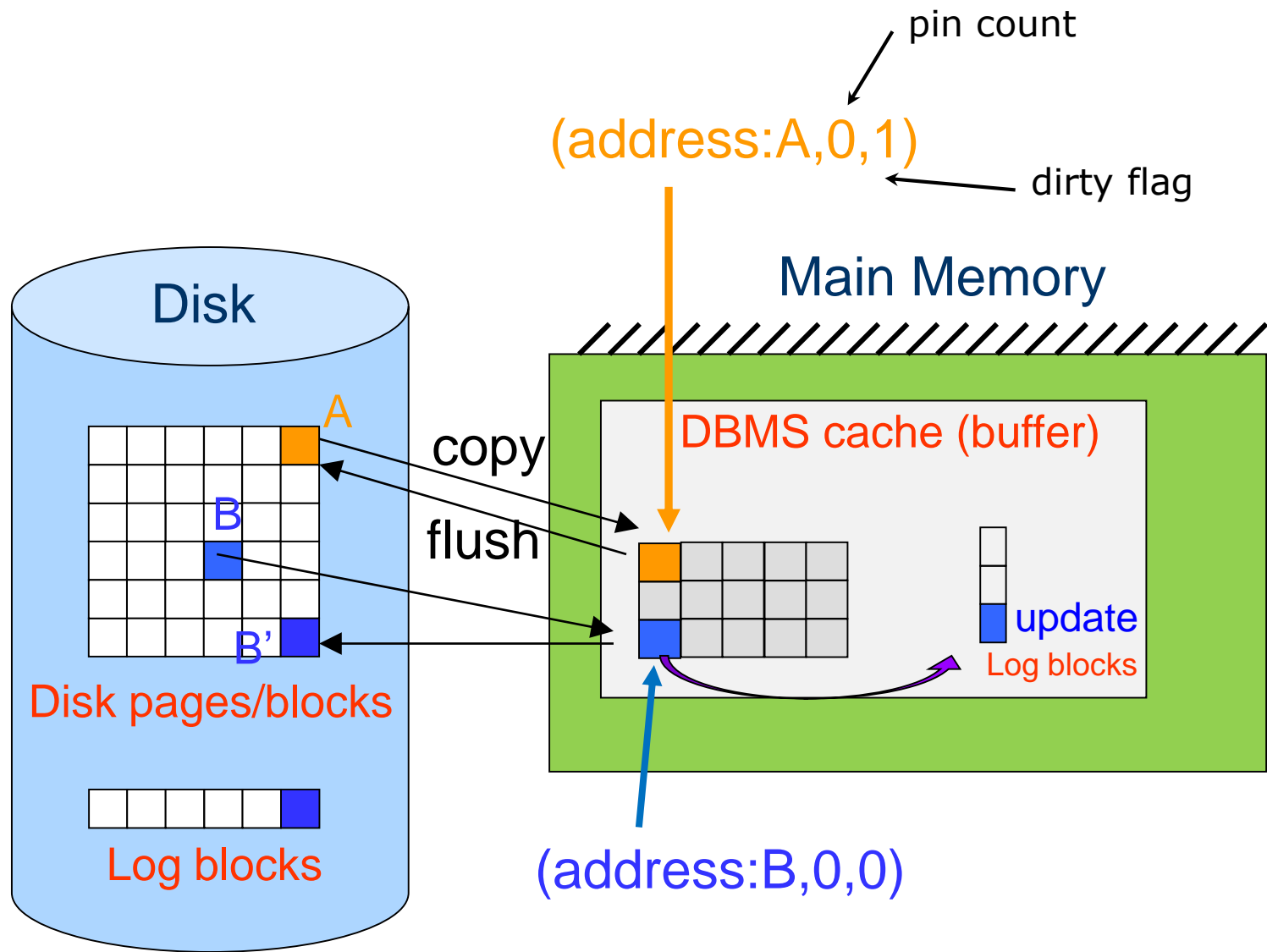


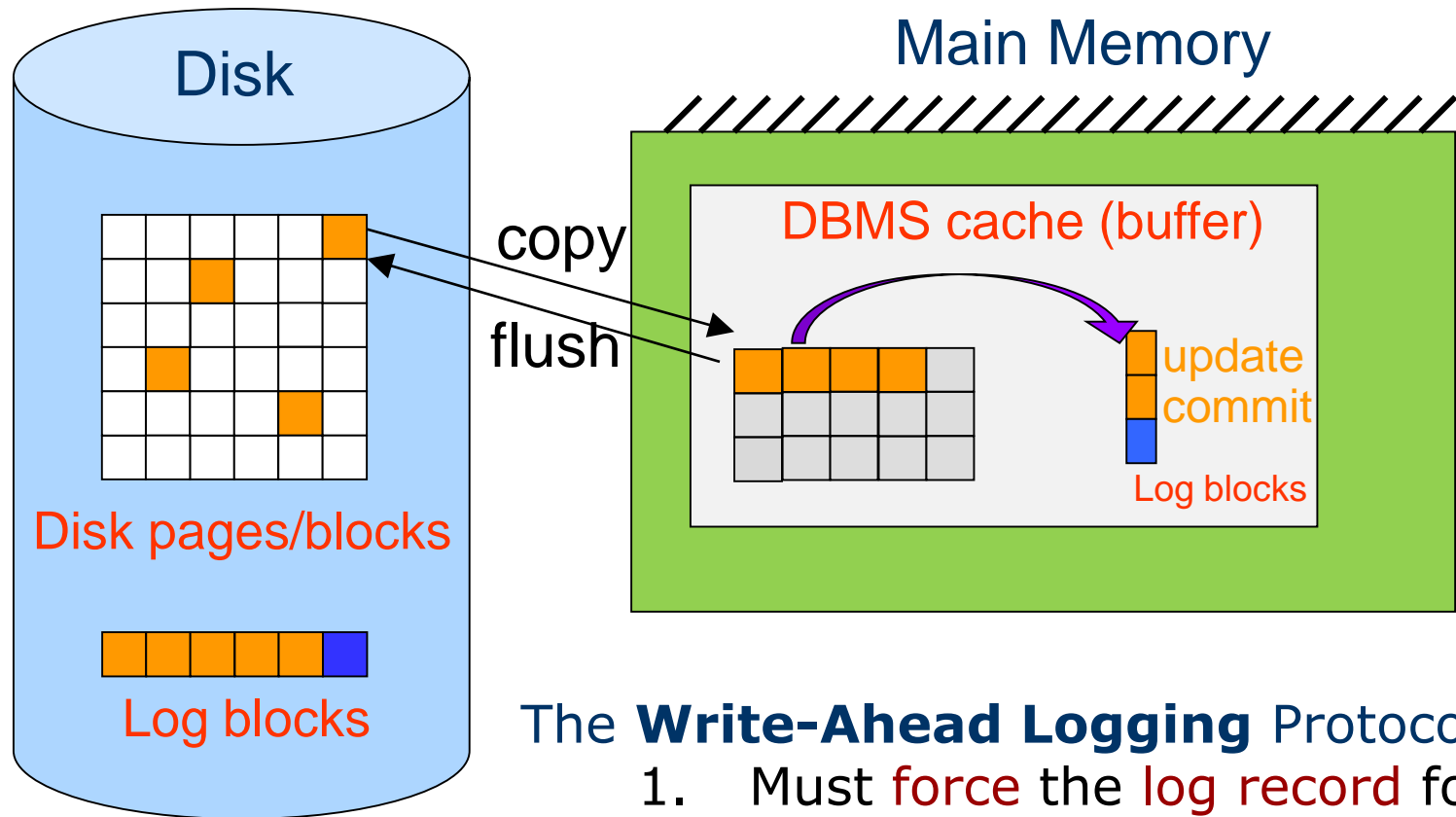
# COURSE 4

## Database Recovery 2

# Data Update

- **Immediate Update:** As soon as a data item is modified in cache, the disk copy is updated.
- **Deferred Update:** All modified data items in the cache is written either after a transaction ends its execution or after a fixed number of transactions have completed their execution.
- **In-place update:** The disk version of the data item is overwritten by the cache version.
- **Shadow update:** The modified version of a data item does not overwrite its disk copy but is written at a separate disk location.

