

CS162 Operating Systems and Systems Programming Lecture 18 TCP's Flow Control, Transactions

October 31, 2012
Ion Stoica
<http://inst.eecs.berkeley.edu/~cs162>

Goals of Today's Lecture

- TCP flow control (continued)
- Transactions (ACID semantics)

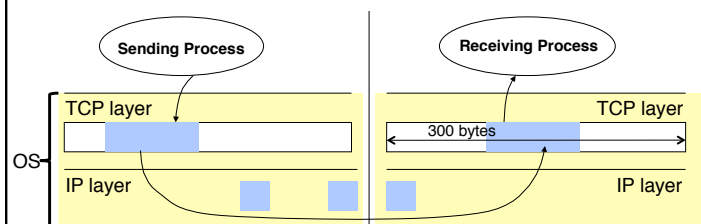
Note: Some slides and/or pictures in the following are adapted from lecture notes by Mike Franklin.

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.2

TCP Flow Control



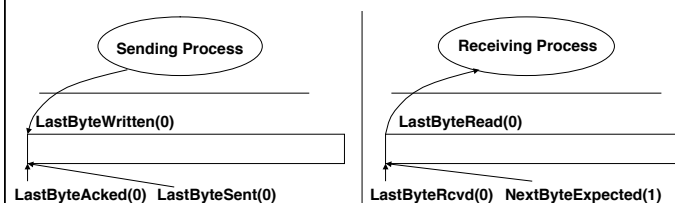
- Example assumptions:
 - Maximum IP packet size = 100 bytes
 - Size of the receiving buffer (MaxRcvBuf) = 300 bytes
- Recall, ack indicates the **next expected byte** in-sequence, not the last received byte
- Use circular buffers

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.3

TCP Flow Control



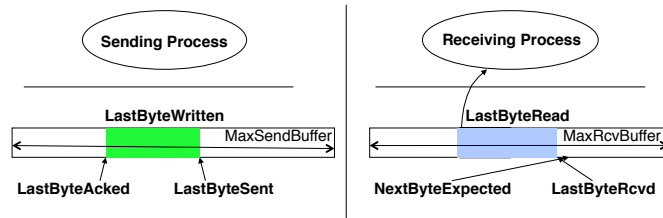
- LastByteWritten: last byte written by sending process
- LastByteSent: last byte sent by sender to receiver
- LastByteAcked: last ack received by sender from receiver
- LastByteRcvd: last byte received by receiver from sender
- NextByteExpected: last **in-sequence** byte expected by receiver
- LastByteRead: last byte read by the receiving process

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.4

TCP Flow Control



- AdvertisedWindow: number of bytes TCP receiver can receive

$$\text{AdvertisedWindow} = \text{MaxRcvBuffer} - (\text{LastByteRcvd} - \text{LastByteRead})$$

- SenderWindow: number of bytes TCP sender can send

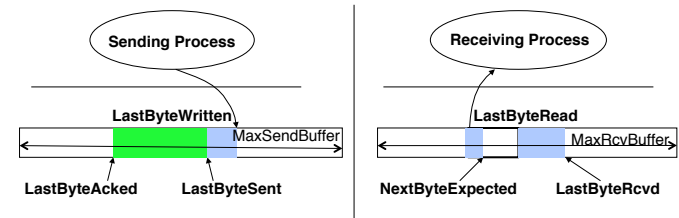
$$\text{SenderWindow} = \text{AdvertisedWindow} - (\text{LastByteSent} - \text{LastByteAcked})$$

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.5

TCP Flow Control



- Still true if receiver missed data....

$$\text{AdvertisedWindow} = \text{MaxRcvBuffer} - (\text{LastByteRcvd} - \text{LastByteRead})$$

- WriteWindow: number of bytes sending process can write

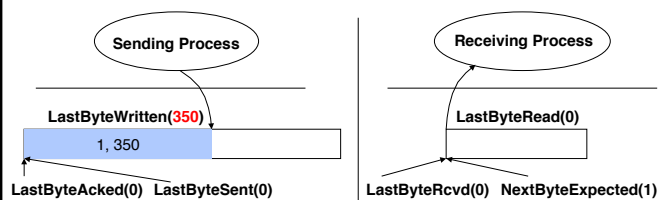
$$\text{WriteWindow} = \text{MaxSendBuffer} - (\text{LastByteWritten} - \text{LastByteAcked})$$

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.6

TCP Flow Control



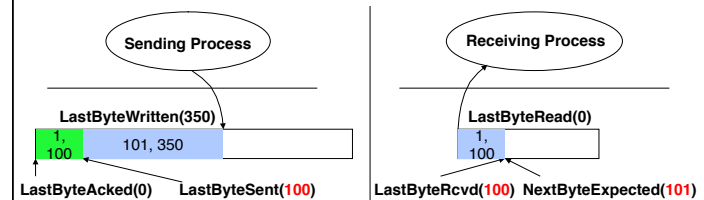
- Sending app sends 350 bytes
- Recall:
 - We assume IP only accepts packets no larger than 100 bytes
 - MaxRcvBuf = 300 bytes, so initial Advertised Window = 300 bytes

