

# CS/ECE 4457

## Computer Networks: Architecture and Protocols

### Lecture 17 Switch/Router Architecture

**Qizhe Cai**



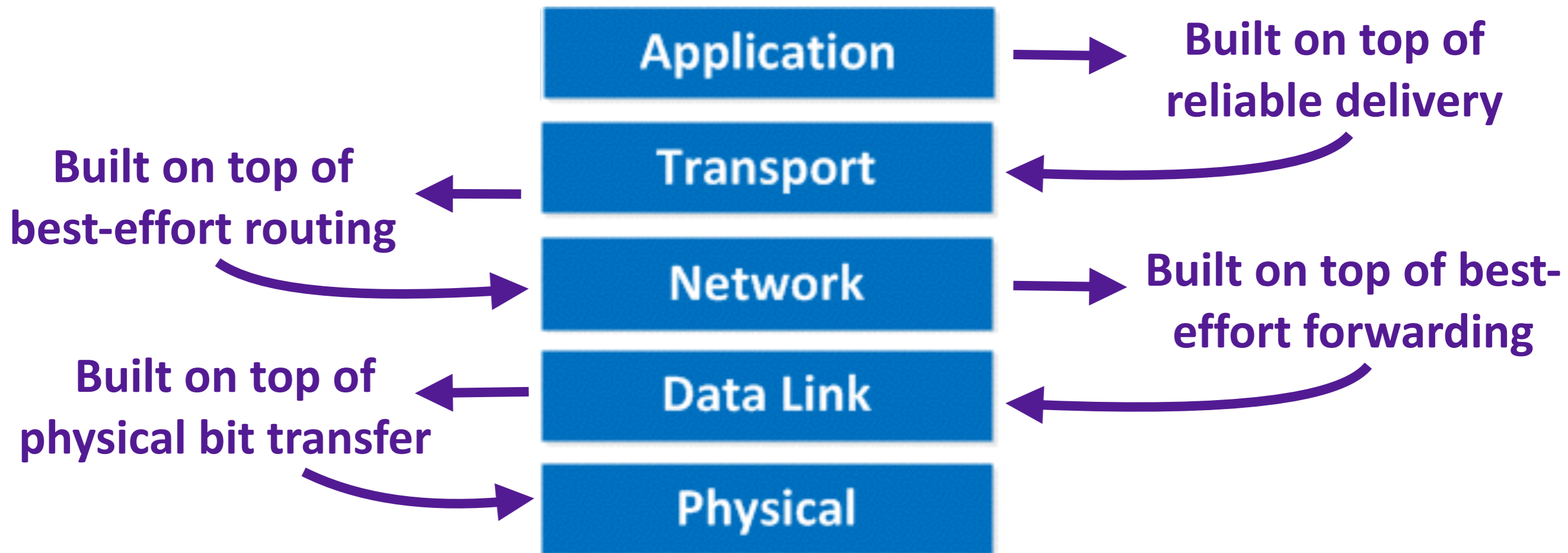
# Announcements

- **Exam 2 solutions released**

# Goals for Today's Lecture

- Understand switch/router architecture

# Where are we?

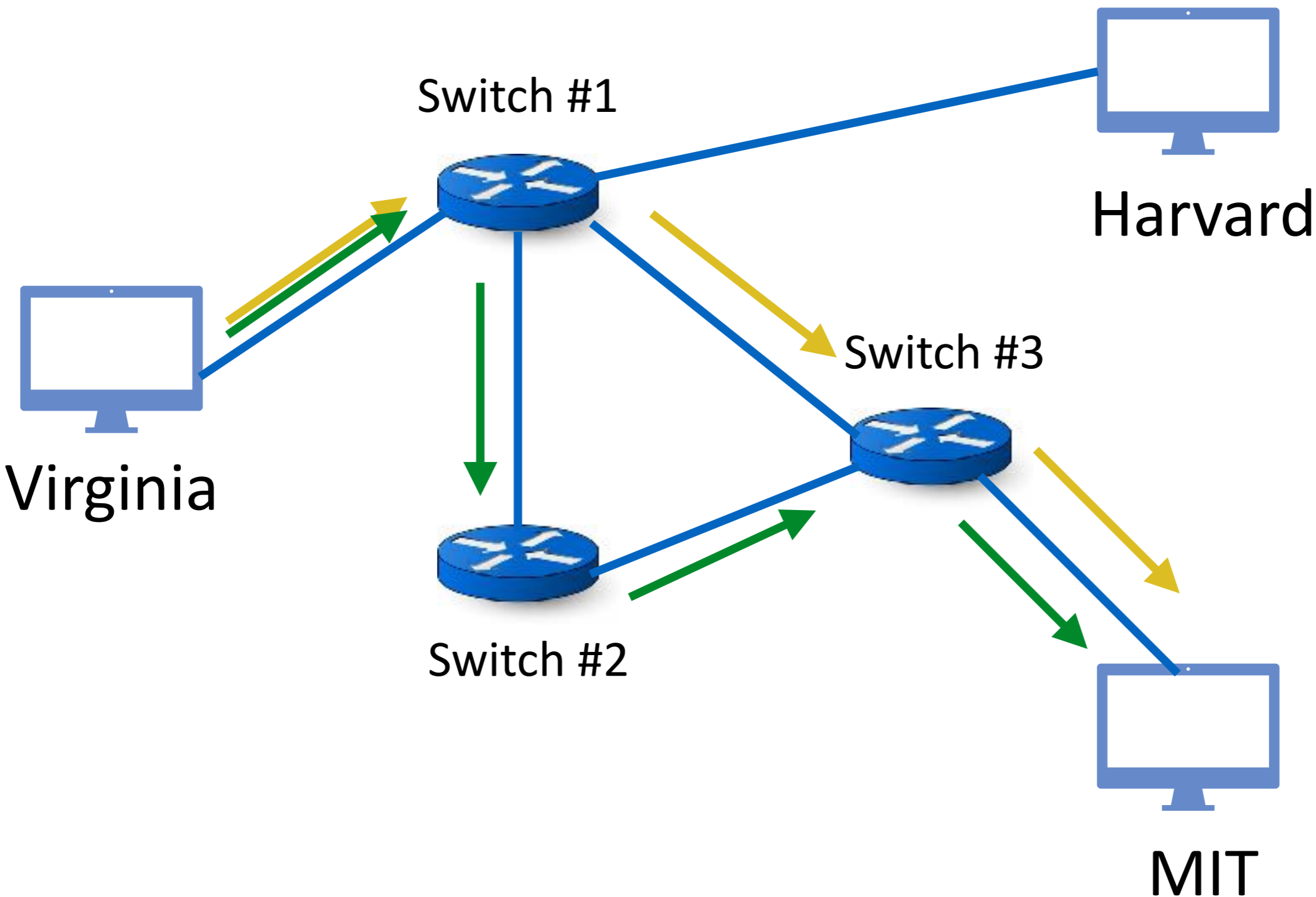


# **Switch/Router Architecture**

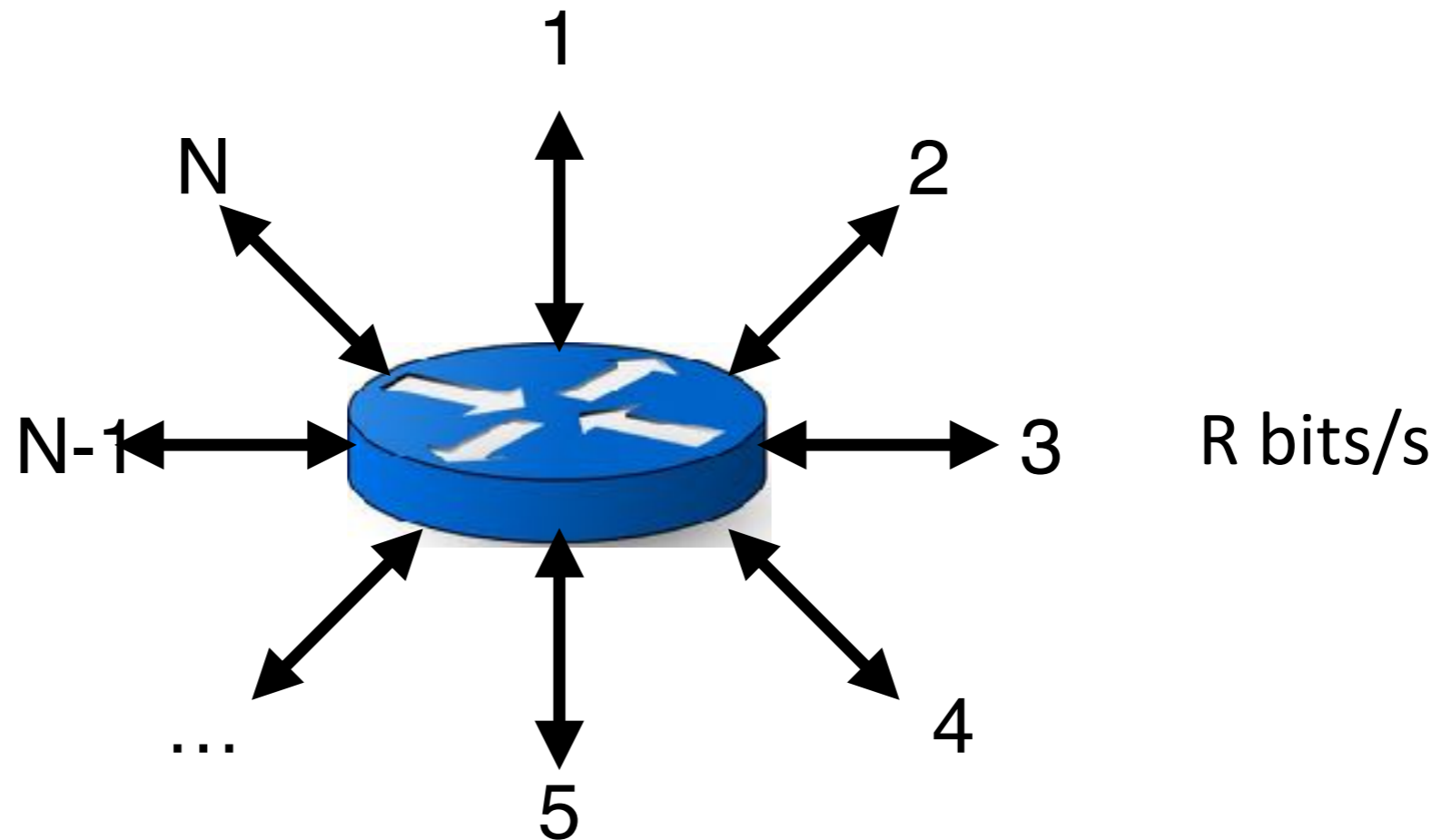
# IP Routers and Switches (used interchangeably today)

- Core building block of Internet infrastructure
- \$120B+ industry
- Vendors: Cisco, Huawei, Juniper, Alcatel-Lucent (account for >90%)

