

# **EEL-4930/5934**

# **Principles of Computer System Design**

Lecture Slides 21  
Textbook Chapter 9  
Atomicity Logs

# Introduction

- Addressing performance implications associated with all-or-nothing actions
- Key idea:
  - Optimize for the common case (reads, writes)
  - Recovery must be possible, but it is the uncommon case – not the main focus of performance optimization

# Atomicity Logs

- How to benefit from
  - All-or-nothing of journal storage
  - Performance of cell storage
- Write data twice
  - Journal storage (Log) – atomicity
  - Cell storage (install) – fast reads
- Don't you just double the work?

# Logging

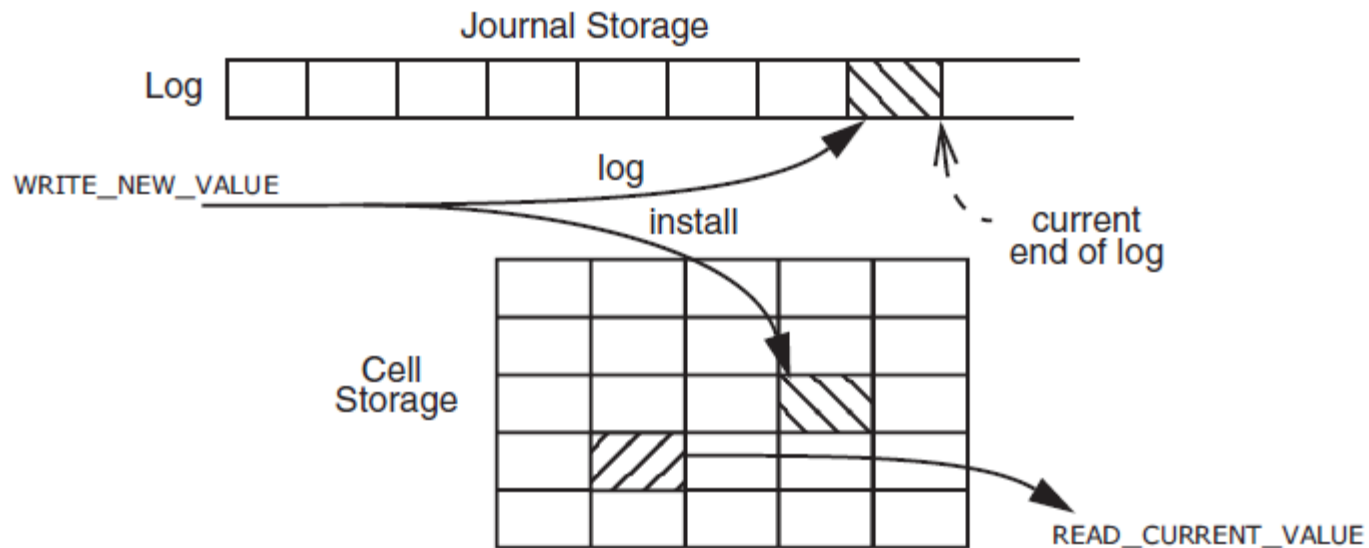
- Optimize writing of multiple objects by mapping to single log
- Key benefit: writes are sequential
  - Single pointer for all updates
  - Dedicated disk storage for logging: minimize seek times, significantly improving performance
- Drawback
  - Histories of different objects are interleaved
  - Longer lists - lookup is slower; recovery however is uncommon

# Cell storage

- “Installing” a new version is an overwrite
- Faster subsequent reads
- No concern about overwriting a single copy, if it’s been already logged
  - Cache that allows bypassing overhead of traversing log on reads
  - Also avoids seek overhead by keeping the arm in place at log’s end
- Authoritative copy always in log
  - Cell storage – can always be reconstructed from log

# Primitives

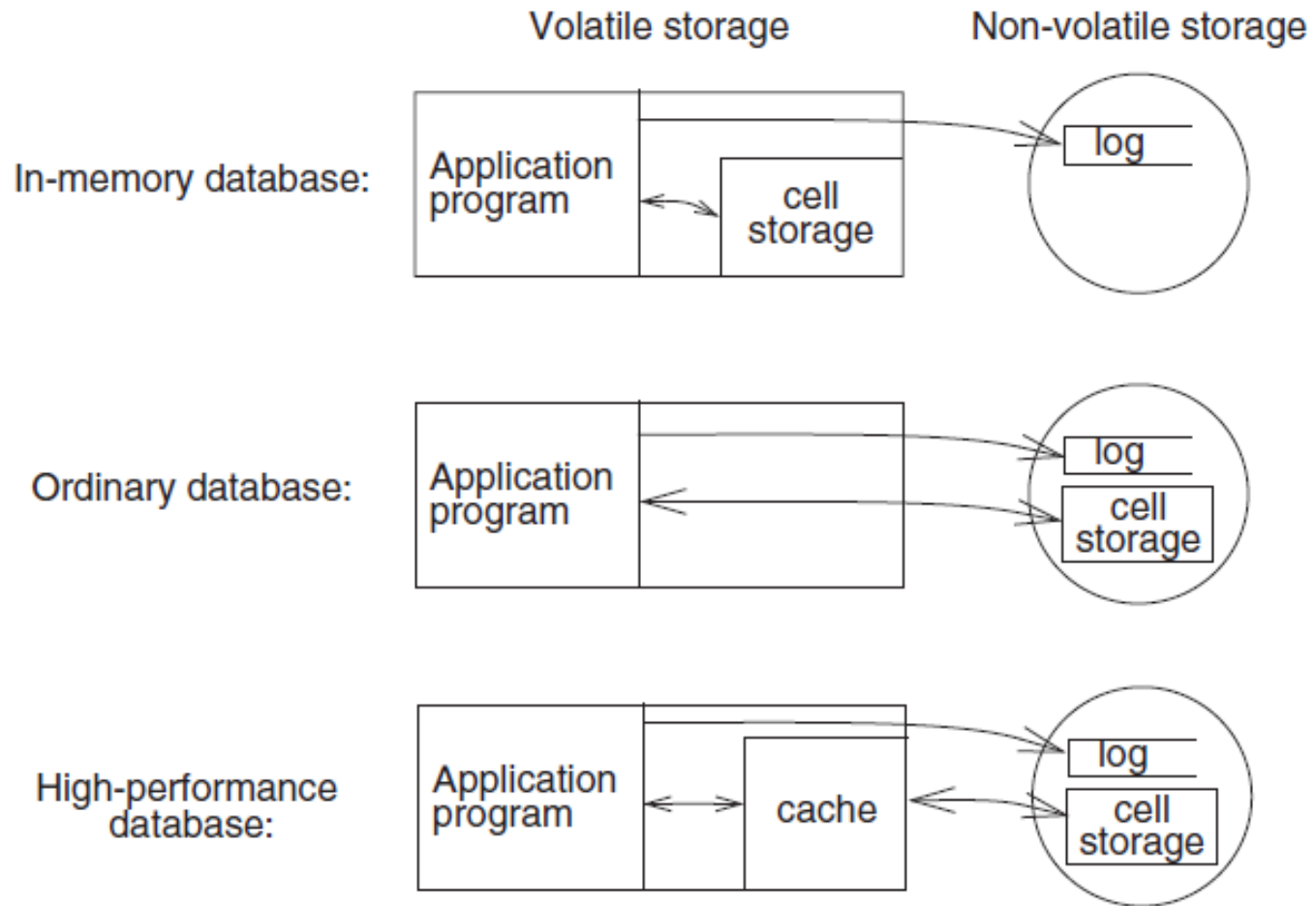
- WRITE\_NEW\_VALUE()
  - LOG
    - Append entry to end of log
  - INSTALL
    - Overwrite cell storage with new value



