# SQL

(5) SELECT [Distinct] col1, col2, agg(col3)
(1) FROM table1, table2,
(2) WHERE <predicate> AND <predicate> OR
(3) GROUP BY <column list>
(4) HAVING <predicate>
(6) ORDER BY <columns> [DESC] [ASC]
(7) LIMIT <integer>;
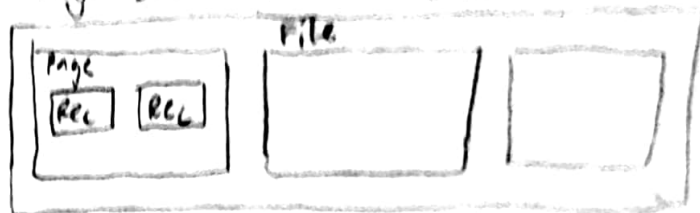
> Everything in SELECT must be in GROUP-BY or is an aggregate

> Must contain only GROUP-BY columns of aggregating functions

## SQL String Comparison
- Old School SQL    starts with B
  WHERE S.name LIKE 'B_o' ← return Bob
  _ = any single char ; % = 0+ chars
- Standard Regular Expressions
  WHERE S.name ~ 'B.*' ← returns Bob, McBob
  . = any char ; * = repeat (0+ instances of previous)

## SQL Join Variants

From table1
  [INNER | NATURAL | {LEFT | RIGHT | FULL}
  OUTER ] JOIN table2
  ON <qualification list>
- INNER: join tables where "ON" qualification
- NATURAL: join tables for pair of attributes
  w/ same name

| FROM t1, t2 | = | FROM t1 INNER JOIN t2 |
| WHERE t1.id=t2.id | | ON t1.id = t2.id |

FROM t1 NATURAL JOIN t2

(assume only matching column names = "id")

- LEFT/RIGHT/FULL OUTER JOIN
  FROM t1 [ ] OUTER JOIN
  ON t1.id = t2.id
  - LEFT: if t1.id has no match, t2.id is NULL
     t1 ⬤ t2
  - Right: if t2.id has no match, t1.id is NULL
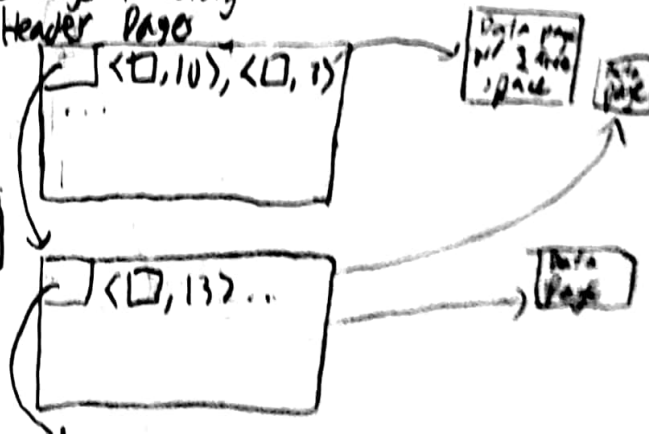     t1 ⬤ t2
  - FULL: LEFT and RIGHT both nulls  t1 ⬤ t2

---

No aggregates!

# Relational Table Properties
- Schema is fixed
  - unique attribute names
  - atomic (primitive) types
- Tables are not ordered (sets, multisets)
- Tables are Flat (no nested tables)
  - First Normal Form

## Disk Representation    [pageID, location on page]
- DB File: collection of pages
- Page: collection of records

File


## Unordered Heap Files
Records placed arbitrarily across pages
1. File as Doubly Linked List



Pages w/ full space
Pages with free space

2. Page Directory
Header Page



## Unordered Heap Files vs Sorted File

| 2.5 | 1.6 | 4.7 | 3.10 | 8.9 | vs. | 1.2 | 3.4 | 5.6 | 7.8 | 9.10 |

| | Heap | Sorted | Assume |
|---|---|---|---|
| Scan all records | B·D | B·D | append page |
| Equality search | 0.5·B·D | (log B)·D | to end of heap files |
| Range search | B·D | (log B + pages)·D | Sorted files |
| Insert | 2·D | (log B * B)·D | are packed |
| Delete | (0.5·B)·D | (log B * B)·D | |