# Recap: Routers Forward Packets

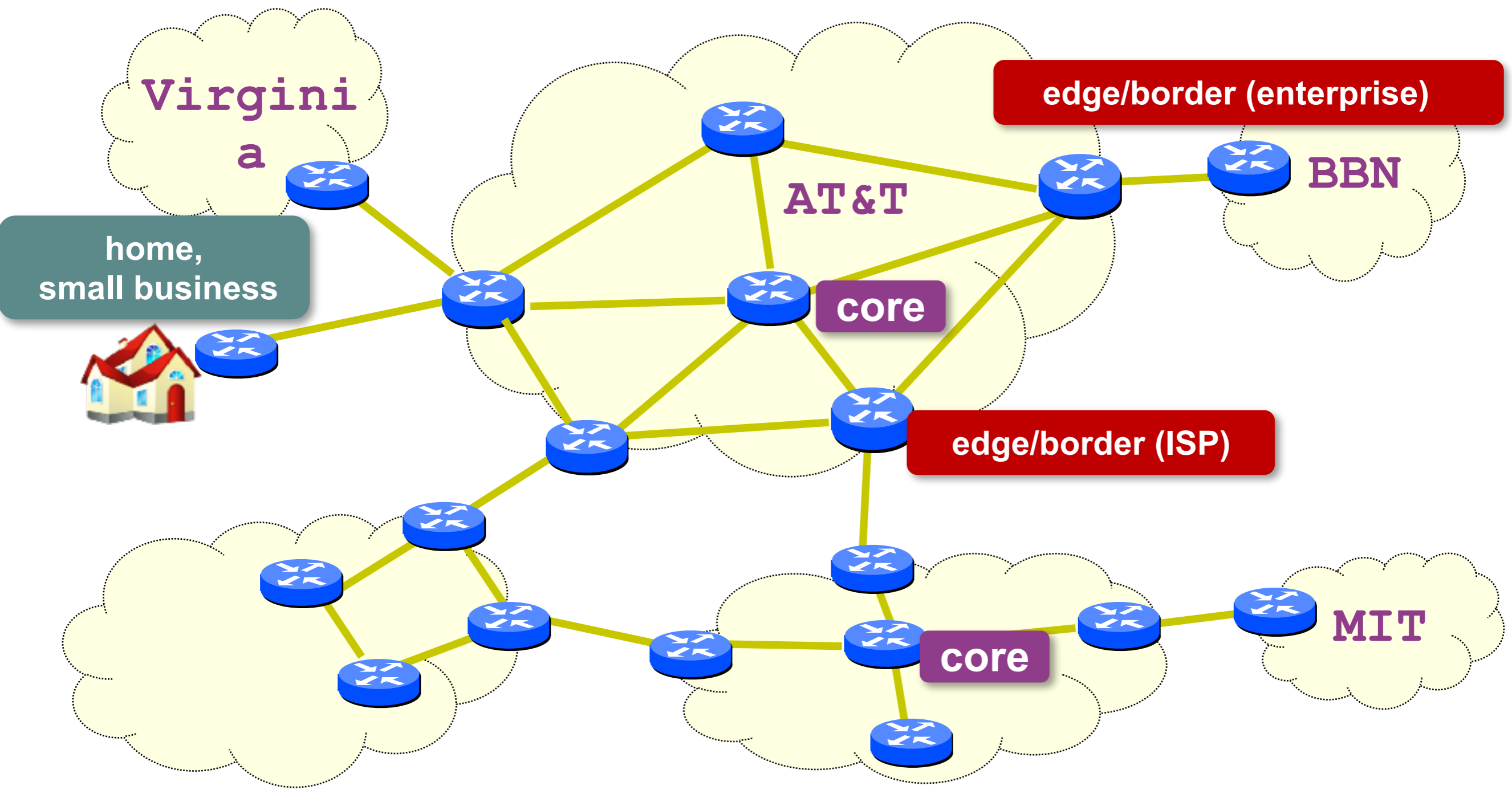


# Router Definitions



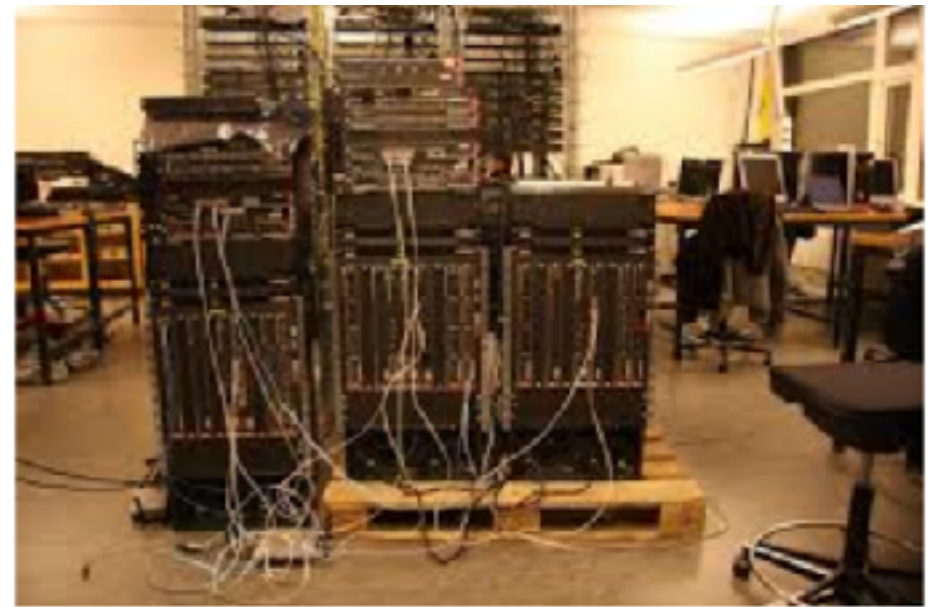
- $N$  = No. Of external router ports
- $R$  = bandwidth (“line rate”) of a port
- Router capacity =  $N \times R$