# Summary

- Writing twice can *improve* performance
  - Reduce seek overhead on log operations for WRITE\_NEW\_VALUE
  - Faster READ\_CURRENT\_VALUE – no need to traverse log
- Writing twice *maintains* all-or-nothing nature of operations
  - Log is written all-or-nothing, as we did before
  - Cell storage is just a copy

# Logging strategies





# Issues

- Failures can happen between writes
  - Log should be authoritative copy
  - Cell storage may not be properly installed
  - Note that cell storage may be non-volatile
- Enforce ordering
  - Log first, install after logged
  - Write-ahead log
- Run recovery procedure after crash
  - Recall check\_and\_repair

# Logging protocols

- Basic element – *log record*
- Record type (e.g. CHANGE)
- ID of all-or-nothing action performing update
- Action that, if performed, installs intended value in cell storage (“redo”)
- Action that, if performed, reverses the effect in cell storage (“undo”)

# Redo/Undo actions

- Redo:
  - Part of the log record that is invoked when a recovery procedure needs to install committed data into cell storage
  - Crash that happens before a committed value gets to be installed
- Undo:
  - Part of the log record that is invoked when an action aborts after installing into cell storage (possibly during recovery procedure)

# Primitives

- LOG:
  - Appends log record to log storage
  - Must itself be all-or-nothing
  - Again, can bootstrap from single-sector primitive (ALL\_OR\_NOTHING\_PUT)
- WRITE\_NEW\_VALUE:
  - LOG then INSTALL (PUT)
- READ\_CURRENT\_VALUE:
  - Read (GET) from cell storage

# Primitives

- **NEW\_ACTION:**
  - Build upon LOG primitive
  - Logs a BEGIN record that establishes the new action's identity
- Action's writes log CHANGE records through its pre-commit phase
- **COMMIT/ABORT:**
  - Log OUTCOME records
  - Commit point – becomes authoritative outcome or action

# Account transfer revisited

```
1 procedure TRANSFER (debit_account, credit_account, amount)
2   my_id ← LOG (BEGIN_TRANSACTION)
3   dbvalue.old ← GET (debit_account)
4   dbvalue.new ← dbvalue.old - amount
5   crvalue.old ← GET (credit_account, my_id)
6   crvalue.new ← crvalue.old + amount
7   LOG (CHANGE, my_id,
8     "PUT (debit_account, dbvalue.new)",           //redo action
9     "PUT (debit_account, dbvalue.old)" )          //undo action
10  LOG ( CHANGE, my_id,
11    "PUT (credit_account, crvalue.new)"           //redo action
12    "PUT (credit_account, crvalue.old)" )          //undo action
13  PUT (debit_account, dbvalue.new)                 // install
14  PUT (credit_account, crvalue.new)                 // install
15  if dbvalue.new > 0 then
16    LOG ( OUTCOME, COMMIT, my_id)
17  else
18    LOG (OUTCOME, ABORT, my_id)
19    signal("Action not allowed. Would make debit account negative.")
20  LOG (END_TRANSACTION, my_id)
```

# Notes

- INSTALL of a value
  - Must occur after LOG of a CHANGE
  - But not necessarily after LOG of OUTCOME
- BEGIN, CHANGE, OUTCOME, END
  - Most general
  - Implementations may combine (first CHANGE means BEGIN; OUTCOME +END)

# Log records

...	<i>type</i> : CHANGE <i>action_id</i> : 9979 <i>redo_action</i> : PUT( <i>debit_account</i> , \$90) <i>undo_action</i> : PUT( <i>debit_account</i> , \$120)	<i>type</i> : OUTCOME <i>action_id</i> : 9974 <i>status</i> : COMMITTED	<i>type</i> : CHANGE <i>action_id</i> : 9979 <i>redo_action</i> : PUT( <i>credit_account</i> , \$40) <i>undo_action</i> : PUT( <i>credit_account</i> , \$10)
← older log records			newer log records →