## Page Layout

- FIXED record lengths, UNPACKED records

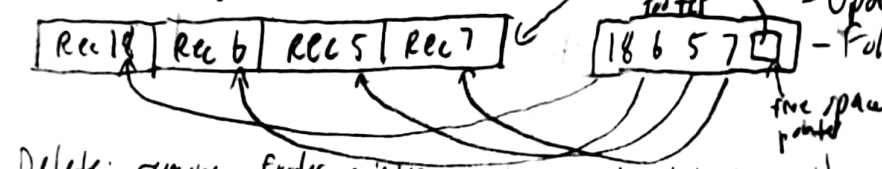| Header | 1 | 0 | 1 | 0 | Rec | X | Rec | X | Rec |
|--------|---|---|---|---|-----|---|-----|---|-----|

Insert: find first empty slot in bitmap ⌃
Delete: find record and clear bit in Header ☺
 - No organization needed
- FIXED record lengths, PACKED RECORDS

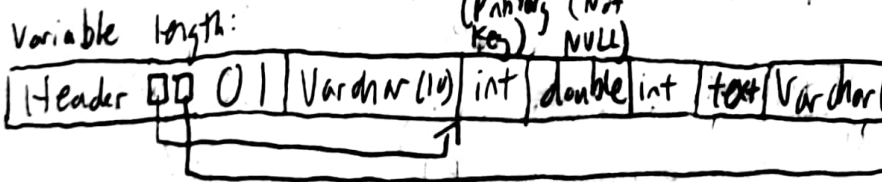| Header | Rec | Rec | Rec | Rec |
|--------|-----|-----|-----|-----|

Insert: just append to end ☺
Delete: Scan for record, delete, reorganize ⌃
  all records (b/c pack)

- VARIABLE record lengths, UNPACKED records
  (Slotted page)

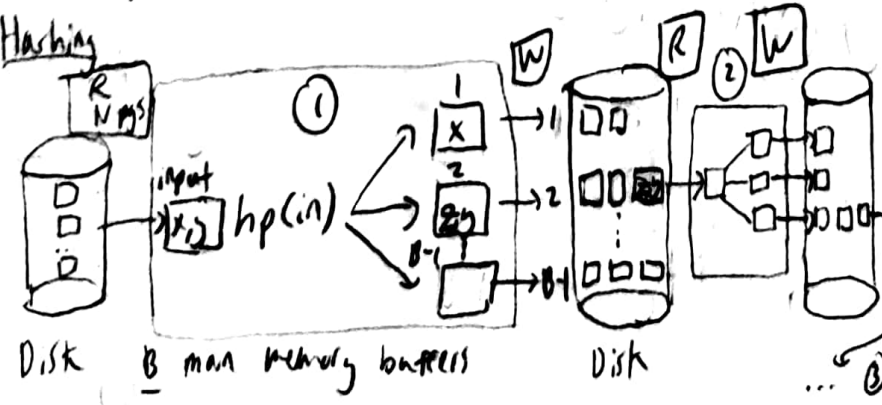| Rec 18 | Rec 6 | REC 5 | Rec 7 | ... | 18 6 5 7 □ |
|--------|-------|-------|-------|-----|-----|

footer
free space pointer

Delete: remove footer pointer (set record slot to null)
Insert: First empty pointer slot of free space pointer
NOTE: footer stores one int and one pointer for EACH
 record length

- Fixed length:

| 4 | 8 | 1 | 4 | 7 |
|---|---|---|---|---|

int double bool     int char(7)   → store 24 bytes

- Variable length:

(Primary Key)  (Not NULL)

| Header | □□ | ( ) | Varchar(10) | int | double | int | text | Varchar(B) |
|--------|----|-----|-------------|-----|--------|-----|------|------------|

