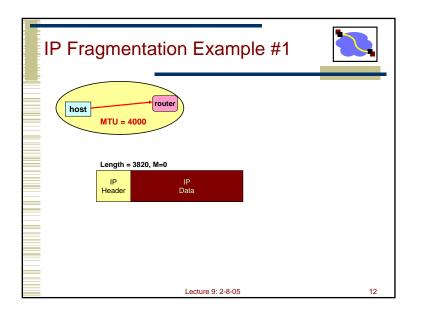
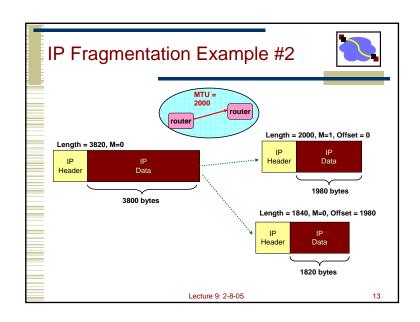
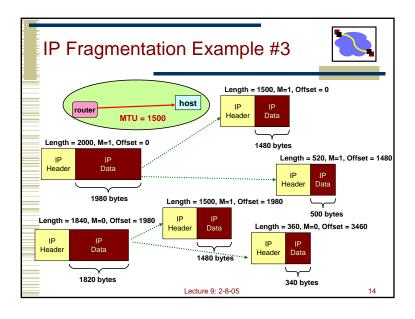


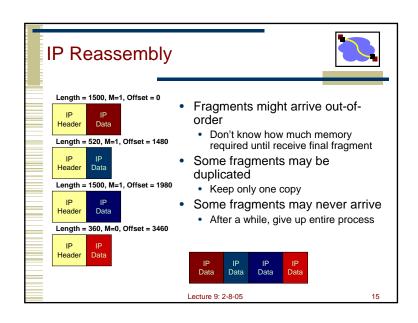
Peassembly Where to do reassembly? End nodes or at routers? End nodes Avoids unnecessary work where large packets are fragmented multiple times If any fragment missing, delete entire packet Dangerous to do at intermediate nodes How much buffer space required at routers? What if routes in network change? Multiple paths through network All fragments only required to go through destination

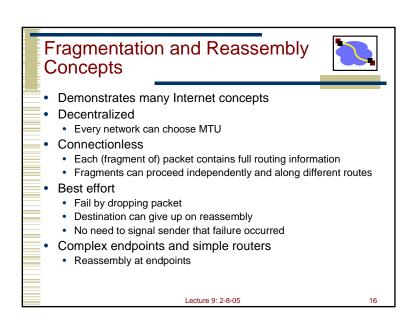
Fragmentation Related Fields Length Length of IP fragment Identification To match up with other fragments Flags Don't fragment flag More fragments flag Fragment offset Where this fragment lies in entire IP datagram Measured in 8 octet units (13 bit field)











Fragmentation is Harmful



- Uses resources poorly
 - Forwarding costs per packet
 - · Best if we can send large chunks of data
 - · Worst case: packet just bigger than MTU
- · Poor end-to-end performance
 - Loss of a fragment
- Path MTU discovery protocol → determines minimum MTU along route
 - · Uses ICMP error messages
- · Common theme in system design
 - · Assure correctness by implementing complete protocol
 - · Optimize common cases to avoid full complexity

Lecture 9: 2-8-05

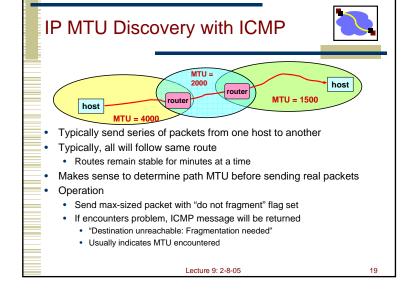
Internet Control Message Protocol (ICMP)