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.7

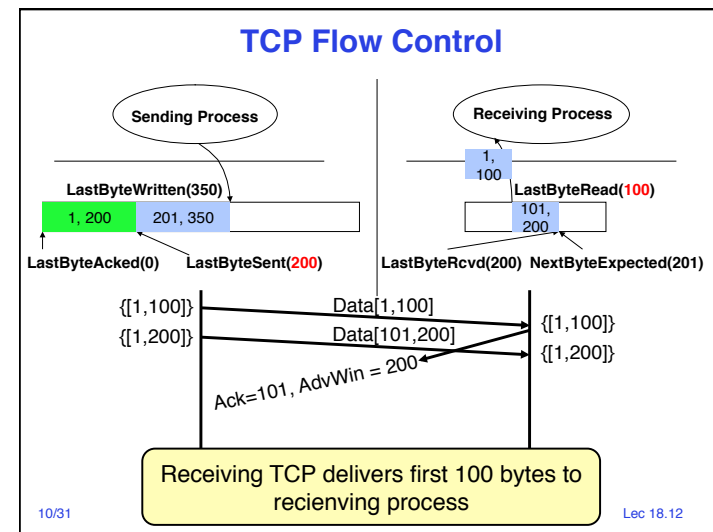
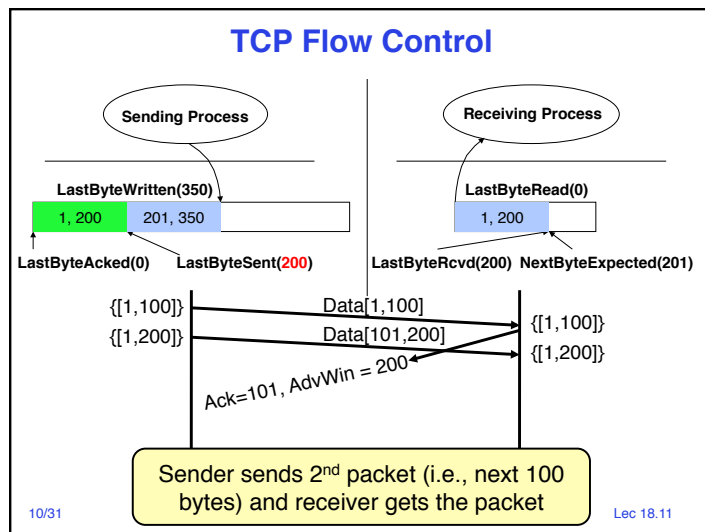
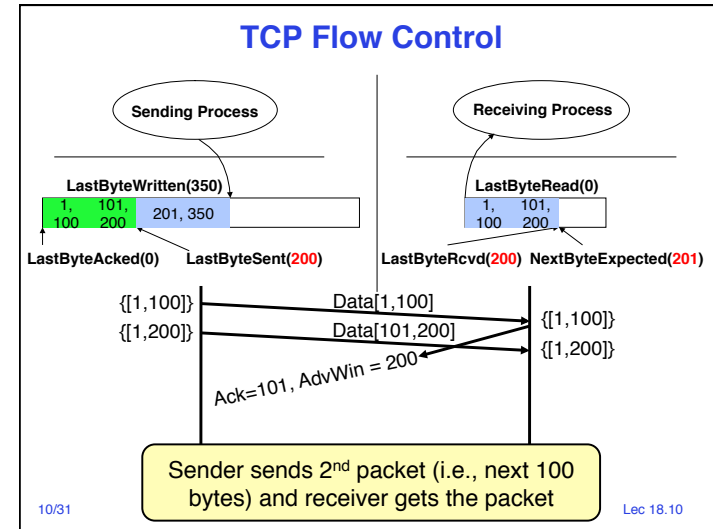
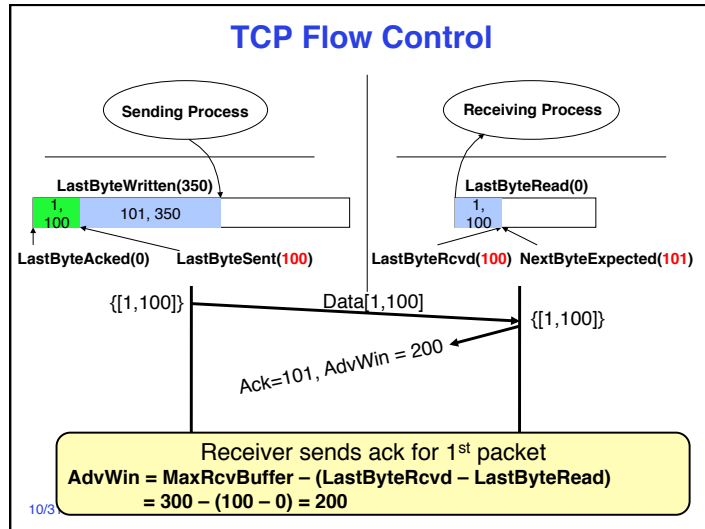
TCP Flow Control