- 2 pointers to variable length attribute
- bitmap (2 bits) where int and text are NULL
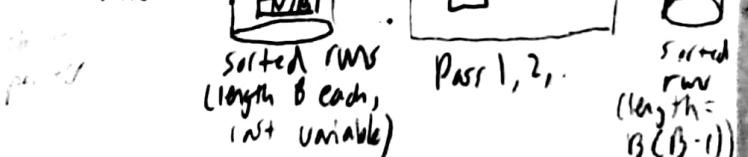
## Hashing



Disk    B main memory buffers         Disk

## Index Files (B+ trees) – acts similar to B-tree

- Property: # entries per node, $d = order \leq entries$
  $$\# leaves = (2d+1)^h \quad 2d$$
- Alternative 1: record contents stored in the leaf node itself
- Alternative 2: Leaf nodes: $\langle K, rid \text{ of matching record} \rangle$
- Alternative 3: Leaf nodes: $\langle K, \text{list of rids of records} \rangle$

- Property 1: leaf node entries: $d = order \leq entries \leq 2d$
- Property 2: all leaves same dist from root
- Property 3: inner node w/ k keys have k+1 children
- Leaf split     insert(8*)



not kept if inner node

- Bulk Load
 - Input sorted record keys
   (1*, 2*, 3*, ...)
 - Fill leaf page to fill factor > d
 - Update inner pages to full
 - Follow leaf and inner page split rules
   Sorting

- General External Merge Sort



Disk pass 0       sorted runs        Pass 1,2,       sorted
                  (length B each,                    run
                   int variable)                     (length =
                                                      B(B-1))

 - N pages to sort, B buffer pages
 - Number of passes: $1 + \lceil \log_{B-1} \lceil N/B \rceil \rceil$
 - Total I/Os = (I/Os per pass) × (# passes)
   $$= 2N \times (1 + \lceil \log_{B-1} \lceil N/B \rceil \rceil)$$
 - memory: $B(B-1)$ after 2 passes
   ①

Divide: split pages into B-1 partitions
  using hp hash function
  ②

Conquer: Repeat for big partitions
  (> B pages), rehash w/ new
  hash function $hp_i$

③ When partition fits in memory
  (<= B pages), read + write to
  check (Build Pass)

JOINS [R] = # pages in R, $p_R$ = # records per page in R, |R| = # records ~~per~~

**• Simple Nested Loop join on θ:**
```
for record r in R:          } cost = scan R pages +
   for record s in S:       }   scan S pages per record
      if θ(r,s):            }   = [R] + |R|[S]
         add <r,s> to result buffer
```

**• Block Nested Loop Join**
```
for rchunk of B-2 pages of R:   } cost =
   for spage of S:             } [R] + ⌈[R]/(B-2)⌉[S]
      add matching tuples to res buffer
```

**• Sort-Merge Join**
① Sort R and S by join key
② Join: merge-scan sorted partitions, emit matching tuples
```
mark = NULL
while r in range:
   if not mark:
      while r<S, advance r
      while r>S, advance s
      mark = S   # start of block of s
   if r == s:
      result = <r,S>
      advance s
      return result
   reset s to mark
   advance r
   mark = NULL
```

r →

| sid | sname |
|-----|-------|
| 22  | dustin |
| 28  | allen |
| ?1  | john |
| 31  | joe |
| 44  | mama |
| 58  | rusty |

S →

| sid | bid |
|-----|-----|
| 28  | 10001 |
| 28  | 1 0  |
| ?1  | 1010 |
| 31  | 1000l |
| 42  | 0 l0l1 |
| 58  | 0 0110 |