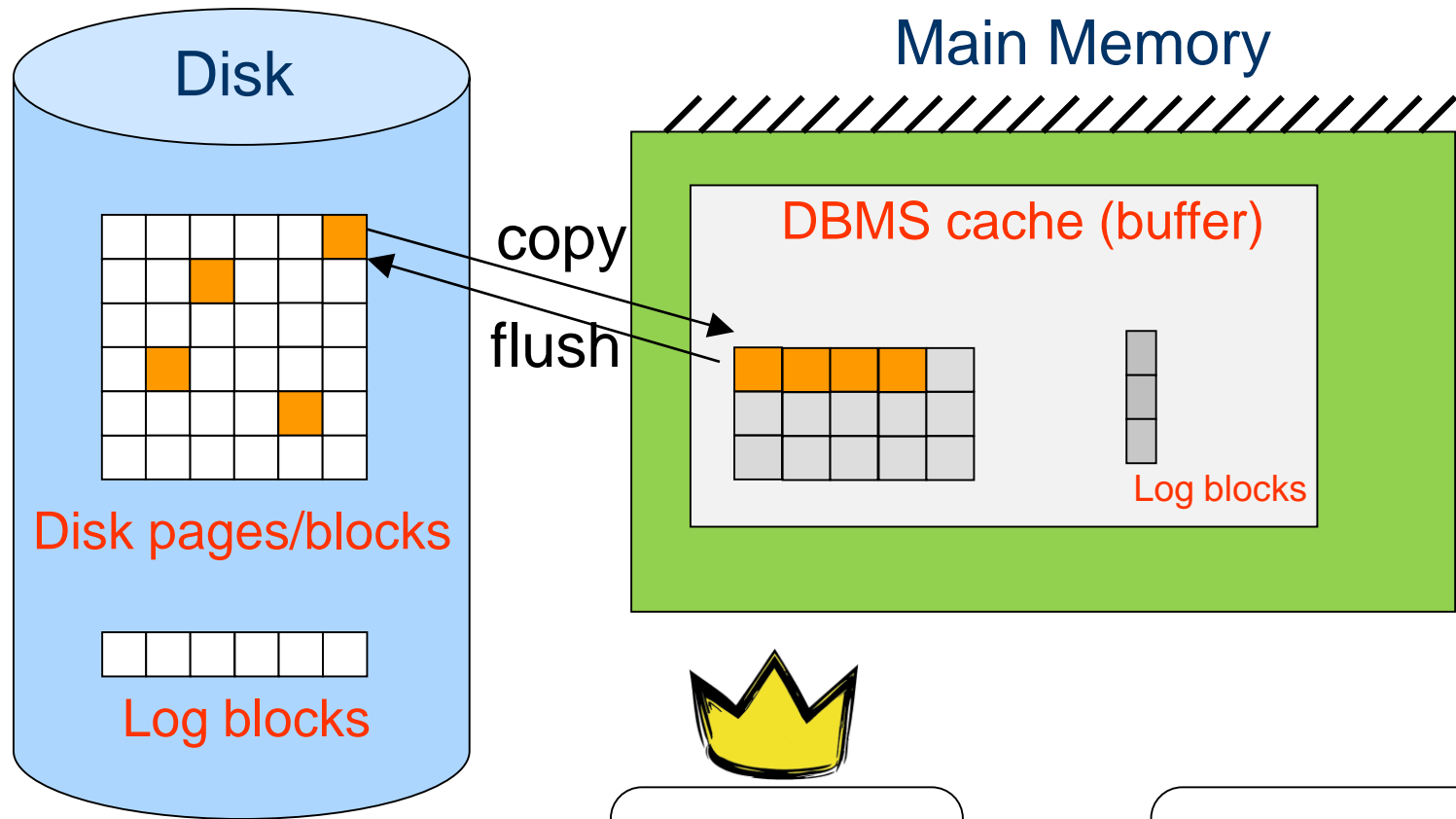




### The **Write-Ahead Logging** Protocol:

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 transaction **before commit**.

