



**EPL646 – Advanced Topics in Databases**

**Lecture 10**

**Crash Recovery: Undo, Redo, Undo/Redo  
Logging and Recovery**

**Chapter 17: Database Systems: The Complete Book**

**Garcia-Molina, Ullman, Widom, 2ED**

**Demetris Zeinalipour**

<http://www.cs.ucy.ac.cy/~dzeina/courses/epl646>

# Recovery: Outline



- Recovery: Definitions, Purpose, Failure Reasons, ACID Properties and Responsibilities
- **Three Types of Recovery**
  - **Undo** (uncommitted) Logging
  - **Redo** (committed) Logging
  - **Undo/Redo** Logging (not discussed)
- Checkpointing and Nonquiescent Checkpointing

# Database Recovery

## Ανάκαμψη σε Βάσεις Δεδομένων



- **Purpose of Database Recovery (Σκοπός Ανάκαμψης)**
  - To bring the database into the last **consistent state**, which existed **prior to the failure**.
  - To **preserve transaction properties** (**Atomicity**, Consistency, Isolation and **Durability**), especially the bold properties.
- **Example:**
  - A **system crashes** before a **fund transfer** transaction commits its execution,
  - Either **one** or **both** accounts may have an **incorrect value**.
  - Thus, the database must be **restored to the state** before the transaction modified **any** of the accounts.

# Failure Reasons of Transactions

(Λόγοι Σφάλματος Δοσοληψιών)



Type of Crash	Prevention
Wrong data entry	Constraints and Data cleaning
Disk crashes	Redundancy: e.g. RAID, archive
Fire, theft, bankruptcy...	Buy insurance, Change jobs...
System failures: e.g. power, OS, RAM	<b>Our focus:</b> <b>DATABASE</b> <b>RECOVERY</b>

Most frequent

# Review: The ACID properties

## (Επανάληψη: Οι ιδιότητες ACID)



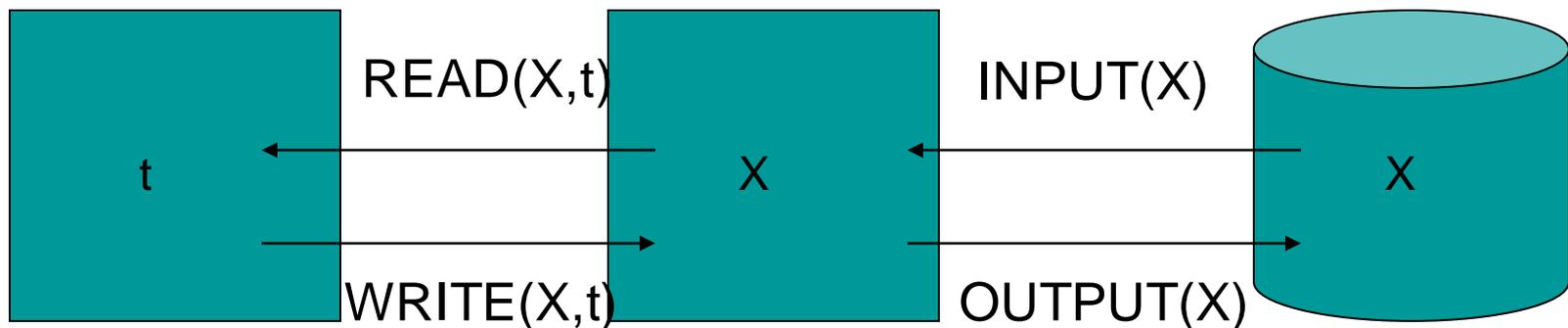
- **Atomicity (Ατομικότητα):** All actions in the Xact happen, or none happen (Responsibility: **Recovery Manager**).
- **(Semantic) Consistency (Συνέπεια):** If each Xact is consistent, and the DB starts consistent, it ends up consistent (Responsibility: **User ... using constraints**).
- **Isolation (Απομόνωση):** Execution of one Xact is isolated from that of other Xacts (Responsibility: **Concurrency Control Manager**).
- **Durability (Μονιμότητα):** If a Xact commits, its effects persist (Responsibility: **Recovery Manager**).
- The **Recovery Manager** guarantees **Atomicity & Durability**

# System Model & Definitions

## (Μοντέλο Συστήματος & Ορισμοί)



- Assumption: the database is composed of reads/writes over some **elements**
  - Usually 1 element = 1 block
  - Can be smaller (=1 record) or larger (=1 relation)
- Symbols Utilized
  - X database object
  - t local (program) variable



program (memory)

memory buffer  
(aka Buffer Manager)

