

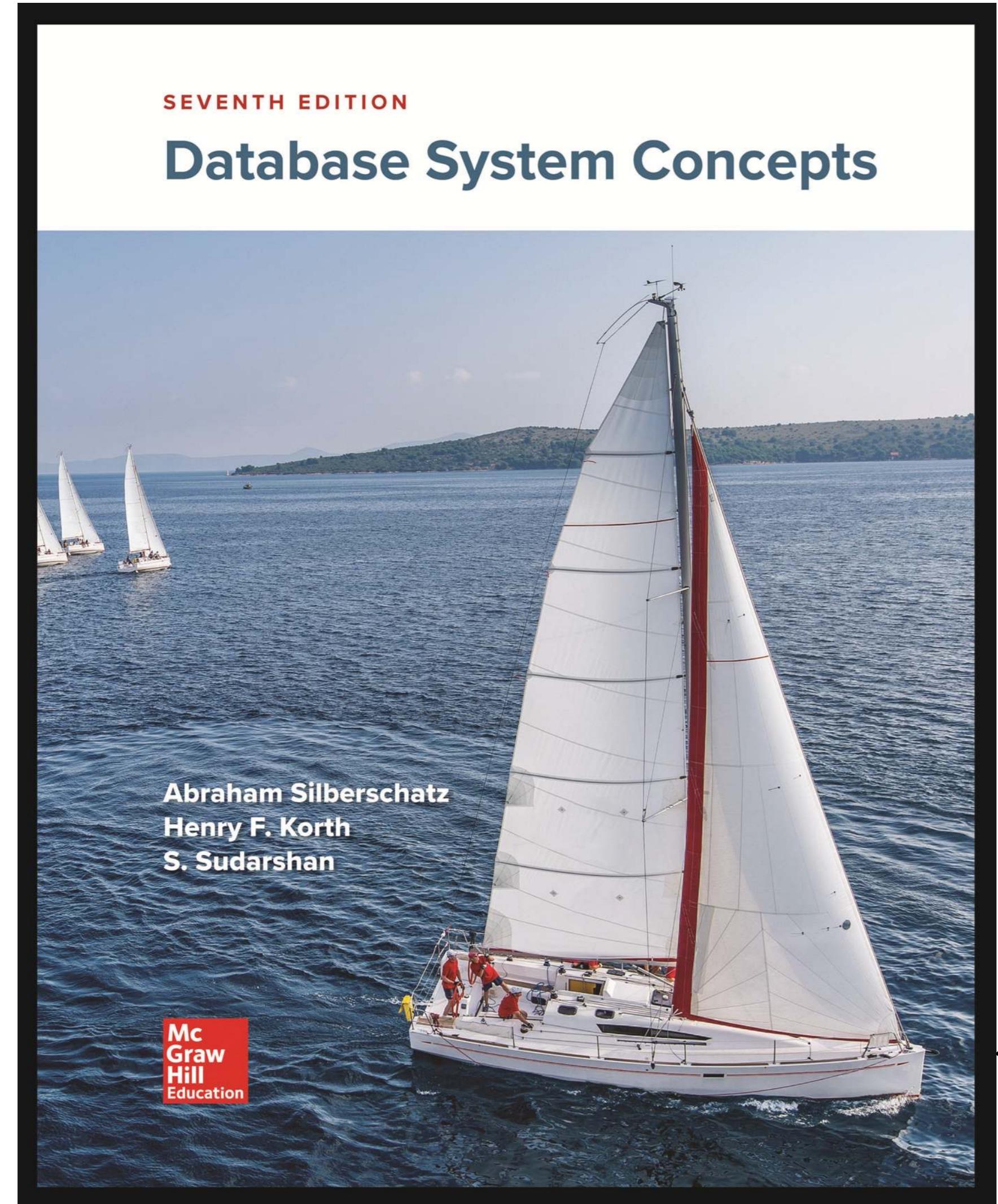
# IF3140 – Sistem Basis Data Recovery System

SEMESTER I TAHUN AJARAN 2024/2025



KNOWLEDGE & SOFTWARE ENGINEERING





# Sumber

Silberschatz, Korth, Sudarshan:  
“Database System Concepts”, 7<sup>th</sup>  
Edition

