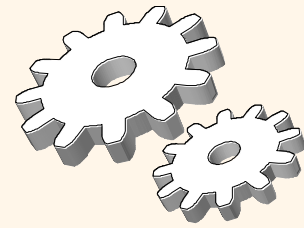


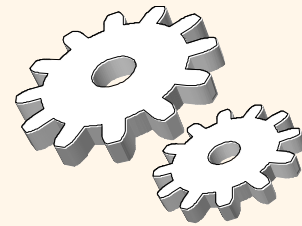
Crash Recovery

Chapter 18



Review: The ACID properties

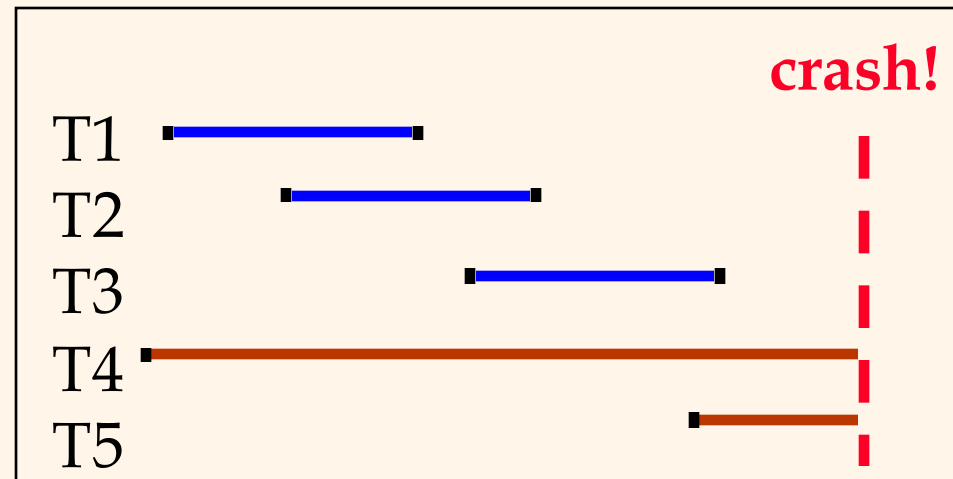
- ❖ **A**tomicity: All actions in the Xact happen, or none happen.
- ❖ **C**onsistency: If each Xact is consistent, and the DB starts consistent, it ends up consistent.
- ❖ **I**solation: Execution of one Xact is isolated from that of other Xacts.
- ❖ **D**urability: If a Xact commits, its effects persist.
- ❖ The **Recovery Manager** guarantees Atomicity & Durability.

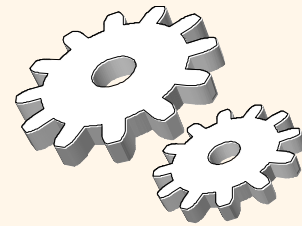


Motivation

- ❖ Atomicity:
 - Transactions may abort (“Rollback”).
- ❖ Durability:
 - What if DBMS stops running? (Causes?)

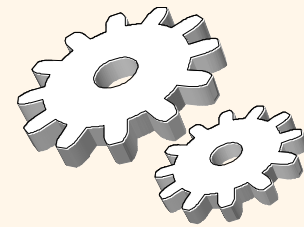
- Desired Behavior after system restarts:
 - T1, T2 & T3 should be durable.
 - T4 & T5 should be aborted (effects not seen).





Assumptions

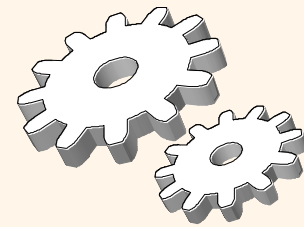
- ❖ Concurrency control is in effect.
 - **Strict 2PL**, in particular.
- ❖ Updates are happening “in place”.
 - i.e. data is overwritten on (deleted from) the disk.
- ❖ A simple scheme to guarantee Atomicity & Durability?



Handling the Buffer Pool

- ❖ **Force** every write to disk?
 - Poor response time.
 - But provides durability.
- ❖ **Steal** buffer-pool frames from uncommitted Xacts?
 - If not, poor throughput.
 - If so, how can we ensure atomicity?

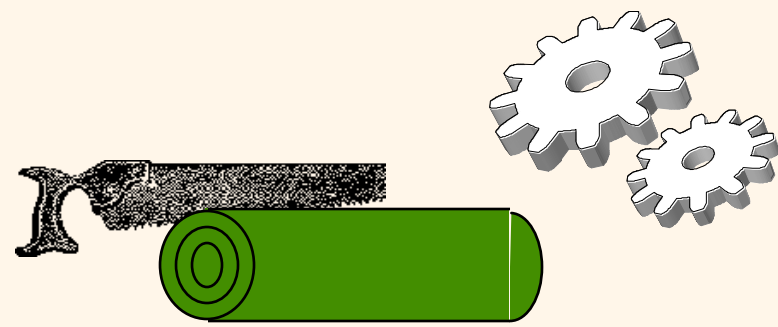
	No Steal	Steal
Force	Trivial	
No Force		Desired



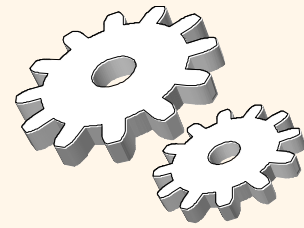
More on Steal and Force

- ❖ **STEAL** (why enforcing Atomicity is hard)
 - *To steal frame F*: Current page in F (say P) is written to disk; some Xact holds lock on P.
 - What if the Xact with the lock on P aborts?
 - Must remember the old value of P at steal time (to support **UNDO**ing the write to page P).
- ❖ **NO FORCE** (why enforcing Durability is hard)
 - What if system crashes before a modified page is written to disk?
 - Write as little as possible, in a convenient place, at commit time, to support **REDO**ing modifications.

Basic Idea: Logging



- ❖ Record REDO and UNDO information, for every update, in a *log*.
 - Sequential writes to log (put it on a separate disk).
 - Minimal info (diff) written to log, so multiple updates fit in a single log page.
- ❖ Log: An ordered list of REDO/UNDO actions
 - Log record contains:
 - <XID, pageID, offset, length, old data, new data>
 - and additional control info (which we'll see soon).

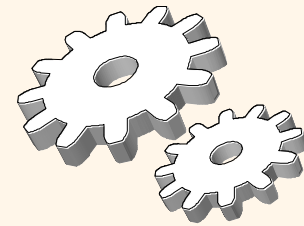
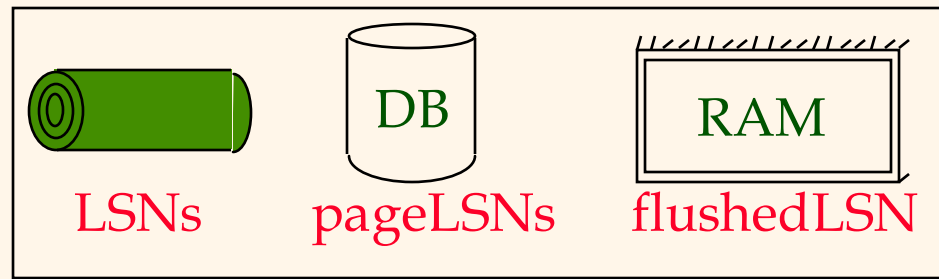


Write-Ahead Logging (WAL)

- ❖ The Write-Ahead Logging Protocol:
 - ① Must **force** the **log record** for an update before the corresponding **data page** gets to disk.
 - ② Must **write all log records** for a Xact before commit.
- ❖ #1 guarantees Atomicity.
- ❖ #2 guarantees Durability.

- ❖ Exactly how is logging (and recovery!) done?
 - We'll study the ARIES algorithms.

WAL & the Log



❖ Each log record has a unique **Log Sequence Number (LSN)**.

- LSNs always increasing.

❖ Each data page contains a **pageLSN**.

- The LSN of the most recent *log record* for an update to that page.

