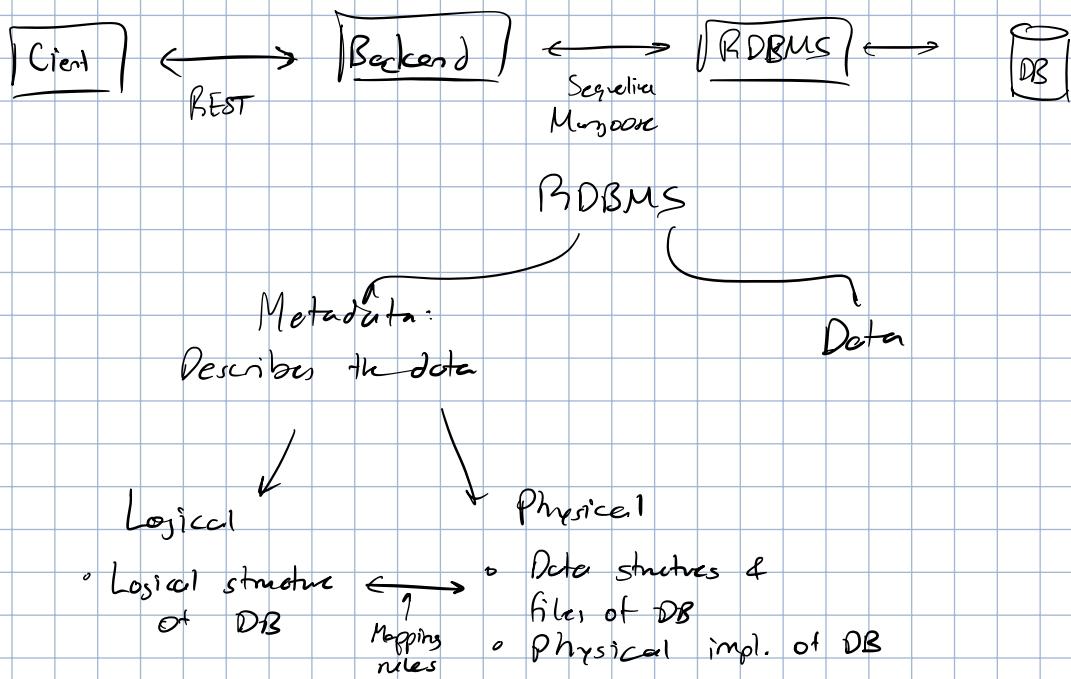


MODULE 1: OVERVIEW OF DATA MANAGEMENT

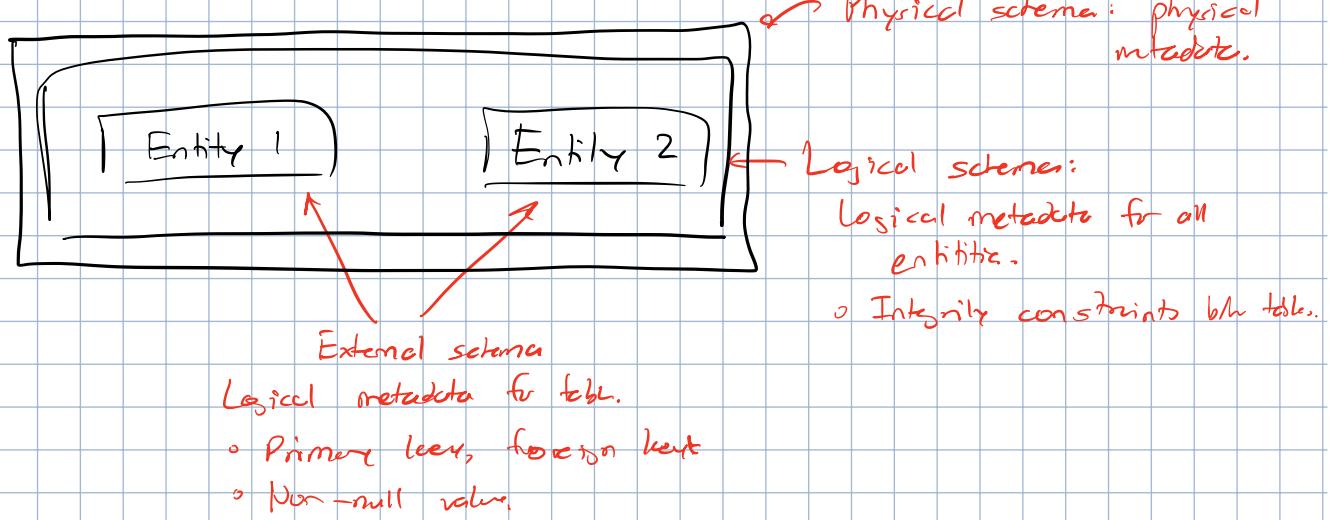
Components of RDBMS

RDBMS: relational DB management system

Typically, architecture is software:



Metadata is composed of schema:



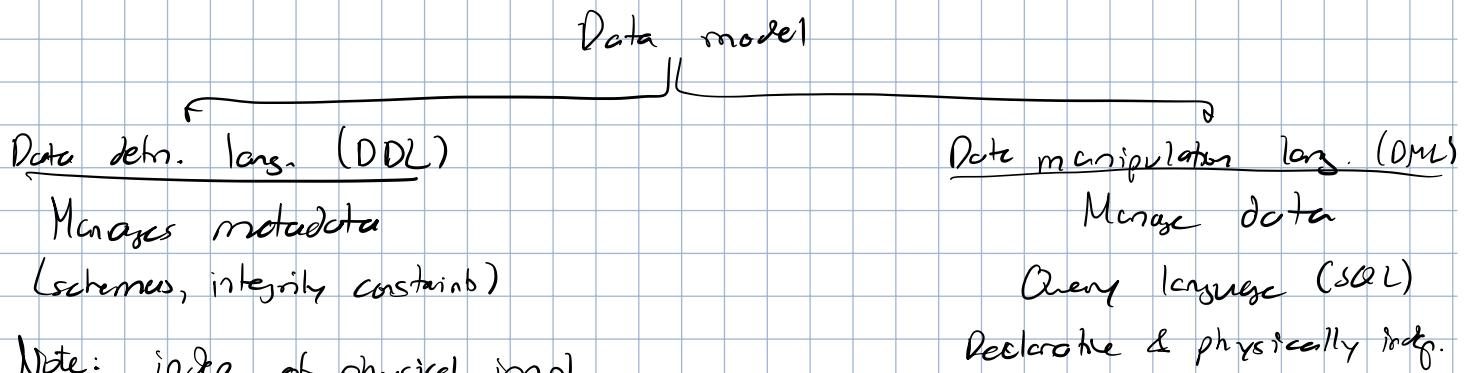
To access DB, we use queries. Criteria for good querying:

1. Indep. of physical structure.
2. Indep. of logical metadata (reject if violates logical schema)
3. Declarative: shouldn't tell RDBMS how to perform.

RDBMS creates model of data \Rightarrow data model

Model: nature of metadata, how do we retrieve & revise data

How to specify the model?



Each schema is not sharing logic

RDBMS Requirements:

#1: RDBMS should be handle concurrent transactions

- o Transaction: indivisible DML requests

- o RDBMS should be ACID compliant

- A: atomic (not dividing transactions, all/nothing)

- C: Consistent (always in a good state)

- I: Isolation (transactions not interfering w/ each other)

- D: Durable (recover data)

#2: RDBMS should be under access control

#3: RDBMS should be able to maintained & monitored

To query, we used predicate logic (SQL)

MODULE 2: RELATIONAL MODEL

Signatures & Relational Calculus

Tables are defined via signatures:

NAME / (column1, column2, col3...)

Data in table is in tuples:

A

col 1	col 2	...	col n
x	y		x
y	z		y
z	c		z

→ tuple

We can specify attributes & tuples via set comprehension & predicate logic.

Ex:// EMP table.

1. Pairs of employees who work for same boss

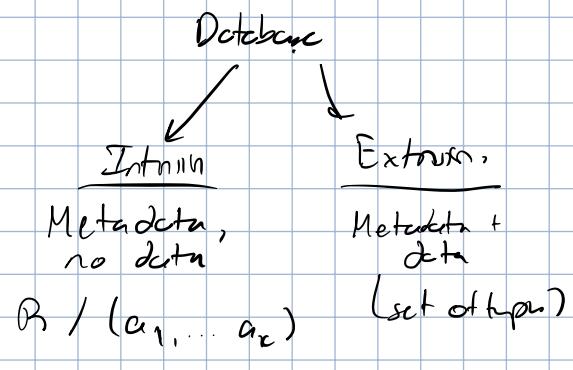
$$\{ (x_1, x_2) \mid \exists y_1, y_2, z. \text{EMP}(x_1, y_1, z) \wedge \text{EMP}(x_2, y_2, z) \}$$

Syntactic sugar.

$$\{ (x_1, x_2) \mid \exists z. \text{EMP}(x_1, -, z) \wedge \text{EMP}(x_2, -, z) \}$$

Database in R.C.

- Universe:
- Set of all rules
 - Equality
 - Constants



All relations in D.B.
(Tables)

Signature: all relations & # of cols in database:

$$\varphi = \{ R_1 / L_1, R_2 / L_2, \dots \} \leftarrow \begin{matrix} \# \text{ of cols in } R_i \\ \uparrow \end{matrix}$$

$$R_1 (col_1, col_2, \dots), \\ R_2 (col_1, col_2, \dots), \\ \vdots$$

Just tells the schema / relations, nothing about data

Instance: actual extension (data) in relation.