- Chapter 19: Recovery System



KNOWLEDGE & SOFTWARE ENGINEERING

# *Objectives*

Students are able to:

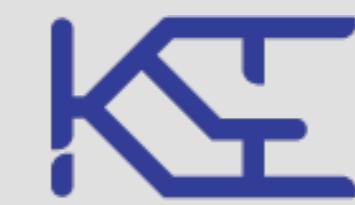
- Explain kinds of failures that can happen in a database system
- Explain stable storage and recovery task
- Explain when and why rollback is needed and how logging assures proper rollback
- Write a backup plan for a database



KNOWLEDGE & SOFTWARE ENGINEERING

# *Outline*

- Failure Classification
- Storage Structure
- Recovery and Atomicity
- Log-Based Recovery
- Shadow Paging
- Remote Backup Systems



KNOWLEDGE & SOFTWARE ENGINEERING

# Failure Classification

## Transaction failure

→ ada aspek logic dr transaksi yg membuat transaksi tidak bisa dilanjut

- **Logical errors:** transaction cannot complete due to some internal error condition
- **System errors:** the database system must terminate an active transaction due to an error condition (e.g., deadlock)

## System crash

→ tdk akan menyebabkan secondary storage rusak, yg hilang hny data yg ada di dlm main Memory

- **Fail-stop assumption:** non-volatile storage contents are assumed to not be corrupted by system crash
- Database systems have numerous integrity checks to prevent corruption of disk data

## Disk failure

→ kerusakan pd secondary storage, jd ada bagian data yg hilang

- Destruction is assumed to be detectable: disk drives use checksums to detect failures



## *Failure Example*

Suppose transaction  $T_i$  transfers \$50 from account A to account B

subtract  
50 from A

add  
50 to B

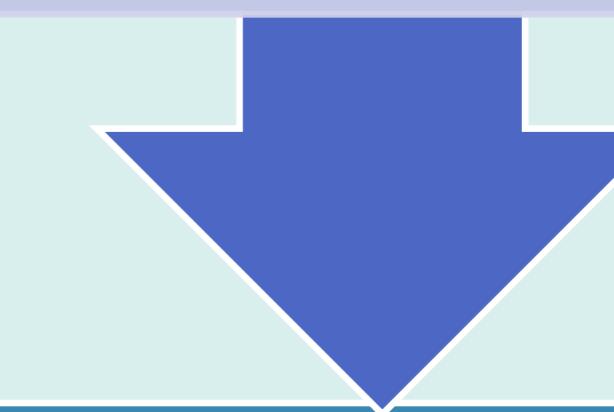


KNOWLEDGE & SOFTWARE ENGINEERING

## *Recovery Algorithm*

Actions taken during normal  
transaction processing

to ensure enough information exists to recover  
from failures



Actions taken after a failure

to recover the database contents to a state that ensures atomicity, consistency and durability

# Storage Structure



- main memory  
- cache memory

## Volatile storage:

- Does not survive system crashes

semua datanya hilang

- disk - flash memory  
- tape - non volatile RAM

- Survives system crashes
- But may still fail, losing data → kalau disk crash

## Stable storage:

- A mythical form of storage that survives all failures
- Approximated by maintaining multiple copies on distinct nonvolatile media

↳ duplicate data dr bbrp  
non-volatile storage

# ***Stable-Storage Implementation***

- Maintain multiple copies of each block on separate disks
  - copies can be at remote sites to protect against disasters such as fire or flooding.
- Block transfer can result in
  - Successful completion
  - Partial failure
  - Total failure
- Protecting storage media from failure during data transfer
  - Write the information onto the first → tulis informasi ke copy yg pertama physical block.
  - Then, write the same information onto the second physical block. → tulis informasi ke bagian kedua (copyannya)
  - The output is completed only after the second write successfully completes.