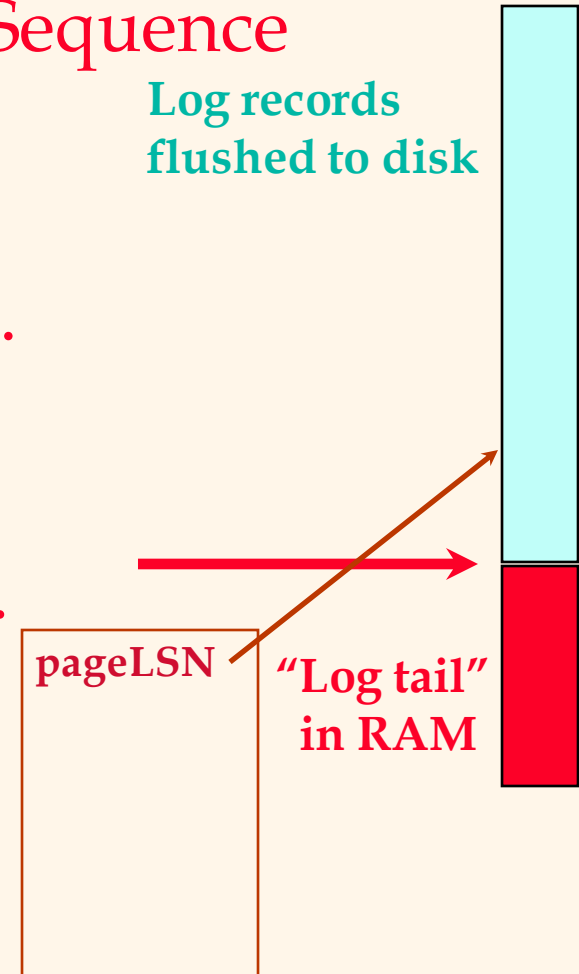
❖ System keeps track of **flushedLSN**.

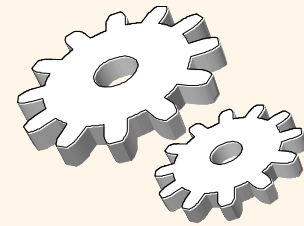
- The max LSN flushed so far.

❖ WAL: *Before* a page is written,

- $\text{pageLSN} \leq \text{flushedLSN}$

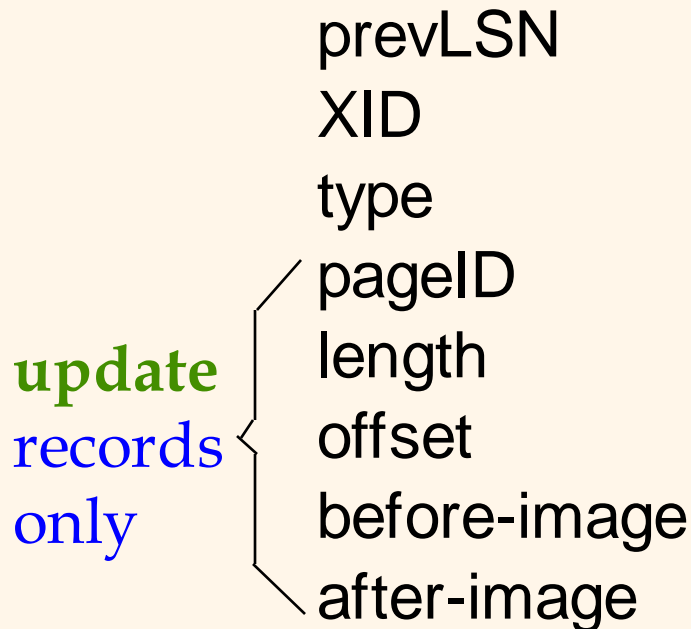
Log records
flushed to disk





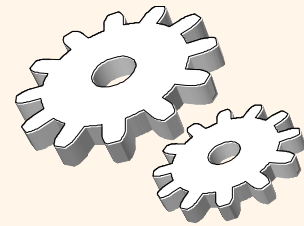
Log Records

LogRecord fields:



Possible log record types:

- ❖ **Update**
- ❖ **Commit**
- ❖ **Abort**
- ❖ **End** (signifies end of commit or abort)
- ❖ **Compensation Log Records (CLRs)**
 - for UNDO actions



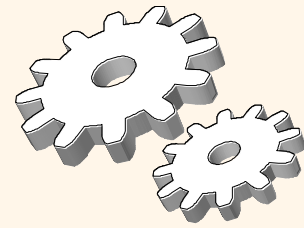
Other Log-Related State

❖ Transaction Table:

- One entry per active Xact.
- Contains **XID**, **status** (running/committed/aborted), and **lastLSN**.

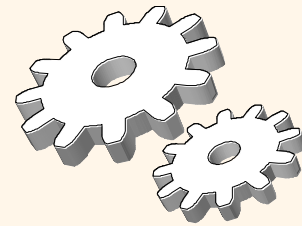
❖ Dirty Page Table:

- One entry per dirty page in buffer pool.
- Contains **recLSN** -- the LSN of the log record which *first* caused the page to be dirty.



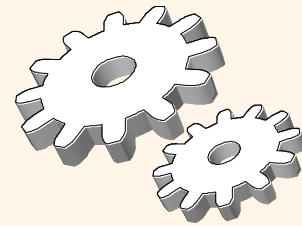
Normal Execution of an Xact

- ❖ Series of **reads & writes**, followed by **commit or abort**.
 - We will assume that write is atomic on disk.
 - In practice, additional details to deal with non-atomic writes.
- ❖ **Strict 2PL**.
- ❖ **STEAL, NO-FORCE** buffer management, with **Write-Ahead Logging**.



Checkpointing

- ❖ Periodically, the DBMS creates a checkpoint, in order to minimize the time taken to recover in the event of a system crash. Write to log:
 - **begin_checkpoint** record: Indicates when chkpt began.
 - **end_checkpoint** record: Contains current *Xact table* and *dirty page table*. This is a **'fuzzy checkpoint'**:
 - Other Xacts continue to run; so these tables accurate only as of the time of the **begin_checkpoint** record.
 - No attempt to force dirty pages to disk; effectiveness of checkpoint limited by oldest unwritten change to a dirty page. (So it's a good idea to periodically flush dirty pages to disk!)
 - Store LSN of chkpt record in a safe place (*master* record).



The Big Picture: What's Stored Where



LogRecords

prevLSN
XID
type
pageID
length
offset
before-image
after-image



Data pages

each
with a
pageLSN

master record



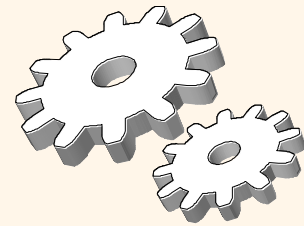
Xact Table

lastLSN
status

Dirty Page Table

recLSN

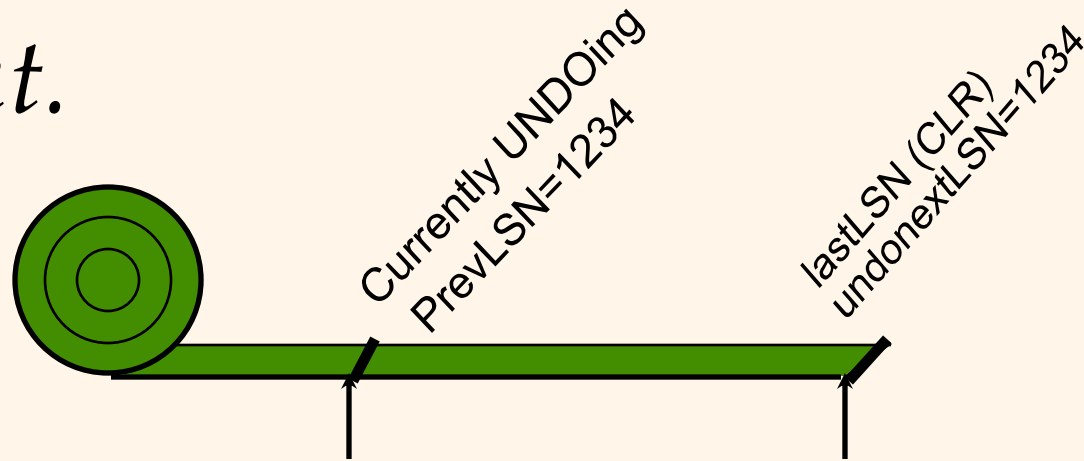
flushedLSN



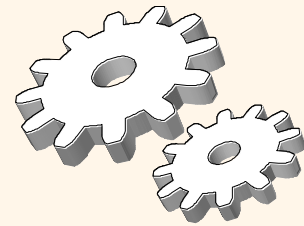
Simple Transaction Abort

- ❖ For now, consider an explicit abort of a Xact.
 - No crash involved.
- ❖ We want to “play back” the log in reverse order, UNDOing updates.
 - Get **lastLSN** of Xact from Xact table.
 - Can follow chain of log records backward via the **prevLSN** field.
 - Before starting UNDO, write an **Abort log record**.
 - For recovering from crash during UNDO!