# Networks and Routers

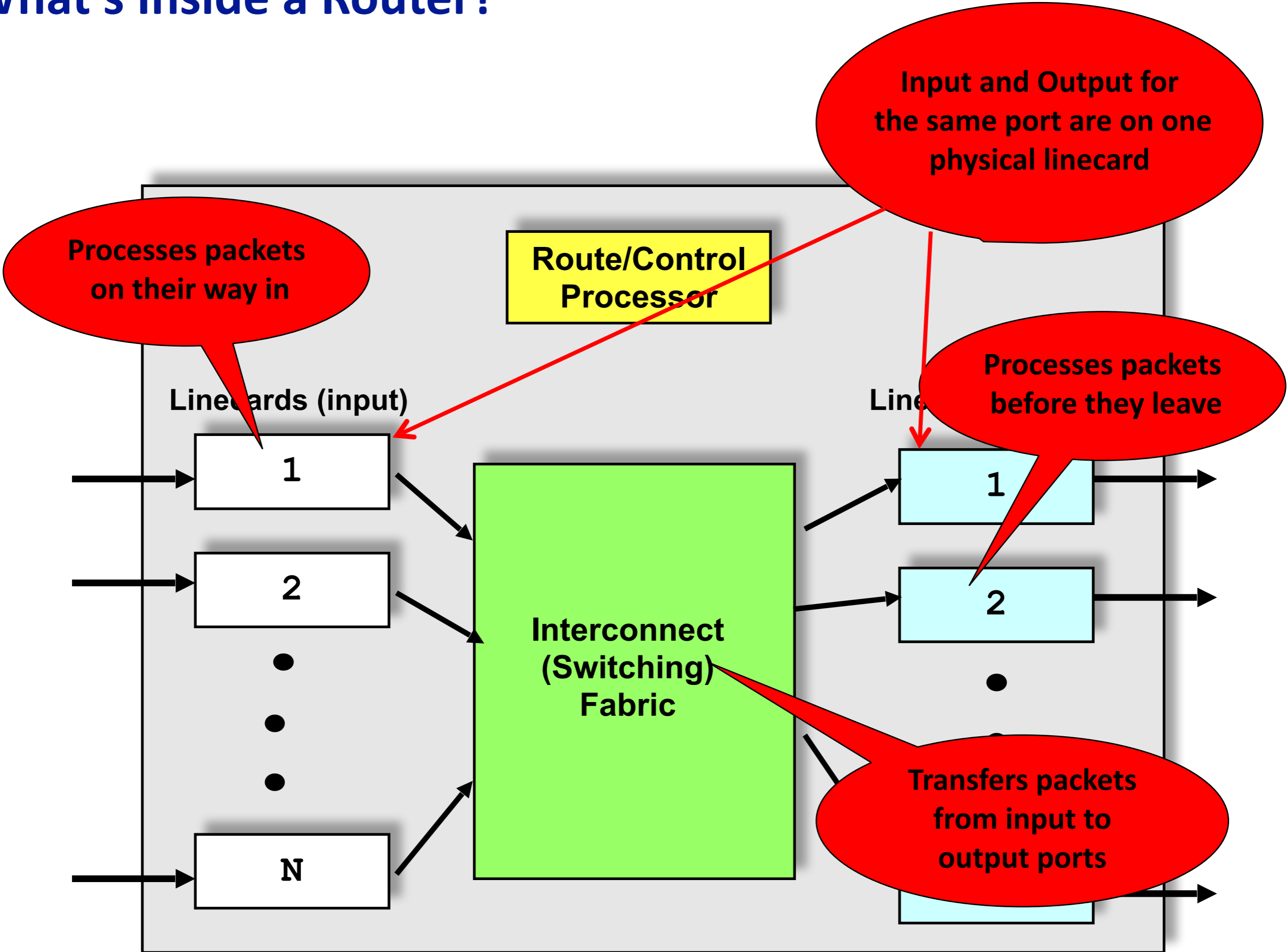


# Examples of Routers (core)

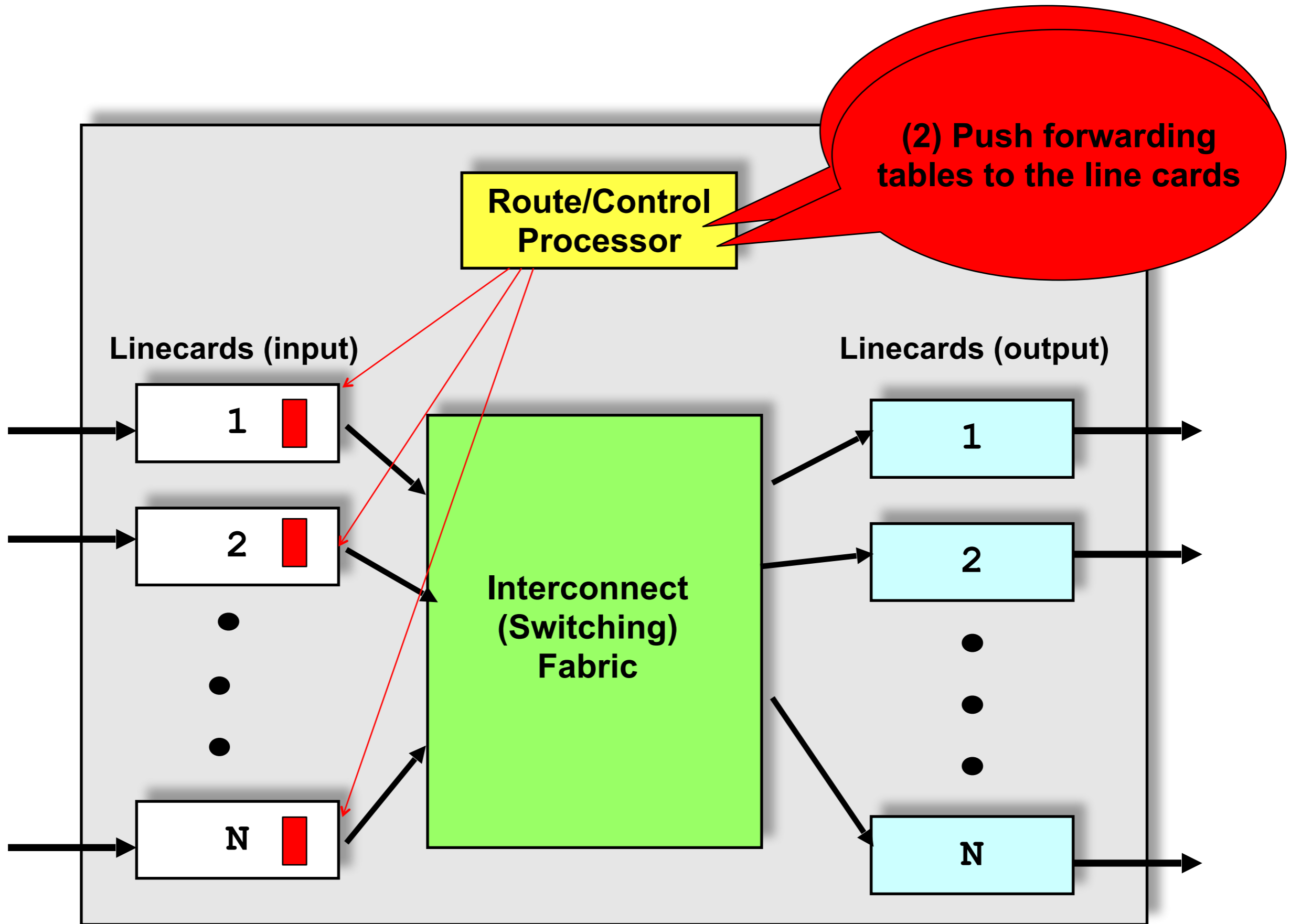
- **Core:** Cisco CRS
  - R = 10/40/100 Gbps
  - NR = 922 Tbps
  - Netflix: 0.7 GB/hr (1.5Mb/s)
