

# CSE 344 - Introduction to Database Management

Datalog, RC, Cost Estimation

# Datalog Preview

# Datalog, introduction

- A subset of Prolog
- Using non-recursive datalog with negation for the purposes of this class
- Not implementing datalog (Use DLV if you're curious)
- Used in big data applications (Google page rank algorithm)

Actor(id, fname, lname)  
Casts(pid, mid)  
Movie(id, name, year)

# Datalog, a brief overview

- A datalog *rule*:  $Q1(y) :- \text{Movie}(x,y,1940)$
- Basically a query
- A datalog *fact*:  $\text{Actor}(7920, \text{'Tom'}, \text{'Hanks'})$
- Basically a tuple
- Like RC, uses an *unnamed perspective*, meaning attributes are defined by position rather than name

Actor(id, fname, lname)

Casts(pid, mid)

Movie(id, name, year)

## A more complex rule

- Q2(b) :- Actor(z, 'Tom', 'Hanks'), Casts(z, a), Movie(a, b, 1995)

# A more complex rule

Actor(id, fname, lname)  
Casts(pid, mid)  
Movie(id, name, year)

- Q2(b) :- Actor(z, 'Tom', 'Hanks'), Casts(z, a), Movie(a, b, 1995)
- Names of movies released in 1995 that Tom Hanks was cast in

SQL Alternative:

```
SELECT m.name
FROM Actor a, Casts c, Movie m
WHERE a.id = c.pid
AND c.mid = m.id
AND a.fname = 'Tom'
AND a.lname = 'Hanks'
AND year = 1995
```

Facts:

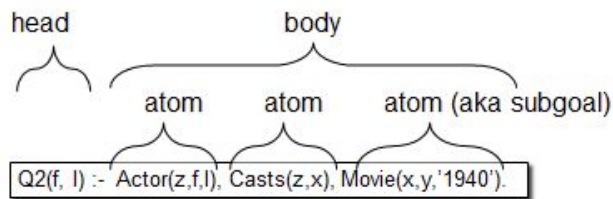
Q2('Apollo 13')  
Q2('Toy Story')

# Datalog continued

Actor(id, fname, lname)  
Casts(pid, mid)  
Movie(id, name, year)

- *Program*: A collection of rules
- *External database relations*: Input relations (Actors, Casts, Movies)
- *Internal database relations*: Output relations (Q1, Q2, Q3, B1, B2, Q4)

## Datalog: Terminology



f, l = head variables  
x, y, z = existential variables

# Safety

- Two unsafe queries:
  - $U1(x,y) :- \text{Movie}(x,z,1994)$
  - $U2(x) :- \text{Movie}(x,z,1994), \text{not Casts}(u,x)$
- No recursion, for example:
  - $T(x,y) :- E(x,y)$
  - $T(x,z) :- E(x,y), T(y,z)$



# Example problem

- Creating a datalog program that uses negation

Consider the following database schema:

*Neighbors(name1, name2, duration)*

*Colleagues(name1, name2, duration)*

Write datalog query that returns all neighbors who do not have any colleagues in common:

# Example problem

- Creating a datalog program that uses negation

Consider the following database schema:

*Neighbors(name1, name2, duration)*

*Colleagues(name1, name2, duration)*

Write datalog query that returns all neighbors who do not have any colleagues in common:

**NonAnswers(n1, n2) :- Neighbors(n1, n2, -),  
Colleagues(n1, c, -), Colleagues(n2, c, -)**

**A(n1, n2) :- Neighbors(n1, n2, -), NOT  
NonAnswers(n1, n2)**

# Cost Estimation Revisited

# Cost Parameters

- **Cost = I/O + CPU + Network BW**
  - We will focus on I/O
- **Parameters:**
  - **$B(R)$**  = # of blocks (i.e., pages) for relation R
  - **$T(R)$**  = # of tuples in relation R
  - **$V(R, a)$**  = # of distinct values of attribute a
    - When  **$a$**  is a key,  **$V(R, a) = T(R)$**
    - When  **$a$**  is not a key,  **$V(R, a)$**  can be anything  $< T(R)$
- Where do these values come from?
  - DBMS collects **statistics** about data on disk

# Estimating Cost

We have 3 relations:

```
Student(sid, name, age, addr)  Book(bid,  
title, author)  Checkout(sid, bid, date)
```

We want to run this query:

```
SELECT S.name  
FROM Student S, Book B, Checkout C  WHERE S.sid =  
C.sid  
AND B.bid = C.bid  
AND B.author = 'Vladimir Putin'  AND S.age > 11  
AND S.age < 20
```

S(sid, name, age, addr)

B(bid, title, author)

C(sid, bid, date)

# Assumptions

Student: S, Book:B, Checkout: C

Sid, bid foreign key in C referencing S and B resp.