COST:
best: Sort R + Sort S + [R] + [S]
worst: Sort R + Sort S + |R|[S]
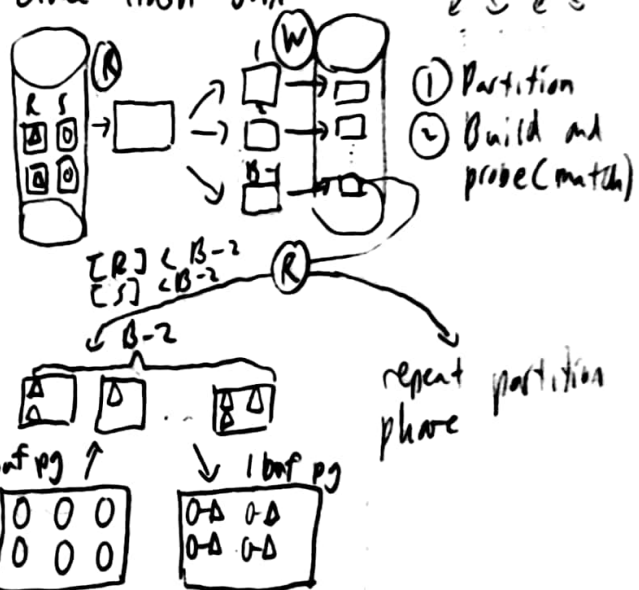
**• Pages Nested Loop Join**
```
for rpage in R:        } cost = [R] + [R][S]
   for spage in S:     }
      for rtuple in rpage
         for stuple in spage.
            if θ(r,s): add <r,s> result buffer
```

**• Index Nested Loops Join**
```
for each tuple r in R:
   for each tuple s in S, r==s:
      add <r,s> to res buffer   ⟵ lookup(r)
cost
= [R] + |R| (cost to find matching tuples in S)
```

**• Grace Hash Join**
① Partition
② Build and probe (match)

$[R] \leq B-2$
$[S] \leq B-2$

↓ B-2

repeat partition phase

1 buf pg ↑    ↓ 1 buf pg

---

**Relational Algebra**

Projection (π) ≡ SELECT
Selection (σ cond) ≡ WHERE
Renaming (ρ temp2(R2.sid → sid1, S1.sid → sid2) (R1×S1))
   output table name / old col → new col / old table

Union (U) ≡ tuples in r1 OR r2
Set-difference (−) ≡ tuples in r1 but not in r2
Cross-Product (×) ≡ each r1 row paired w/ each r2 row
Intersection (∩) ≡ S1 − (S1 − S2)
Theta Join ($\bowtie_\theta$) ≡ join on logical expression θ
Natural Join ($\bowtie$) ≡ equi-join on matching col names
Left Outer Join ($\bowtie$)
Right Outer Join ($\bowtie$)
Full Outer Join ($\bowtie$)
Group By ($\gamma_{age, count(*)>2}$) ≡ group by age, having count > 2

**Iterators**
• Materializing: save iterator as list on disk (costs I/O)
• Streaming: don't need to save on disk (SELECT, ×)

**Pushdown: equivalent statements**
• Selection (σ)
   cascade: $\sigma_{c_1 \wedge c_2 \wedge ... \wedge c_n}(R) \equiv \sigma_{c_1}(\sigma_{c_2}(...(\sigma_{c_n}(R))...))$
   commute: $\sigma_{c_1}(\sigma_{c_2}(R)) \equiv \sigma_{c_2}(\sigma_{c_1}(R))$
• Projections
   cascade: $\pi_{a_1}(R) \equiv \pi_{a_1}(\pi_{a_2}(\pi_{a_3, a_4 ...}(R)))$
• Cartesian Product
   associative: $R \times (S \times T) \equiv (R \times S) \times T$
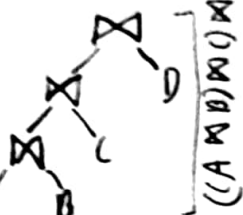   commutative: $R \times S \equiv S \times R$

# Query Optimization

## ① Plan Space
- Try diff queries (equivalent relational algebra)
- Try different types of joins (brute, hash, block)
- Heuristics
  - Left-deep trees only
  - Avoid cartesian products
  - Selection/projection pushdown



$((A \bowtie B) \bowtie C) \bowtie D$

## ② Cost Estimation (cost = #I/O + CPU factor * #tuples)
- Catalog: containing tuple statistics (updated periodically)
- Selectivity (sel) = |output|/|input| of $\sigma$ term
  - $\sigma$ term1 ∧ term2 ∧ term3
  - reflects impact of $\sigma$ term