# ABORT Procedure

- An action may install values in cell storage that need to be restored
  - Not necessarily by thread that started action (e.g. thread deadlocks)
  - Can't leave an action pending indefinitely
- ABORT procedure
  - Scan log backwards for aborted action
  - Perform logged “undo” action for every CHANGE log record
  - Stop at BEGIN record
- Expensive (but uncommon)

# ABORT procedure

```
1  procedure ABORT (action_id)
2      starting at end of log repeat until beginning
3          log_record ← previous record of log
4          if log_record.id = action_id then
5              if (log_record.type = OUTCOME)
6                  then signal ("Can't abort an already completed action.")
7              if (log_record.type = CHANGE)
8                  then perform undo_action of log_record
9              if (log_record.type = BEGIN)
10                 then break repeat
11  LOG (action_id, OUTCOME, ABORTED)           // Block future undos.
12  LOG (action_id, END)
```

# Recovery Procedure

- Must be executed after any crash and before application uses data
- Bring cell storage contents in accordance to log
- System can crash during recovery procedure itself
  - Idempotent recovery
- Cell storage in volatile memory simplifies recovery

# Recovery for in-memory db

```
1  procedure RECOVER () // Recovery procedure for a volatile, in-memory database.
2      winners ← NULL
3      starting at end of log repeat until beginning
4          log_record ← previous record of log
5          if (log_record.type = OUTCOME)
6              then winners ← winners + log_record                // Set addition.

7      starting at beginning of log repeat until end
8          log_record ← next record of log
9          if (log_record.type = CHANGE)
10             and (outcome_record ← find (log_record.action_id) in winners)
11             and (outcome_record.status = COMMITTED) then
12                 perform log_record.redo_action
```

Backward scan looks for records with OUTCOME

Forward scan applies redo actions

As if all all-or-nothing actions that committed before the crash had run to completion, and as if those that didn't commit (aborted, pending) never existed

# Recovery for non-volatile cell

- Cell storage in volatile memory simplifies recovery
  - All cell storage data is lost on a crash; start from clean slate
- However, the database may be too large to fit in memory
- In general, need cell storage to also be in non-volatile storage
  - Performance enhancement – a cache in volatile memory

# Recovery for non-volatile cell

- The additional complexity now is that installs persist crashes
  - There may be all-or-nothing actions that installed changes, but were still pending when the system crashed
  - Those changes will persist – recovery needs to undo them
- Furthermore, reinstalling the whole database in cell storage from the log upon recovery would be too slow for large databases

# Modified recovery procedure

- Scan backwards, but instead of looking for “winners” (actions that had an outcome), build a set of “losers”
  - Those that were pending
  - I.e., on backward scan, if the first record found for an action is *not* END, it is a loser
- Perform undo action for loser set
  - Mark them as ENDED so to not undo in future recoveries
- Then, forward scan and perform redo of committed actions

# Recovery

```
1  procedure RECOVER ()// Recovery procedure for non-volatile cell memory
2      completeds ← NULL
3      losers ← NULL
4      starting at end of log repeat until beginning
5          log_record ← previous record of log
6          if (log_record.type = END)
7              then completeds ← completeds + log_record           // Set addition.
8          if (log_record.action_id is not in completeds) then
9              losers ← losers + log_record           // Add if not already in set.
10             if (log_record.type = CHANGE) then
11                 perform log_record.undo_action

12     starting at beginning of log repeat until end
13         log_record ← next record of log
14         if (log_record.type = CHANGE)
15             and (log_record.action_id.status = COMMITTED) then
16                 perform log_record.redo_action

17     for each log_record in losers do
18         log (log_record.action_id, END)           // Show action completed.
```



# Improvement

- Scanning forward and redoing every install – slow
- Cannot know if all installs have been done (even by committed action) until an END record is found
  - Write-ahead log protocol specifies all installs must be done when an END is logged
- Any committed action that has logged END has completed its installs
  - Can modify code from previous page to avoid those installs

# Rollback logging

- Can avoid the need for redoing any installs if application is required to perform all of its installs before it logs an OUTCOME record
- Recovery procedure needs to only undo installs of losers, and skip forward scan

# Rollback recovery

```
1  procedure RECOVER ()                // Recovery procedure for rollback recovery.
2      completeds ← NULL
3      losers ← NULL
4      starting at end of log repeat until beginning        // Perform undo scan.
5          log_record ← previous record of log
6          if (log_record.type = OUTCOME)
7              then completeds ← completeds + log_record    // Set addition.
8          if (log_record.action_id is not in completeds) then
9              losers ← losers + log_record                // New loser.
10             if (log_record.type = CHANGE) then
11                 perform log_record.undo_action

12  for each log_record in losers do
13      log (log_record.action_id, OUTCOME, ABORT)        // Block future undos.
```

# Before-or-after

- With logs, when an action installs data in cell storage, it is visible
  - Intermediate steps would be visible to other concurrent threads
- Our assumption - single thread
  - Data used later in time (after commit) or, on a crash, after recovery process
- Supporting multiple threads requires dealing with before-or-after atomicity

# Simple serialization

- Assign monotonically-increasing transaction IDs
- Transaction  $n+1$  does not read/write data until transaction  $n$  has committed or aborted
- Enforces that *transactions* run in serial order
  - Threads can still execute concurrently outside transaction boundaries

# Simple serialization

- Simple and correct; conservative
  - Transactions  $n$ ,  $n-1$  may work on completely independent data, yet one blocks the other
  - If higher performance is needed, there are alternatives – more efficient and complex

```
1 procedure BEGIN_TRANSACTION ()  
2   id ← NEW_OUTCOME_RECORD (PENDING)           // Create, initialize, assign id.  
3   previous_id ← id - 1  
4   wait until previous_id.outcome_record.state ≠ PENDING  
5   return id
```