When query using predicate logic: diff. answers depending on valuation

↳ Valuation: $\Theta: \{ \text{coll} \rightarrow \overbrace{x_1, x_2, \dots, x_n}^{\infty} \}$

Valuations are diff. for diff. instances.

↳ Formally:

Query: $\{ (\text{coll}, \text{col2}, \dots, \text{colj}) \mid \varphi \}$

Result: $\{ \Theta(\text{col1}), \Theta(\text{col2}), \dots \mid \text{DB}, \Theta \models \varphi \}$

↳ some cond. about relations.

To specify query in rel. calculus, use set comprehension.

↳ Ex:// Who are bosses that have 2 employees.

$$\{ b \mid \exists e_1, e_2 : \text{EMP}(e_1, -, b) \wedge \text{EMP}(e_2, -, b) \wedge \neg(e_1 = e_2) \}$$

Integrity Constraints

Defn: conditions that are true for all relations & instances

Specify via predicate logic.

Ex://

1. Every boss manages a unique department

$$\forall b, d_1, d_2 : \text{EMP}(-, d_1, b) \wedge \text{EMP}(-, d_2, b) \Rightarrow d_1 = d_2$$

Why do we want integrity constraints:

- ↳ Keys: primary key, foreign key
 - ↳ Disjointness
 - ↳ Coverage
- } Enforce behaviors
in DB.

View: create table out of integrity constraint

$$\forall x_1, x_2, \dots, x_n : R(x_1, \dots, x_n) \hookleftarrow \varphi$$

Vier

All tuples satisfying φ in R

Predicate logic

query

Database

φ : ~~single~~

Σ : integrity constraint.

Safety & Finiteness

In any DB, 2 conditions:

1. Extension is finite

2. Queries are safe: answer is finite if DB is finite

Ex: // $\{(x, y) \mid x = 0\} \Rightarrow \text{infinite}$

Domain independent queries: all ans. of qns. only found in extension.

↳ Way to create domain indep. query: range restriction.

$\{(x_1, \dots, x_n) \mid \psi\}$ \nwarrow range restricted (grammar format)

MODULE 3: INTRO TO SQL

Relational Schemata & Conjunctive Queries



SQL is preferred lang. for DDL / DML requests.

① Create tables:

```
create table NAME (
    col-name TYPE not null,
```

:
primary key (col1, col2, ...) \Rightarrow Uniquely defines tuples
foreign key (col-k) references TABLE (col-m),

) ↳ Col-k = col-m

② Query:

```
SELECT col1, col2, ..., colk
FROM TABLE
WHERE  $\psi$ 
```

$\left.\right\} \rightarrow$ Be able to
Rel. quer \Leftrightarrow SQL

Note: can specify multiple tables to query from \Rightarrow Cartesian product

↳ Implication: if 1 table is empty, Cart. product is empty
⇒ no query.

(3) Aliasing.

column name of query
↓
SELECT coll AS cool-name . . .

Set Operations & First Order Queries

Always consider query as a table \Rightarrow set!

Set operations:

1. Union: $Q_1 \cup Q_2$ (all tuples in Q_1 & Q_2)
2. Difference: $Q_1 \setminus Q_2$ (all tuples in Q_1 not in Q_2)
3. Intersection: $Q_1 \cap Q_2$ (all tuples in both Q_1 & Q_2)

Q_1, Q_2
share same
attributes.

Name queries:

WITH query-name AS (query)

Put \Rightarrow query-name as table in another query

Can inline queries.

SELECT . . .
FROM (SELECT . . .)

Nested Queries

SELECT . . .
FROM . . .
WHERE condition (SELECT . . .)

↑ |
 | Q

1. <attr> in Q (query returns list of values)

2. <attr> not in Q

||

3. <attr> op ($>$, $<$, $=$) SOME Q

4. ||

ALL Q (maximum, minimum)
constraint.

5. EXISTS Q

6. NOT EXISTS &

Ex://

Select distinct pubid
from article
where endpage-startpage >= all ()
 select endpage-startpage
 from article
)

} Maximum range.

Nested queries can use variables defined in outer query

Ex:// select * from note r

where not exist (

 select *

 from note (s)

 where s.publication = r.publication
 and s.author < r.author

)

Can't use to output result. Use to check conditions.

Aggregations

Goal: calculate attribute for groups

Syntax: ④ SELECT $\frac{x_1 \dots x_n}{\sum}$ agg1 (attr)

① FROM ...

② WHERE ____ \Rightarrow filtering

③ GROUP BY $x_1 \dots x_n$ \Rightarrow Creates groups / combo
 ↳ Calculate agg (.) / combo

Sometimes, filter out groups based on aggregation

↳ HAVING: filters on aggregation.

SELECT ___, agg1 (attr) as agg

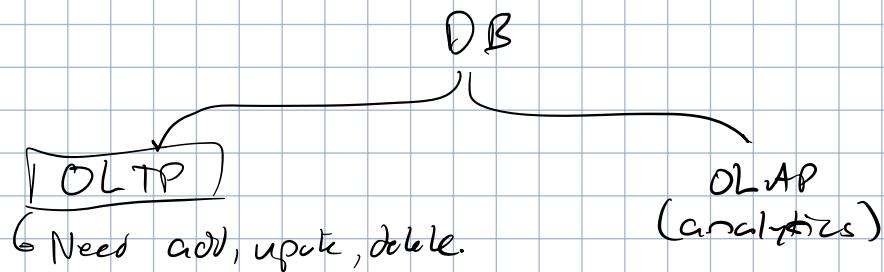
FROM _____

WHERE _____

GROUP BY _____

HAVING agg condition

Transactions & Database Updates



Commands:

① INSERT: $\text{INSERT INTO } \langle T \rangle \text{ VALUES } (\dots)$

In order of
attribut of T

Ex://

T

attr1,	attr2	attr3
—	—	—
val1	val2	val3

Add row:

$\dots \text{VALUES } (val_1, val_2, val_3)$

$\text{VALUES } (\dots) \Leftrightarrow \text{query}$

② DELETE: $\text{DELETE FROM } \langle T \rangle \text{ WHERE } \langle \text{cond} \rangle$

③ UPDATE: $\text{UPDATE } \langle T \rangle \text{ SET } (\text{attr1} = val_1, \dots)$

WHERE $\langle \text{cond} \rangle$

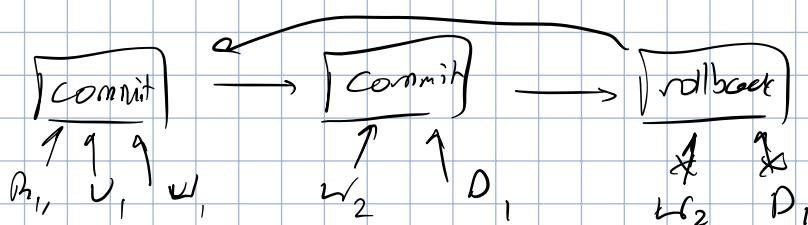
↳ Update rule

↳ Selected tuples to update

Transactions: reads & writes in transactions

↳ Records will not be modified at same time

↳ Way to save changes / remove changes: commit / rollback



MODULE 4: ADVANCED SQL

General Integrity Constraints & Views

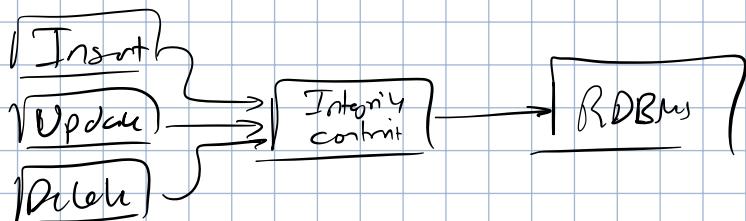
Integrity constraint

o Syntax:

CREATE ASSERTION <name>

CHECK <cond>

↳ SQL query



Ex://

Create unique-pubid

Check (not exist)

Or going to use subquery expr.

Select * from publication p1, publication p2,
where p1.pubid = p2.pubid and p1.title <> p2.title
))

Be able to convert RC integrity \leftrightarrow SQL integrity constraints.

View: like a table defines on certain query

o Syntax:

CREATE VIEW <name> AS <query>

SQL

When is view updated:

1. View is dependent on 1 table
2. View has simple attribute (no crazy expressions for columns)
3. No aggregations / distinct
4. No nested views
5. No set ops.

Multiset Semantics

Relation actually has a count on # of tuples

$\{ d : \exists c. \text{DEPT}(c, d) \}$

What we see:

CS
CS
MATH
:
:

To get distinct tuples

RC:

{ val : elim \forall }

Ex: // $\{ d : \text{elim } \exists c. \text{DEPT}(c, d) \}$

↳ CS |
 MATH |
 : |

SQL:

SELECT DISTINCT --
FROM __
WHERE __

Multiset operations:

Previous set ops don't allow duplicate tuples (behaves like set)

- 1. $Q_1 \cup Q_2$
 - 2. $Q_1 \cap Q_2$
 - 3. $Q_1 - Q_2$
- } Same as set ops, but multiple tuples are allowed

Null Values

Null value

Inapplicable

Column doesn't apply to tuple

↓
Remedied by adding tables

Unknown values

Tuple could have col. value but it doesn't

SOLN

SOLN

Default

How does SQL handle NULLS:

1. Expressions: will have value of NULL if it involves NULL

SELECT col1, col2, col4 + cols \Rightarrow If $col4 = \text{NULL}$ /
 $col5 = \text{NULL} \Rightarrow$
 $col4 + col5 = \text{NULL}$

2. Comparisons: unknown if dealing w/ NULL

SELECT _____

FROM _____

WHERE col1 < col2 \Rightarrow Unknown if col1 or col2 = NULL

All comparisons have 3 valuations: T, F, unknown.