Abort, cont.

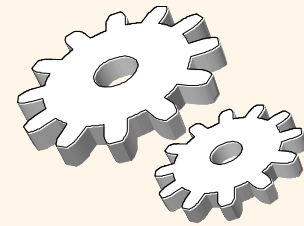


- ❖ To perform UNDO, must have a lock on data!
 - No problem!
- ❖ Before restoring old value of a page, write a CLR:
 - You continue logging while you UNDO!!
 - CLR has one extra field: **undonextLSN**
 - Points to the next LSN to undo (i.e. the prevLSN of the record we're currently undoing).
 - CLR's *never* Undone (but they might be Redone when repeating history: guarantees Atomicity!)
- ❖ At end of UNDO, write an "end" log record.

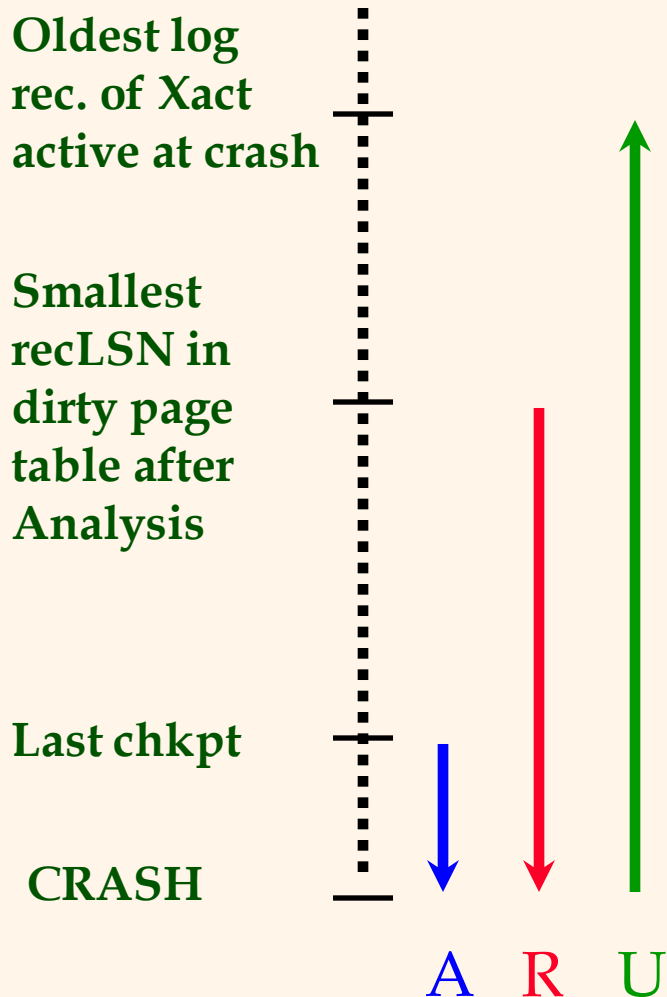


Transaction Commit

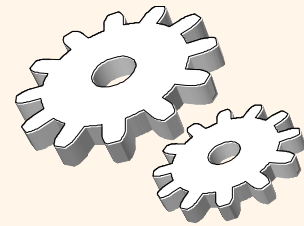
- ❖ Write **commit** record to log.
- ❖ All log records up to Xact's **lastLSN** are flushed.
 - Guarantees that **flushedLSN** \geq **lastLSN**.
 - Note that log flushes are sequential, synchronous writes to disk.
 - Many log records per log page.
- ❖ Commit() returns.
- ❖ Write **end** record to log.



Crash Recovery: Big Picture

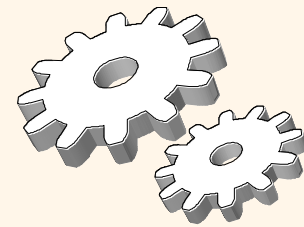


- Start from a **checkpoint** (found via **master** record).
- Three phases. Need to:
 - Figure out which Xacts committed since checkpoint, which failed (**Analysis**).
 - **REDO** *all* actions.
 - (repeat history)
 - **UNDO** effects of failed Xacts.



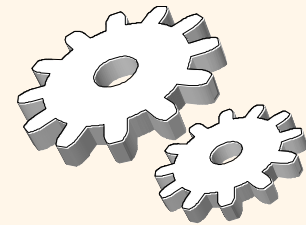
Recovery: The Analysis Phase

- ❖ Reconstruct state at checkpoint.
 - via **end_checkpoint** record.
- ❖ Scan log forward from checkpoint.
 - **End** record: Remove Xact from Xact table.
 - **Other records**: Add Xact to Xact table, set **lastLSN=LSN**, change Xact status on **commit**.
 - **Update** record: If P not in Dirty Page Table,
 - Add P to D.P.T., set its **recLSN=LSN**.



Recovery: The REDO Phase

- ❖ We *repeat History* to reconstruct state at crash:
 - Reapply *all* updates (even of aborted Xacts!), redo CLR's.
- ❖ Scan forward from log rec containing smallest *recLSN* in D.P.T. For each CLR or update log rec *LSN*, REDO the action unless:
 - Affected page is not in the Dirty Page Table, or
 - Affected page is in D.P.T., but has *recLSN* > *LSN*, or
 - *pageLSN* (in DB) \geq *LSN*.
- ❖ To REDO an action:
 - Reapply logged action.
 - Set *pageLSN* to *LSN*. No additional logging!



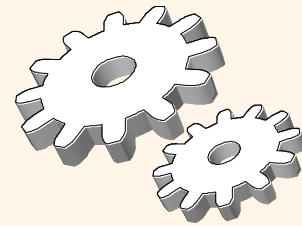
Recovery: The UNDO Phase

ToUndo = { l | l a lastLSN of a “loser” Xact }

Repeat:

- Choose largest LSN among ToUndo.
- If this LSN is a CLR and undonextLSN == NULL
 - Write an End record for this Xact.
- If this LSN is a CLR, and undonextLSN != NULL
 - Add undonextLSN to ToUndo
- Else this LSN is an update. Undo the update, write a CLR, add prevLSN to ToUndo.

Until ToUndo is empty.