$\sigma_{col=1}$ ⟹ , sel $= \frac{1}{5}$

| Predicate | Sel | | Assume |
|---|---|---|---|
| $c=v$ | 1/distinct c | | $|c|$ |
| $c=v$ | 1/10 | | |
| $c_1 = c_2$ | 1/MAX(distinct c1, distinct c2) | $|c_1|, |c_2|$ | |
| $c_1 = c_2$ | 1/distinct $c_i$ | $|c_i|$ only | |
| $c_1 = c_2$ | 1/10 | | |
| $c < v$ | $(v-\min(c))/(\max(c)-\min(c)+1)$ | $\max(c), \min(c)$ | |
| $c > v$ | $(\max(c)-v)/(\max(c)-\min(c)+1)$ | $\max(c), \min(c)$ | |
| $c < v$ | 1/10 | | |
| $c > v$ | 1/10 | |  ⎱ $c=$ int |
| $c <= v$ | $(v-\min(c))/(\max(c)-\min(c)+1)+\frac{1}{|c|}$ | $\max(c), \min(c)$ | |
| $c >= v$ | $(\max(c)-v)/(\max(c)-\min(c)+1)+\frac{1}{|c|}$ | $\max(c), \min(c)$ | |
| $c <= v$ | 1/10 | | |
| $c >= v$ | 1/10 | | |
| $c >= v$ | $(\max(c)-v)/(\max(c)-\min(c))$ | $\max(c), \min(c)$ | ⎱ $c=$ float |
| $c >= v$ | 1/10 | | |
| $c <= v$ | $(v-\min(c))/(\max(c)-\min(c))$ | $\max(c), \min(c)$ | |
| $c <= v$ | 1/10 | | |
| $P_1 \cap P_2$ | $S(P_1) * S(P_2)$ | independent terms | |
| $P_1 \cup P_2$ | $S(P_1)+S(P_2)-S(P_1)*S(P_2)$ | independent terms | |
| NOT $p$ | $1-S(p)$ | | |

Lossiness: can we construct original relation?
- Decompose R into A and B
  Decomp is lossless iff $F^+$ contains:
  $X \cap Y \to X$  BCNF = always lossless
  $X \cap Y \to Y$

## ③ Search Algorithm (bot joins)
pass 1: find minimum cost access method
(index scan, full scan, etc.)
for each table, and interesting order
on interesting index

pass i:
- consider only: left deep and not cartesian product (unless all are cartesian)
- advance only: cheapest cost plan every subset group of relations and interesting order

### Interesting Orders
- ORDER BY attributes
- GROUP BY attributes
- downstream join attributes

## Functional Dependencies
- FD: $X \to Y$ ($X$ determines $Y$)
- Superkey: $X$ is a superkey of $R$ if $X \to [\text{all attributes of } R]$
- Candidate key: minimal super key of $R$ (no subset is a superkey)
- Armstrong's Axioms ($X, Y, Z$ are attribute sets)
  - Reflexivity: if $Y \subseteq X$, $X \to Y$
  - Augmentation: if $X \to Y$, $XZ \to YZ$ (not other way)
  - Transitivity: if $X \to Y$, and $Y \to Z$, then $X \to Z$
  - Union: if $X \to Y$ and $X \to Z$, then $X \to YZ$
  - Decomposition: if $X \to YZ$, then $X \to Y$ and $X \to Z$
- Closure: set of FDs $F^+$ s.t. $F \to F^+$
  - ex: $\{A \to B, A \to C\} \to \{A \to B, A \to C, (A \to (B\}$
- Attribute closure: set of attributes $X^+$ s.t. $X \to X^+$ is in $F^+$
- Closure algorithm: to find $X^+$
  ```
  closure = X
  repeat until closure does not change:
    for U → V in F:
      if U ⊆ closure: closure = closure ∪ V
  ```

## Normalization: BCNF Decomposition
- R is in BCNF if:
  for every FD $X \to A$ in R, $A \subseteq X$ or $X$ is superkey
- Decompose relation R not in BCNF into multiple relations in BCNF:
  ```
  for each X → Y in F+:
    if X → Y violates BCNF:
      Decompose R into (R−X+) ∪ X, X+
  ```

# Transactions and Concurrency

Transactions follow rules

- **Atomicity**: All operations in a transaction happens or no ops happen
- **Consistency**: Data starts and ends consistent
- **Isolation**: Execution of txns looks like running one at a time
- **Durability**: if txn commits, effects persist

Non-conflict swapping:

if two consecutive operations from two txns can be swapped if non-conflicting.