Database disk



# Example Used in these Slides

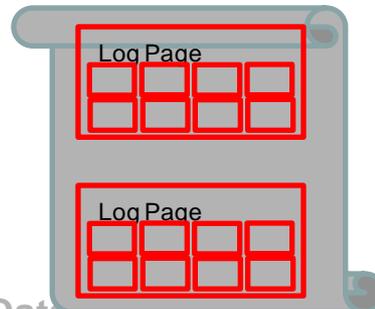
**READ(A,t); t := t\*2;WRITE(A,t); READ(B,t); t := t\*2;WRITE(B,t)**

Action	t	Mem A	Mem B	Disk A	Disk B
INPUT(A)		8		8	8
READ(A,t)	8	8		8	8
t:=t*2	16	8		8	8
WRITE(A,t)	16	16		8	8
INPUT(B)	16	16	8	8	8
READ(B,t)	8	16	8	8	8
t:=t*2	16	16	8	8	8
WRITE(B,t)	16	16	16	8	8
OUTPUT(A)	16	16	16	16	8
OUTPUT(B)	16	16	16	16	16

# System Failures and Logs



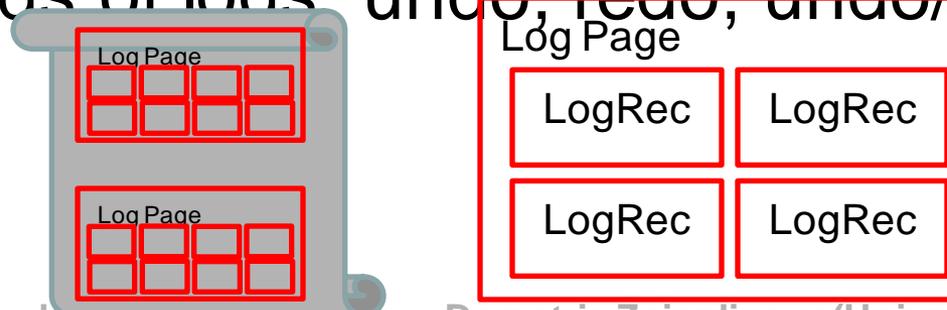
- Each transaction (program) has an ***internal state***
- When system crashes, internal state is lost
  - Don't know which parts executed and which didn't
- Remedy: use a **log**
  - A **sequential file** that keeps track of all **transaction operations** that **affect the values** of database **items**.
  - Log is maintained both on disk and buffer (will see next) where, how, when.



# Transaction Log (Κατάστιχο Δοσοληψιών)



- An **append-only** file containing **log records**
- **Note:** multiple transactions run concurrently, log records are interleaved
- After **system crash**, use log to do either or both:
  - **Undo** updates of **uncommitted** xacts (**from end**)
  - **Redo** updates of **committed** xacts (that have not been OUTPUTTED to disk) (**from start**)
- Three kinds of logs: undo, redo, undo/redo



# Recovery A: Undo Logging



## Log records

- **<START T>**
  - transaction T has begun
- **<COMMIT T>**
  - T has committed
- **<ABORT T> (aka. Rollback)**
  - T has aborted
- **<T,X,v>**
  - T has updated element X, and its old value was v

# Undo-Logging Example

## (Log-then-output, Commit later)



Action	t	Mem A	Mem B	Disk A	Disk B	Log
						<START T>
READ(A,t)	8	8		8	8	
t:=t*2	16	8		8	8	<i>old value</i>
WRITE(A,t)	16	16		8	8	<T,A,8>
READ(B,t)	8	16	8	8	8	
t:=t*2	16	16	8	8	8	
WRITE(B,t)	16	16	16	8	8	<T,B,8>
<b>FLUSH LOG</b>				(now all log records appear on disk)		
<small>Takes time</small> <b>OUTPUT(A)</b>	16	16	16	16	8	
<small>Takes time</small> <b>OUTPUT(B)</b>	16	16	16	16	16	<b>!!! CRASH =&gt; UNDO OUTPUT(A)</b>
						<COMMIT T>
<b>FLUSH LOG</b>				(now COMMIT T appears on disk as well)		

# Undo-Logging Rules



The following two rules must be obeyed:

**U1 (Log then Output)** : If **T** modifies **X**, then  $\langle T, X, v \rangle$  must be written to disk **BEFORE** **X** is output to disk  
 *$\pi.\chi.$ , Write to Log:  $\langle T, A, \delta \rangle$  **then** Flush Log **then**  
Write to Disk: OUTPUT A*

**U2 (Commit Later)**: If **T** commits, then  $\langle \text{COMMIT } T \rangle$  must be written to disk only after all changes by **T** are output to disk  
 *$\pi.\chi.$ , Write to Disk: OUTPUT A **then**  
Write to Log: COMMIT A **then** Flush Log*

– Hence: **OUTPUTs** are done early, before the transaction commits

# Recovery with Undo Log (Abort Uncommitted)



**General Idea:** Undo Uncommitted Transactions (see next slide for details)

...