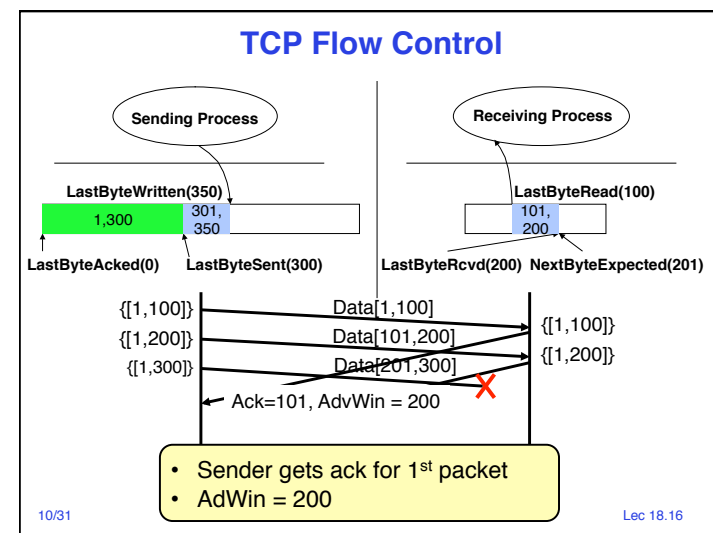
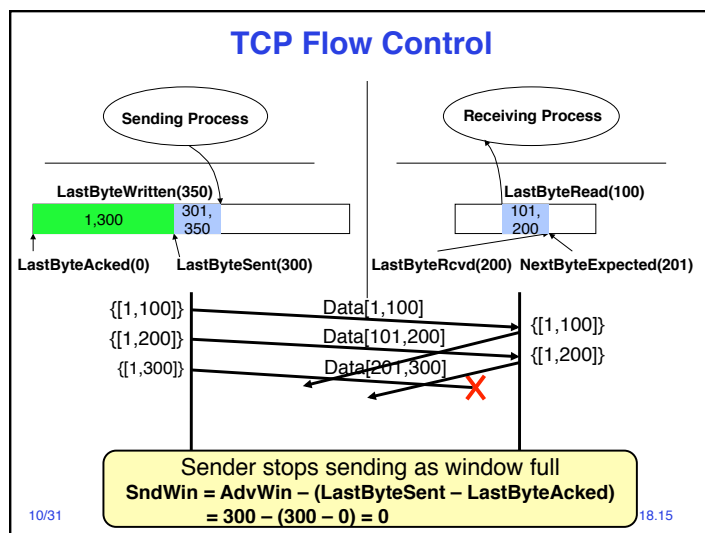
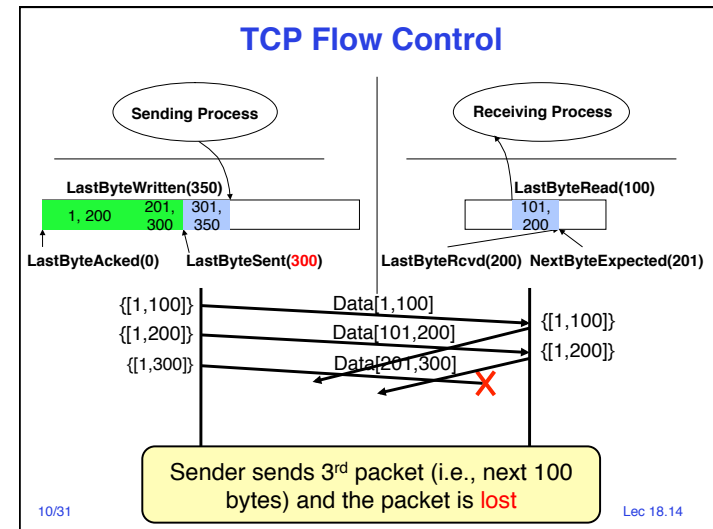
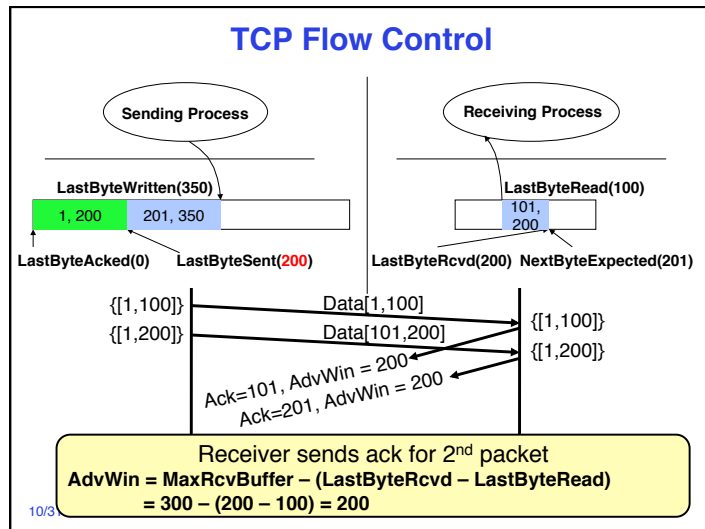