Example of Recovery



Xact Table

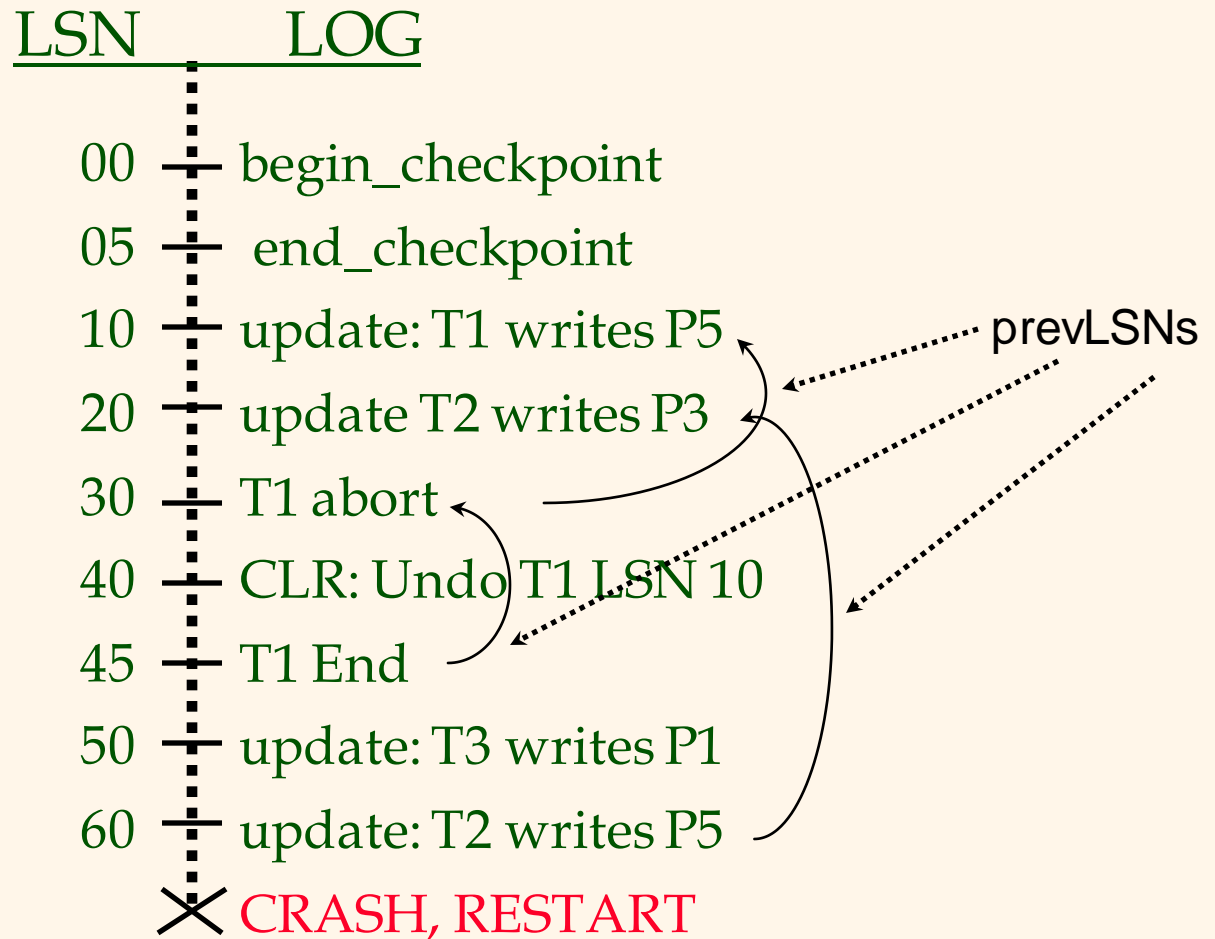
lastLSN
status

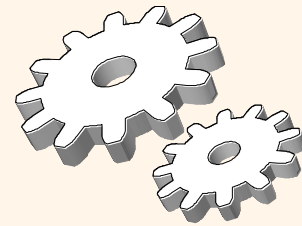
Dirty Page Table

recLSN

flushedLSN

ToUndo





Example: Crash During Restart!



Xact Table

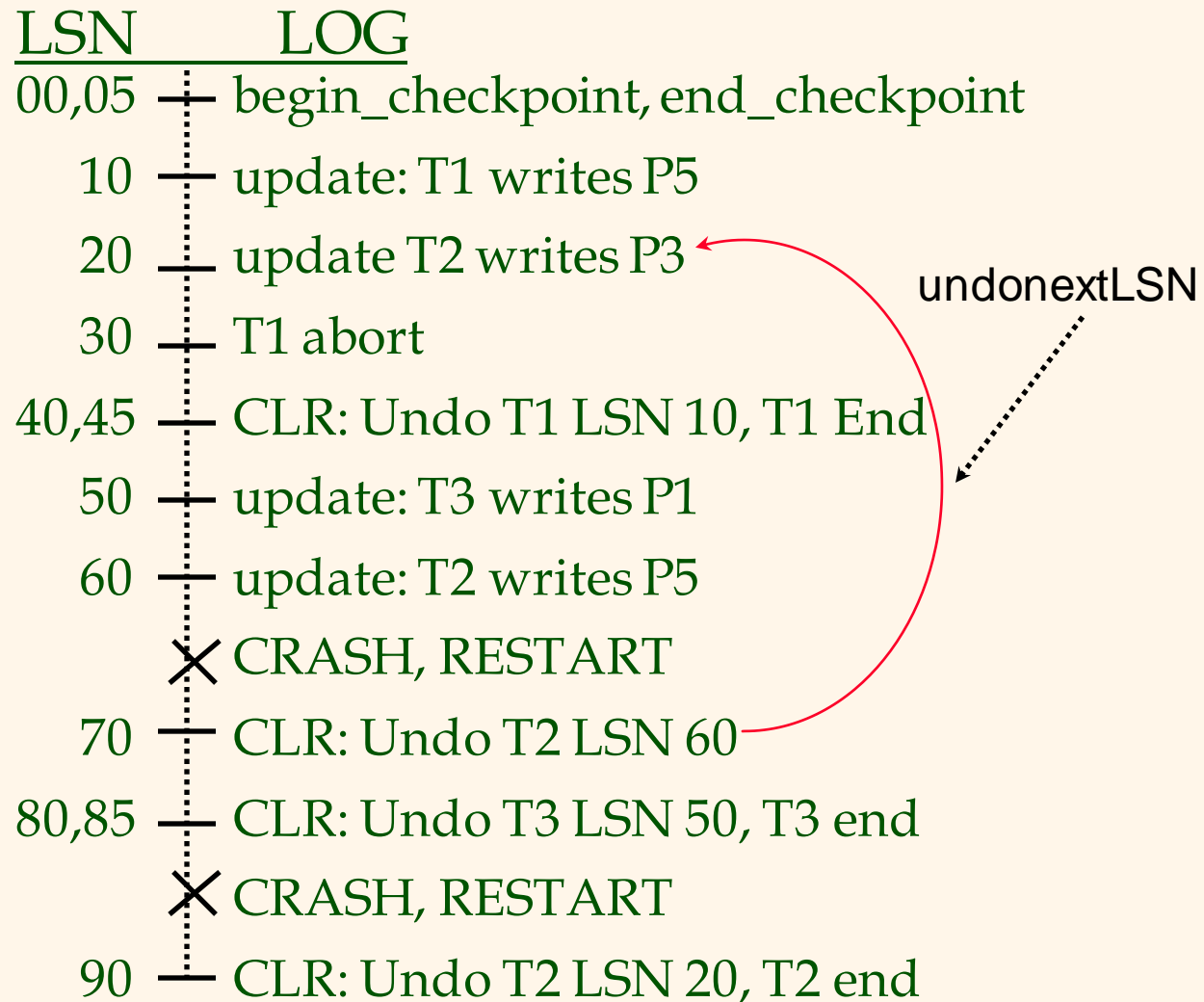
lastLSN
status

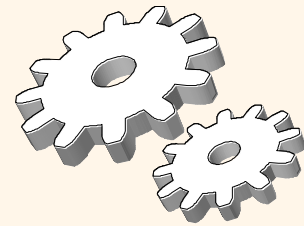
Dirty Page Table

recLSN

flushedLSN

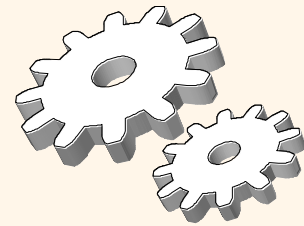
ToUndo





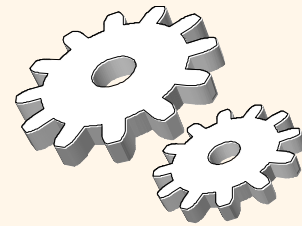
Additional Crash Issues

- ❖ What happens if system crashes during Analysis? During REDO?
- ❖ How do you limit the amount of work in REDO?
 - Flush asynchronously in the background.
 - Watch “hot spots”!
- ❖ How do you limit the amount of work in UNDO?
 - Avoid long-running Xacts.



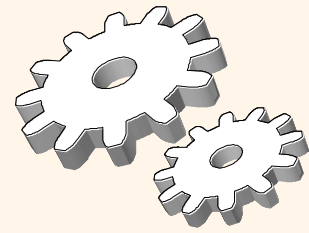
Summary of Logging/Recovery

- ❖ **Recovery Manager** guarantees Atomicity & Durability.
- ❖ Use WAL to allow STEAL/NO-FORCE w/o sacrificing correctness.
- ❖ LSNs identify log records; linked into backwards chains per transaction (via prevLSN).
- ❖ pageLSN allows comparison of data page and log records.



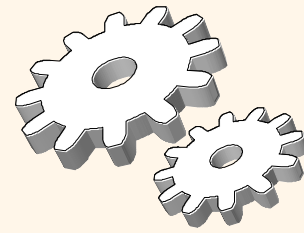
Summary, Cont.

- ❖ **Checkpointing:** A quick way to limit the amount of log to scan on recovery.
- ❖ Recovery works in 3 phases:
 - **Analysis:** Forward from checkpoint.
 - **Redo:** Forward from oldest recLSN.
 - **Undo:** Backward from end to first LSN of oldest Xact alive at crash.
- ❖ Upon Undo, write CLR's.
- ❖ Redo “repeats history”: Simplifies the logic!