...

<T6,X6,v6>

...

...

<**START T5**>

<START T4>

<T1,X1,v1>

<**T5,X5,v5**>

<**T4,X4,v4**>

<**COMMIT T5**>

<T3,X3,v3>

<T2,X2,v2>



crash

**Question 1:**

Which updates are undone ?

**Question 2:**

How far back do we need to read in the log ?

# Recovery with Undo Log



After system's crash, run Recovery Manager

- **Step 1.** Process the Log from **END** and decide for each **transaction T** whether it is completed or not (how? See next slide)
  - <START T>....<COMMIT T>.... = COMPLETE
  - <START T>....<ABORT T>..... = COMPLETE
  - <START T>..... = INCOMPLETE
- **Step 2.** Undo all modifications by incomplete transactions (see next slide)

# Recovery with Undo Log



## How does the Recovery Manager classify Xacts?

- Read log **from the end**; cases:
  - **<COMMIT T>**: mark T as completed
  - **<ABORT T>**: mark T as completed
  - **<T,X,v>**: if (T is completed) then  
Ignore  
else // incomplete  
Write  $X = \text{oldest}(v)$  to disk (i.e. reset X to its initial value, higher in the log)
  - **<START T>**: ignore

# Crashing During Recovery with Undo Log



- **System Crash during Recovery with Undo Log. What happens?**

- We can repeat all actions from scratch for a second time, no harm is done.

- **Why?**

- All undo commands are idempotent (*ταυτοδύναμες*)

- **T1(A); crash; T1(A); T2(A)**

Generates the same result with:

- **T1(A); T2(A)**

because « $\langle T_i, A, v \rangle$  holds the previous value  $v$  of object  $A$  (wouldn't apply if  $v$  was holding the difference from previous value)

# Recovery with Undo Log



## When do we stop reading the log ?

- We cannot stop until we reach the **beginning of the log file**

## Why?

- Think about `<START T1>` on first line without `<COMMIT|ABORT T1>`

## So Recovery is not very practical!

- Better idea: use checkpointing

# Checkpointing (Σημείο Έλεγχου)



**Idea: Checkpoint the database periodically.**

## How?

- Stop accepting new transactions
- Wait until all current transactions complete
- Write a <CKPT> log record
- Resume transactions

# Undo Recovery with Checkpointing



During recovery,  
Can stop at first  
<CKPT>

...

...

<T9,X9,v9>

...

...

(all completed)

**<CKPT>**

<START T2>

<START T3>

<START T5>

<START T4>

<T1,X1,v1>

<T5,X5,v5>

<T4,X4,v4>

<COMMIT T5>

<T3,X3,v3>

<T2,X2,v2>

} other transactions

} transactions T2,T3,T4,T5

crash

# Nonquiescent Checkpointing (Μη-αδρανές Checkpointing)



- **Problem with checkpointing:** What happens with long running xacts ? => database freezes during checkpoint
- Would like to checkpoint while database is operational

- **Idea: Nonquiescent Checkpointing** (μη-αδρανές σημείο έλεγχου):
  - Write a **<START CKPT(T1,...,Tk)>**  
where **T1,...,Tk** are **all** active transactions
  - Continue normal operation
  - When all of **T1,...,Tk** have completed, write **<END CKPT>**

# Undo Recovery with Nonquiescent Checkpointing



During recovery,  
Can stop at first  
<CKPT>

(provided an  
<END CKPT>  
is seen)

New Xacts might have  
begun!

...

...

*// Start Monitoring active  
Xacts T4, T5, T6*

**<START CKPT T4, T5, T6>**

...

*//Wait until all active commit  
or abort (do not prohibit other  
from starting)*

...

...

*Write when T4, T5, T6  
complete*

**<END CKPT>**

earlier transactions plus  
T4, T5, T6

T4, T5, T6, plus  
later transactions

After END CKPT has been written, all records prior START CKPT can be deleted

later transactions

crash

# Recovery B: Redo Logging



**Problems of UNDO Logging:** It requires to **OUTPUT** the data before **COMMIT**. Thus, in order to offer **Strict Schedules** we must write all in-memory data to disk.