Clustered index on C(bid, sid)

There are 10,000 Student records stored on 1,000 pages.

There are 50,000 Book records stored on 5,000 pages.

There are 300,000 Checkout records stored on 15,000 pages.

There are 8,000 unique students who have an entry in Checkout

There are 10,000 unique books that are referenced in Checkout

There are 500 different authors.

$8 \leq \text{student age} \leq 23$

$V(B, \text{author}) = 500$

$T(S) = 10,000$

$B(S) = 1,000$

$S(\text{sid}, \text{name}, \text{age}, \text{addr})$

$V(C, \text{sid}) = 8000$

$T(B) = 50,000$

$B(B) = 5,000$

$B(\text{bid}, \text{title}, \text{author})$

$V(C, \text{bid}) = 10000$

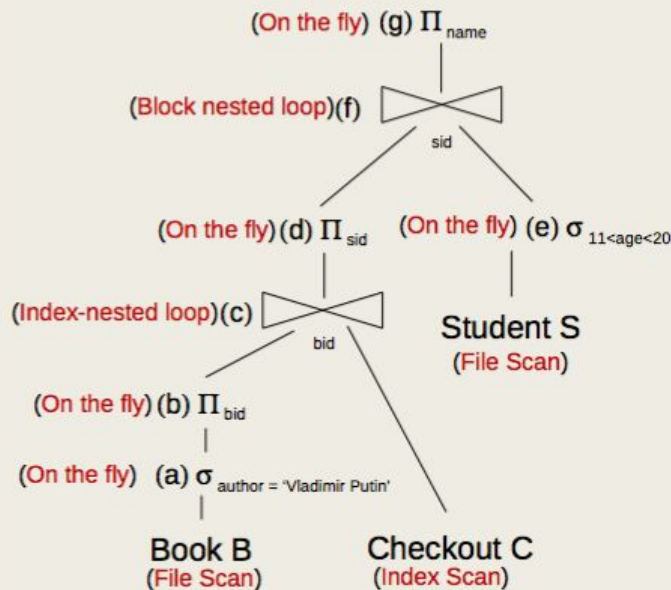
$T(C) = 300,000$

$B(C) = 15,000$

$C(\text{sid}, \text{bid}, \text{date})$

$8 \leq \text{age} \leq 23$

# Selectivity



■ (a)  $\sigma_{\text{author} = \text{'Vladimir Putin'}}$

■ (c) Join predicate bid

■ (e)  $\sigma_{11 < \text{age} < 20}$

■ (f) Join predicate sid

■ (a)  $1 / V(B, \text{author})$   
 $= 1 / 500$

■ (c)  $1 / \max(V(B, \text{bid}), V(C, \text{bid}))$   
 $= 1 / \max(50000, 10000)$   
 $= 1 / 50000$

■ (e)  $(\# \text{ ages covered}) / (\# \text{ possible ages})$   
 $= 8 / 16$   
 $= 1 / 2$

■ (f)  $1 / \max(V(C, \text{sid}), V(S, \text{sid}))$   
 $= 1 / \max(8000, 10000)$   
 $= 1 / 10000$