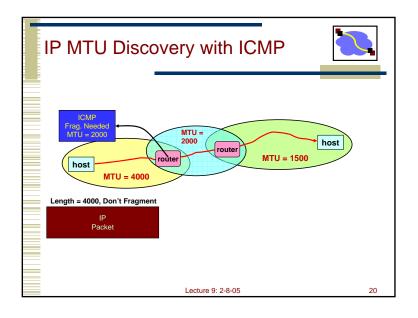


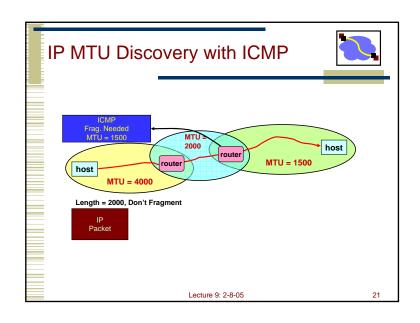
- Short messages used to send error & other control information
- Examples
- Ping request / response
 - Can use to check whether remote host reachable
 - · Destination unreachable
 - · Indicates how packet got & why couldn't go further
- Flow control
 - · Slow down packet delivery rate
- Redirect
 - · Suggest alternate routing path for future messages
- · Router solicitation / advertisement
 - · Helps newly connected host discover local router
- Timeout
 - Packet exceeded maximum hop limit