## Log records

- $\langle \text{START } T \rangle$  = transaction  $T$  has begun
- $\langle \text{COMMIT } T \rangle$  =  $T$  has committed
- $\langle \text{ABORT } T \rangle$  =  $T$  has aborted
- $\langle T, X, v \rangle$  =  $T$  has updated element  $X$ , and its new value is  $v$

# Redo-Logging Example

## (Commit-then-Output)

IREAD: INPUT & READ



Action	T	Mem A	Mem B	Disk A	Disk B	Log
						<START T>
IREAD(A,t)	8	8		8	8	
t:=t*2	16	8		8	8	
WRITE(A,t)	16	16		8	8	<T,A,16>
IREAD(B,t)	8	16	8	8	8	
t:=t*2	16	16	8	8	8	
WRITE(B,t)	16	16	16	8	8	<T,B,16>
						<COMMIT T>
<b>FLUSH LOG</b>	(now all log records appear on disk). The next OUTPUT commands can in practice come much later					
<del>OUTPUT(A)</del>	16	16	16	16	8	<b>!!! CRASH =&gt;</b>
<del>OUTPUT(B)</del>	16	16	16	16	16	<b>REDO COMMITTED not OUTPUTED!</b>

Other Xacts might start using A, B (in a “strict” way) without pulling them from disk (i.e., directly from buffer manager) . Not applicable in Undo Logging (where data needs to be first OUTPUT to disk (before COMMIT))

# Redo-Logging Rule (Write-Ahead Rule)



**R1 (Commit-then-Output):** If **T** modifies **X**, then both  $\langle \mathbf{T}, \mathbf{X}, \mathbf{v} \rangle$  and  $\langle \mathbf{COMMIT T} \rangle$  must be written to disk before **X** is written to disk

i.e., *Write to Log: COMMIT A*

*Write to Disk: OUTPUT A*

- Hence: OUTPUTs are done late
- ***This rule also known as Write-Ahead-Rule i.e., Write Ahead Log (WAL), i.e., writing log records ahead of actual data records.***

# Recovery with Redo Log (Repeat Committed)



<START T1>  
<T1,X1,v1>  
<START T2>  
<T2, X2, v2>  
<START T3>  
<T1,X3,v3>  
<COMMIT T2>  
<T3,X4,v4>  
<T1,X5,v5>  
...  
...

## Question 1:

Which updates are redone ?

## Question 2:

How far back do we need to read in the log down ?

crash

# Recovery with Redo Log (Repeat Committed)



After system's crash, run recovery manager

- Step 1. Process the Log from **START** and decide for each **transaction T** whether it is completed or not (how? See next slide)
  - **<START T>....<COMMIT T>.... = COMPLETE**
  - **<START T>....<ABORT T>..... = COMPLETE**
  - **<START T>..... = INCOMPLETE**
- Step 2. Read **log from the beginning, redo all updates of committed transactions (not outputed)**
  - **Do not repeat the uncommitted ones (let the application that initiated them worry about repeating them)**

# Recovery with Redo Log



## How does the Recovery Manager classify Xacts?

- Read log from end :
  - **<COMMIT T>**: mark T as completed
  - **<ABORT T>**: mark T as completed
- Read log from the start; cases:
  - **<T,X,v>**: if (T is **incomplete**) then // opposite to UNDO  
Ignore  
else // complete & committed  
Write **X=newest(v)** to disk
  - **<START T>**: ignore
- For each incomplete T write an **<Abort T>** record to the end of the log and **flush the log.**

# Nonquiescent Checkpointing with REDO Logging



- Write a `<START CKPT(T1,...,Tk)>`  
where **T1,...,Tk** are all active transactions
  - *Same with UNDO Checkpointing*
- **Flush to disk all blocks of committed transactions (*dirty blocks*), while continuing normal operation**
  - This solves the problem: DIRTY blocks go to disk
- When all blocks have been written, write `<END CKPT>`
  - *Same with UNDO Checkpointing*

# Redo Recovery with Nonquiescent Checkpointing



**Step 1:** look for the last  
<END CKPT>

**Step 2:** find  
corresponding  
<START CKPT ...>

```
...
<START T1>
...
<COMMIT T1>
...
<START T4>
...
<START CKPT T4, T5, T6>
...
All OUTPUTs of T1 are
now known to be on disk
...
<END CKPT>
...
// now T4 committed in
log but not OUTPUTed.
...
<START CKPT T9, T10>
...
```

**Step 3:** redo from  
the earliest start of a  
transaction listed in  
<START CKPT ...>,  
i.e. T4, T5, T6

(transactions  
committed  
earlier can be  
ignored, e.g., T1., as  
these were  
OUTPUTed by prior  
START..END block)