If comparison involves unknown \Rightarrow <cond> IS UNKNOWN

If filtering on nulls \Rightarrow <val> IS NULL

3. Set ops: use special values for duplicates

4. Aggregates: don't count NULLS

\hookrightarrow COUNT (*): will count all tuples returned by query, regardless if row has NULL value.

Joins

FROM T₁ JOIN T₂ ON <cond> ↗ Boolean reduction
between table column

Join types: FULL, LEFT, RIGHT, INNER (default)



If LEFT: join if right col is NULL
If RIGHT: || left ||

Ordering and Limits

Query results can be non-deterministic \Rightarrow has no order

To order:

SELECT
FROM
WHERE
ORDER BY col1, col2, ... colk [DESC/ASC]

To limit results:

;

↓

LIMIT C, [offset 0] \Rightarrow Will take all vals from 0 \rightarrow 0 + c

Triggers and Authorization

Trigger: SQL command that happens on event

o Syntax:

```
CREATE TRIGGER <name>
AFTER <event> ON <table/view>
[REFERENCING OLD AS <old-alias>]
[REFERENCING NEW AS <new-alias>]
FOR EACH ROW [WHEN <condition>] BEGIN ATOMIC
    <DML>
end
```

o Ex:// Create a trigger to delete all rows with author id in note upon delete.

```
create trigger delete-author
after delete on author
referencing old as a
for each row begin atomic
    delete from note where note.author = a.old
end
```

o Integrity constraint checking after triggers

Foreign keys:

B/c auto / cascade deletion are super common, can create foreign reference behaviors on table creation.

o Syntax:

```
FOREIGN KEY (<from-attr>)
REFERENCES <table> [<to-attr>]
[ON DELETE / UPDATE <action>]
```

↳ CASCADE,
SET NULL
;

Authorization:

Data control lang. } { GRANT / REVOKE <role> TO <user>
GRANT / REVOKE <action> ON <obj> TO <user/role>
↓
CONNECT, ALTER, SELECT ↓ Table view

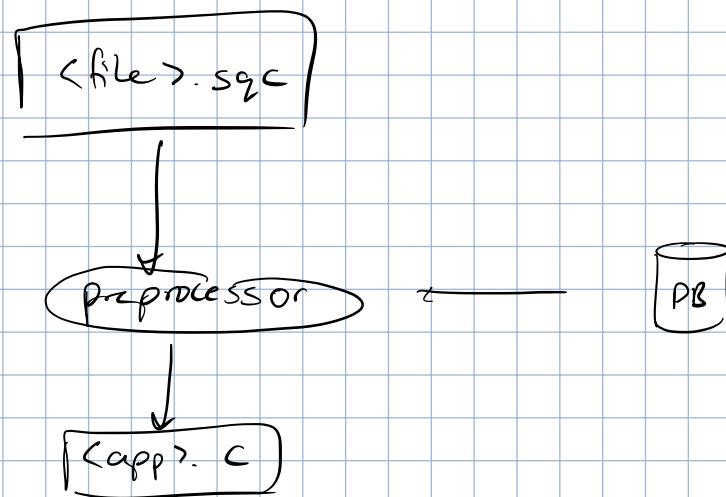
MODULE S: APPLICATION PROGRAMMING

Obj: how to make C-based apps to interact w/ DB.

Embedded SQL in C



General architecture:



Code structure:

```

#include SQL
int main () {
    declarations;
    connect to db;
    work;
    error handling;
    commit/abort & disconnect;
}
  
```

① Include & declarations

1. All .sqlc files start like so:

```

#include <stdio>
:
EXEC SQL INCLUDE SCRLCA
  
```

2. Declare host variables:

Host vars: like a param. in SQL - Can hold values that can be used in SQL & in C!

EXEC SQL INCLUDE SQLCA;

EXEC SQL BEGIN DECLARE SECTION;

C-data-type varName;
:
:

EXEC SQL END DECLARE SECTION;

How to use host variables?

| : hostVar

② Error handling:

```
int main () {
```

EXEC SQL WHENEVER SQLERROR GO TO end;

||

NOT FOUND

||

||

SQL WARNING

↳

NOT FOUND is useful even if not to move to end!

CODE

end:

check_error ("Error: ", &sqlca) \Rightarrow output error

EXEC SQL WHENEVER SQLERRM CONTINUE. \Rightarrow skip own error

EXEC SQL ROLLBACK \Rightarrow Revert

EXEC SQL CONNECT root; \Rightarrow root DB connect.

exit (1)

③ Connecting to DB:

```
int main () {
```

:

:

EXEC SQL CONNECT TO :dB;

Host var that contains DB name.

④ Query:

1 row result: EXEC SQL query > INTO : hostVar
INDICATOR :ind

Indicator < 0 \Rightarrow host variable is null!

What if you have multiple rows in query? Cursor!

① Declare

EXEC SQL DECLARE <name> CURSOR FOR <query>

② Error handling for cursor

EXEC SQL WHENEVER NOT FOUND GO TO <label>

③ To query:

EXEC SQL OPEN <cursorName>

for(;;) {

 EXEC SQL FETCH cursorName INTO :hostVar1,

:hostVar2

 :

}

④ Close:

EXEC SQL CLOSE <cursorName>

⑤ Compile:

.sql \rightarrow 'db2 prep <sql file>' \rightarrow C file \rightarrow compile
 \downarrow
run

Stored Procedures

Functions written in RDBMS

Syntax:

CREATE FUNCTION <name> (var <sql-data-type>)

RETURNS <sql-data-type>

LANGUAGE SQL

RETURN <query that uses var>

Usage:

SELECT <name> (someval), . . . From TABLE

MODULE 6: DYNAMIC APP. PROGRAMMING

Dynamic Embedded SQL.

Q: How do we process dynamic queries?

Main tool:

EXEC SQL PREPARE stmt FROM :strig;

SQL query

parsed in query

What if we want to parameterize query?

↳ Soln:

1. Reconstruct ASCII of query w/o param \Rightarrow recompile

2. Specify param inclusion in query string:

"SELECT . . . FROM '?' . . ."

To execute, must have a param.

How to execute:

If query does not return tuples (DML changes, transactions)

EXEC SQL EXECUTE stmt USING :var1, :var2, . . .

Prepared SQL query

Have vars for params.

If query has tuples: cursors

(1) Declare: exact same

(2) Open:

EXEC SQL OPEN <cursorName> USING :var1, :var2, . . .

(3) Fetch & close: exact same

Helpful tool: descriptor

Stores metadata for a query.

To use:

(1) Inclue \Rightarrow EXEC SQL INCLUDE SQLDA

② Declared way:

```
struct sql_da * <struct_name>  
int main() { . . . }
```

③ Struct info:

```
init_da(<struct_name>, 1);
```

```
'  
'
```

```
EXEC SQL DESCRIBE <query> INTO :<struct_name>
```

④ Access info:

```
<struct_name> → fields → . . .
```

ODBC

Refer to slides, note of my time.

MODULE 7: ENTITY RELATIONSHIP DATA MODEL

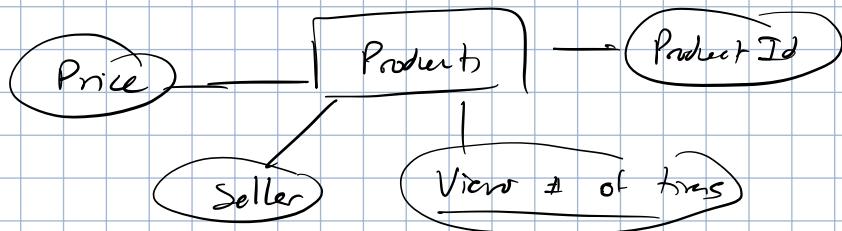
ER Modelling

Entity: set of things that are similar to each other



- Entity set: set of entities

Attributes: properties of entities



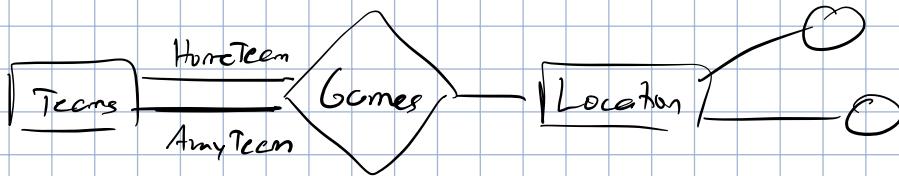
- Each attr has domain: set of all possible attr values.

Relationships: associations b/w entities.