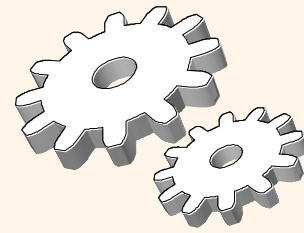
Hash-Based Indexes

Chapter 11



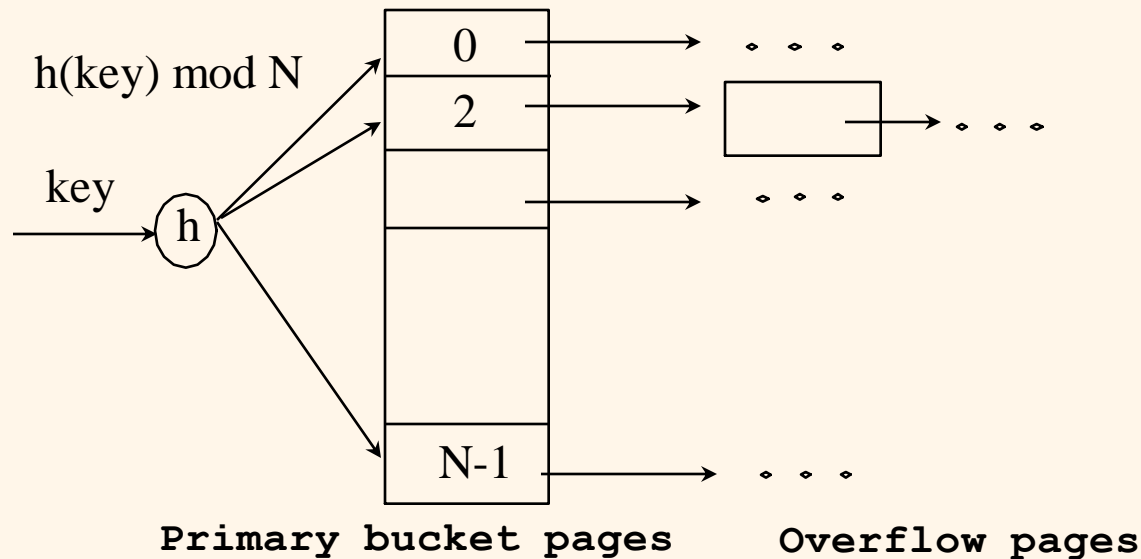
Introduction

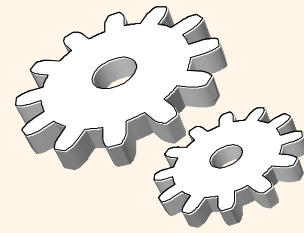
- ❖ *As for any index, 3 alternatives for data entries \mathbf{k}^* :*
 - Data record with key value \mathbf{k}
 - $\langle \mathbf{k}, \text{rid of data record with search key value } \mathbf{k} \rangle$
 - $\langle \mathbf{k}, \text{list of rids of data records with search key } \mathbf{k} \rangle$
 - Choice orthogonal to the *indexing technique*
- ❖ *Hash-based* indexes are best for *equality selections*.
Cannot support range searches.
- ❖ Static and dynamic hashing techniques exist;
trade-offs similar to ISAM vs. B+ trees.



Static Hashing

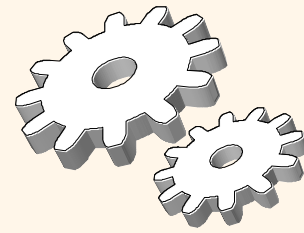
- ❖ # primary pages fixed, allocated sequentially, never de-allocated; overflow pages if needed.
- ❖ $h(k) \bmod M =$ bucket to which data entry with key k belongs. ($M = \#$ of buckets)





Static Hashing (Contd.)

- ❖ Buckets contain *data entries*.
- ❖ Hash fn works on *search key* field of record r . Must distribute values over range $0 \dots M-1$.
 - $h(key) = (a * key + b)$ usually works well.
 - a and b are constants; lots known about how to tune h .
- ❖ **Long overflow chains** can develop and degrade performance.
 - *Extendible* and *Linear Hashing*: Dynamic techniques to fix this problem.

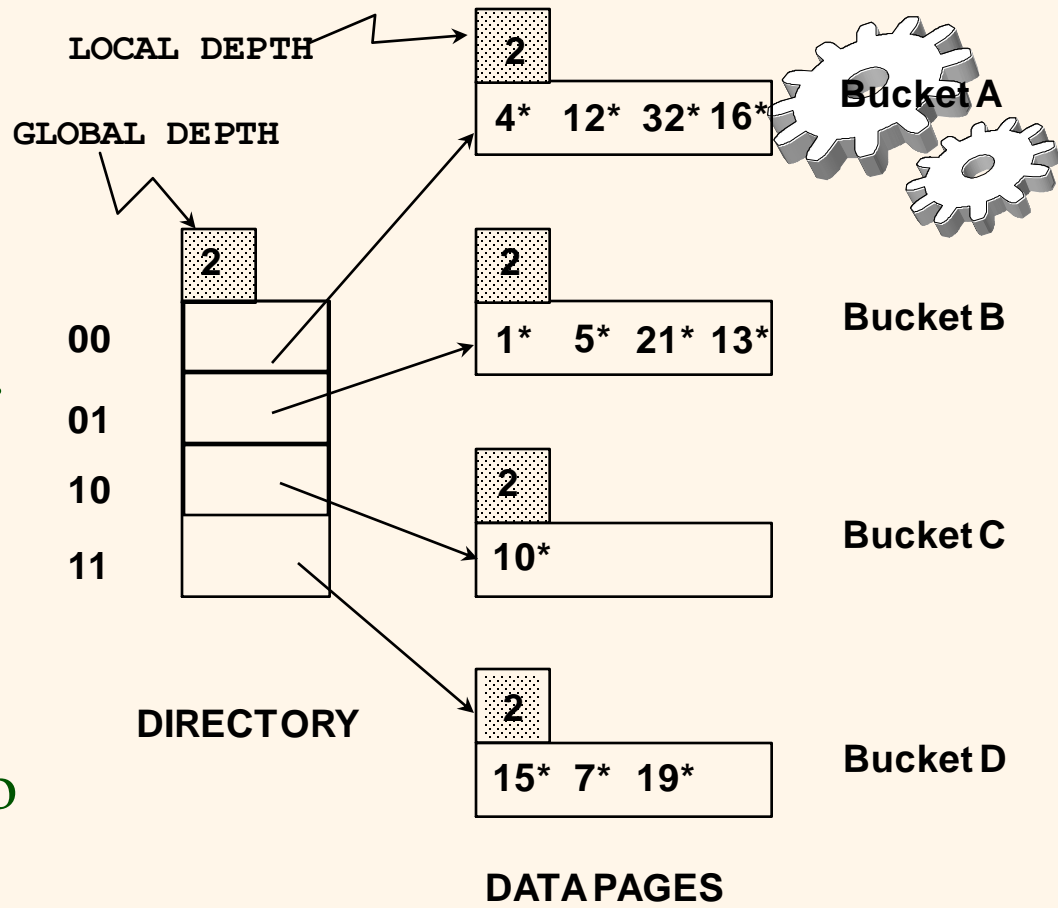


Extendible Hashing

- ❖ Situation: Bucket (primary page) becomes full. Why not re-organize file by *doubling* # of buckets?
 - Reading and writing all pages is expensive!
 - Idea: Use directory of pointers to buckets, double # of buckets by *doubling the directory*, splitting just the bucket that overflowed!
 - Directory much smaller than file, so doubling it is much cheaper. Only one page of data entries is split.
No overflow page!
 - Trick lies in how hash function is adjusted!