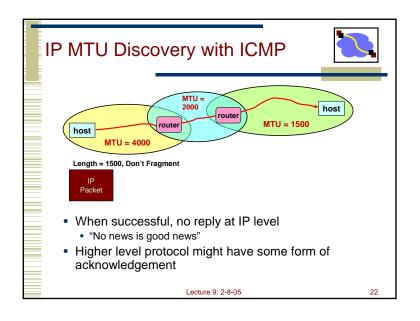
Lecture 9: 2-8-05

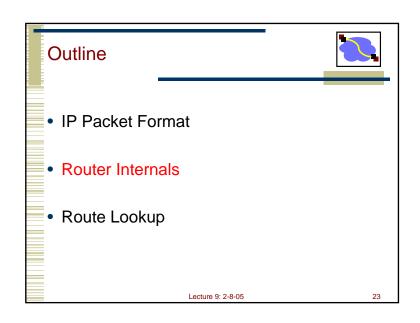
18

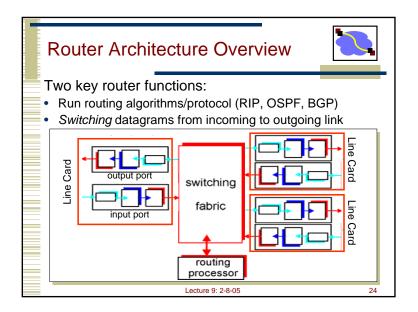


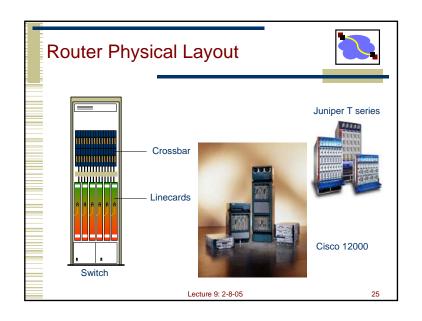


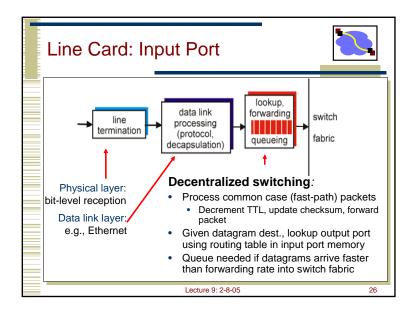


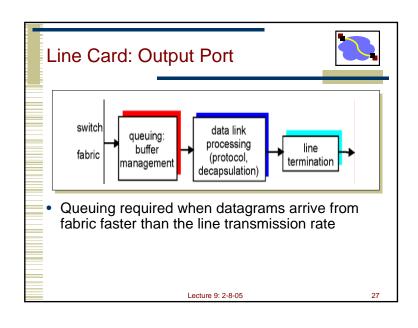


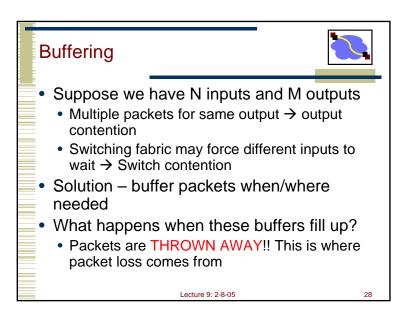










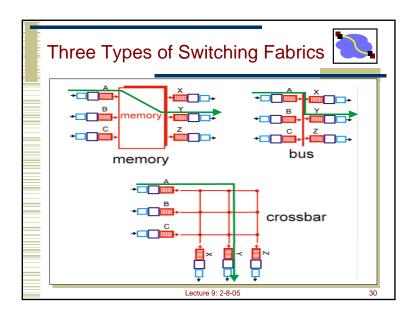


Network Processor

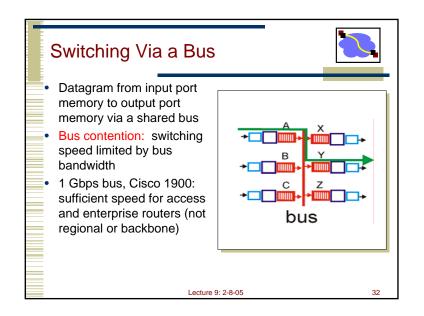


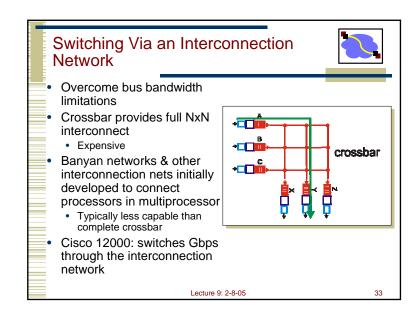
- Runs routing protocol and downloads forwarding table to forwarding engines
- Performs "slow" path processing
 - ICMP error messages
 - IP option processing
 - Fragmentation
 - · Packets destined to router

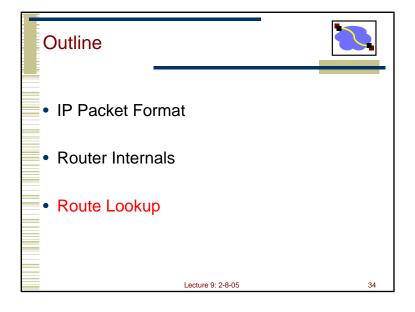
Lecture 9: 2-8-05

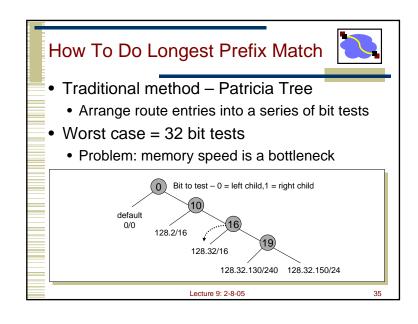


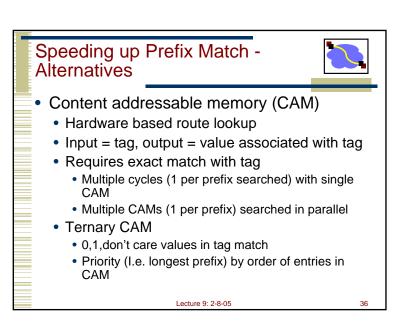
Switching Via a Memory First generation routers → looked like PCs Packet copied by system's (single) CPU • Speed limited by memory bandwidth (2 bus crossings per datagram) Memory Output Input Port Port System Bus Modern routers • Input port processor performs lookup, copy into memory Cisco Catalyst 8500 Lecture 9: 2-8-05











Speeding up Prefix Match - Alternatives



- Route caches
 - Packet trains → group of packets belonging to same flow
 - Temporal locality
 - · Many packets to same destination
- Other algorithms
 - Routing with a Clue [Bremler-Barr Sigcomm 99]
 - Clue = prefix length matched at previous hop
 - · Why is this useful?