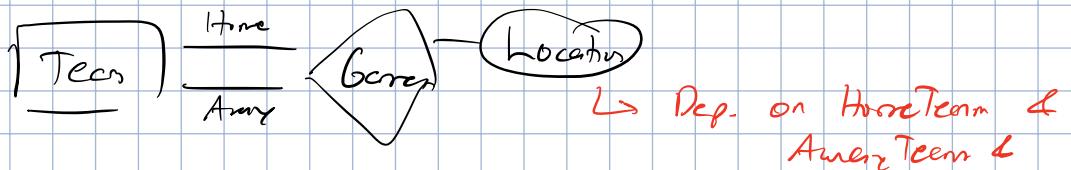


- Note: entities can take on diff. roles in rel.

↳ Self joins



- Can add attr. to relationships but dep. of entities



Integrity Constraints

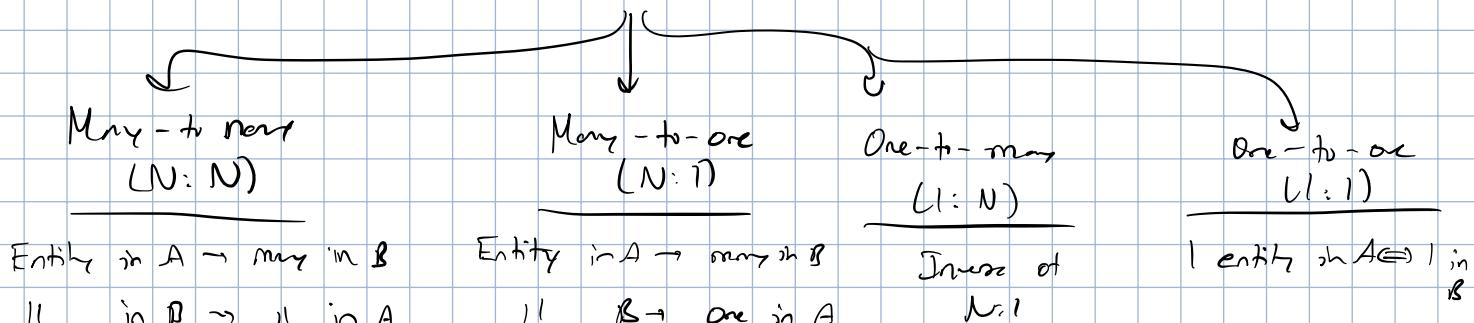
① Primary key

Underline attr. that is primary key



② Binary relationships:

Q: How many rel. does 1 entity have w/ another?



This is putting limits on rel. #, can have less!

To show:

↳ If \rightarrow : 1 : 1 or arrowhead.

Ex: // Employees ↔ Dept. is a 1:N rel.



③ Existence dependencies

x is existent dependent on $y \Rightarrow x$ cannot exist without y

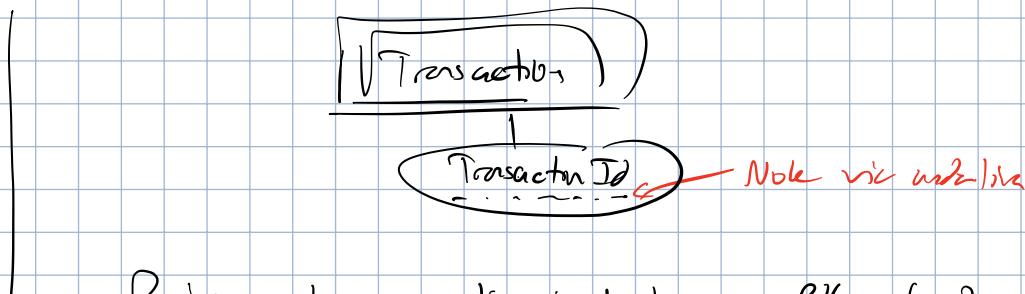
Notation: double bars on subord. & rel.



Trans. is exist deg. on Account

Weak entity set: subord entity

↳ Discriminatr: what dominant entity, what is the PK other?



Primary key: discriminator + PK of dominant

Precursive! Dominant could also be weak.

(4) General cardinality constraint.

Specific limit on # of rel. Give lower & upper bound.

Ex. 11 Students can take 3-5 courses. Courses must have 6-100 students.



To put this int RC.

$\chi = \{$ entity-name / self,
 ;
 ;
 ; ref. to own table
 attribute-name / (self, attribute-name) \Rightarrow Adds on
 ;
 ; attr. val on
 relationship-name / (el, e2, . . .)
 ;

Use SQL queries to write integ.- constraints

Ex:// 1. AID is PK of AUTHOR:

$\wedge \text{AUTHOR}(e_1)$

$$\forall e_1, e_2, v. \text{AUTHOR}(e_1) \wedge \text{AID}(e_1, v) \wedge \text{AID}(e_2, v) \wedge \\ \Rightarrow e_1 = e_2$$

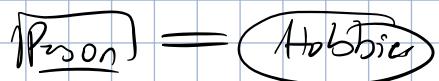
Extended ER modeling

① Structured attr

- Composite attr: attribute w/ other attributes:



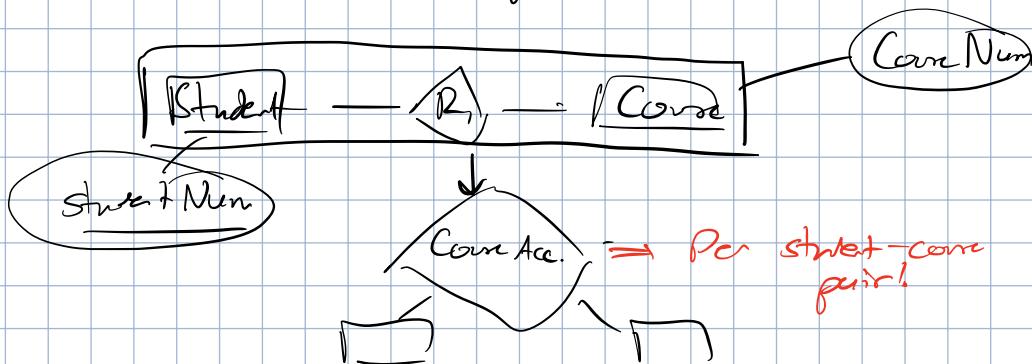
- Multivalued attr: 1 instance = several instances of attr.



② Aggregation:

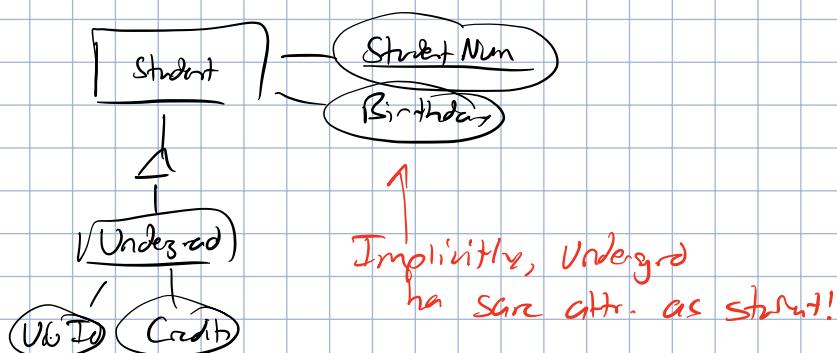
Reity (group) entity & relationships \Rightarrow super entity

Model rel. according to super entity:



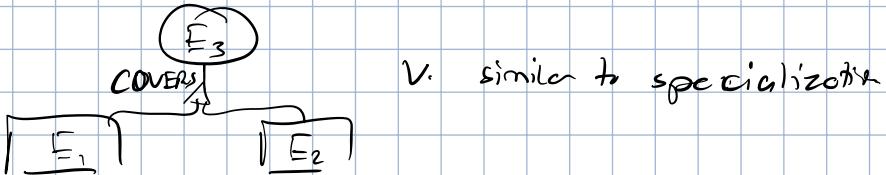
③ Specialization

Polymerism: entity of 1 set is an entity of another specialized set



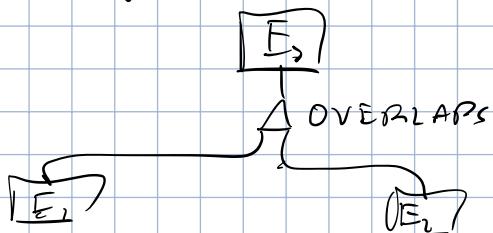
↑
Implicitly, Undergrad
has same attr. as student!

④ Generalization : opposite of specialization.



By default, E_1 is disjoint from E_2

To not have disjointness



How to use in RC?

Specialization:

$$\forall e. \text{specialization}(e) \Rightarrow \text{generalization}(e)$$

Generalization:

$$\forall e. \text{special1}(e) \Rightarrow \text{generalize}(e)$$

$$\forall e. \text{special2}(e) \Rightarrow \text{generalize}(e)$$

$$\forall e. s(e) \Rightarrow s1(e) \vee s2(e)$$

$$\text{Disjoint: } \forall e. s1(e) \Rightarrow \neg s2(e)$$

Aggregation:

1. Add new relations:

rel-ent (self)

el-ent (self, e1)

!

en-ent (self, eN)

2. Add intg. constraint = rule @ slides

3. Make rel. a rich.

MODULE 8: LOGICAL MAPPING

Extended ER \longrightarrow DB schema

Basic ER Diagram Mappings

Entity \rightarrow table. Attr \rightarrow cols. Relationship \rightarrow table/col

① Entity:

Entity set w/ attr. $a_1 \dots a_n \Rightarrow$ table with cols $a_1, \dots a_n$

\therefore 1 entity \rightarrow 1 row in table.