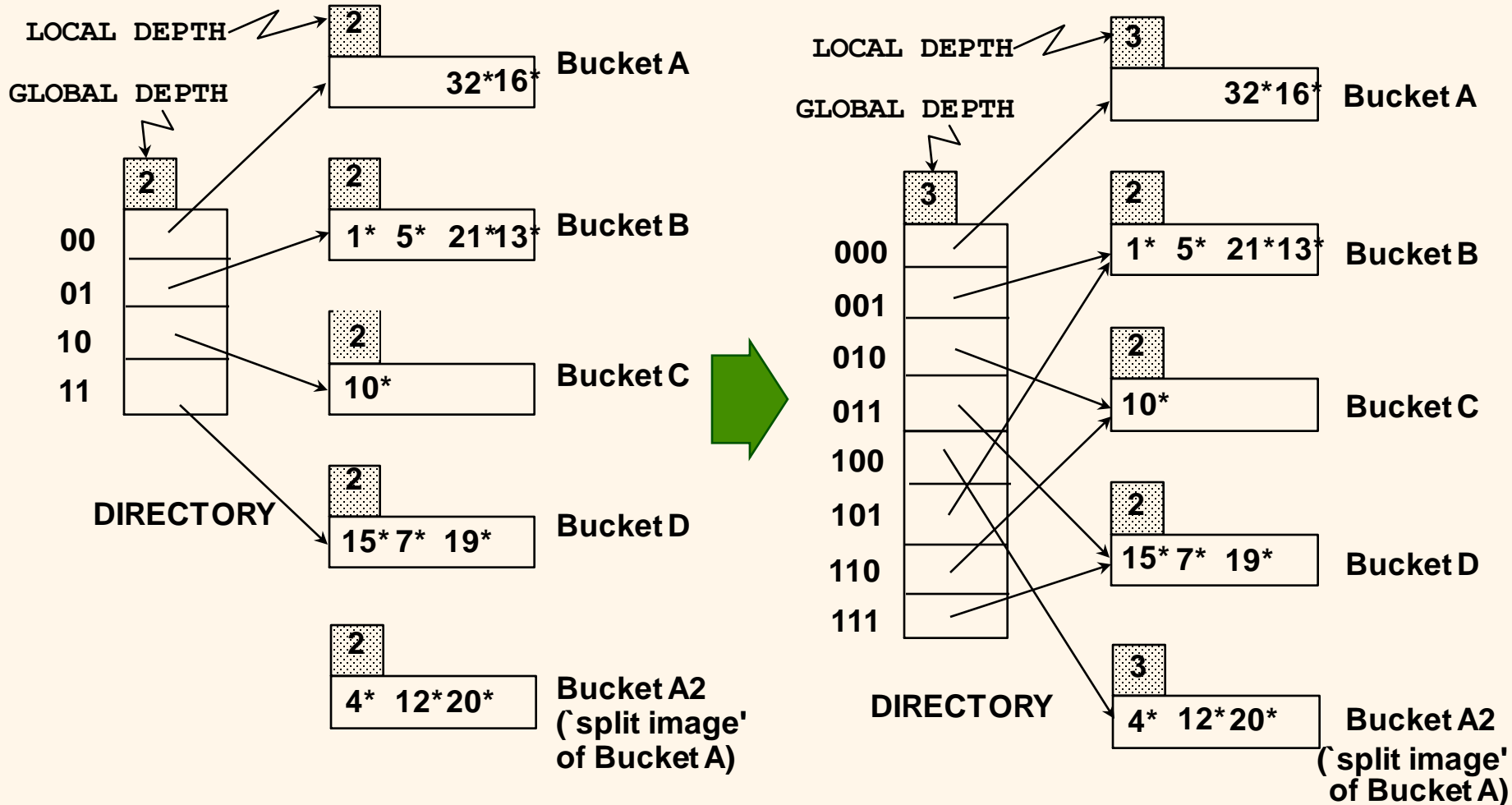
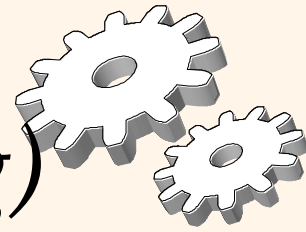
Example

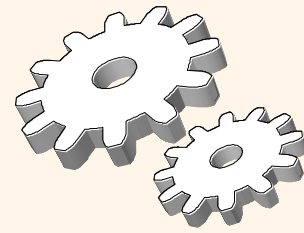
- ❖ Directory is array of size 4.
- ❖ To find bucket for r , take last '*global depth*' # bits of $\mathbf{h}(r)$; we denote r by $\mathbf{h}(r)$.
 - If $\mathbf{h}(r) = 5 = \text{binary } 101$, it is in bucket pointed to by 01.



- ❖ **Insert:** If bucket is full, *split* it (allocate new page, re-distribute).
- ❖ If necessary, double the directory. (As we will see, splitting a bucket does not always require doubling; we can tell by comparing *global depth* with *local depth* for the split bucket.)

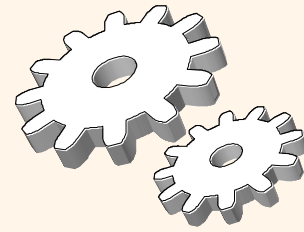
Insert $h(r)=20$ (Causes Doubling)





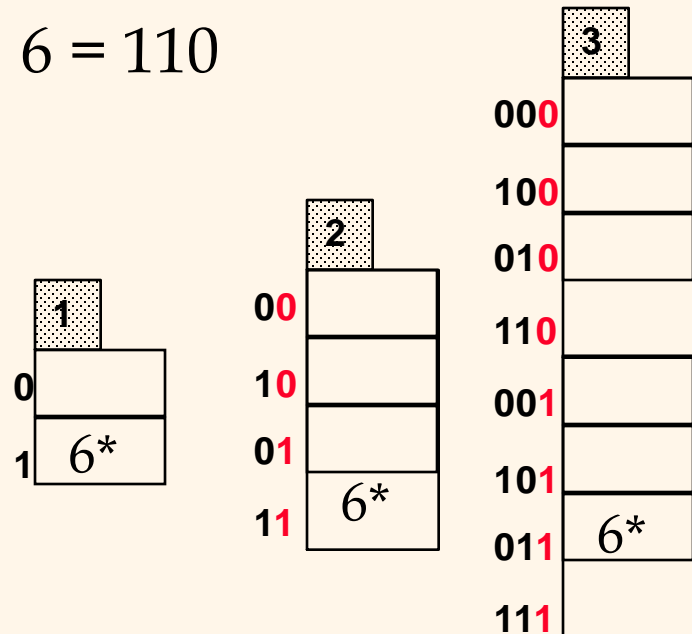
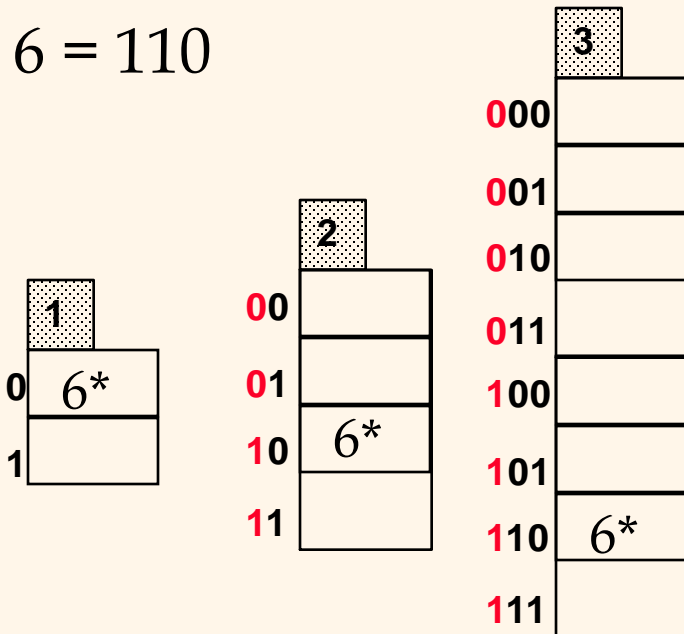
Points to Note

- ❖ 20 = binary 10100. Last 2 bits (00) tell us r belongs in A or A2. Last 3 bits needed to tell which.
 - *Global depth of directory*: Max # of bits needed to tell which bucket an entry belongs to.
 - *Local depth of a bucket*: # of bits used to determine if an entry belongs to this bucket.
- ❖ When does bucket split cause directory doubling?
 - Before insert, *local depth* of bucket = *global depth*. Insert causes *local depth* to become $>$ *global depth*; directory is doubled by *copying it over* and 'fixing' pointer to split image page. (Use of least significant bits enables efficient doubling via copying of directory!)



Directory Doubling

Why use least significant bits in directory?
⇔ Allows for doubling via copying!