  - ~600 million concurrent Netflix users
- **Edge (ISP):** Cisco ASR
  - R = 1/10/40 Gbps
  - NR = 120 Gbps
- **Edge (enterprise):** Cisco 3945E
  - R = 10/100/1000 Mbps
  - NR < 10 Gbps



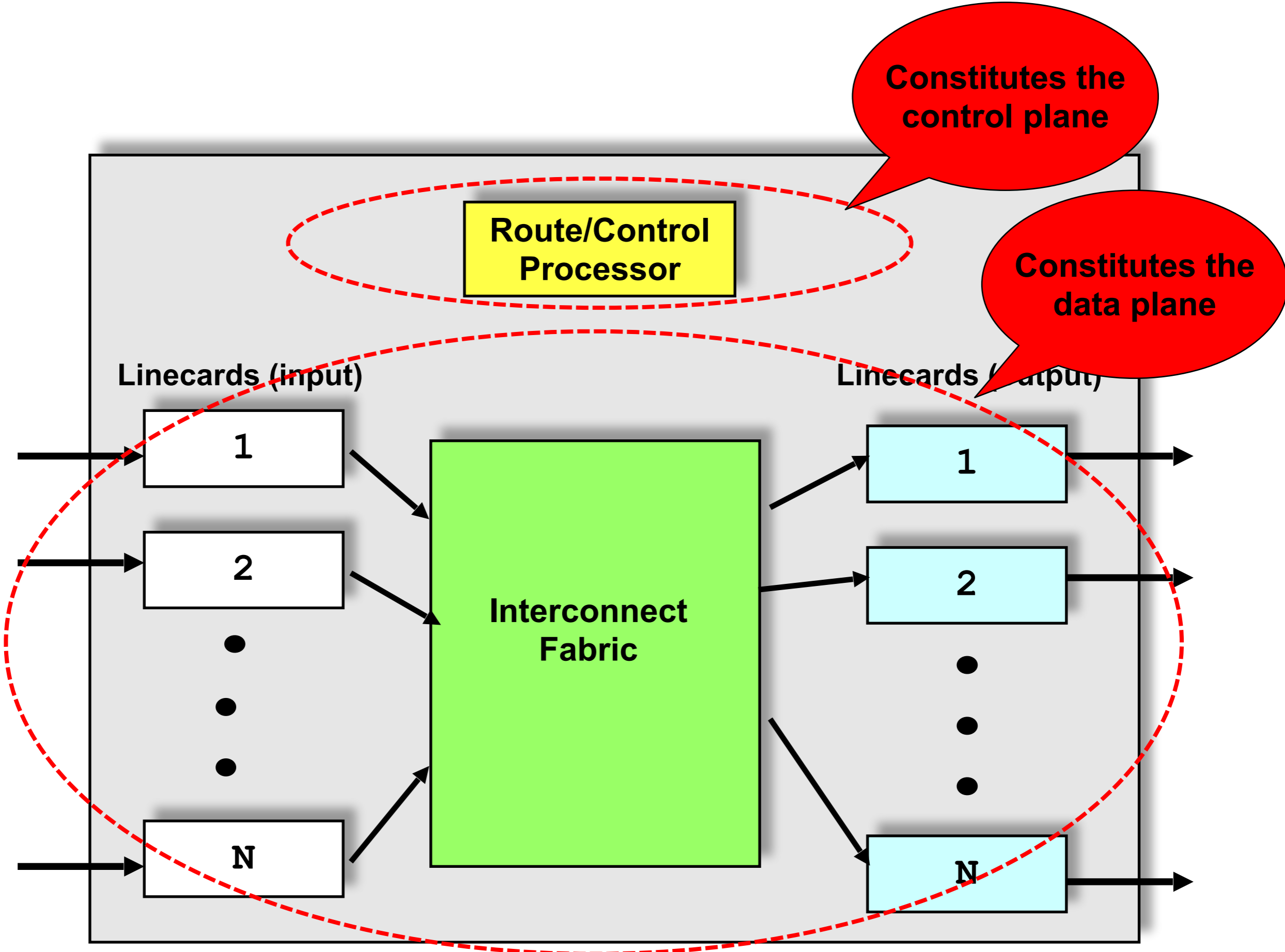
# What's Inside a Router?