Additional integrity constraints:

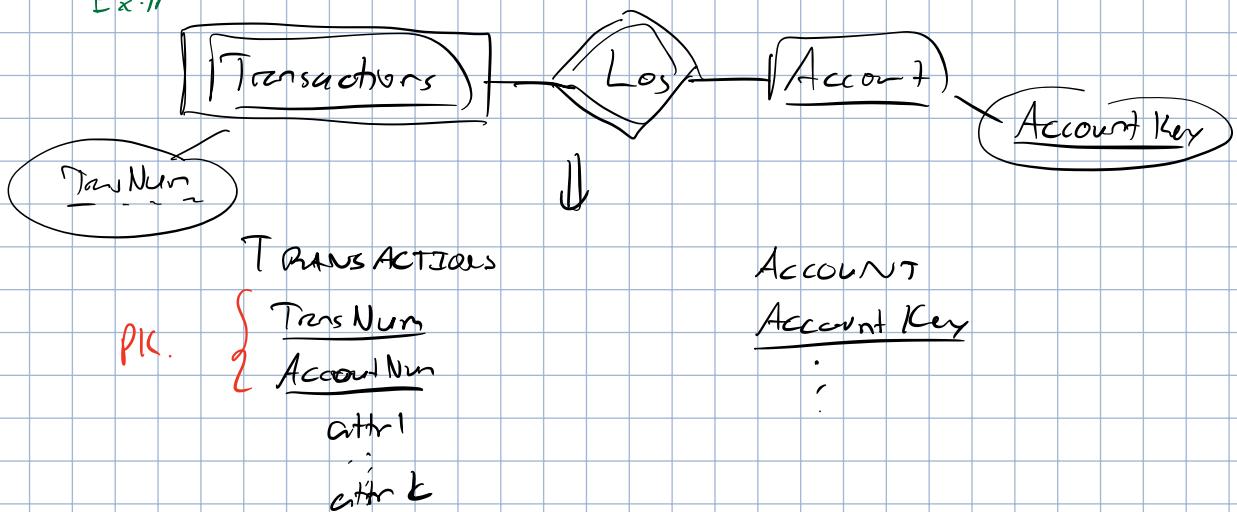
① Primary key: primary attr \Rightarrow primary key.

② Existence dep:

Together \Rightarrow PK of "weak" table

Weak entity set \Rightarrow table w/ discriminator, dominant PK, attr.

Ex://



② Relationships:

Case #1: Relation. for weak entity.

\hookrightarrow Do nothing.

Case #2: (1,1) cardinality constraint

1. Add relationship attributes to components
2. Add PK of component to other components.

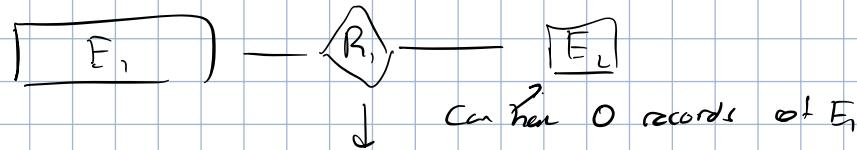
Case #3: Else:

1. Create new table
2. PK of table = PKs of components

3. Attr. of rel. entity \Rightarrow cols of table

Exception: (0, 1) cardinality

\hookrightarrow PK of table = PK of table with possibly 0 of other records



R , PK = PK of E_2 (not PK of $E_1 + E_2$)

Extended ER Mapping

① Aggregation.

Any relation that uses aggregation \Rightarrow new table

\hookrightarrow PK: can include PK of the component tables in aggregation.

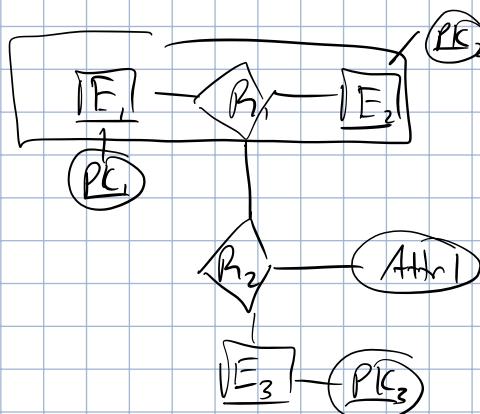


Table for:

E_1, E_2, E_3, R_1, R_2

PKs:

$R_2 \rightarrow$ circle do PK of E_3

$R_2 \rightarrow$ $PK = PK_1 + PK_2$

PK of aggregation.

FKs:

R_1 : FKs to E_1 and E_2

R_2 : FKs to R_1 not E_1 / E_2

\hookrightarrow Looking at E_1, E_2, R as circular.

Bottom line \Rightarrow same as before, outside relations reference agg. relations, not entities.

② Specialization

Specialization is 'weak entity' of generalized table

\hookrightarrow Child table w/ PK being PK of generalized tab.

③ Generalization

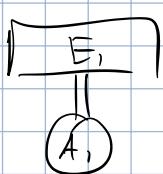
Multiple specialization. Take PK of generalized table

If we have disjoint tables \Rightarrow add integrity constraints.

Tables qualified for views if no FK reference it

\hookrightarrow All attributes must be defined on child tables instead.

④ Multivalued attributes:



\Rightarrow Table E,
PK₁, A-PK
Table A₁:
A-PK₁, PK₂

Triple Store

View that contains all data in DB.

Triple store:

subject	property	object
---------	----------	--------

To change DB to accommodate TS:

\hookrightarrow Add OID to PK \Rightarrow URIs for each table (OID = unique id)

Triple store:

(Select OID as subject, 'in' as prop., 'table-name' as obj. from table)

UNION

(
;

(Select OID as subject, 'attr1' as prop., A1 as obj. from table)

UNION

(
;
} All tables & attributes.

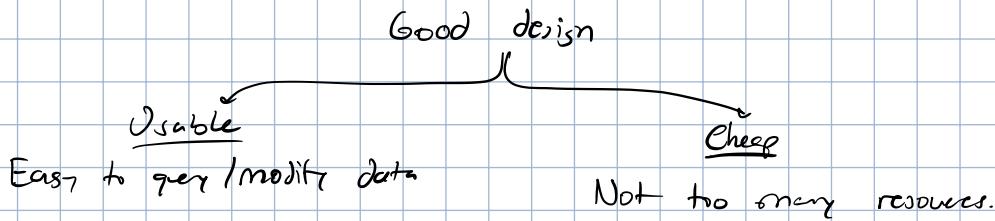
Query on regular DB \Rightarrow Query on triple store

Useful for graph DBs



MODULE 9: DATA DEPENDENCIES & NORMAL FORM

Objective: Create good design for DB



Usually: need to decompose tables for ease of query & modification

Data Dependencies

$R / (A_1 \dots A_k)$. $X \subseteq \text{attr.}$, $Y \subseteq \text{attr.}$ If you can determine values of Y from values of X : $X \rightarrow Y$ (X functionally det. Y)

Ex: II EmpProj / (SIN, PNum, Hours, EName, PName, PLoc, Allorace)

Example of functional dep:

1. SIN \rightarrow EmployeeName
2. PNum \rightarrow PName, PLoc
3. PLoc, Hours \rightarrow Allorace

$\underbrace{\text{Hours}}$
Hours both

If this + identify dependencies in table gives dep:

- ① See when col. values match
 ↳ If col 1 always col 2 match \rightarrow col 1 \rightarrow col 2

Functional dependency closure:

F is functional dep. F^+ : list of all dependencies implied by F .

How to derive more functional dep from existing deps? Armstrong Axiom!

Armstrong Axioms:

1. Reflexivity: $Y \subseteq X \Rightarrow X \rightarrow Y$
2. Augmentation: $X \rightarrow Y \Rightarrow XZ \rightarrow YZ$
3. Transitivity: $X \rightarrow Y$, $Y \rightarrow Z \Rightarrow X \rightarrow Z$

4. Union: $X \rightarrow Y, X \rightarrow Z \Rightarrow X \rightarrow YZ$
5. Decomposition: $X \rightarrow YZ \Rightarrow X \rightarrow Y$

} Need to prove to us.

Ex: // $F = \{ SIN \rightarrow ENam; PNum \rightarrow PName, PLoc; PLoc, Hours \rightarrow Allowance, SIN, PName \rightarrow Hours \}$

Determine $SIN, PNum \rightarrow Allowance$.

Use SE212 principles

1. $SIN, PNum \rightarrow Hours (\in F)$
2. $PNum \rightarrow PName, PLoc (\in F)$
3. $PLoc, Hours \rightarrow Allowance (\in F)$
4. $SIN, PNum \rightarrow PNum (\in F)$
5. $SIN, PNum \rightarrow PName, PLoc$ (transitivity on 2, 4)
6. $SIN, PNum \rightarrow PLoc$ (decomp. on 5)
7. $SIN, PNum \rightarrow PLoc, Hours$ (union on 1, 6)
8. $SIN, PNum \rightarrow Allowance$ (transitivity on 3, 7)

Keys:

1. Superkey: set of cds that determine entire relation.
2. Candidate key: minimal (least amount of keys) superkey.

↳ Primary key: selection for DB

How to know if you have superkey? Want to compute closure of other.