# Opportunities for Concurrency

Object		value of object at end of transaction					
		1	2	3	4	5	6
Object ↓ A  B  C  D	A	0	+10		+12		0
	B	0	-10	-6		-12	-2
	C	0		-4		+2	
	D	0			-2		
	outcome						
	record state	Committed	Committed	Committed	Aborted	Committed	Pending

transaction

- 1: initialize all accounts to 0
- 2: transfer 10 from *B* to *A*
- 3: transfer 4 from *C* to *B*
- 4: transfer 2 from *D* to *A* (aborts)
- 5: transfer 6 from *B* to *C*
- 6: transfer 10 from *A* to *B*

# Opportunities for Concurrency

Object		value of object at end of transaction					
		1	2	3	4	5	6
↓ A  B  C  D	A	0	+10		+12		0
	B	0	-10	-6		-12	-2
	C	0		-4		+2	
	D	0			-2		
	outcome	Committed	Committed	Committed	Aborted	Committed	Pending
	record state						

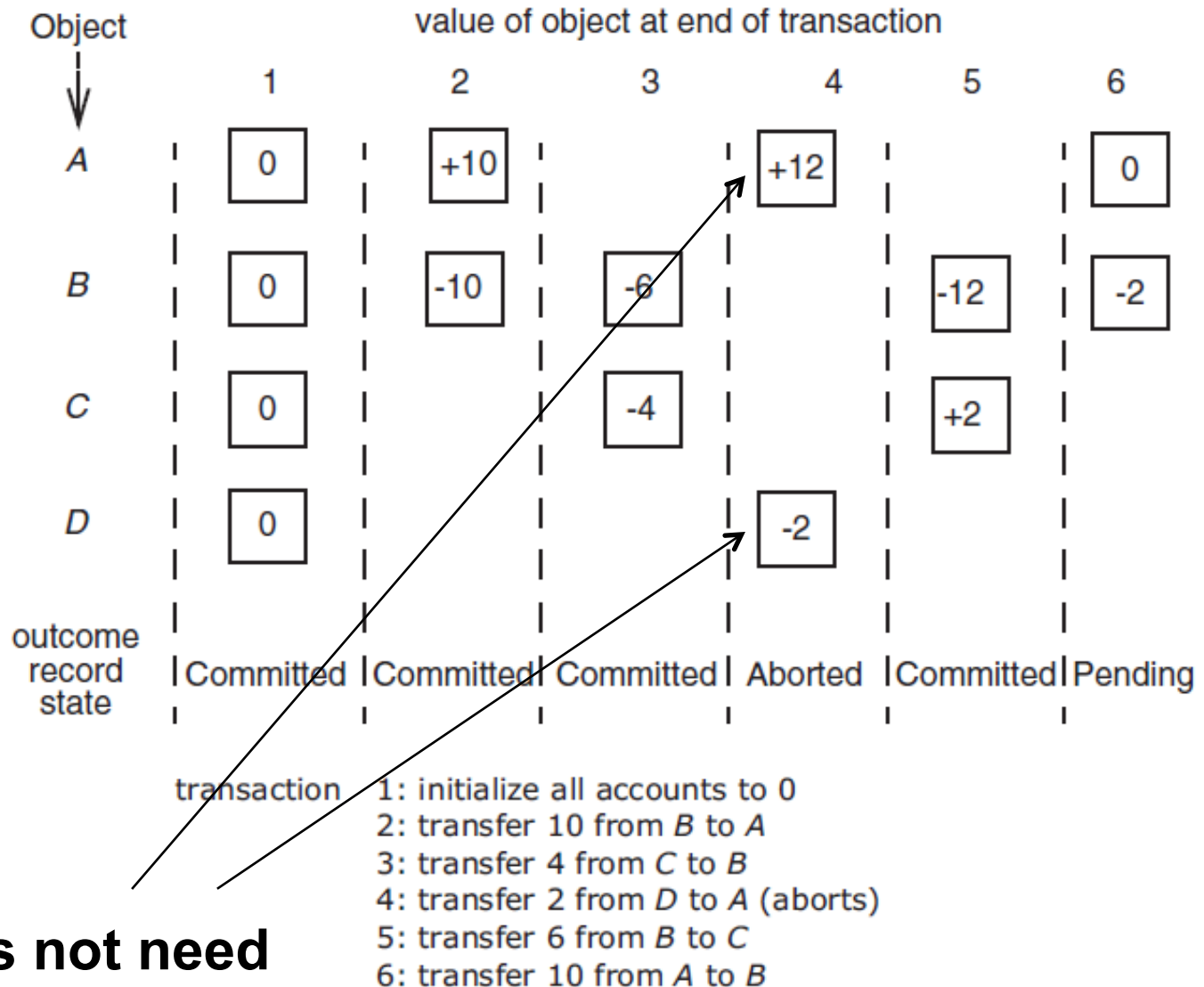
transaction

- 1: initialize all accounts to 0
- 2: transfer 10 from B to A
- 3: transfer 4 from C to B
- 4: transfer 2 from D to A (aborts)
- 5: transfer 6 from B to C
- 6: transfer 10 from A to B