Lecture 9: 2-8-05

37

Important Concepts



- Base-level protocol (IP) provides minimal service level
 - · Allows highly decentralized implementation
 - Each step involves determining next hop
 - · Most of the work at the endpoints
- ICMP provides low-level error reporting
- IP routers
 - Architecture
 - · Optimized for common case processing
 - Complex/expensive lookup algorithms (especially in comparison to ATM fixed length lookup)

Lecture 9: 2-8-05

38

Important Concepts



- IP forwarding → global addressing, alternatives, lookup tables
- IP addressing → hierarchical, CIDR,
- IP service → best effort, simplicity of routers
- IP packets → header fields, fragmentation,
 ICMP

Lecture 9: 2-8-05

05

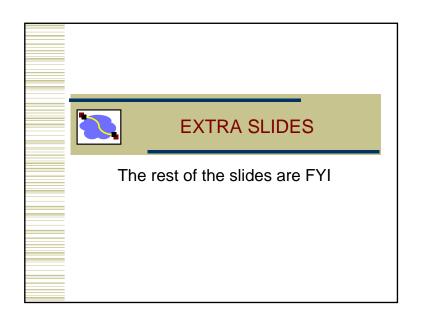
Next Lecture

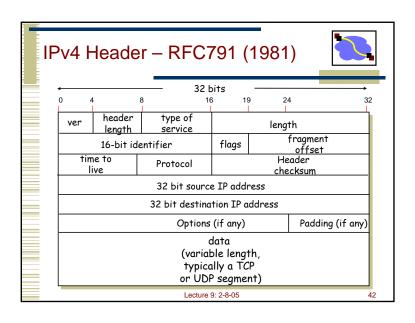


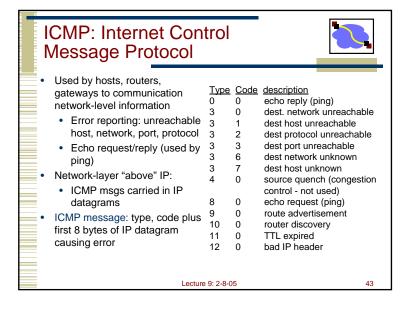
- How do forwarding tables get built?
- Routing protocols
 - · Distance vector routing
 - · Link state routing

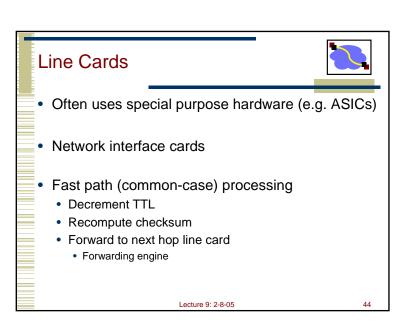
Lecture 9: 2-8-05

05









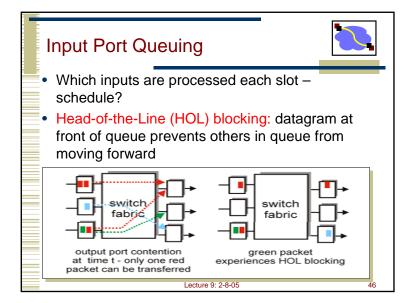
Switch Buffering



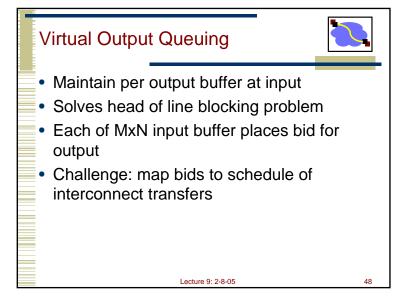
45

- · 3 types of switch buffering
 - Input buffering
 - Fabric slower than input ports combined → queuing may occur at input queues
 - Can avoid any input queuing by making switch speed = N x link speed
 - Output buffering
 - Buffering when arrival rate via switch exceeds output line speed
 - Internal buffering
 - Can have buffering inside switch fabric to deal with limitations of fabric

Lecture 9: 2-8-05



Output Port Queuing Output Port Contention of Time I Scheduling discipline chooses among queued datagrams for transmission Can be simple (e.g., first-come first-serve) or more clever (e.g., weighted round robin) Lecture 9: 2-8-05 Available of the provided round robin of the pr



Speeding up Prefix Match



- Cut prefix tree at 16/24/32 bit depth
 - Fill in prefix tree entries by creating extra entries
 - Entries contain output interface for route
 - Add special value to indicate that there are deeper tree entries
 - Only keep 24/32 bit cuts as needed
- Example cut prefix tree at 16 bit depth
 - Only 64K entries

Lecture 9: 2-8-05