# Protecting Storage Media From Failure



## Find inconsistent blocks:

Expensive solution: Compare the two copies of every disk block

Better solution: Record in-progress disk writes on non-volatile storage

## Recovery:

If either copy is detected to have an error, overwrite it

If both have no error but are different, overwrite the second block

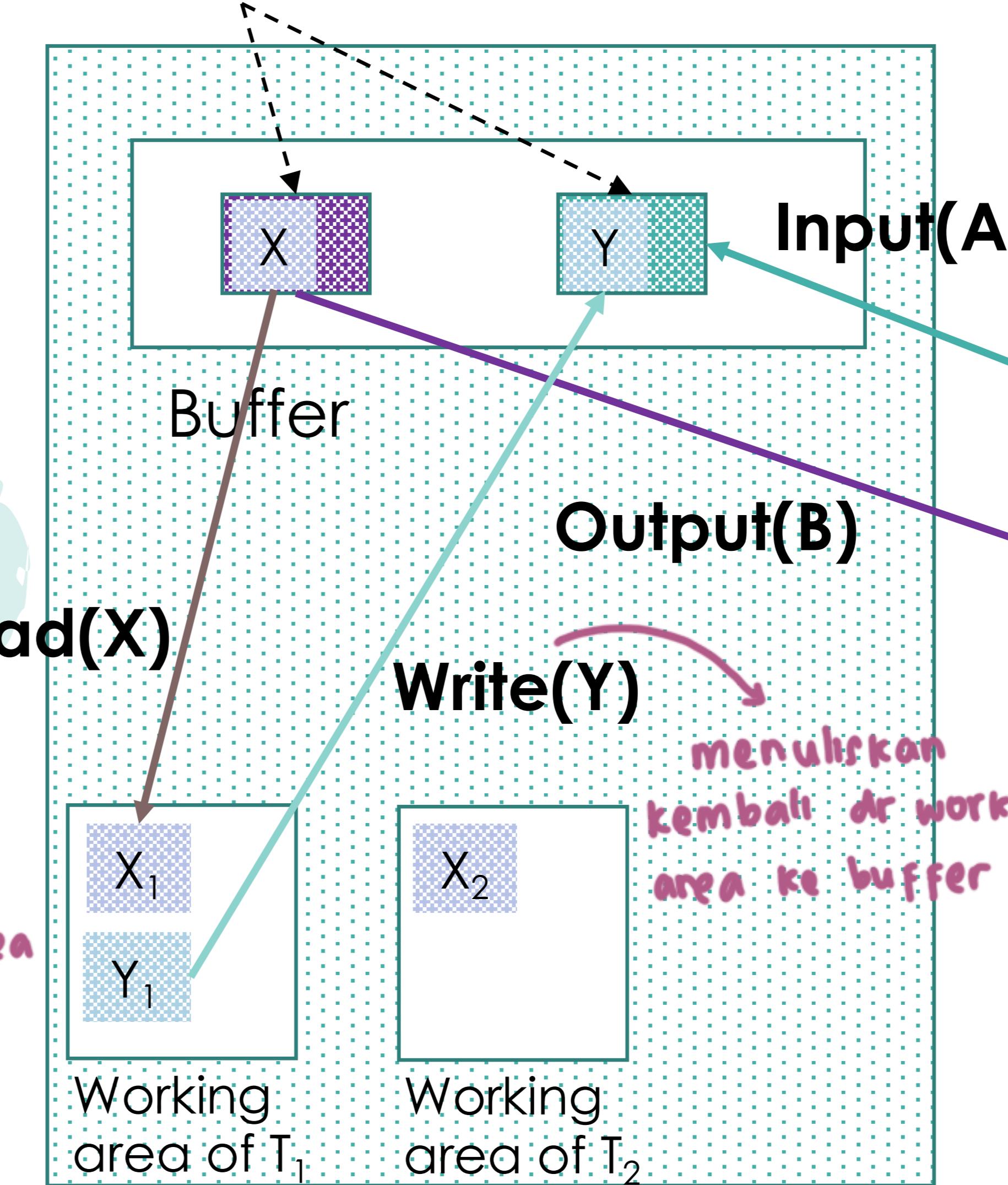
yg dipake primary block

input & output sifatnya seamless

# Data Access

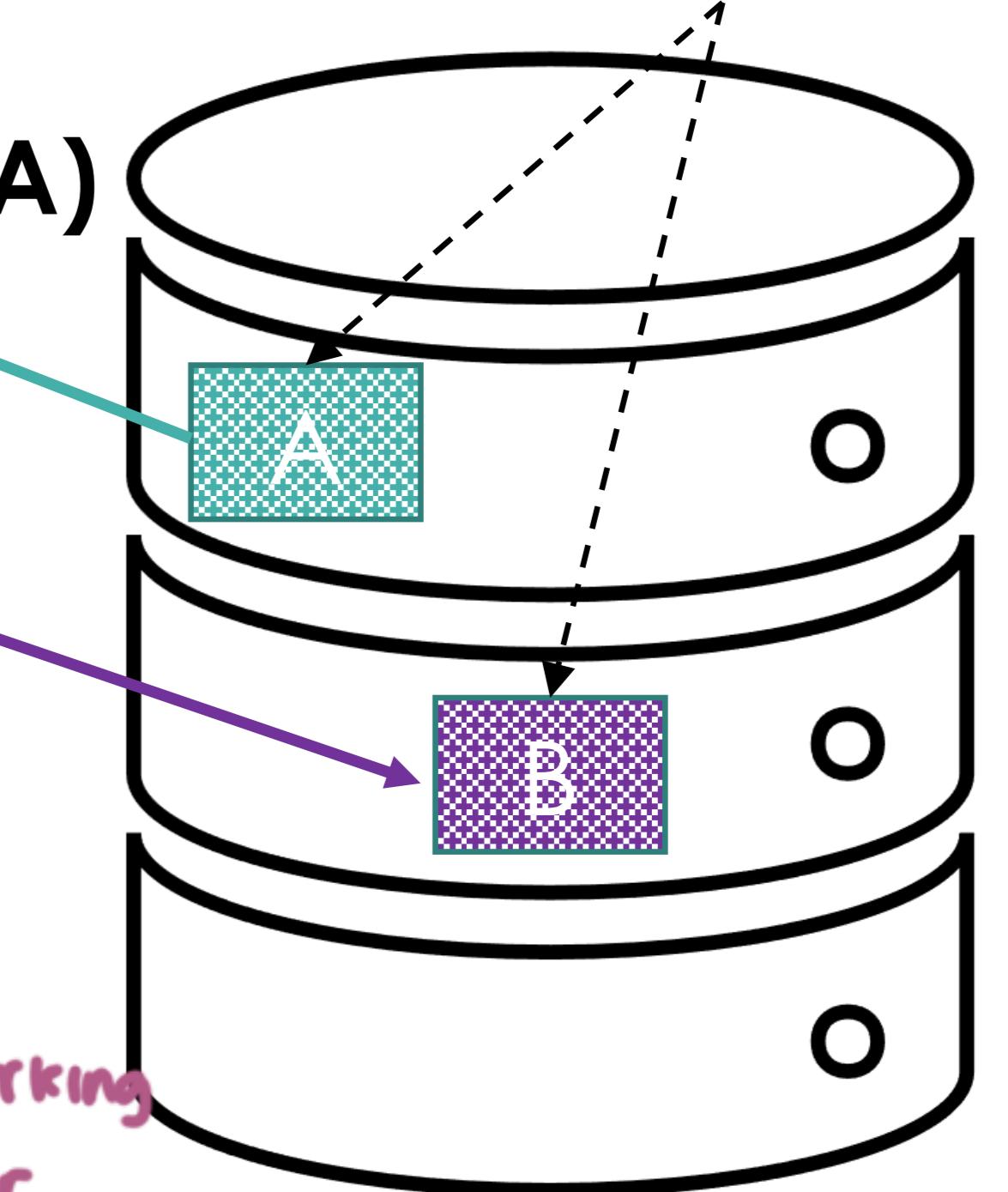
menyalin dr  
buffer ke working area

Buffer Blocks

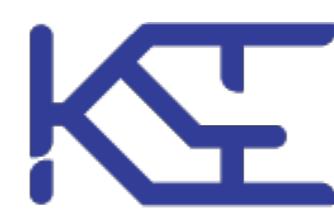


Memory

Physical Blocks



Disk  
(secondary storage)