Sender sends first packet (i.e., first 100 bytes) and receiver gets the packet

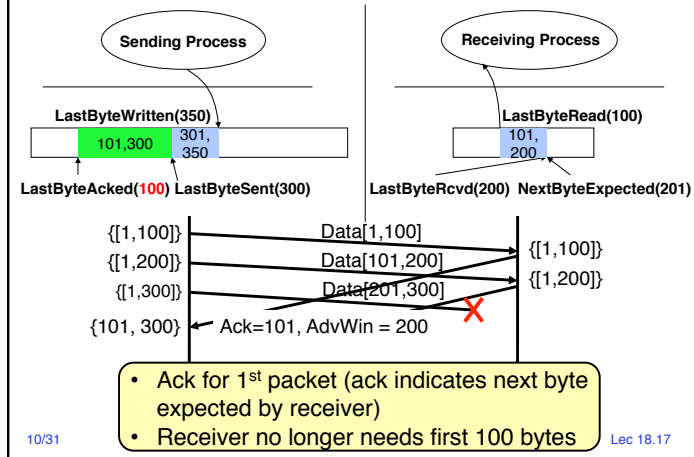
10/31

Lec 18.8

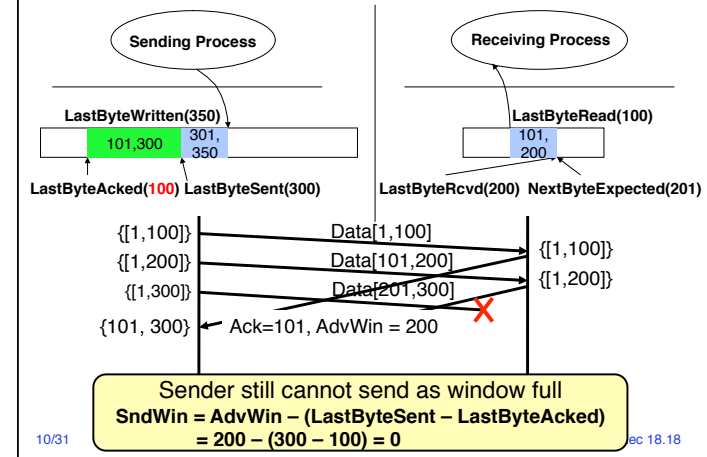




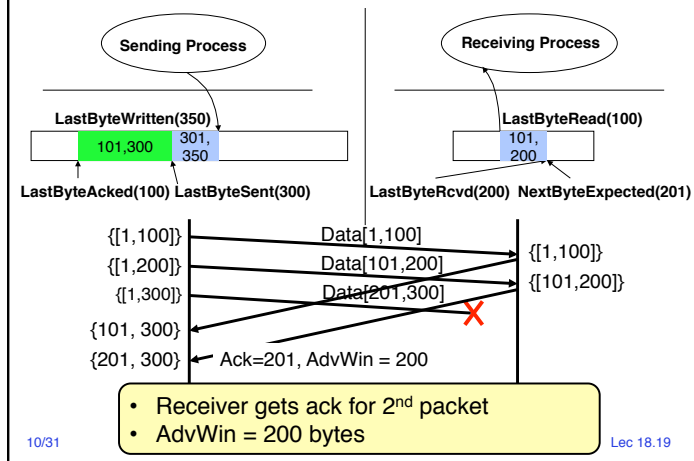
TCP Flow Control