- **STEAL** (why enforcing *Atomicity* is hard)
  - *To steal frame F*: Current page in F (say P) is written to disk; some transaction holds lock on P.
    - What if the transaction 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.



Performance

vs



Atomicity & Durability

# 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.

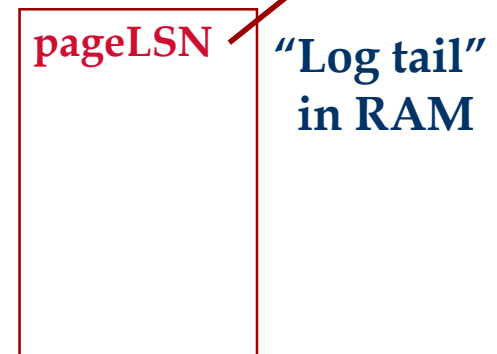
- System keeps track of **flushedLSN**.

- 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:

LSN  
prevLSN  
TransID  
type  
update  
records  
only { pageID  
length  
offset  
before-image  
after-image

## Possible log record types:

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



# Compensation Log Record (CLR)

- Written just before change recorded in one update log record is undone
- Contains a field called **undoNextLSN**
  - LSN of next log record to be undone for the transaction that wrote the update record
  - Set to prevLSN of the update log record
- Indicate which actions have already been undone
- Prevent undoing same action twice

# Other Log-Related State

## ■ Transaction Table:

- One entry per active transaction.
- 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 a Transaction

- 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
- In-place update.

# The Big Picture: What's Stored Where



## LogRecords

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



## Data pages

each  
with a  
pageLSN

## Master record



## Transaction Table

lastLSN  
status

## Dirty Page Table

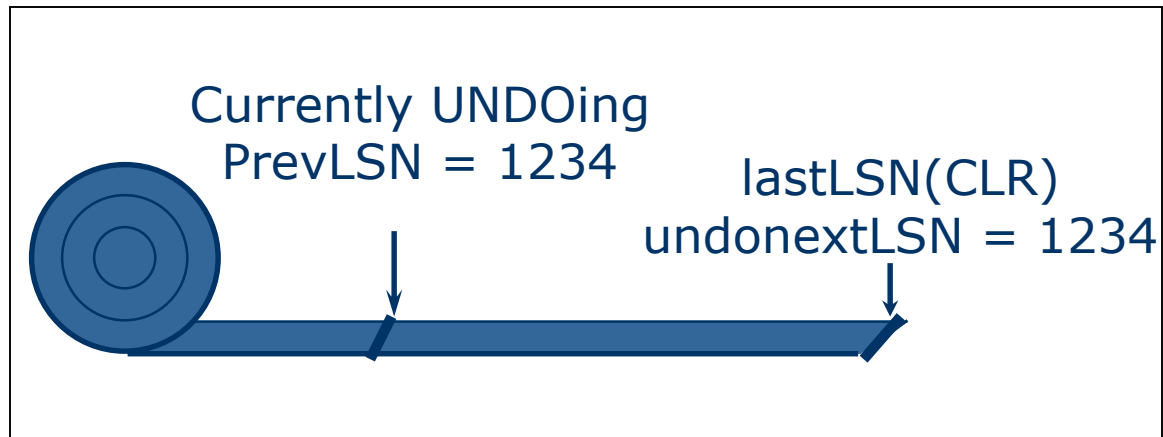
recLSN

## flushedLSN

# Simple Transaction Abort

- For now, consider an explicit abort of a transaction (*no crash involved*).
- We want to “play back” the log in reverse order, *UNDO*ing updates.
  - Get **lastLSN** of transaction from trans. 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!
- 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).
  - CLRs *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 transaction'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.