# *Recovery and Atomicity*

**LOG-BASED RECOVERY MECHANISMS  
IN DETAIL**  
**LESS USED ALTERNATIVE: SHADOW-  
COPY AND SHADOW-PAGING**

Suppose transaction  $T_i$  transfers \$50 from account A to account B

subtract  
50 from A

add  
50 to B



KNOWLEDGE & SOFTWARE ENGINEERING

# Log-Based Recovery

Transaction  $T_{10}$  transfers \$50 from account A to account B  
Value of A = 1000 and B = 2000

## A log

- A sequence of **log records**
- The **log** is kept on stable storage

mengandung catatan/informasi updating yg terjadi thd data yg terdapat di database

$\langle T_{10} \text{ start} \rangle$

$\langle T_{10}, A, 1000, 950 \rangle$  subtract 50 from A

$\langle T_{10}, B, 2000, 2050 \rangle$  add 50 to B

$\langle T_{10} \text{ commit} \rangle$

## When transaction $T_i$ starts

- $\langle T_i \text{ start} \rangle$  log record

## Before $T_i$ executes **write**(X)

- $\langle T_i, X, V_1, V_2 \rangle$  log record

old val X ← → new val X

## When $T_i$ finishes

- $\langle T_i \text{ commit} \rangle$  log record

## Immediate database modification

semua write thd item diaplikasikan langsung ke buffer block, bs aja dia ke output ke secondary storage pdhl blm commit

## Deferred database modification

modifikasi basis data ditunda saat transaksi akan commit

semua perubahan pd item dilakukan di private area



yg kesimpulannya yg kesimpulan ke secondary storage hny updated data yg sdh committed

# Immediate Database Modification

Update log record must be written before database item is written

Output of updated blocks:

- can take place at any time
- order can be different from the order in which they are written

Transaction Commit:

- A transaction is said to have committed when its commit log record is output to stable storage
- Writes performed by a transaction may still be in the buffer when the transaction commits, and may be output later

sudah tercatat kalau  
transaksi ini commit

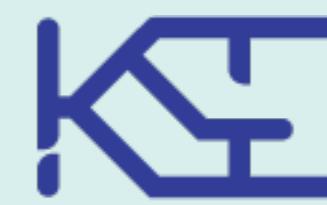


KNOWLEDGE & SOFTWARE ENGINEERING

# *Immediate Database Modification Example*

$T_0$  : transfer 50 dr A ke B

$T_1$  : tarik dana 100 dr C



KNOWLEDGE & SOFTWARE ENGINEERING

Log	Write	Output
$\langle T_0 \text{ start} \rangle$ $\langle T_0, A, 1000, 950 \rangle$ $\langle T_0, B, 2000, 2050 \rangle$		$A = 950$ $B = 2050$
$\langle T_0 \text{ commit} \rangle$ $\langle T_1 \text{ start} \rangle$ $\langle T_1, C, 700, 600 \rangle$		$C = 600$ <div style="background-color: #00AEEF; color: white; padding: 5px; border-radius: 10px;"> <math>B_C</math> output before <math>T_1</math> commits         </div>
$\langle T_1 \text{ commit} \rangle$	$B_B, B_C$ <div style="background-color: #00AEEF; color: white; padding: 5px; border-radius: 10px;"> <math>B_A</math>  <math>B_A</math> output after <math>T_0</math> commits         </div>	

# Deferred Database Modification

- The **deferred-modification** scheme performs updates to buffer/disk only at the time of transaction commit
  - Simplifies some aspects of recovery
  - But has overhead of storing local copy



KNOWLEDGE & SOFTWARE ENGINEERING

# Concurrency Control and Recovery

All concurrent transactions share a single disk buffer and a single log

A buffer block can have data items updated by one or more transactions

We assume strict two-phase commit

i.e., the updates of uncommitted transactions should not be visible to other transactions

semua lock dihold sampai commit

Log records of different transactions may be interspersed in the log



# Recovering from Failure

Transaction  $T_i$  needs to be **undone**

- Contains the record  $\langle T_i \text{ start} \rangle$ ,
- But does not contain either the record  $\langle T_i \text{ commit} \rangle$  or  $\langle T_i \text{ abort} \rangle$ .