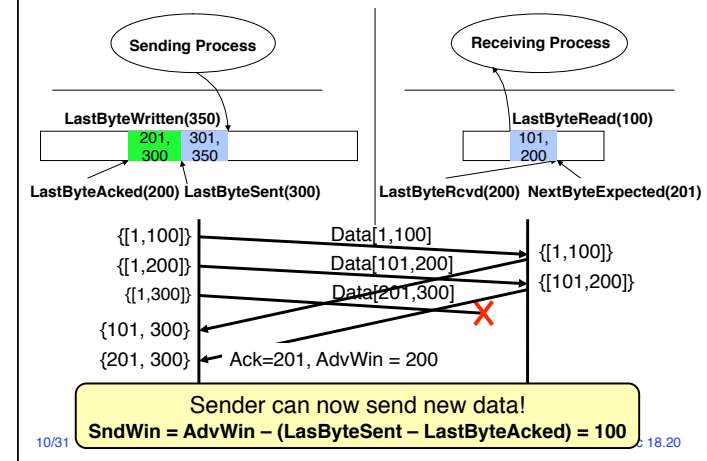
TCP Flow Control

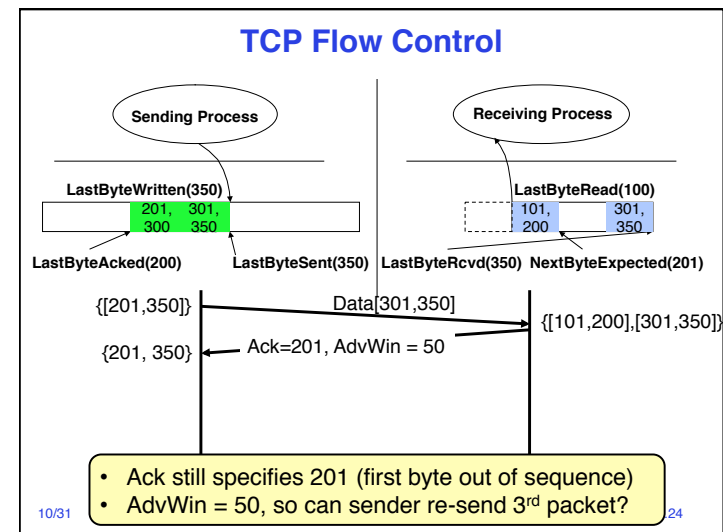
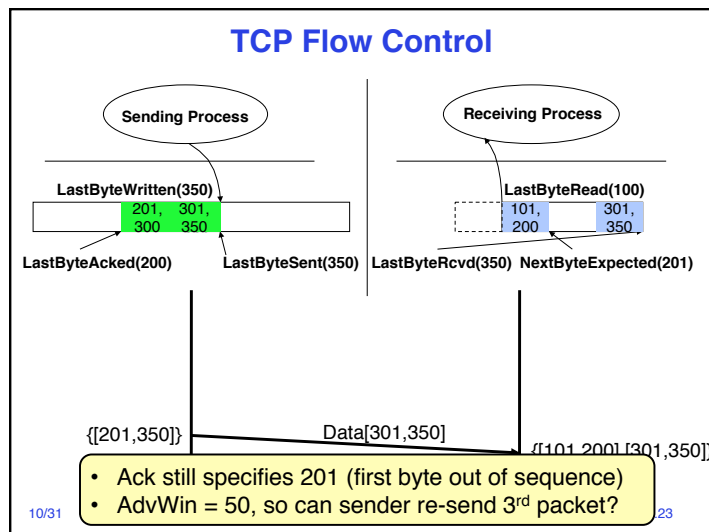
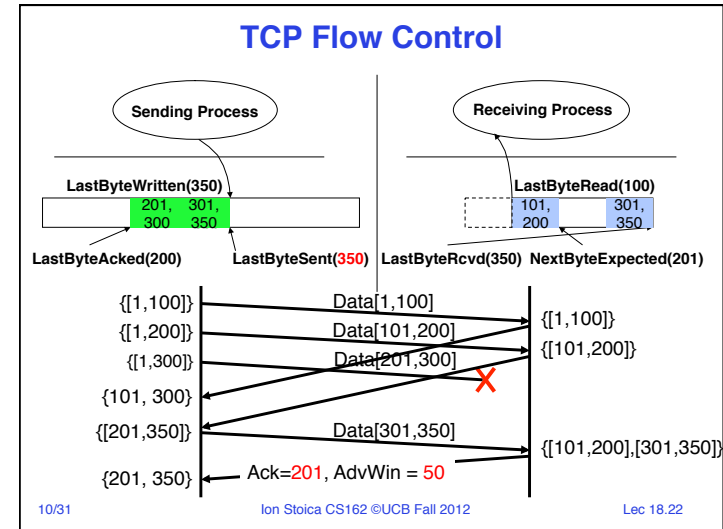
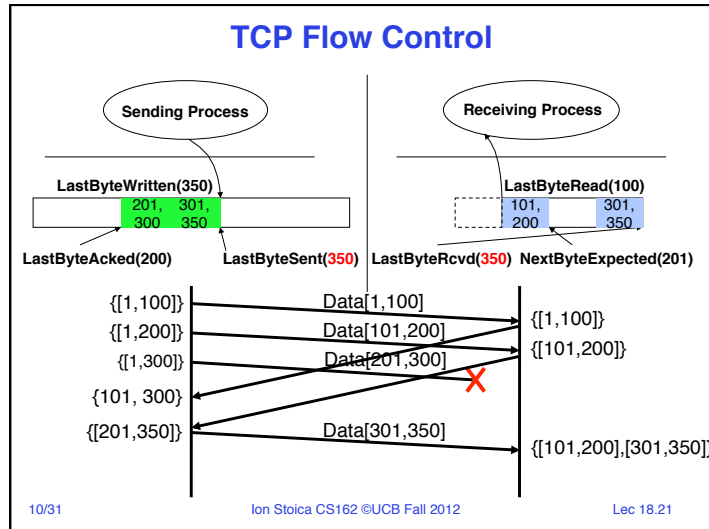