$T_1: R(A) \cdots W(A)$ } reverse ok OR $W(A) \cdots W(A)$

$T_2: \qquad W(A)$

↳ Conflicting

**Conflict equivalent**: Schedules w/ same txn operations, different orderings (non-conflict swap)

**Serializable**: a schedule $S$ is equivalent to a serial schedule

**Conflict Serializable**: schedule $S$ is conflict serializable to some other serial schedule
($S \rightarrow$ Serial)

**Dependency graph**:
- one node per transaction
- If an operation in $T_i$ conflicts w/ an operation in $T_j$ and $T_i$ comes first, add edge from $T_i$ to $T_j$.
- Schedule is conflict serializable IFF dependency graph is acyclic

| $T_1$ | $R(A)$ $R(B)^*$ | | $R(A)$ |
|---|---|---|---|
| $T_2$ | | $R(B)$ $W(B)^*$ | |

$(T_1) \rightarrow (T_2)$

**View Serializable**: same as conflict serializable but calls $W(A) \cdots W(A)$, $R(A) \cdots W(A)$ non-conflicting, only $W(A) \cdots R(A)$ is conflicting
- $(S \rightarrow VS)$

## Simple Locking

| | S | X |
|---|---|---|
| S | ✓ | X |
| X | X | X |

on a resource

- An S lock lets a txn read a resource
  - many txns can hold S locks on a resource at once
- An X lock lets a txn modify a resource
  - No other txn can have any type of lock while a txn has X lock
- If a txn can't get a lock it wants, it blocks and waits until another txn releases the conflicting lock

## Deadlocks

$T_1: R(A), W(B)$ } $T_1$ waits for $T_2$'s lock on
$T_2: R(B), W(A)$ } $B$, while $T_2$ waits for $T_1$'s lock on $A$

- **Priority**: older txns = higher priority
- **wait-die**: if $T_i$ wants $T_j$'s lock
  - if $T_i$ higher priority, wait for $T_j$ to release
  - if $T_i$ lower priority, $T_i$ aborts
- **wound-wait**: if $T_i$ wants $T_j$'s lock
  - if $T_i$ higher priority, $T_j$ aborts ("wound")
  - if $T_i$ lower priority, $T_j$ waits
- **"waits-for" graph**
  - One node for each txn
  - If $T_j$ waits for $T_i$, edge from $T_j$ to $T_i$
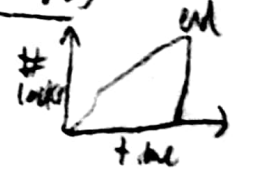  - Cycle means deadlock

## 2 Phase Locking (2PL)

- A txn may not acquire a lock after releasing any lock
- Conflict serializable guaranteed

[graph: # locks held vs time, rising then falling]

## Strict 2-Phase Locking (Strict 2PL)

- Same as 2PL except all locks released at end of txn
- Avoids cascading aborts

[graph: # locks held vs time, rising then dropping at end]

## Multi-granularity Locking

| | NL | IS | IX | SIX | S | X |
|---|---|---|---|---|---|---|
| NL | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| IS | ✓ | ✓ | ✓ | ✓ | ✓ | X |
| IX | ✓ | ✓ | ✓ | X | X | X |
| SIX | ✓ | ✓ | X | X | X | X |
| S | ✓ | ✓ | X | X | X | X |
| X | ✓ | X | X | X | X | X |

- To get S or IS lock on a node, need IS or IX on parent node
- To get X, IX, or SIX on a node, need IX or SIX on parent node
- Assume strict 2PL unless interfere w/ txn's new operation

# Recovery

**STEAL:** txns can flush pages w/ uncommitted updates
Pro: Can maximize use of buffer pages
Con: Must undo if program crashes before commit
**FORCE:** force dirty pages to disk before commit ⟶
**NO-STEAL:** txns cannot flush pages w/ uncomm changes
**NO-FORCE:** dirty pages allowed to not be forced on disk upon commit