Attribute Closure: list of all columns that can be det. from X & F

↳ Compute $X^+ (X, F)$

- ① $X^+ := X$
- ② While exists $Y \rightarrow Z \in F \wedge Y \subseteq X^+ \wedge Z \not\subseteq X^+$
 - ↳ Add $X^+ := X^+ \cup Z$
- ③ Return X^+

This also correctly creates our attribute closure! Thus in $O(n)$ time.

If X is superkey $\Leftrightarrow \text{Compute } X^+(X, F) = \text{all attr.}$

Table Decomposition.

Goals:

1. Lossless
2. Dependency preserving
3. Normal form

① Lossless:

Defn: able to reconstruct original tables from decomposed without additional tuples (spurious)

How to determine if decomp. is lossless?

Decomposition of $\{R_1, R_2\}$ from R . This is lossless iff

$$R_i \cap R_j \rightarrow R_i \quad \text{or} \quad R_i \cap R_j \rightarrow R_j$$

↳ Common attr. or superkeys of either table in decomp.

Ex:// $R = \{\text{Student}, \text{Assignment}, \text{Group}, \text{Mark}\}$

$F = \{\text{Student}, \text{Assignment} \rightarrow \text{Group}, \text{Mark}\}$

$R_1 = \{\text{Student}, \text{Group}, \text{Mark}\}$

$R_2 = \{\text{Assignment}, \text{Mark}\}$

Is this lossless?

$R_1 \cap R_2 = \{\text{Mark}\} \Rightarrow$ This is not a superset of
 R_1 or $R_2 \Rightarrow$ not lossless.

② Dependency Preservation

All original F should hold without inter-relational dependency clashes

Ex:// $R / (\text{Proj}, \text{Dept}, \text{Div})$

$F = \{\text{Proj} \rightarrow \text{Dept}; \text{Dept} \rightarrow \text{Div}; \text{Proj} \rightarrow \text{Div}\}$

Decompositions:

$$D_1 = \{ R_1 / (\text{Proj}, \text{Dept}), R_2 / (\text{Dept}, \text{Div}) \}$$

Is this good decomp?

① Lossless? Yes ($R_1 \cap R_2$ is superkey of R_2)

② Dependency preserving?

a) Look at each dep. of original

b) Check if dep. can hold if table in isolation

c) If dep. not hold in isolation, check if dependencies remains can be inferred

R_1 has $\text{Proj} \rightarrow \text{Dept}$ which leads to $\text{Proj} \rightarrow \text{Div}$ since $\text{Dept} \rightarrow \text{Div}$ Good

$$D_2 = \{ R_1 / (\text{Proj}, \text{Dept}), R_2 / (\text{Proj}, \text{Div}) \}$$

① Lossless? Yes

② Dep. preserving? No

$\hookrightarrow \text{Dept} \rightarrow \text{Div}$ is missing.

③ Normal form

Normalization \Rightarrow reducing as much redundancy as possible

Ls If redundant \Rightarrow chance anomalies (multiple records of same type, L) one change not cascading to others)

Normal forms

BCNF

3NF

A) BCNF:

Defn: R is in BCNF wrt F iff

$\forall (X \rightarrow Y) \in F^+$ and $XY \in R$, then

1. $X \rightarrow Y$ is trivial ($Y \subseteq X$)

2. X is superkey of R .

OR

Ex:// 1. $R / (A, B, C)$. $F = \{AB \rightarrow C, C \rightarrow B\}$.

Is R in BCNF?

a) $AB \rightarrow C$: AB is superkey of R , so it's good.

b) $C \rightarrow B$: not trivial, $C \neq$ superkey \Rightarrow not in BCNF

How to change rel. into BCNF?

↪ Compute BCNF (R, F):

1. Result := $\{R\}$

2. As long as $\exists R_i \in \text{Result}$ and some dep- $X \rightarrow Y \in F^+$

that violates BCNF

Replace R_i w/ $R_i - (Y - X)$

Add $X \cup Y$ to Result

3. Return result.

This also returns BCNF, lossless decompos. of R .

Ex:// Decompose $R / (A, B, C, D, E)$, $F = \{E \rightarrow CD, BC \rightarrow A, D \rightarrow B, A \rightarrow E\}$.
into BCNF:

Use Compute BCNF algo:

① $E \rightarrow CD$, E is not superkey of R .

Result = $\{R_1 / (A, B, E), R_2 / (E, C, D)\}$

②

1. $E \rightarrow CD$, $ECD \not\subseteq R_1$, so no chkd.

$ECD \subseteq R_2$, E is superkey

2. $BC \rightarrow A$, $BCA \not\subseteq R_1$, R_2

3. $D \rightarrow B$, $DB \not\subseteq R_1, R_2$

4. $A \rightarrow E$: $AE \subseteq R_1$. Is A a superkey?

Yes! Prove via Armstrong's Axioms ($A \rightarrow BE$)

③ Retn:

$\{R_1 / (A, B, E), R_2 / (E, C, D)\}$

Problems:

1. Non-deterministic: depends on order of dep. in algo.

2. Not always dep- present

(B) 3NF

Defn: R is in 3NF wrt $F \Leftrightarrow \forall X \rightarrow Y \in F^+ \wedge X \cdot Y \subseteq R$, then:

1. $X \rightarrow Y$ is trivial ($Y \subseteq X$)
2. X is superkey
3. Each attr. of $Y - X$ is candidate key of R .

Why better than BCNF?

Both lossless & dependency-present decomp.

To create 3NF, we need a minimal cover

↳ Minimal covr: least amt of FD necessary.

F is minimal if:

1. Every RHS of D.P is singular ($A \rightarrow^* BC$, $A \rightarrow^* B$)
2. For no $X \rightarrow A$ is $F - \{X \rightarrow A\}$ equiv. \rightarrow to $F - \{X \rightarrow A\}$
3. For no $X \rightarrow A$ & $Z \subset X$ is $F - \{X \rightarrow A\} \cup \{Z \rightarrow A\}$ equiv.

How to compute: repeat algo until no updates

1. Replace $X \rightarrow YZ$: $X \rightarrow Y$, $X \rightarrow Z$

2. Remove $X \rightarrow A$ if $A \in$ Compute X^+ ($X, F - \{X \rightarrow A\}$)

↳ Can reach A w/out $X \rightarrow A$!

3. Remove A from X in $X \rightarrow B$ if $B \in$ Compute X^+ ($X - \{A\}, F$)

↳ $AB \rightarrow C$. If C can be reached w/ only B , no need to have A .

To create 3NF decomp, follow this algo:

Compute 3NF (R, F):

1. Result := \emptyset , G := minimal cover of F

2. $\forall X \rightarrow Y \in G$:

Result := Result \cup ($X \cdot Y$) } Updating w/
Concat.

3. If no $R_i \in$ Result s.t. R_i contains candidate key for R

a) Compute candidate key K for R

b) Result := Result \cup (K)

4. Return Result.

Ex:// R/(A, B, C, D, E, F). $F = \{F \rightarrow CDE, DE \rightarrow BC, E \rightarrow C, C \rightarrow F\}$
 Construct 3NF decomp. of R.

① Minimal cover of F:

1. First step: $\{F \rightarrow C, F \rightarrow D, F \rightarrow E, DE \rightarrow B, DE \rightarrow C, E \rightarrow C, C \rightarrow F\}$
2. 2nd step: $\{F \rightarrow D, F \rightarrow E, DE \rightarrow B, DE \rightarrow C, E \rightarrow C, C \rightarrow F\}$
3. 3rd step: $\{F \rightarrow D, F \rightarrow E, E \rightarrow B, E \rightarrow C, C \rightarrow F\}$

② Update result initial:

Result: $\{FD, FE, EB, EC, CF\}$

③ Check for candidate keys:

- ① Check candidate key (look back at minimal cover work)
- ② AF is candidate key.
- ③ Add candidate key in Result if not included.

Result: $\boxed{\{FD, FE, EB, EC, CF, AF\}}$

Additional Dependencies & Normal Forms

① Multivalued dependency:

Multivalued attributes: can take multiple to one entity \Rightarrow multiple rows.

Ex:// CAR:

Model Num	Brand	Color
A	Toyota	Black
A	Toyota	Blue
A	Toyota	Green

Multivalued dep: 2 attributes are independent of each, but both dep. on common column.