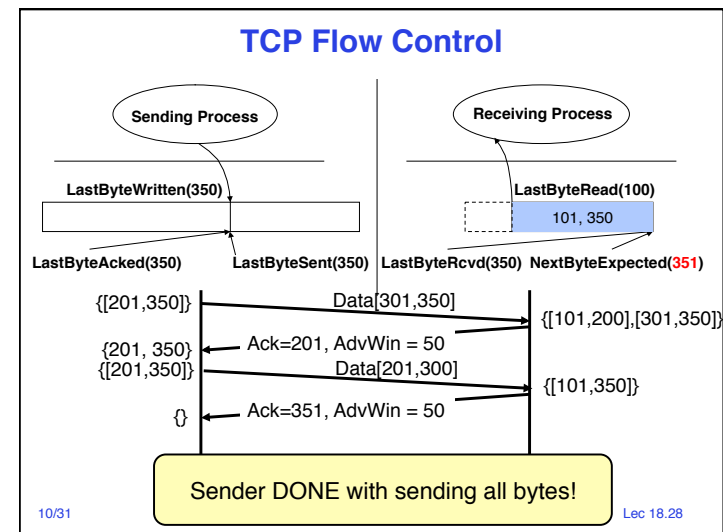
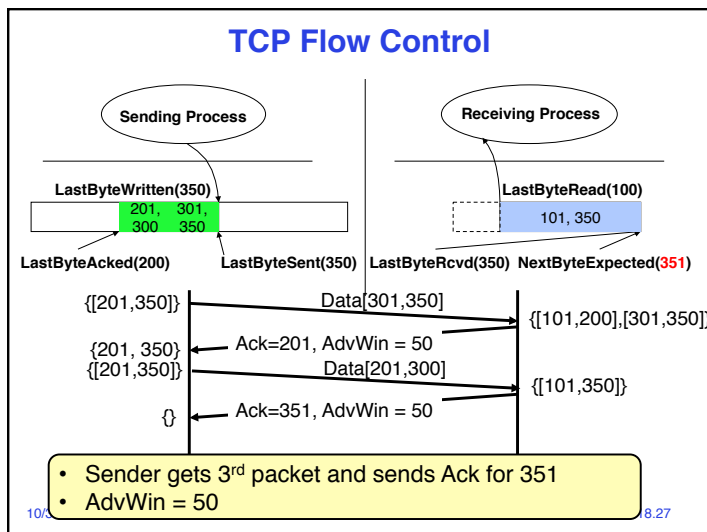
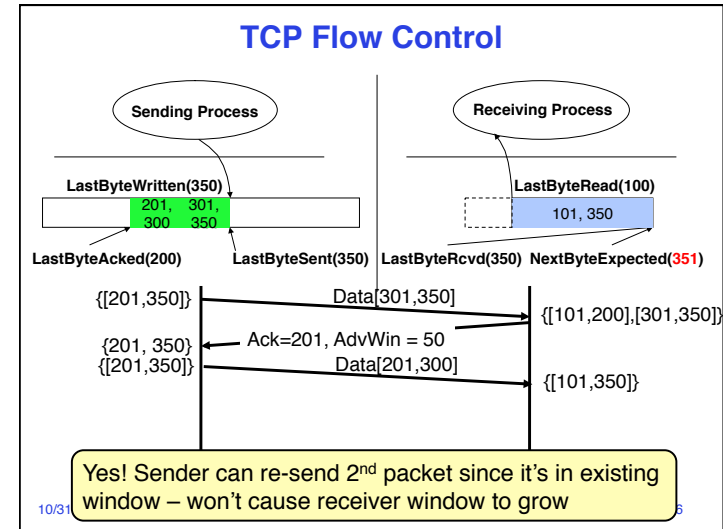
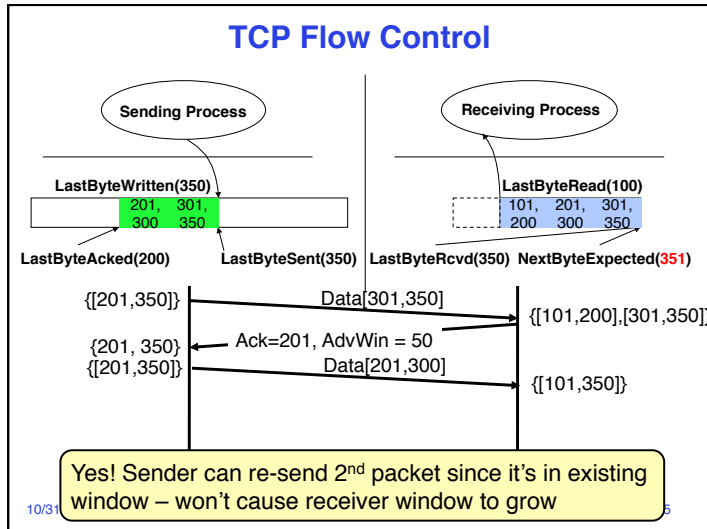
TCP Flow Control



TCP Flow Control







Quiz 18.1: Flow-Control

- Q1: True ☐ False ☐ Flow control is responsible with detecting packet losses and retransmissions
- Q2: True ☐ False ☐ Flow control always allows a sender to resend a lost packet
- Q3: True ☐ False ☐ With TCP, the receiving OS can deliver data to the application out-of-sequence (i.e., with gaps)
- Q4: True ☐ False ☐ Flow control makes sure the sender doesn't overflow the receiver