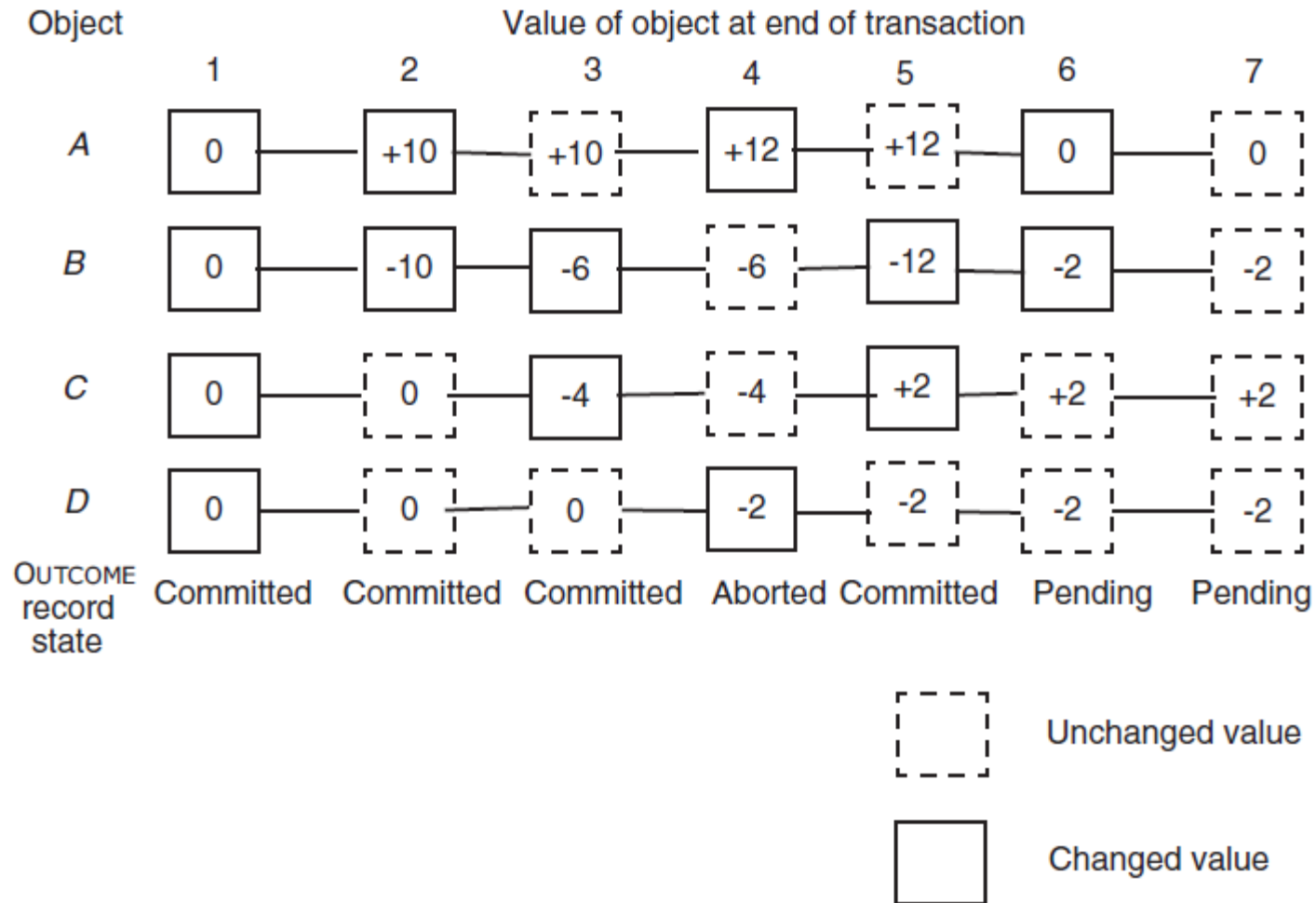
3 must wait for 2



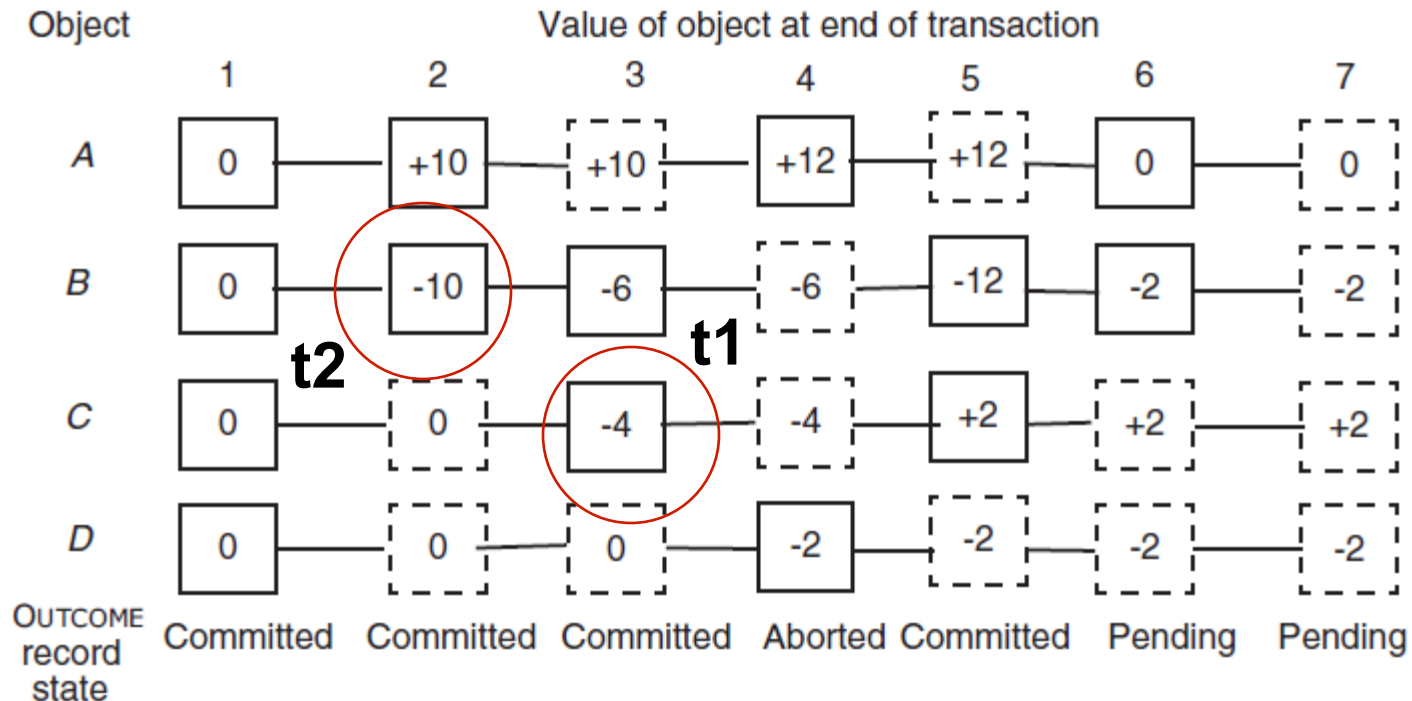
# Opportunities for Concurrency



# Opportunities for Concurrency

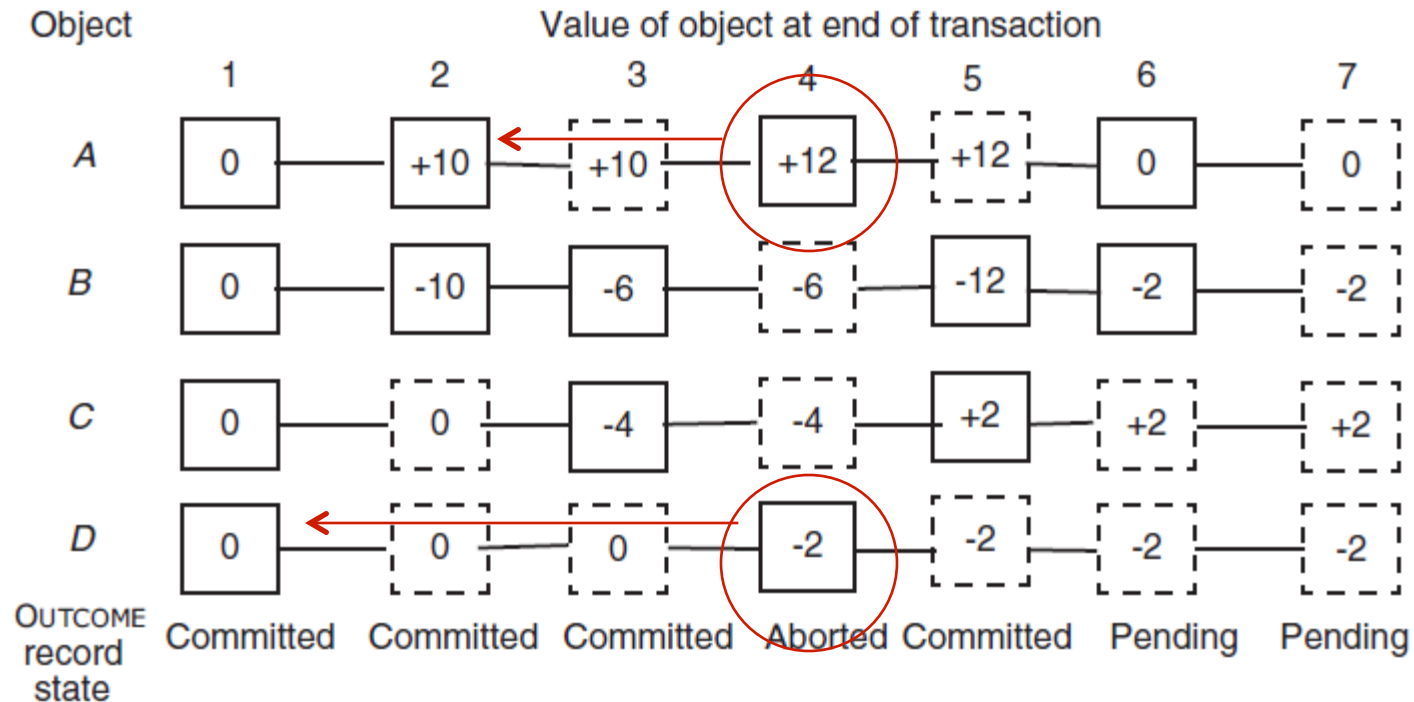


# Opportunities for Concurrency



**Real-time in which values are placed in boxes does not need to be strictly serial – as long as the results are the same as one would get with serial execution**

# Opportunities for Concurrency



**Transaction 4 “looking backwards”  
needs to wait for 1, 2 to commit  
(not 3)**

# Mark Point Discipline

- Recall version history/journal discussion
  - READ\_CURRENT\_VALUE should not expose any pending, non-committed value
    - All values should be exposed at once at commit
    - No intermediate values visible (expect for self)
      - READ\_MY\_CURRENT\_VALUE
- Additional requirement:
  - A thread must *wait* until a pending value from an earlier thread that it needs to read is committed

# Mark Point Discipline

- READ\_CURRENT\_VALUE
  - Skip “later” transaction (larger IDs)
    - Concurrent execution - later threads may be creating their own versions
  - Wait for pending values from “earlier” transactions (smaller IDs)
- How does a transaction know which pending values to wait for?
  - Explicitly programmed in the application

# Mark Point Discipline

- Each transaction explicitly creates new, pending versions of each object it intends to modify
  - “Marking” an object
  - WRITE\_NEW\_VALUE: NEW\_VERSION (mark), WRITE\_VALUE (update)
- Transaction also announces that it has done so (for *all* objects)
  - The mark point
  - MARK\_POINT\_ANNOUNCE

# Mark Point Discipline

- Enforce mark point discipline at the beginning of a transaction
- BEGIN\_TRANSACTION – wait until preceding transaction has reach mark point (or no longer pending)
- By design – transactions must reach mark point in transaction ID order
  - But they can overlap



# BEGIN / MARK

```
1 procedure BEGIN_TRANSACTION ()
2   id ← NEW_OUTCOME_RECORD (PENDING)
3   previous_id ← id - 1
4   wait until (previous_id.outcome_record.mark_state = MARKED)
5     or (previous_id.outcome_record.state ≠ PENDING)
6   return id
```