## ARIES (STEAL, NO FORCE)

- Writes
  UPDATE log record for every write
- Commit:
  Write COMMIT log record, flush log to disk, "committed" is status, release txn locks, write END log record, remove from txn table
- Abort (see UNDO)
  write ABORT log record, write CLR record for each UPDATE the txn logged in reverse order, release txn locks, write END log record, remove from txn table
- Assume strict 2-phase locking and write-ahead logging
- Data structures
  - Transaction Table: information on currently active transaction (running, committing, aborting)
  - Dirty page table: dirty pages in memory not yet written to memory ⟶

① **Analysis:** identify changes not written to disk and active txn at time of crash
- Start at last BEGIN-CHECKPOINT and scan all following log records
- For each record:
  if record is UPDATE/CLR:
   if txn not in ATT: add it (Txn ID, Status, lastLSN)
   if txn in ATT: update lastLSN
  if record is ABORT/COMMIT: update ATT (Status, lastLSN)
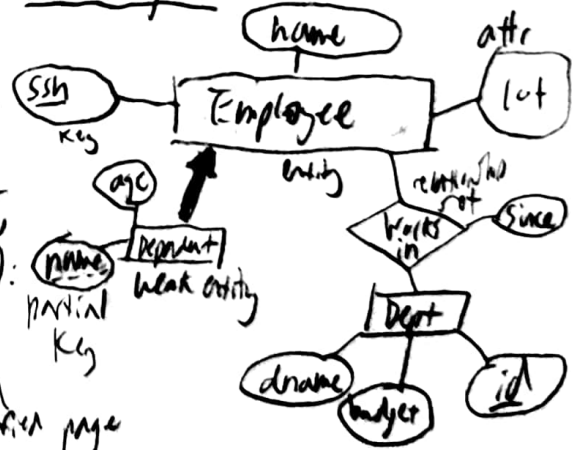  if record is END: remove transaction from ATT

② **REDO:** Repeat all actions to restore database state to what it was at the time of the crash from the smallest recLSN
- For each UPDATE/CLR beginning from min (all recLSN in DPT):
  REDO operation UNLESS any of the following is true:
   - affected page not in DPT
   - recLSN of page in DPT > LSN
   - pageLSN of page (in DB disk) ≥ LSN    pageLSN = LSN of last record that modified page

③ **UNDO:** undo the actions of transactions that did not commit
- For all "loser" (active) transactions:
  undoLSN = lastLSN
  while undoLSN != NULL:              undo = add record to CLR
   if record at undoLSN is UPDATE:
    undo record w/ undoNextLSN = record's prevLSN
    undoLSN = prevLSN
  add END record for this txn to CLR

(left margin bracket: Abort)

---

## Write ahead logging

- txn not committed until log records for changes written to disk
- changes must be logged before data modified on disk

### Undo Logging (Atomicity)    updated element (dirty page)    old value
⟨Start T⟩ ⟨Commit T⟩ ⟨Abort T⟩ ⟨T, X, v⟩
- If T commits, then FLUSH(x) must be written to disk before ⟨Commit T⟩
- If T modifies X, then ⟨T, X, v⟩ written to log before FLUSH(X)

### Redo Logging (Durability)    updated element (dirty pg)    new value
⟨Start T⟩ ⟨Commit T⟩ ⟨Abort T⟩ ⟨T, X, v⟩
- If T modifies X, both ⟨T, X, v⟩ and ⟨Commit T⟩ must be logged before FLUSH(X)

| Undo | Redo |
|------|------|
| flush(X) | Commit(txn) |
| Flush(Y) | flush(X) |
| commit(txn) | flush(Y) |

### ATT

| Txn | LastLSN | Status |
|-----|---------|--------|
| T1 | 10 | running |
| T2 | 20 | committing |
| T3 | 40 | running |
| T4 | 50 | Aborting |

### DPT

| Page | recLSN |
|------|--------|
| P1 | 0 |
| P3 | 10 |

## ER diagrams



- Key constraint ⟶ participates in at most one
- Participation constraint ▬▬ participates in at least one
- Key constraint with total participation exactly one ➡
- Non-key partial participation 0 or more