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.29

Quiz 18.1: Flow-Control

- Q1: True ☐ False ☒ Flow control is responsible with detecting packet losses and retransmissions
- Q2: True ☒ False ☐ Flow control always allows a sender to resend a lost packet
- Q3: True ☐ False ☒ With TCP, the receiving OS can deliver data to the application out-of-sequence (i.e., with gaps)
- Q4: True ☒ False ☐ Flow control makes sure the sender doesn't overflow the receiver

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.30

Summary: Reliability & Flow Control

- Flow control: three pairs of producer consumers
 - Sending process → sending TCP
 - Sending TCP → receiving TCP
 - Receiving TCP → receiving process
- AdvertisedWindow: tells sender how much **new** data can the receiver buffer
- SenderWindow: specifies how many more bytes can sender sent
 - Depends on AdvertisedWindow and on data sent since sender got AdvertisedWindow

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.31

5min Break

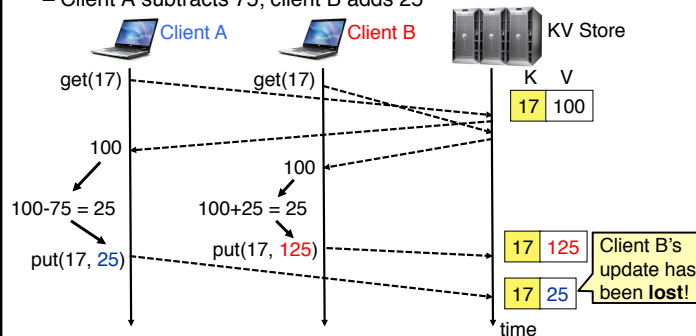
10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.32

Need for Transactions

- Example: assume two clients updating same value in a key-value (KV) store at the same time
 - Client A subtracts 75; client B adds 25



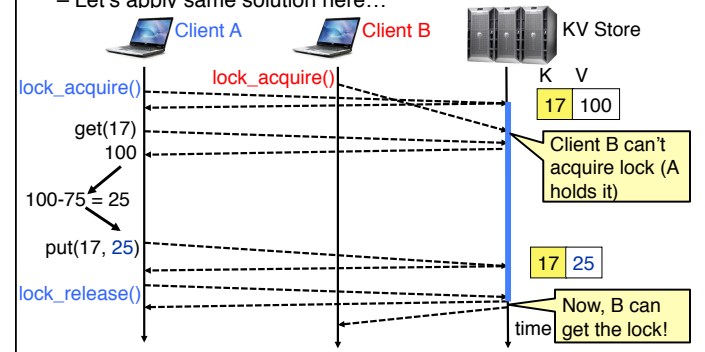
10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.33

Solution?

- How did we solve such problem on a single machine?
 - Critical section, e.g., use locks
 - Let's apply same solution here...



10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.34

Discussion

- How does client B get the lock?
 - Pooling: periodically check whether the lock is free
 - KV storage system keeps a list of clients waiting for the lock, and gives the lock to next client in the list
- What happens if the client holding the lock crashes?
- Network latency might be higher than update operation
 - Most of the time in critical section spent waiting for messages
- What is the lock granularity?
 - Do you lock every key? Do you lock the entire storage?
 - What are the tradeoffs?

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.35

Better Solution

- Interleave reads and writes from different clients
- Provide the same semantics as clients were running one at a time
- Transaction** – database/storage system's abstract view of a user program, i.e., a sequence of reads and writes

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.36

"Classic" Example: Transaction

BEGIN; --BEGIN TRANSACTION

```
UPDATE accounts SET balance = balance -  
100.00 WHERE name = 'Alice';
```

```
UPDATE branches SET balance = balance -  
100.00 WHERE name = (SELECT branch_name  
FROM accounts WHERE name = 'Alice');
```

```
UPDATE accounts SET balance = balance +  
100.00 WHERE name = 'Bob';
```

```
UPDATE branches SET balance = balance +  
100.00 WHERE name = (SELECT branch_name  
FROM accounts WHERE name = 'Bob');
```

COMMIT; --COMMIT WORK

Transfer \$100 from Alice's account to Bob's account

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.37

The ACID properties of Transactions

- **Atomicity:** all actions in the transaction happen, or none happen
- **Consistency:** transactions maintain data integrity, e.g.,
 - Balance cannot be negative
 - Cannot reschedule meeting on February 30
- **Isolation:** execution of one transaction is isolated from that of all others; no problems from concurrency
- **Durability:** if a transaction commits, its effects persist despite crashes

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.38

Atomicity

- A transaction
 - might *commit* after completing all its operations, or
 - it could *abort* (or be aborted) after executing some operations
- Atomic Transactions: a user can think of a transaction as always either *executing all its operations*, or *not executing any operations at all*
 - Database/storage system *logs* all actions so that it can *undo* the actions of aborted transactions

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.39

Consistency

- Data follows integrity constraints (ICs)
- If database/storage system is consistent before transaction, it will be after
- System checks ICs and if they fail, the transaction rolls back (i.e., is aborted)
 - A database enforces some ICs, depending on the ICs declared when the data has been created
 - Beyond this, database does not understand the semantics of the data (e.g., it does not understand how the interest on a bank account is computed)

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.40

Isolation

- Each transaction executes as if it was running by itself
 - It cannot see the partial results of another transaction
- Techniques:
 - Pessimistic – don't let problems arise in the first place
 - Optimistic – assume conflicts are rare, deal with them *after* they happen

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.41

Durability

- Data should survive in the presence of
 - System crash
 - Disk crash → need backups
- All committed updates and only those updates are reflected in the database
 - Some care must be taken to handle the case of a crash occurring during the recovery process!