- $A \rightarrow B$, $A \rightarrow C \Rightarrow A \rightarrow BC$

Specific axioms for MVD

Dependency basis

Lossless decomp. or MVD: $\{R_1, R_2\}$ is lossless:

$$(R_1 \cap R_2) \rightarrow (R_1 - R_2) \text{ or } (R_1 \cap R_2) \rightarrow R_2 - R_1$$

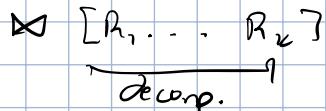
② 4NF:

R is 4NF $\Leftrightarrow x \Rightarrow y \in F^+ \wedge xy \in R$, then either:

1. $x \rightarrow y$ is trivial ($y = x$, or $xy = R$)
 2. x is superior of R

Exactly same as BCNF, so we can also try construct 4NF

③ Join dependency:



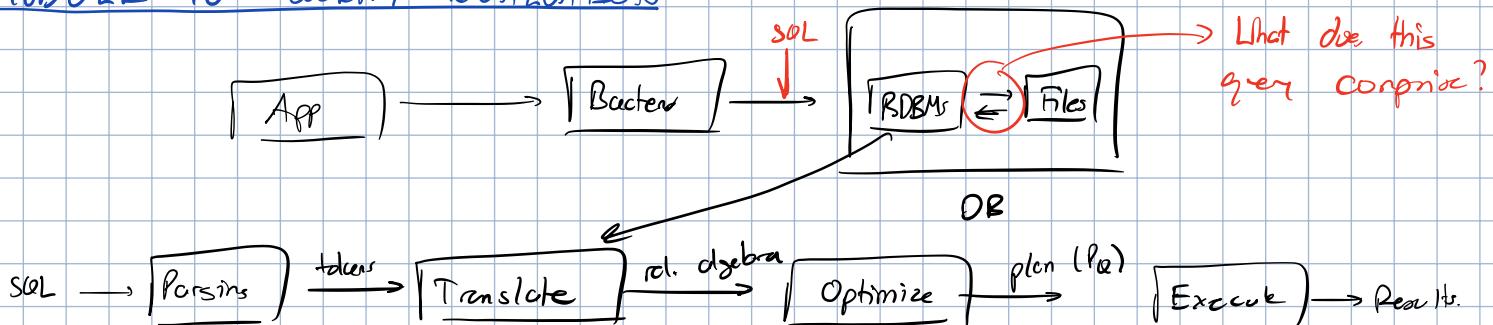
2 constraints:

- $\forall x_i \cdot R_i(x_i) \Leftrightarrow \exists y \cdot \underline{R}(x_i, y) \quad \forall i, 1 \leq i \leq k \quad \text{↳ Number cols. in } \underline{\underline{L}} - 3$
decomp Original rel.
 - $\forall x \cdot R(x) \Leftrightarrow \underline{R}_1(x_1) \wedge \dots \wedge \underline{R}_k(x_k)$

Meaning? we can reconstruct original relation w/ decomps.

This is generalization of MVD

MODULE 10: QUERY EVALUATION



Note: rel. algebra \neq rel. calculus.

Relational Algebra



Unary Operators :

① Selection operator:

$\sigma_{\text{condition}} \text{ (relation)}$

Ex.: // Account / (anum, type, balance, bank, bnum, cnt)

Q: Find all accounts of type A

$\sigma_{\#2 = 'A'} \text{ (ACCOUNT)}$

Column to search (1-index)

Q: Accounts w/ balances > \$1000

$\sigma_{\#3 > 1000} \text{ (ACCOUNT)}$

Condition can be any Boolean condition

↪ Column-to-col. comparison is used a lot!

② Projection:

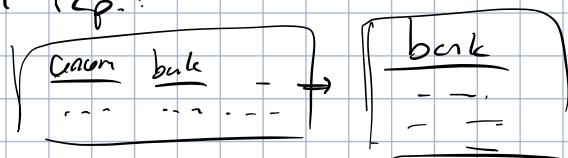
$\Pi_{\#i, \#j, \dots} \text{ (relation)}$

Basically 'select' columns.

Q: Want bank num for all accounts w/ balances > \$1000

$\Pi_{\#4} (\sigma_{\#3 > 1000} \text{ (ACCOUNT)})$

Actual rep.:



③ Duplication elimination:

elim (relation) \Rightarrow removes all duplicate tuples.

Q: What distinct bank nums

elim $\Pi_{\#4} (\sigma_{\#3 > 1000} \text{ (ACCOUNT)})$

Binary operators:

① Cross product:

$R_1 \times R_2$

Take 1 tuple in $R_1 \Rightarrow$ pair w/ all tuples in R_2 & repeat for all tuples in R_2

Adding on more columns.

Ex:// Accounts w/ balances > \$1000

$$\sigma_{\#3 = \#7} (\text{ACCOUNTS} \times 1000)$$

Added 1 column to ACCOUNTS w/ value 1000

② Union:

Exact same as set union

③ Difference:

|| set diff.

To evaluate rel. algebra:

Eval (IQ, DB)
↳ instance to query over.

Codd's Theorem:

Query in RC \Rightarrow Query in RA.

Chainin multiple operators \Rightarrow pipelining! Efficient

↳ No need to store intermediate values.

To make RA more efficient, use index scans.

↳ Indices on table are usually in Btree format

Ex:// PROF / (pnum, lname, dept)

Primary index of PROF is a Btree

RID - P	pnum	lname	dept
- -	- -	- -	—
- -	- -	—	—
- -	- -	—	—
- -	- -	—	—

↳ Unique id / PK in table \Rightarrow Really fast for search.

Why is this important?

↳ Vor index scans rather than rel. scans.

Ex:// 1. Get last name & dept for prof # 19.

Unoptimized:

$$\pi_{\#2, \#3} (\sigma_{\#1=19} (\text{PROF}))$$

Optimized:

$$\pi_{\#2, \#3} (\sigma_{\#2=19} (\text{PROF - PRIMARY}))$$

↳ Why? PROF-PRI~~IMARY~~, pnum, lno...~

Often times, e.g. put secondary indices to make query faster.

Nested join:

$$R_1 \times \sigma_{\#i=\#j \text{ l}} (R_2)$$

Joining R_1 on R_2 s.t. Col j of R_1 = Col i on R_2

Query Optimization

Idea: construct all plans \Rightarrow calculate cost \Rightarrow select cheapest plan

Difficult & unsatisfiable, but reasonable approx. can be done.

To estimate costs, we look @ # of pages need to be queried from Bstore.

Formulas:

① Search on clustered primary index:

$$\text{Cost} = 2 + \frac{\# \text{ of tuples to search}}{\text{blocking factor of index}}$$

② Search on clustered secondary index:

$$\text{Cost} = 2 + \# \text{ of tuples.}$$

③ Entire index scan:

$$\text{Cost} = \frac{|\text{relation}|}{\text{blocking factor}}$$

Note: try to reduce amount of int. results & subexpr. costs

MODULE II: TRANSACTION AND RECOVERY MANAGEMENT

Concurrency Control

Problem: multiple users reading & writing records @ same time.

Recall transactions!

↪ Indivisible DB request: reads & writes.

Schedule ops. w/ transactions to prevent integrity issues.

Notation:

$T_i = r_i[x_1] w_i[x_B] r_i[x_2] w_i[x_B] \dots$

↑ reads object x_A
Transaction # : Read Units

$\{T_1, T_2, \dots, T_r\} \Rightarrow \boxed{\text{RDBMS}} \Rightarrow \{r_1[\dots], w_1[\dots], r_2[\dots], \dots\}$

Schedule

Transactions look as if sequential, but can actually mix & match.

Gold standard for schedules: serializable schedule

↪ Ens result is same as sequentially executing on data!

Ex: // $S_a = w_1[x_2] r_2[x_2] w_1[x_3] c_1, r_2[x_3] c_2$

(
 ↑
 commit (1 is done)

Equiv: $S_a = \underbrace{w_1[x_2] w_1[x_3]}_{T_1} \underbrace{c_1, r_2[x_2] r_2[x_3], c_2}_{T_2}$

Is this equiv? Yes!

Scheduling conflict:

2 Ops conflict if:

1. Ops belong to diff. transaction.
2. Access same obj
3. One of them is write.

Not necessarily bad \Rightarrow only bad if happens in non-serializable sched.

Conflict serializable sched.: have conflicts, but actually serializable.

Q: How to tell if schedule is conflict serializable?

Ans: Use topological sorting

① Create a graph w/ directed edges.

Nodes: transaction

Edges: T_1, T_2 exist if T_1 & T_2 conflict on some op
& T_1 's op schedule is 1st

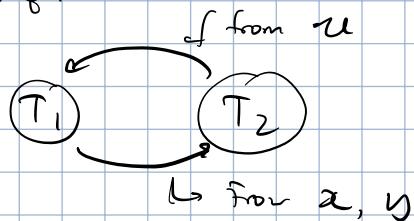
② Check if cycles \Rightarrow if so, not serializable

③ If no cycles \Rightarrow use topological sorting to get serialized schedule.

Ex:// 1. $S_1 = \underline{u_1[x]}, \underline{u_1[y]}, r_2[u], \underline{u_1[z]}, \underline{u_2[x]}, r_2[y], \underline{u_1[z]}$

Is this conflict serializable?

① Create graph:



② Serializable?

No! Cycle.

2. $S_2 = \underline{u_1[x]}, \underline{u_1[y]}, r_2[u], \underline{u_2[x]}, \underline{r_1[y]}, \underline{u_2[z]}, \underline{u_1[z]}$

① Create graph



② Serializable?

Yes! No cycle

\therefore Do T_1 , then T_2

Unique situations:

① Recoverable sched.

Problem: T_i writes to $x \rightarrow T_j$ reads $x \rightarrow T_j$ commits $\rightarrow T_i$ aborts.

Soln: commits in order of read order.

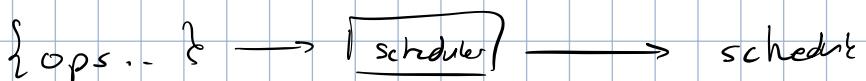
⑦ Cascading aborts

Problem: T_i writes to $x \rightarrow T_j$ reads $x \rightarrow T_i$ aborts

$\hookrightarrow T_i$ aborts $\rightarrow T_j$ is forced to abort \rightarrow cascade of aborts.