# What's Inside a Router?



# What's Inside a Router?



# Input Line Cards: Tasks

- Receive incoming packets (physical layer stuff)
- Update the IP header
  - TTL, Checksum (maybe some other fields)
- Lookup the output port for the destination IP address
- Queue the packet at the switch fabric

# Challenge: Speed!

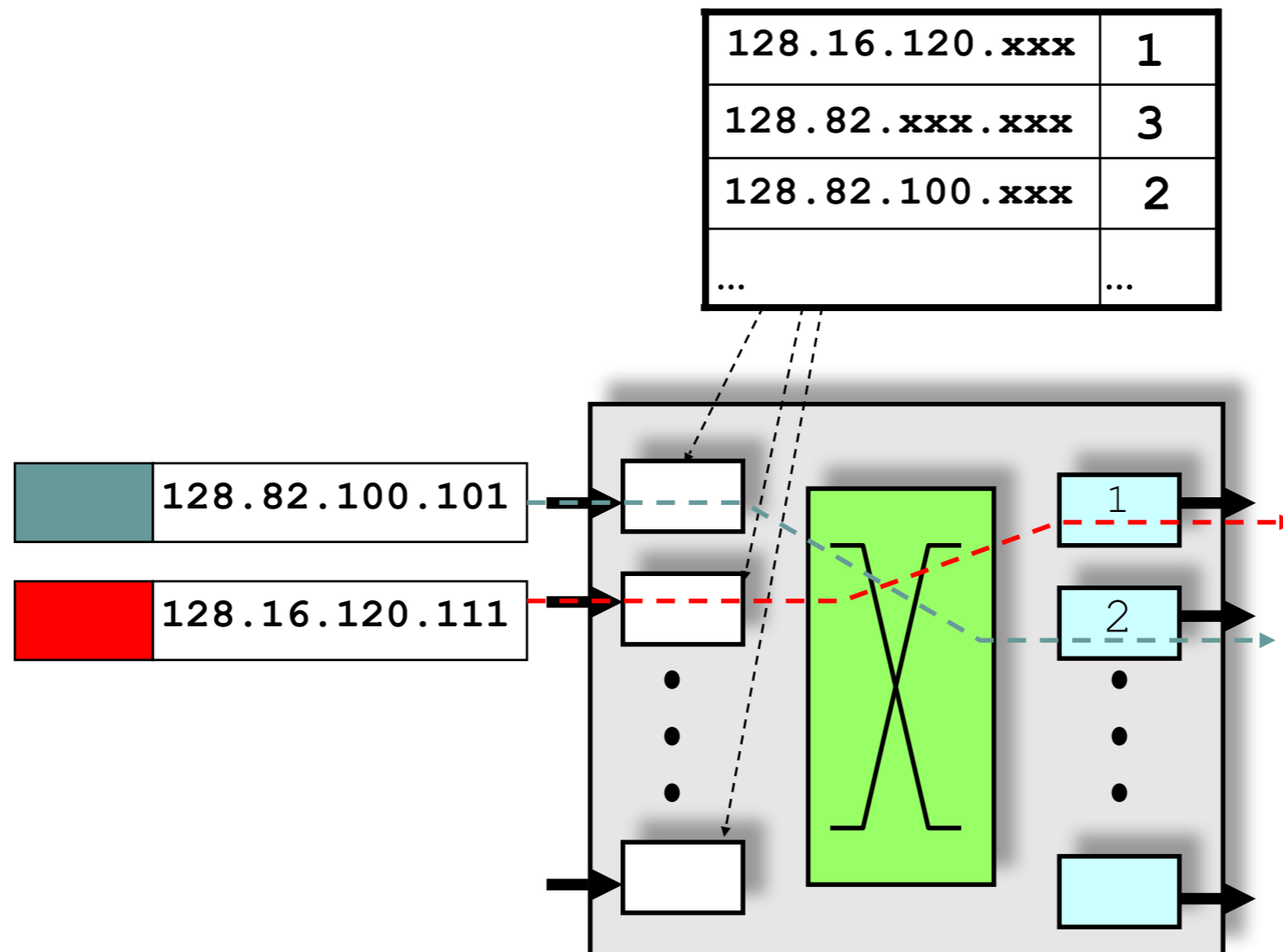
- 100B packets @ 40Gbps => packet every 20 nano secs!
- Typically implemented with specialized hardware
  - ASICs, specialized “network processors”

# Looking up the Output Port

- Upon receiving a packet
  - Inspect the destination IP address in the header
  - Index into the routing/forwarding table
  - If no match, select the **default route**
  - Forward packet out appropriate interface
- **Default route**
  - Configured to cover cases where no matches
  - Allows small tables at edge (w/o routing algorithms)
    - **if it isn't on my subnet, send it to my ISP**

# Scaling the Lookup

- Recall: For scalability, addresses are **aggregated**
- **Longest Prefix match**
  - Find the entry with matching “longest prefix” with destination address



# Finding a Match

- Incoming packet destination: 201.143.7.0

<b>Prefix</b>	<b>Port</b>
201.143.0.0/22	Port 1
201.143.4.0.0/24	Port 2
201.143.5.0.0/24	Port 3
201.143.6.0/23	Port 4

# Finding a Match: Covert to Binary

- Incoming packet destination: 201.143.7.0

11001001	10001111	00000111	11010010
----------	----------	----------	----------

## Routing Table

### 201.143.0.0/22

11001001	10001111	000000 - -	- - - - -
----------	----------	------------	-----------

### 201.143.4.0/24

11001001	10001111	00000100	- - - - -
----------	----------	----------	-----------

### 201.143.5.0/24

11001001	10001111	00000101	- - - - -
----------	----------	----------	-----------

### 201.143.6.0/23

11001001	10001111	0000011-	- - - - -
----------	----------	----------	-----------

# Finding a Match: Covert to Binary

- Incoming packet destination: 201.143.7.0

11001001	10001111	00000111	11010010
----------	----------	----------	----------

## Routing Table

### 201.143.0.0/22

11001001	10001111	000000 - -	- - - - -
----------	----------	------------	-----------

### 201.143.4.0/24

11001001	10001111	00000100	- - - - -
----------	----------	----------	-----------

### 201.143.5.0/24

11001001	10001111	00000101	- - - - -
----------	----------	----------	-----------