Transaction  $T_i$  needs to be **redone**

- Contains the records  $\langle T_i \text{ start} \rangle$
- And contains the record  $\langle T_i \text{ commit} \rangle$  or  $\langle T_i \text{ abort} \rangle$



# Recovering from Failure

Transaction  $T_i$  needs to be **undone**

- Contains the record  $\langle T_i \text{ start} \rangle$ ,
- But does not contain either the record  $\langle T_i \text{ commit} \rangle$  or  $\langle T_i \text{ abort} \rangle$ .

Transaction  $T_i$  needs to be **redone**

- Contains the records  $\langle T_i \text{ start} \rangle$
- And contains the record  $\langle T_i \text{ commit} \rangle$  or  $\langle T_i \text{ abort} \rangle$

**undo( $T_i$ )** -- restores the value of all data items updated by  $T_i$  to their old values, going backwards from the last log record for  $T_i$

- Each time a data item X is restored to its old value  $\forall$  a special log record  $\langle T_i, X, V \rangle$  is written out.
- When undo of a transaction is complete, a log record  $\langle T_i \text{ abort} \rangle$  is written out.

**redo( $T_i$ )** -- sets the value of all data items updated by  $T_i$  to the new values, going forward from the first log record for  $T_i$

- No logging is done in this case



# Immediate DB Modification Recovery Example

<T<sub>0</sub> start>  
 <T<sub>0</sub>, A, 1000, 950>  
 <T<sub>0</sub>, B, 2000, 2050>

Recovery action:

- undo T<sub>0</sub>
- Log <T<sub>0</sub>, B, 2000>
- Restore B to 2000
- Log <T<sub>0</sub>, A, 1000>
- Restore A to 1000
- Log <T<sub>0</sub> abort>



<T<sub>0</sub> start>  
 <T<sub>0</sub>, A, 1000, 950>  
 <T<sub>0</sub>, B, 2000, 2050>  
 <T<sub>0</sub> commit
 <T<sub>1</sub> start>  
 <T<sub>1</sub>, C, 700, 600>

Recovery action:

redo {

- Set A to 950 new value
- Set B to 2050
- Log <T<sub>1</sub>, C, 700>
- Restore C to 700
- Log <T<sub>1</sub> abort>

undo {

<T<sub>0</sub> start>  
 <T<sub>0</sub>, A, 1000, 950>  
 <T<sub>0</sub>, B, 2000, 2050>  
 <T<sub>0</sub> commit
 <T<sub>1</sub> start>  
 <T<sub>1</sub>, C, 700, 600>  
 <T<sub>1</sub> commit

Recovery action:

redo T<sub>0</sub> and redo T<sub>1</sub>

- Set A to 950
- Set B to 2050
- Set C to 600

# Checkpoints

Streamline recovery procedure by periodically performing **checkpointing**

Output all log records currently residing in main memory onto stable storage.

Output all modified buffer blocks to the disk.

Write a log record <checkpoint L> onto stable storage where  $L$  is a list of all transactions active at the time of checkpoint.

All updates are stopped while doing checkpointing

semua perubahan sblm checkpoint udah  
tempat ke secondary storage



# Checkpoints – On Recovery

During recovery we need to consider only the most recent transaction  $T_i$  that started before the checkpoint, and transactions that started after  $T_i$ .

- Scan backwards from end of log to find the most recent <checkpoint L> record
- Only transactions that are in L or started after the checkpoint need to be redone or undone
- Transactions that committed or aborted before the checkpoint already have all their updates output to stable storage.

Some earlier part of the log may be needed for undo operations

- Continue scanning backwards till a record  $<T_i \text{ start}>$  is found for every transaction  $T_i$  in  $L$ .
- Parts of log prior to earliest  $<T_i \text{ start}>$  record above are not needed for recovery and can be erased whenever desired.