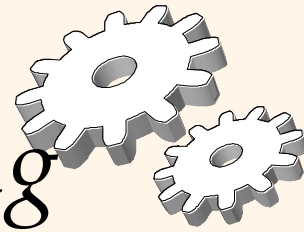


Least Significant

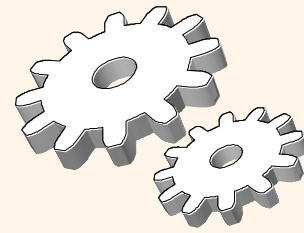
vs.

Most Significant

Comments on Extendible Hashing



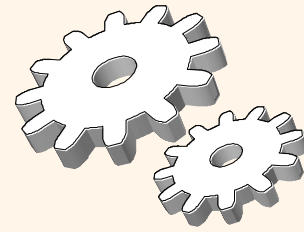
- ❖ If directory fits in memory, equality search answered with one disk access; else two.
 - 100MB file, 100 bytes/rec, 4K pages contains 1,000,000 records (as data entries) and 25,000 directory elements; chances are high that directory will fit in memory.
 - Directory grows in spurts, and, if the distribution of *hash values* is skewed, directory can grow large.
 - Multiple entries with same hash value cause problems!
- ❖ **Delete**: If removal of data entry makes bucket empty, can be merged with 'split image'. If each directory element points to same bucket as its split image, can halve directory.



Linear Hashing

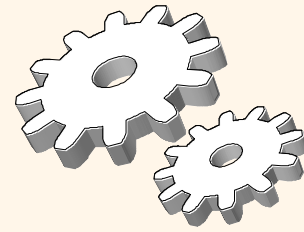
- ❖ This is another dynamic hashing scheme, an alternative to Extendible Hashing.
- ❖ LH handles the problem of long overflow chains without using a directory, and handles duplicates.
- ❖ Idea: Use a family of hash functions $\mathbf{h}_0, \mathbf{h}_1, \mathbf{h}_2, \dots$
 - $\mathbf{h}_i(\text{key}) = \mathbf{h}(\text{key}) \bmod(2^i\mathbf{N})$; \mathbf{N} = initial # buckets
 - \mathbf{h} is some hash function (range is *not* 0 to $\mathbf{N}-1$)
 - If $\mathbf{N} = 2^{d_0}$, for some d_0 , \mathbf{h}_i consists of applying \mathbf{h} and looking at the last d_i bits, where $d_i = d_0 + i$.
 - \mathbf{h}_{i+1} doubles the range of \mathbf{h}_i (similar to directory doubling)

Linear Hashing (Contd.)

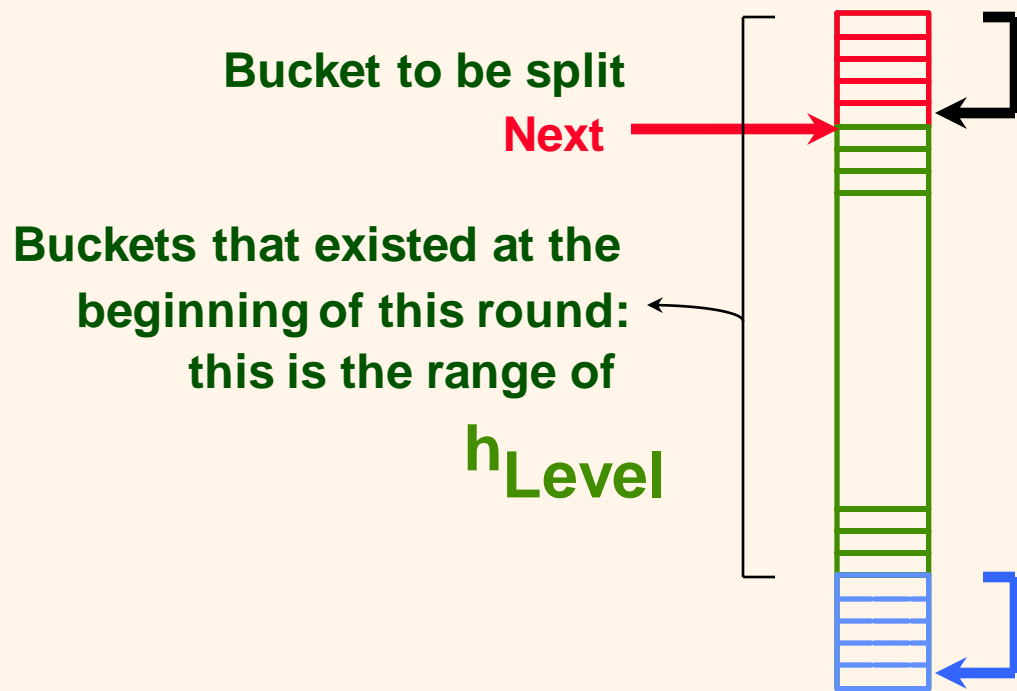


- ❖ Directory avoided in LH by using overflow pages, and choosing bucket to split round-robin.
 - **Splitting proceeds in `rounds`**. Round ends when all N_R initial (for round R) buckets are split. Buckets 0 to *Next-1* have been split; *Next* to N_R yet to be split.
 - **Current round number is *Level***.
 - **Search**: To find bucket for data entry r , find $\mathbf{h}_{Level}(r)$:
 - If $\mathbf{h}_{Level}(r)$ in range `*Next* to N_R ', r belongs here.
 - Else, r could belong to bucket $\mathbf{h}_{Level}(r)$ or bucket $\mathbf{h}_{Level}(r) + N_R$; must apply $\mathbf{h}_{Level+1}(r)$ to find out.

Overview of LH File

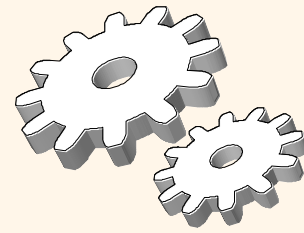


❖ In the middle of a round.



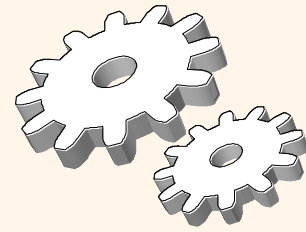
Buckets split in this round:
If h_{Level} (search key value) is in this range, must use $h_{Level+1}$ (search key value) to decide if entry is in 'split image' bucket.

'split image' buckets:
created (through splitting of other buckets) in this round



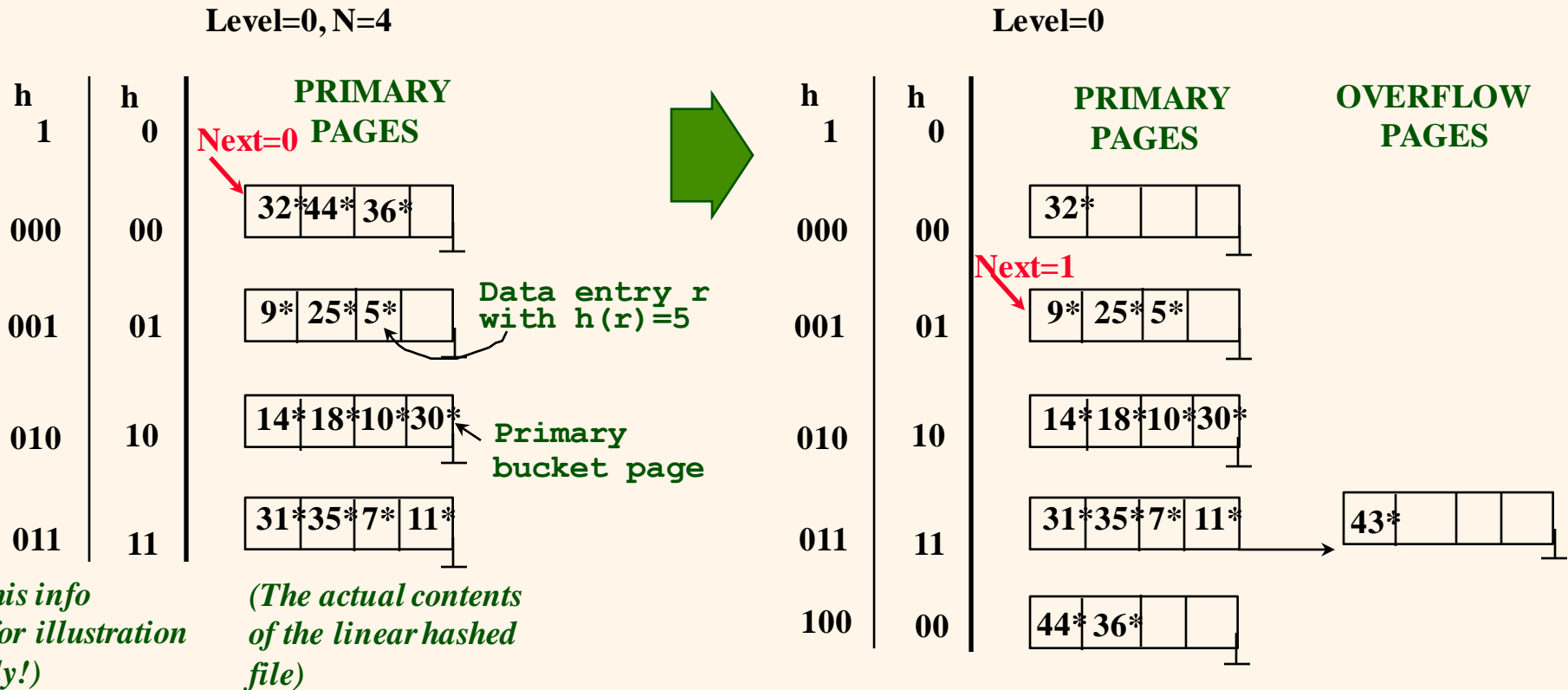
Linear Hashing (Contd.)