$V(B, \text{author}) = 500$

$T(S) = 10,000$

$B(S) = 1,000$

$S(\text{sid}, \text{name}, \text{age}, \text{addr})$

$V(C, \text{sid}) = 8000$

$T(B) = 50,000$

$B(B) = 5,000$

$B(\text{bid}, \text{title}, \text{author})$

$V(C, \text{bid}) = 10000$

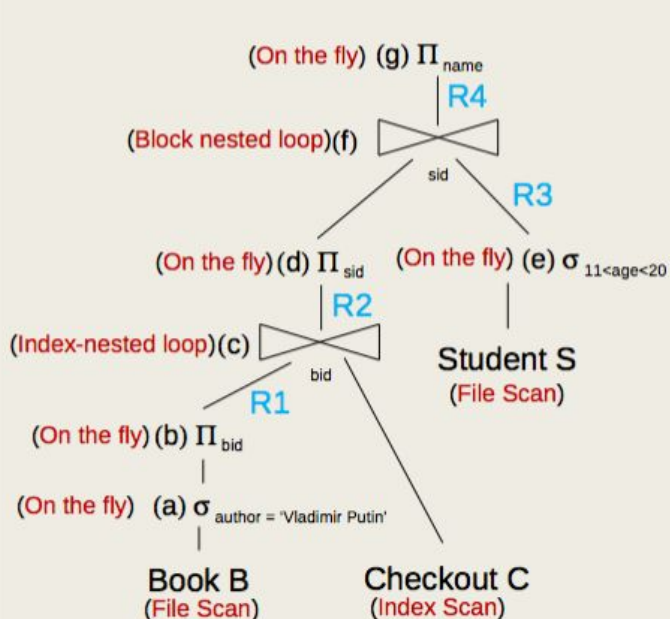
$T(C) = 300,000$

$B(C) = 15,000$

$C(\text{sid}, \text{bid}, \text{date})$

$8 \leq \text{age} \leq 23$

# Cardinality



■  $T(R1) = ?$

■  $T(R2) = ?$

■  $T(R3) = ?$

■  $T(R4) = ?$

■  $T(R1) = T(B) / 500$   
 $= 100$

■  $T(R2) = T(R1) * T(C) / 50000$   
 $= 100 * 300000 / 50000$   
 $= 600$

■  $T(R3) = T(S) / 2$   
 $= 10000 / 2$   
 $= 5000$

■  $T(R4) = T(R2) * T(R3) / 10000$   
 $= 600 * 5000 / 10000$   
 $= 300$



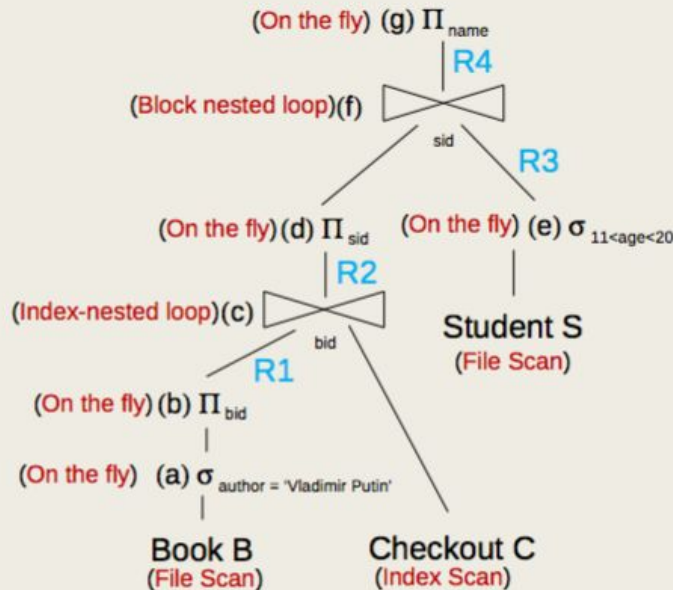
$V(B, \text{author}) = 500$   
 $V(C, \text{sid}) = 8000$   
 $V(C, \text{bid}) = 10000$   
 $8 \leq \text{age} \leq 23$

$T(S) = 10,000$   
 $T(B) = 50,000$   
 $T(C) = 300,000$

$B(S) = 1,000$   
 $B(B) = 5,000$   
 $B(C) = 15,000$

$S(\text{sid}, \text{name}, \text{age}, \text{addr})$   
 $B(\text{bid}, \text{title}, \text{author})$   
 $C(\text{sid}, \text{bid}, \text{date})$

# Cost



- Data not sorted in any way
- Relations can fit in memory
- Compute the cost of each step (a) through (g)
- (a)  $B(B) = 5000$
- (b) 0
- (c)  $T(R1) * B(C) / V(C, \text{bid})$   
 $= 100 * 15000 / 10000 = 150$
- (d) 0
- (e)  $B(S) = 1000$
- (f) 0
- (g) 0
- Total: 6150

# RC/Datalog Worksheet