Karena blm abort / commit



## Example of Checkpoints

$T_1$  can be ignored

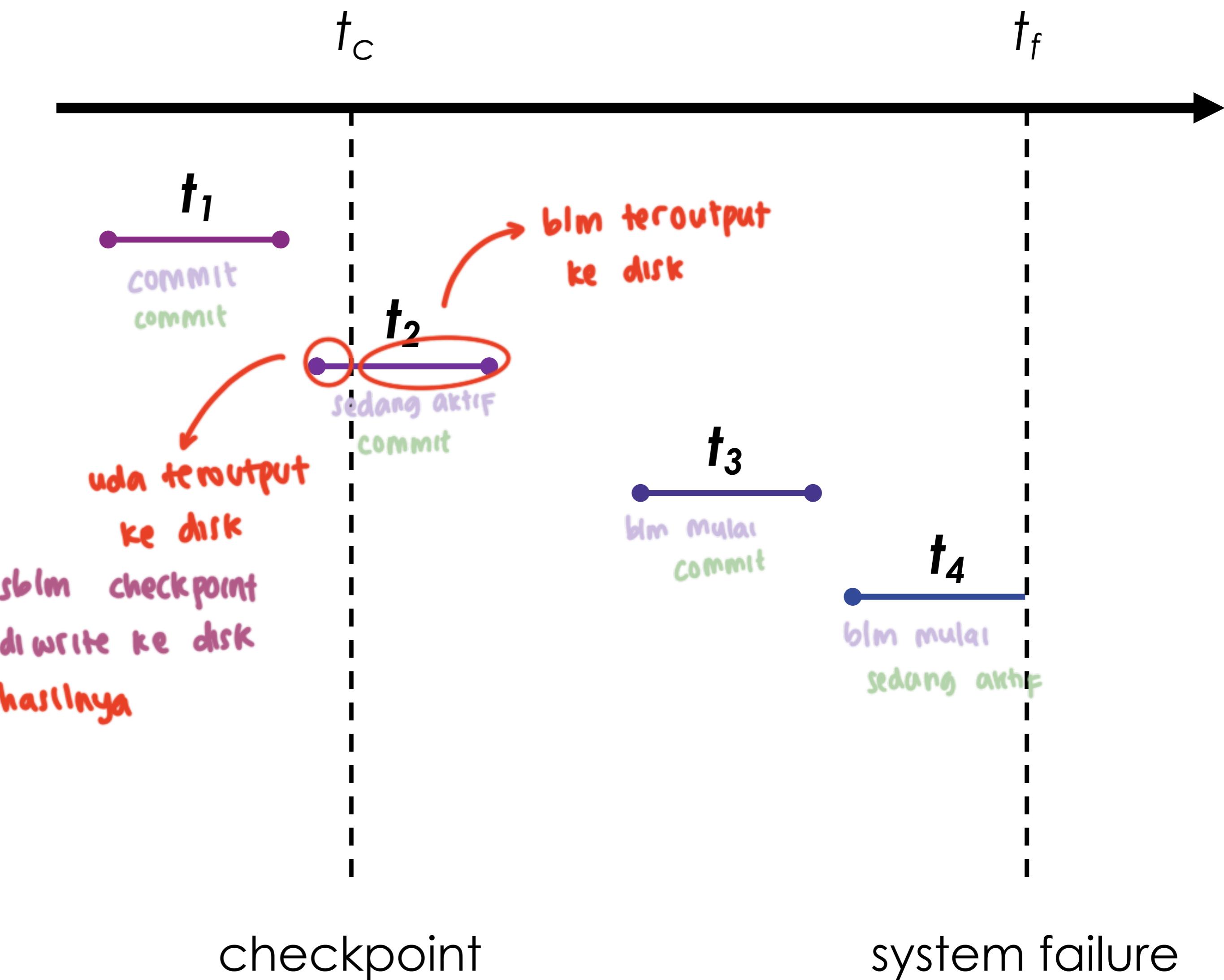
$T_2$  and  $T_3$  redone.