### 201.143.6.0/23

11001001	10001111	0000011-	- - - - -
----------	----------	----------	-----------

# Finding a Match: Covert to Binary

- Incoming packet destination: 201.143.7.0

11001001	10001111	00000 <b>1</b> 11	11010010
----------	----------	-------------------	----------

## Routing Table

**201.143.0.0/22**

11001001	10001111	00000 <b>0</b>	
----------	----------	----------------	--

---

**201.143.4.0/24**

11001001	10001111	00000 <b>1</b> 00	-----
----------	----------	-------------------	-------

**201.143.5.0/24**

11001001	10001111	00000 <b>1</b> 01	-----
----------	----------	-------------------	-------

**201.143.6.0/23**

11001001	10001111	00000 <b>1</b> 1-	-----
----------	----------	-------------------	-------

# Longest Prefix Match

- Incoming packet destination: 201.143.7.0

11001001	10001111	00000 <b>1</b> 11	11010010
----------	----------	-------------------	----------

## Routing Table

**201.143.0.0/22**

11001001	10001111	00000 <b>0</b> -	-----
----------	----------	------------------	-------

**201.143.4.0/24**

11001001	10001111	00000 <b>1</b> 00	-----
----------	----------	-------------------	-------

**201.143.5.0/24**

11001001	10001111	00000 <b>1</b> 01	-----
----------	----------	-------------------	-------

**201.143.6.0/23**

11001001	10001111	00000 <b>1</b> 1-	-----
----------	----------	-------------------	-------

Check an address against all destination prefixes and select the prefixes it matches with and then select the one that has the most bits

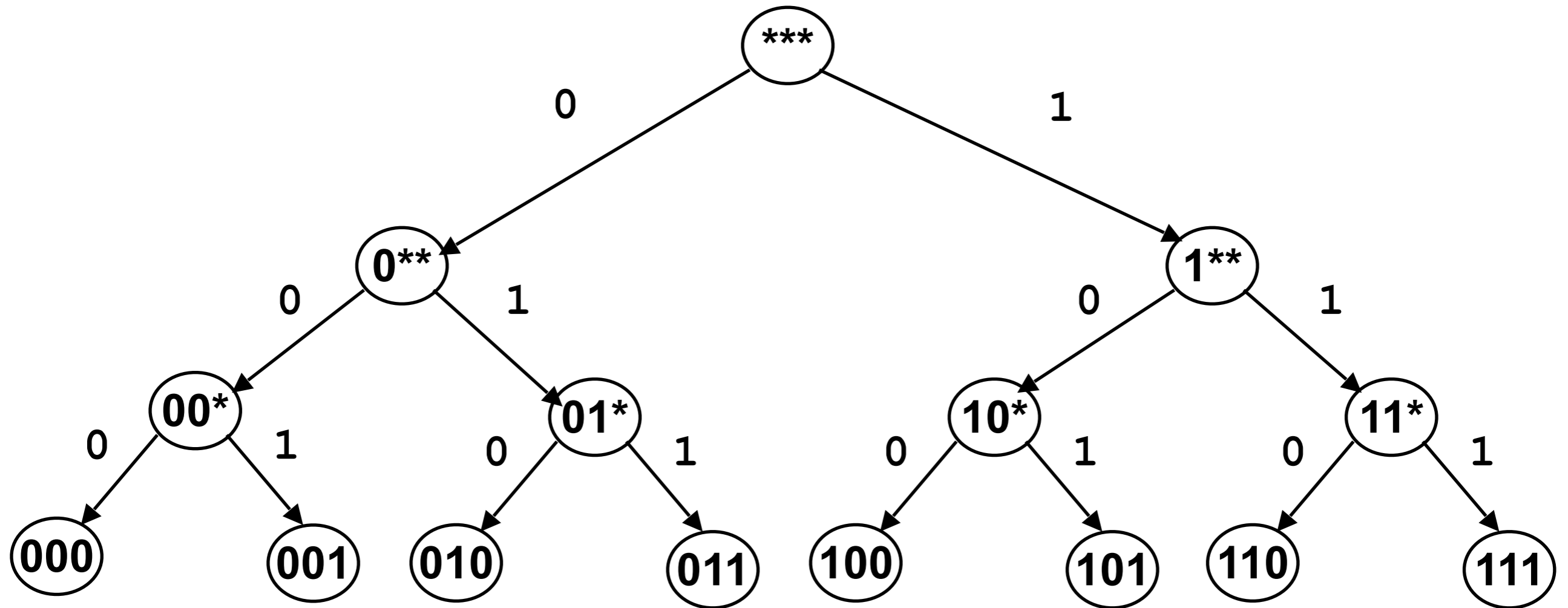
# Finding the Match Efficiently

- Testing each entry to find a match scales poorly
  - Roughly (number of entries)  $\times$  (number of bits)
- Must leverage tree structure of binary strings
  - Set up tree-like data structure
  - Called a **TRIE**
    - We will briefly discuss it; more details in text
    - In case you are interested ....