**Block until previous ID  
marked (or not pending)**

```
7 procedure NEW_OUTCOME_RECORD (starting_state)
8   ACQUIRE (outcome_record_lock)
9   id ← TICKET (outcome_record_sequencer)
10  allocate id.outcome_record
11  id.outcome_record.state ← starting_state
12  id.outcome_record.mark_state ← NULL
13  RELEASE (outcome_record_lock)
14  return id
```

**Bootstrap before-  
or-after**

**Dispense one  
monotonically  
increasing  
outcome record  
at a time**

```
15 procedure MARK_POINT_ANNOUNCE (reference this_transaction_id)
16   this_transaction_id.outcome_record.mark_state ← MARKED
```

**More state  
in outcome record**

# WRITE\_NEW\_VALUE

```
1 procedure NEW_VERSION (reference data_id, this_transaction_id)
2   if this_transaction_id.outcome_record.mark_state = MARKED then
3     signal ("Tried to create new version after announcing mark point!")
4   append new version v to data_id
5   v.value ← NULL
6   v.action_id ← transaction_id
```


**Just a placeholder  
(before mark point)**

```
7 procedure WRITE_VALUE (reference data_id, new_value, this_transaction_id)
8   starting at end of data_id repeat until beginning
9     v ← previous version of data_id
10    if v.action_id = this_transaction_id
11      v.value ← new_value
12    return
13  signal ("Tried to write without creating new version!")
```

**Actual write  
(after mark point)**

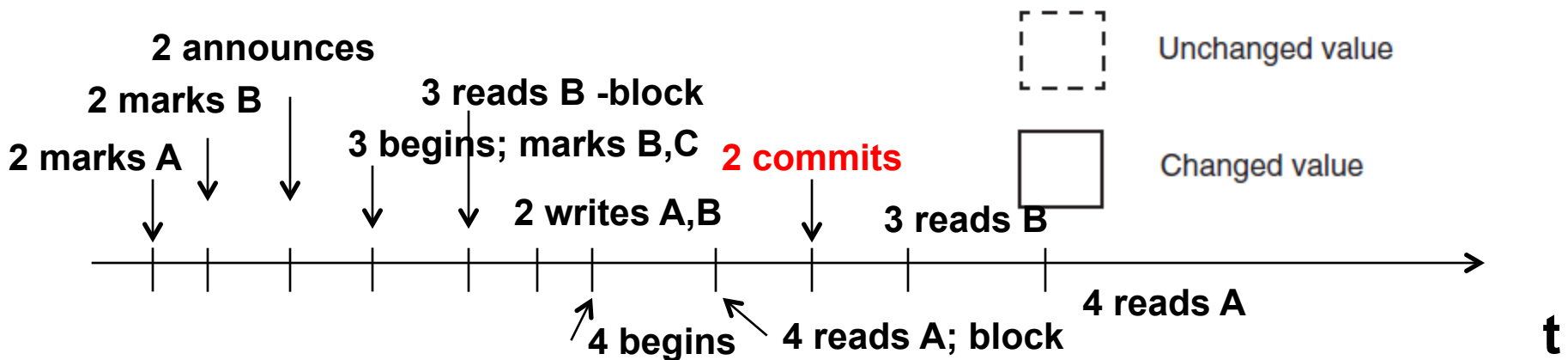
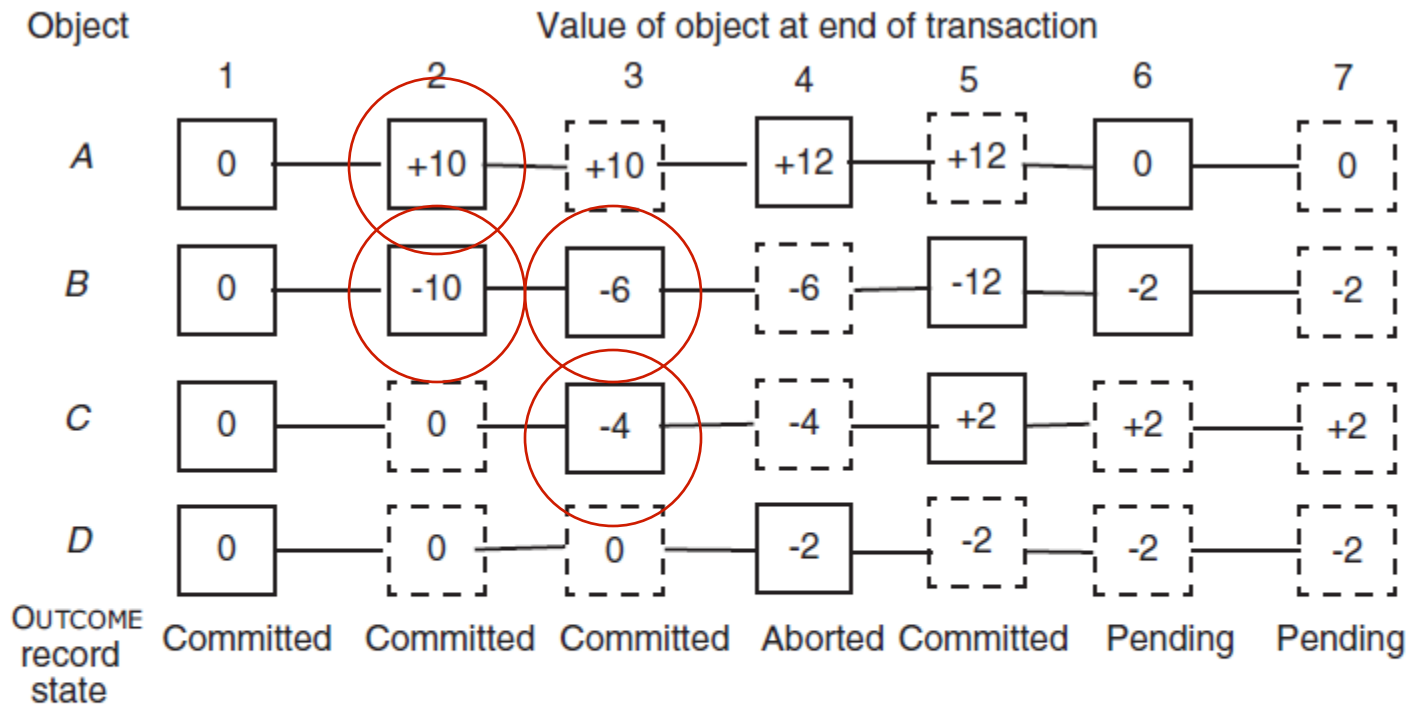
# READ\_CURRENT\_VALUE

```
1 procedure READ_CURRENT_VALUE (data_id, this_transaction_id)
2   starting at end of data_id repeat until beginning
3     v ← previous version of data_id
4     last_modifier ← v.action_id
5     if last_modifier ≥ this_transaction_id then skip v           // Keep searching
6     wait until (last_modifier.outcome_record.state ≠ PENDING)
7     if (last_modifier.outcome_record.state = COMMITTED)
8       then return v.state
9       else skip v                                           // Resume search
10  signal ("Tried to read an uninitialized variable")
```