$T_4$  undone

→ krn  $T_1$  uda commit slm checkpoint  
shg hasilnya uda diwrite ke disk  
tdk ada jaminan hasilnya  
teroutput ke disk

dia aktif saat failure,  
mk hrs diabort

- kondisi pas checkpoint
- kondisi pas failure



# *Recovery Algorithm*

**So far:** we covered key concepts

**Now:** we present the components of the basic recovery algorithm

**Later:** we present extensions to allow more concurrency



KNOWLEDGE & SOFTWARE ENGINEERING

# Recovery Algorithm

bukan krn kegagalan sistem

## Logging (during normal operation)

- $\langle T_i \text{ start} \rangle$  at transaction start
- $\langle T_i, X_j, V_1, V_2 \rangle$  for each update
- $\langle T_i \text{ commit} \rangle$  at transaction end

## Transaction rollback (during normal operation)

- Scan log backwards from the end, and for each log record of  $T_i$  of the form  $\langle T_i, X_j, V_1, V_2 \rangle$ 
  - Perform the undo by writing  $V_1$  to  $X_j$ ,  $X_j$  jadi  $V_1$
  - Write a log record  $\langle T_i, X_j, V_1 \rangle$
  - such log records are called **compensation log records** = redo only log records
- Once the record  $\langle T_i \text{ start} \rangle$  is found stop the scan and write the log record  $\langle T_i \text{ abort} \rangle$



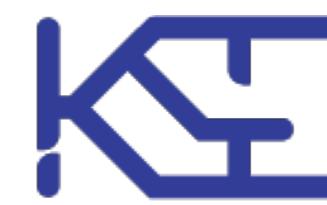
# *Recovery Algorithm – Recovery from Failure*

Redo  
phase

- replay updates of all transactions, whether they committed, aborted, or are incomplete

Undo  
phase

- undo all incomplete transactions

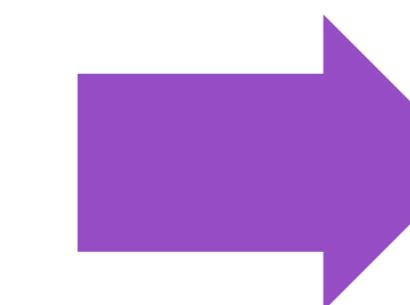


KNOWLEDGE & SOFTWARE ENGINEERING

# Recovery Algorithm – Recovery from Failure

**Redo  
phase**

Find last  
**<checkpoint L>** record,  
and set undo-list to L.



**Undo  
phase**

kala udah selesai redo,  
isinya undo-list hanya  
tran saku : yg incomplete

Scan forward from  
above **<checkpoint L>**  
record

- Whenever a record  $\langle T_i, X_j, V_1, V_2 \rangle$  or  $\langle T_i, X_j, V_2 \rangle$  is found, redo it by writing  $V_2$  to  $X_j$  *(X<sub>j</sub> ja V<sub>2</sub>)*
- Whenever a log record  $\langle T_i \text{ } \underline{\text{start}} \rangle$  is found, add  $T_i$  to undo-list
- Whenever a log record  $\langle T_i \text{ } \underline{\text{commit}} \rangle$  or  $\langle T_i \text{ } \underline{\text{abort}} \rangle$  is found, remove  $T_i$  from undo-list



# Recovery Algorithm – Recovery from Failure



Scan log backwards from end

Whenever a log record  $\langle T_i, X_j, V_1, V_2 \rangle$  is found where  $T_i$  is in undo-list  $\rightarrow$  transaction rollback

- perform undo by writing  $V_1$  to  $X_j$ .
  - write a log record  $\langle T_i, X_j, V_1 \rangle$
- $X_j \leftarrow V_1$

Whenever a log record  $\langle T_i \text{ start} \rangle$  is found where  $T_i$  is in undo-list,

- Write a log record  $\langle T_i \text{ abort} \rangle$
- Remove  $T_i$  from undo-list

Stop when undo-list is empty, i.e.,  $\langle T_i \text{ start} \rangle$  has been found for every transaction in undo-list

After undo phase completes, normal transaction processing can commence

# Example of Recovery



KNOWLEDGE & SOFTWARE ENGINEERING