# Phases of ARIES

= Advanced Recovery and Integrated Extraction System

- Analysis: Scan the log forward (from the most recent *checkpoint*) to identify all transactions that were active, and all dirty pages in the buffer pool at the time of the crash
- Redo: Redoes all updates to dirty pages in the buffer pool, as needed, to ensure that all logged updates are in fact carried out and written to disk.
- Undo: The writes of all transactions that were active at the crash are undone (by restoring the *before value* of the update, which is in the log record for the update), working backwards in the log. (Some care must be taken to handle the case of a crash occurring during the recovery process!)



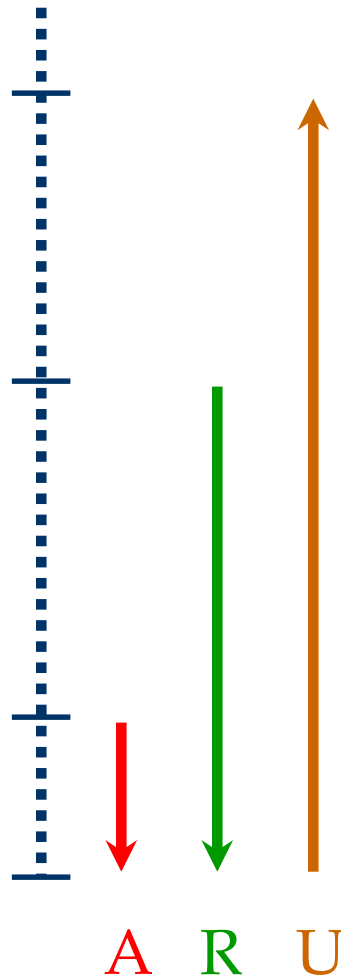
# Crash Recovery: Big Picture

Oldest log  
rec. of trans.  
active at crash

Smallest  
recLSN in  
dirty page  
table after  
Analysis

Last chkpt

CRASH



Start from a checkpoint (found via master record).

Three phases. Need to:

- Figure out which transactions committed since checkpoint, which failed (**Analysis**).
- **REDO** *all* actions (repeat history)
- **UNDO** effects of failed transactions.

# Recovery: The **Analysis** Phase

- Reconstruct state at checkpoint.
  - Set of active transactions and dirty pages
- Scan log forward from checkpoint.
  - **End** record: Remove transaction from transaction table.
  - **Other records**: Add transaction to transaction table, set **lastLSN=LSN**, change transaction 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 transactions!), 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” transaction }

**Repeat:**

Choose largest LSN among ToUndo.

If this LSN is a CLR and undonextLSN == NULL

Write an End record for this transaction.

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



Trans Table

lastLSN

status

Dirty Page Table

recLSN

flushedLSN

ToUndo

LSN	LOG
-----	-----

00	— begin_checkpoint
----	--------------------

05	— end_checkpoint
----	------------------

10	— update: T1 writes P5
----	------------------------

20	— update T2 writes P3
----	-----------------------

30	— T1 abort
----	------------

40	— CLR: Undo T1 LSN 10
----	-----------------------

45	— T1 End
----	----------

50	— update: T3 writes P1
----	------------------------

60	— update: T2 writes P5
----	------------------------

✗ CRASH, RESTART

prevLSNs

# Example of Recovery



Trans Table

lastLSN

status

Dirty Page Table

recLSN

flushedLSN

ToUndo

LSN	LOG
00,05	begin_checkpoint, end_checkpoint
10	update: T1 writes P5
20	update T2 writes P3
30	T1 abort
40,45	CLR: Undo T1 LSN 10, T1 End
50	update: T3 writes P1
60	update: T2 writes P5
	✗ CRASH, RESTART
70	CLR: Undo T2 LSN 60
80,85	CLR: Undo T3 LSN 50, T3 end
	✗ CRASH, RESTART
90,95	CLR: Undo T2 LSN 20, T2 end

undonextLSN

# Additional Crash Issues

- System crash may occur during database recovery:
  - Apply redo and undo to a record either once only, or
  - Make redo and undo as idempotent operations
    - idempotent: Acting as if used only once, even if used multiple times
- 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 transactions.