**Block until outcome  
Record no longer pending  
(e.g. from a previous ID's  
NEW\_VERSION)**

# Opportunities for Concurrency



# Updated transfer function

```
1  procedure TRANSFER (reference debit_account, reference credit_account,  
2                               amount)  
3      my_id ← BEGIN_TRANSACTION ()  
4      NEW_VERSION (debit_account, my_id)  
5      NEW_VERSION (credit_account, my_id)  
6      MARK_POINT_ANNOUNCE (my_id);  
7      xvalue ← READ_CURRENT_VALUE (debit_account, my_id)  
8      xvalue ← xvalue - amount  
9      WRITE_VALUE (debit_account, xvalue, my_id)  
10     yvalue ← READ_CURRENT_VALUE (credit_account, my_id)  
11     yvalue ← yvalue + amount  
12     WRITE_VALUE (credit_account, yvalue, my_id)  
13     if xvalue > 0 then  
14         COMMIT (my_id)  
15     else  
16         ABORT (my_id)  
17         signal("Negative transfers are not allowed.")
```

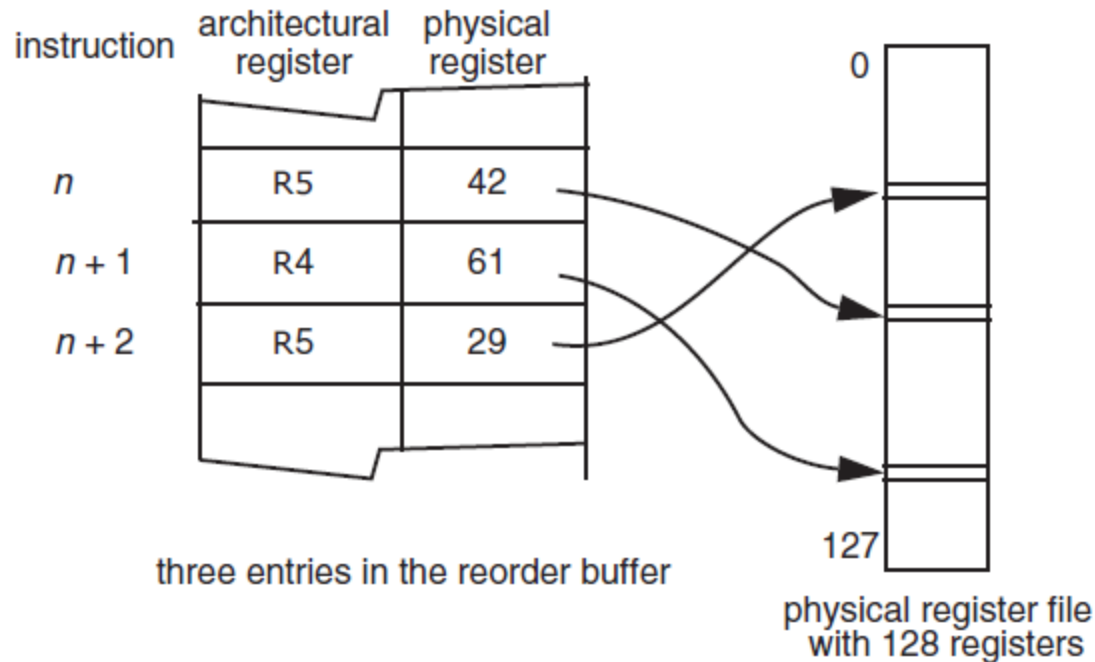
# Example – architecture

- Register renaming
- Legacy architectures have few registers
  - E.g. x86 – 8 registers
- For high instruction-level parallelism
  - Need to overlap many instructions
  - Implementation – 128 physical registers
- Reorder buffer maps architecture register (1 of 8) to physical register (1 out of 128) holding actual value

# Example- architecture

- Instruction issues – assigned next sequential slot in reorder buffer
  - Begin transaction
  - NEW\_OUTCOME\_RECORD (its position in the buffer) and NEW\_VERSION
- Commits
  - Write to physical register
  - WRITE\_NEW\_VALUE and COMMIT

# Register Renaming



**FIGURE 9.34**

Example showing how a reorder buffer maps architectural register numbers to physical register numbers. The program sequence corresponding to the three entries is:

```
 $n$       R5  $\leftarrow$  R4  $\times$  R2      // Write a result in register five.  
 $n + 1$  R4  $\leftarrow$  R5 + R1      // Use result in register five.  
 $n + 2$  R5  $\leftarrow$  READ (117492)  // Write content of a memory cell in register five.
```

Instructions  $n$  and  $n + 2$  both write into register R5, so R5 has two versions, with mappings to physical registers 42 and 29, respectively. Instruction  $n + 2$  can thus execute concurrently with instructions  $n$  and  $n + 1$ .