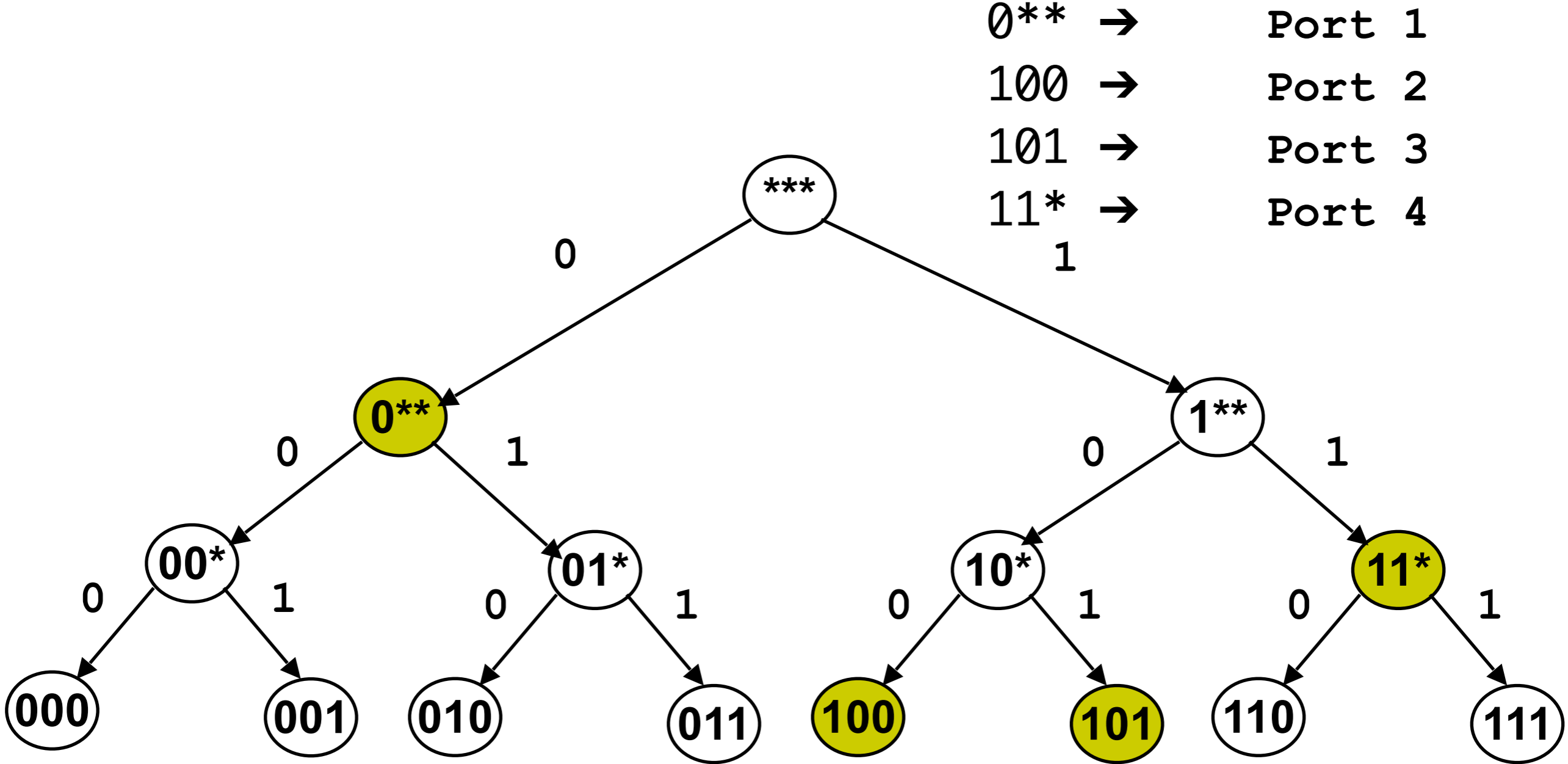
# Consider Four 3-Bit Prefixes

- Just focusing on the bits where all the action is....
- $0^{**}$  → Port 1
- 100 → Port 2
- 101 → Port 3
- $11^*$  → Port 4