Soln: Don't read uncommitted objects

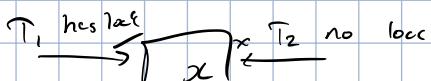
Two Phase Locking



\hookrightarrow can delay, reject, ignore

Lock-based scheduler:

Every transaction access must have object lock before accessing.



- If reading \Rightarrow shared lock
- Multiple shared locks @ same time

- If writing
- Can only have 1 lock

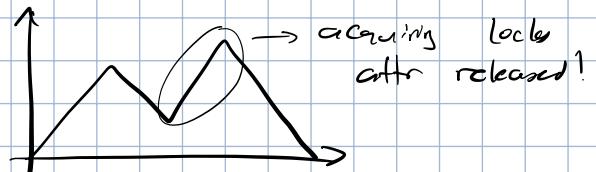
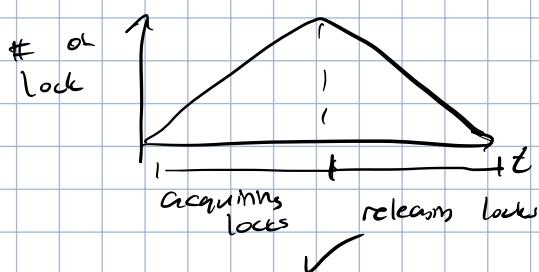
Analogy: blackboard:

\hookrightarrow Everyone can read blackboard @ same time (shared lock)

\hookrightarrow Teacher/student write, cannot read b/c blocking view (exclusive lock)

2 phase locking protocol:

If transaction releases locks, cannot acquire locks again



If schedule is serializable \Rightarrow 2PL schedule possible.

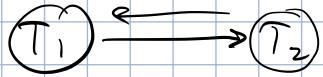
To prevent cascading aborts, use strict 2PL protocol:

- o Cannot release lock until transaction commits/aborts.

Deadlock: T_{i+1} depends on T_i , but T_n depends on T_1

Ex://

$r_1[x], r_2[y], u_2[z], u_1[y]$



Both T_1 & T_2 come first. — deadlock!

How to prevent deadlock?

- ① Use 2PL protocol: locks grants in order of schedule
- ② Detect deadlock via graphs

Other types of locks:

1. Multi-granularity lock: locks depend on size of obj
2. Predicate lock: lock object based on predicate
3. Tree locking: Btree.

Phantom tuple problem:

Issue: get diff # of tuples depending on read & insert & delete order

Ex:// $r_1[x] u_2[z]$

vs.

$u_2[z] r_1[x]$

Soln: lock entire individual predicate lockins.

2PL is too strict \Rightarrow diff. level of concurrency control

Failure Recovery

Module that does recovery must:

1. Atomicity: follow commits/faults
2. Consistent

Schedule \longrightarrow Recovery module \longrightarrow los (append-only)

log structure:

T_i, UNDO, x, old-val
T_i, REND, x, new-val
⋮
T_{i+1}, BEGIN : start of trans.
T_i, COMMIT/ABORT : end of trans.

old-val → new-val

To recover from logs:

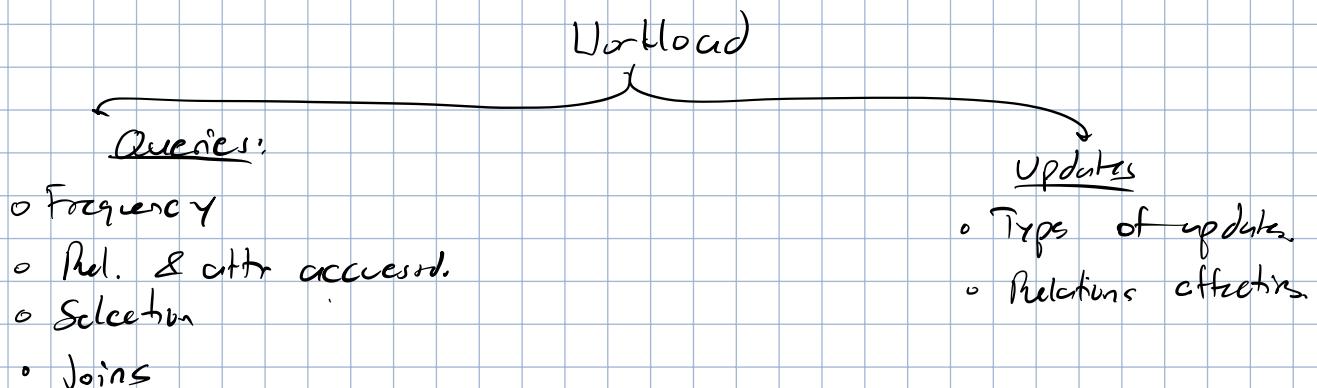
1. Note down all operations from start of failure
2. Note which transactions complete by end.
3. All remaining transactions are aborted \Rightarrow restore to UNDO state.

MODULE 12: DB TUNING

Goal: change phys. design & parameters \rightarrow perf. ↑

L. Might also require logical design changes & DML changes.

To optimize, we need description of workload + bottlenecks.



If we want to optimize, then 2 goals:

- ① Reduce latency
- ② Increase throughput

Physical DB Design

Most powerful design tool \Rightarrow index!

To create indices:

Materialized view:

What is stored in index?

Create view INDEX-NAME as (query)

materialized as DATA-STRUCT. with search key (col)

Btree, heap file, hash file...

index column.

Ex:// 1. Primary index on pnum

Create view PROF-PRIMARY as (select * from prof)

materialized as BTREE with search key (pnum)

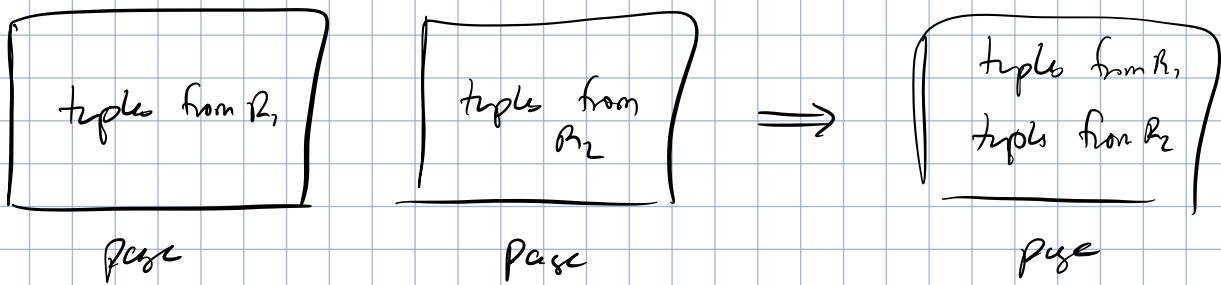
2. Secondary index on lname

Create view PROF-SEC. as (select lname, rid from PROF-PRIMARY)

materialized as BTREE with search key (lname)

To improve perf, we can do co-clustering.

co-clustering.



When to use?

1. Lots of joins b/w R₁ & R₂

2. (1,N) rel.

Pros: speeds up joins.

Cons: sequential reads slower.

Multi-attribute search keys:

Useful in following scenario:

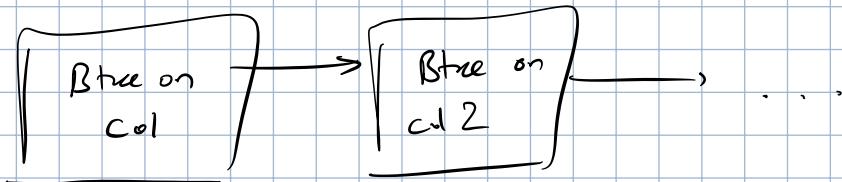
SELECT ...

FROM ...

HAVING col1 = ... , col2 = ... , ...

Order remains
Same across many queries

Physical:



Pro: very fast to do selects on specified order

Con: Not useful if only filters on non-specific order

Join index: just do a materialized view on RIDs of tables to join!

Guidelines:

- Only do indices if benefit > cost of impl.
 - Try + use columns that are used for filtering for indices.
 - Try to choose common columns.

Logical Schema & Query Tuning

Goal: avoid joins as much as possible & return data w few ops.

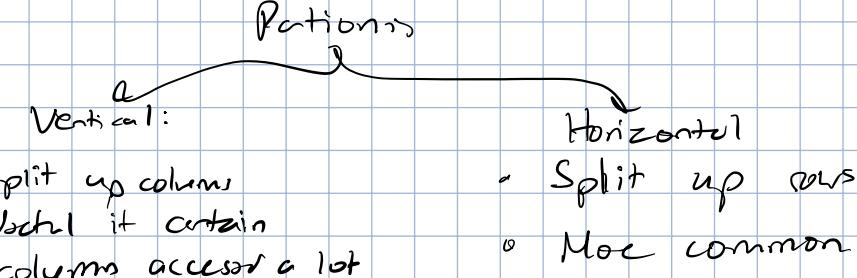
Techniques to grow ions?

① Renormalization: Use a diff. BCNF decomp. that might be better for workload.

② Denormalization: merge tables → redundancy

Why is this good? No need joins!

⑥ Partitioning: split up data across servers \Rightarrow throughput ↑



Techniques to optimize queries:

① Avoid starting, stopping, joining

① Replace subgenu's 25 pins.

③ Make sure plan cost estimation is updated

④ Reduce # of columns & rows returned (WHERE ... useful,
SELECT ... as few
cols as
possible)

⑤ Reorganize transactions & update to avoid hotspots.