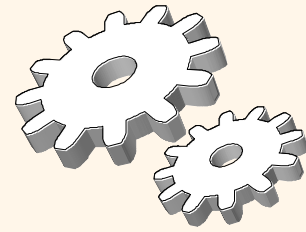
- ❖ **Insert**: Find bucket by applying $h_{Level} / h_{Level+1}$:
 - If bucket to insert into is full:
 - Add overflow page and insert data entry.
 - (*Maybe*) Split *Next* bucket and increment *Next*.
- ❖ Can choose any criterion to `trigger' split.
- ❖ Since buckets are split round-robin, long overflow chains don't develop!
- ❖ Doubling of directory in Extendible Hashing is similar; switching of hash functions is *implicit* in how the # of bits examined is increased.



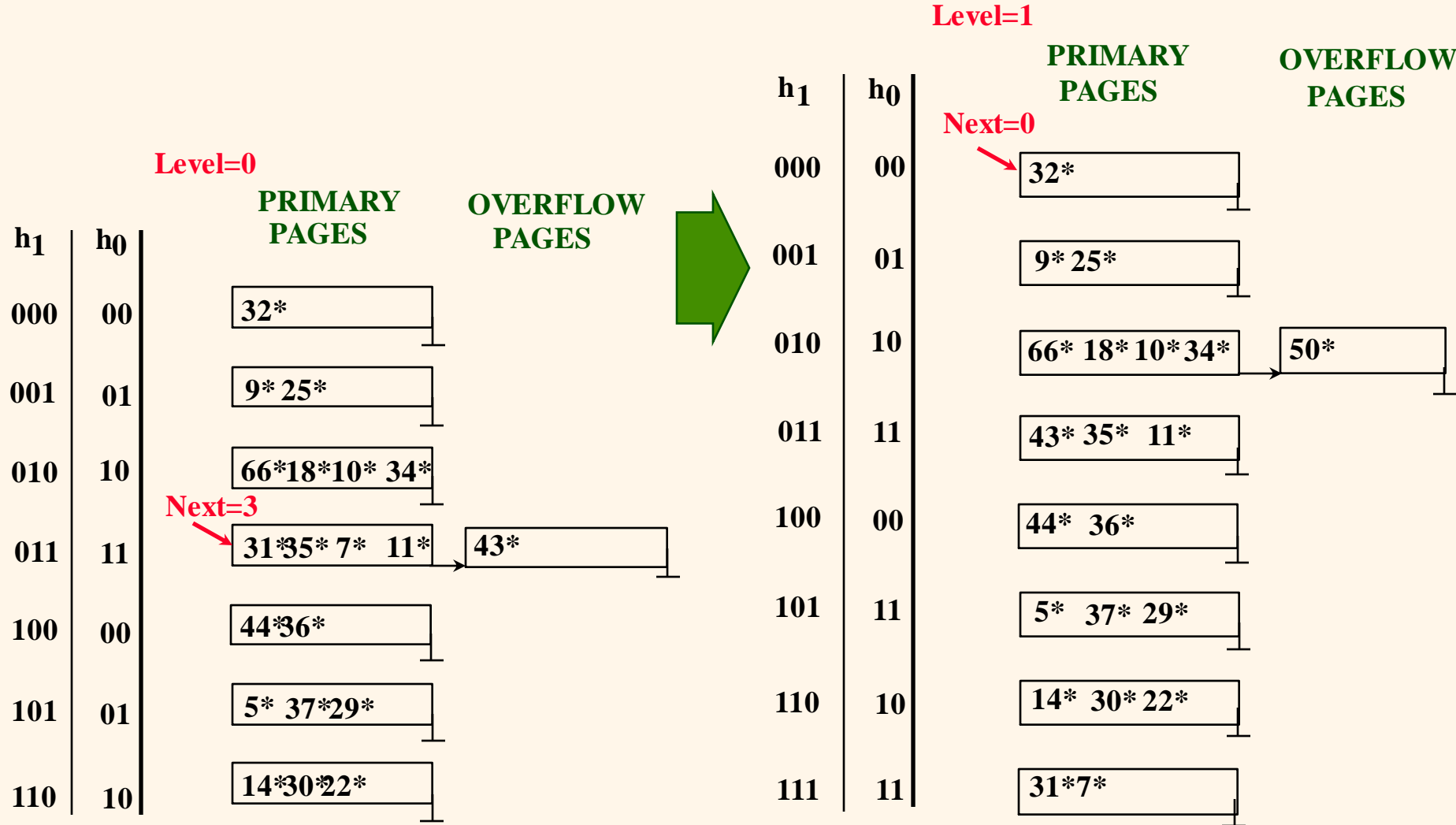
Example of Linear Hashing

- ❖ On split, $h_{\text{Level}+1}$ is used to re-distribute entries.

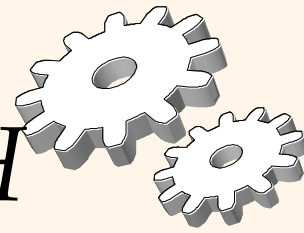




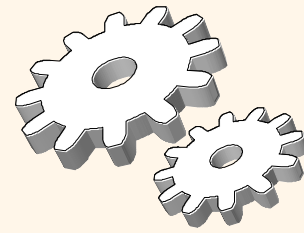
Example: End of a Round



LH Described as a Variant of EH

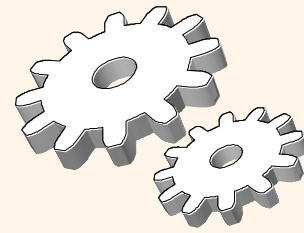


- ❖ The two schemes are actually quite similar:
 - Begin with an EH index where directory has N elements.
 - Use overflow pages, split buckets round-robin.
 - First split is at bucket 0. (Imagine directory being doubled at this point.) But elements $\langle 1, N+1 \rangle$, $\langle 2, N+2 \rangle$, ... are the same. So, need only create directory element N , which differs from 0, now.
 - When bucket 1 splits, create directory element $N+1$, etc.
- ❖ So, directory can double gradually. Also, primary bucket pages are created in order. If they are *allocated* in sequence too (so that finding i 'th is easy), we actually don't need a directory! Voila, LH.



Summary

- ❖ Hash-based indexes: best for equality searches, cannot support range searches.
- ❖ Static Hashing can lead to long overflow chains.
- ❖ Extendible Hashing avoids overflow pages by splitting a full bucket when a new data entry is to be added to it. (*Duplicates may require overflow pages.*)
 - Directory to keep track of buckets, doubles periodically.
 - Can get large with skewed data; additional I/O if this does not fit in main memory.



Summary (Contd.)

- ❖ Linear Hashing avoids directory by splitting buckets round-robin, and using overflow pages.
 - Overflow pages not likely to be long.
 - Duplicates handled easily.
 - Space utilization could be lower than Extendible Hashing, since splits not concentrated on `dense' data areas.
 - Can tune criterion for triggering splits to trade-off slightly longer chains for better space utilization.
- ❖ For hash-based indexes, a *skewed* data distribution is one in which the *hash values* of data entries are not uniformly distributed!