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.42

This Lecture

- Deal with **(I)solation**, by focusing on **concurrency control**
- Next lecture focus on (A)tomicity, and partially on (D)urability

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.43

Example

- Consider two transactions:
 - T1: moves \$100 from account A to account B

```
T1:A := A-100; B := B+100;
```
 - T2: moves \$50 from account B to account A

```
T2:A := A+50; B := B-50;
```
- Each operation consists of (1) a read, (2) an addition/subtraction, and (3) a write
- Example: A = A-100

```
Read(A); // R(A)
A := A - 100;
Write(A); // W(A)
```

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.44

Example (cont' d)

- Database only sees reads and writes

Database View

T1: A:=A-100; B:=B+100; → T1: R(A), W(A), R(B), W(B)

T2: A:=A+50; B:=B-50; → T2: R(A), W(A), R(B), W(B)

- Assume initially: A = \$1000 and B = \$500
- What is the legal outcome of running T1 and T2?
 - A = \$950
 - B = \$550

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.45

Example (cont' d)

T1: A:=A-100; B:=B+100;

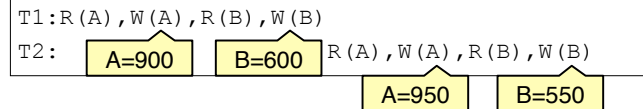
Initial values:

A:=1000

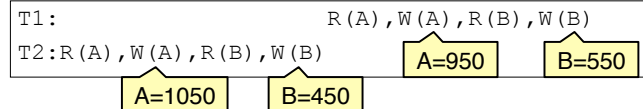
B:=500

T2: A:=A+50; B:=B-50;

- What is the outcome of the following execution?



- What is the outcome of the following execution?



10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.46

Example (cont' d)

T1: A:=A-100; B:=B+100;

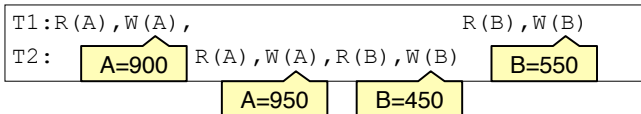
Initial values:

A:=1000

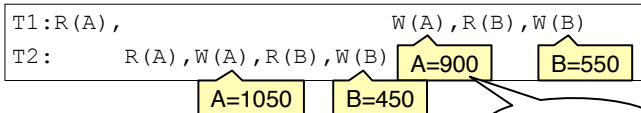
B:=500

T2: A:=A+50; B:=B-50;

- What is the outcome of the following execution?



- What is the outcome of the following execution?



Lost \$50!

10/31

Ion Stoica CS162 ©UCB Fall 2012

Transaction Scheduling

- Why not run only one transaction at a time?
- Answer: low system utilization
 - Two transactions cannot run simultaneously even if they access different data

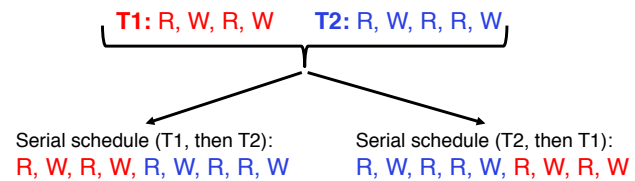
10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.48

Goals of Transaction Scheduling

- Maximize system utilization, i.e., concurrency
 - Interleave operations from different transactions
- Preserve transaction semantics
 - Semantically equivalent to a **serial schedule**, i.e., one transaction runs at a time



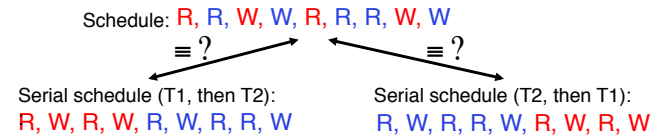
10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.49

Two Key Questions

- 1) Is a given schedule equivalent to a serial execution of transactions?



- 2) How do you come up with a schedule equivalent to a serial schedule?

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.50

Summary

- Transaction: a sequence of storage operations
- ACID:
 - Atomicity: all operations in a transaction happen, or none happens
 - Consistency: if database/storage starts consistent, it ends up consistent
 - Isolation: execution of one transaction is isolated from another
 - Durability: the results of a transaction persists

10/31

Ion Stoica CS162 ©UCB Fall 2012

Lec 18.51