# Tree Structure

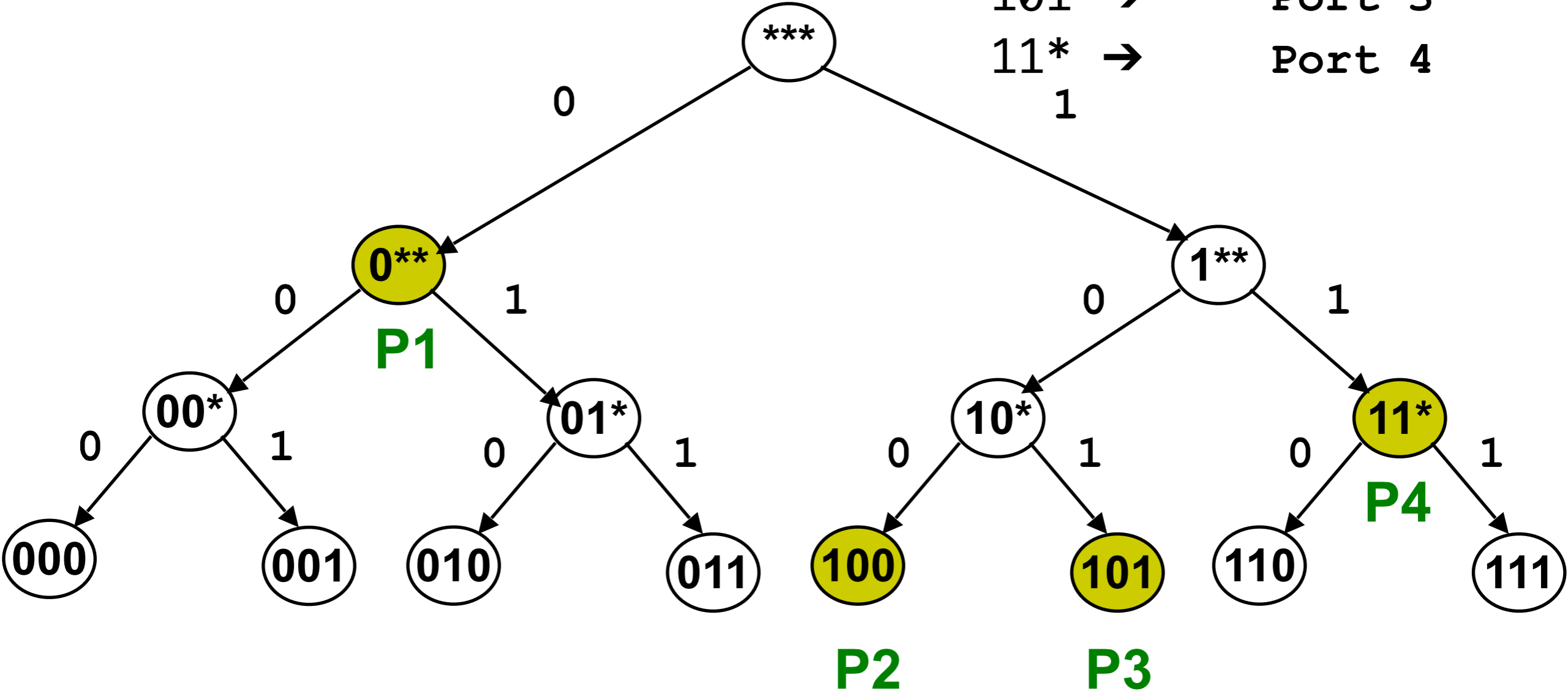


# Walk Tree: Stop at Prefix Entries



# Walk Tree: Stop at Prefix Entries

0\*\* → Port 1  
100 → Port 2  
101 → Port 3  
11\* → Port 4



walking trees takes  $O(\#bits)$

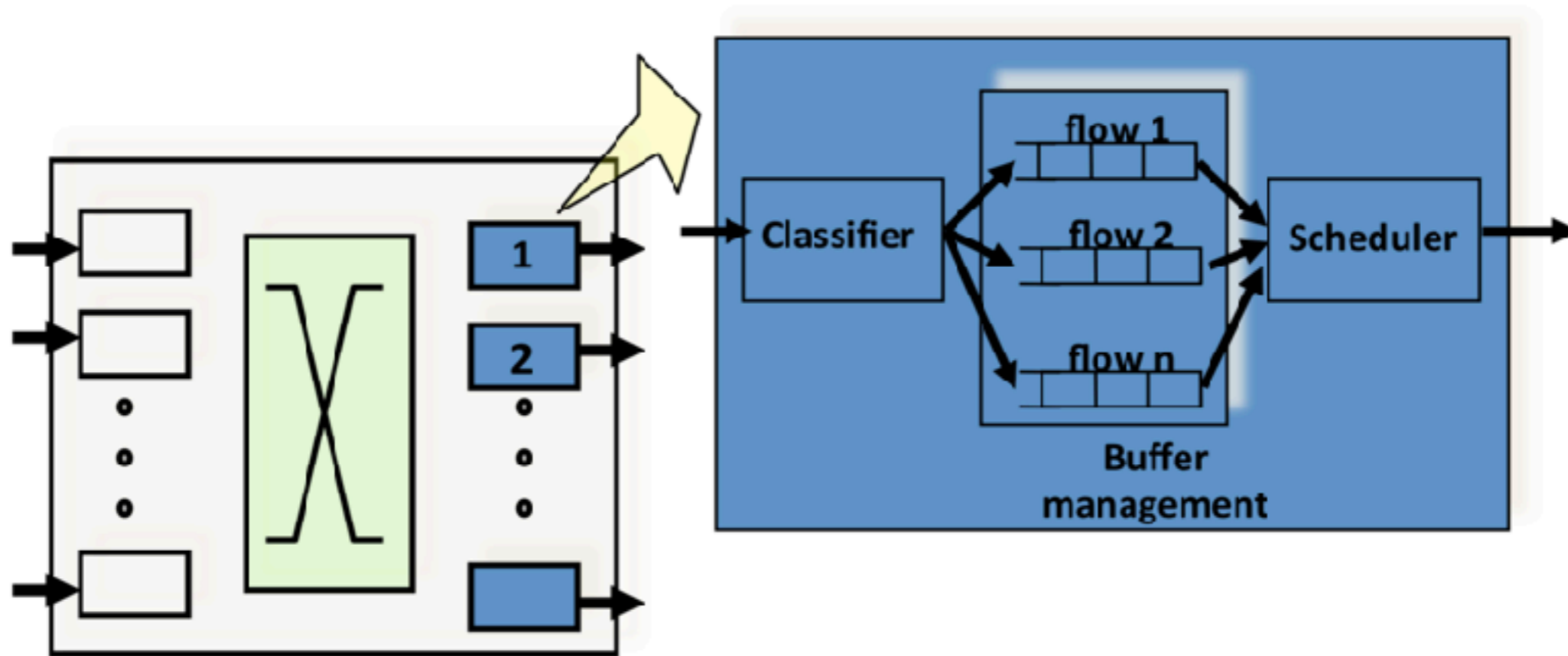
# Longest Prefix Match in Real Routers

- Real routers use far more advanced/complex solutions
  - But what we discussed is the starting point
- With many heuristics and optimizations that leverage real-world patterns
  - Some destinations more popular than others
  - Some ports lead to more destinations
  - Typical fix granularities

# Recap: Input Linecards

- Main challenge is processing speed
  - But what we discussed is the starting point
- Tasks involved
  - Update packet header (easy)
  - Longest prefix match lookup on destinations address (harder)
- Mostly implemented with specialized hardware

# Output Linecard



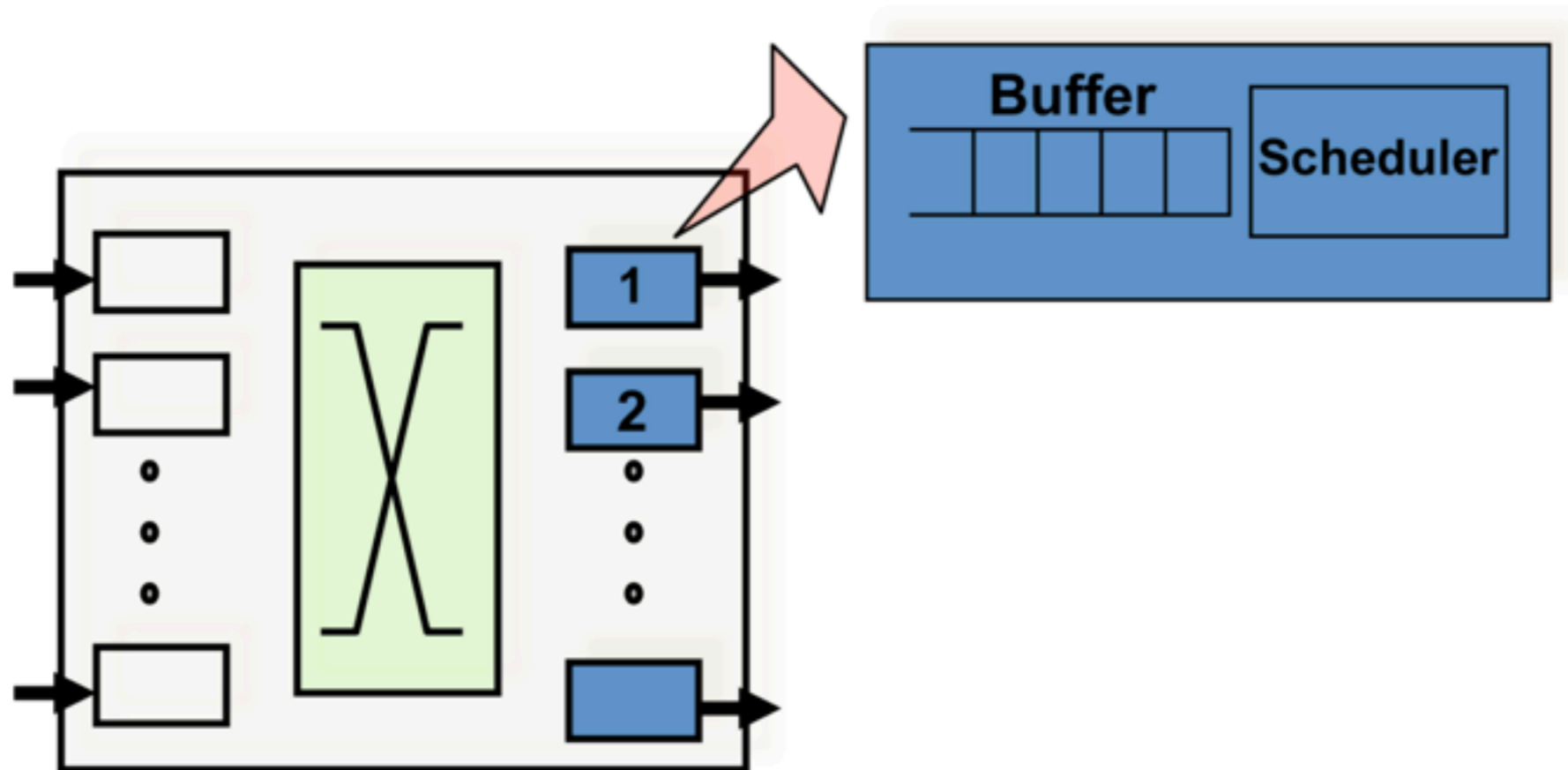
- **Packet Classification:** map each packet to a “flow”
  - Flow (for now): set of packets between two particular endpoints
- **Buffer Management:** decide when and which packet to drop
- **Scheduler:** decide when and which packet to transmit

# Output Linecard

- **Packet Classification:** map each packet to a “flow”
  - Flow (for now): set of packets between two particular endpoints
- **Buffer Management:** decide when and which packet to drop
- **Scheduler:** decide when and which packet to transmit
- Used to implement various forms of policy
  - Deny all e-mail traffic from ISP X to Y (**access control**)
  - Ensure that no more than 50 Mbps are injected from ISP-X (**QoS**)

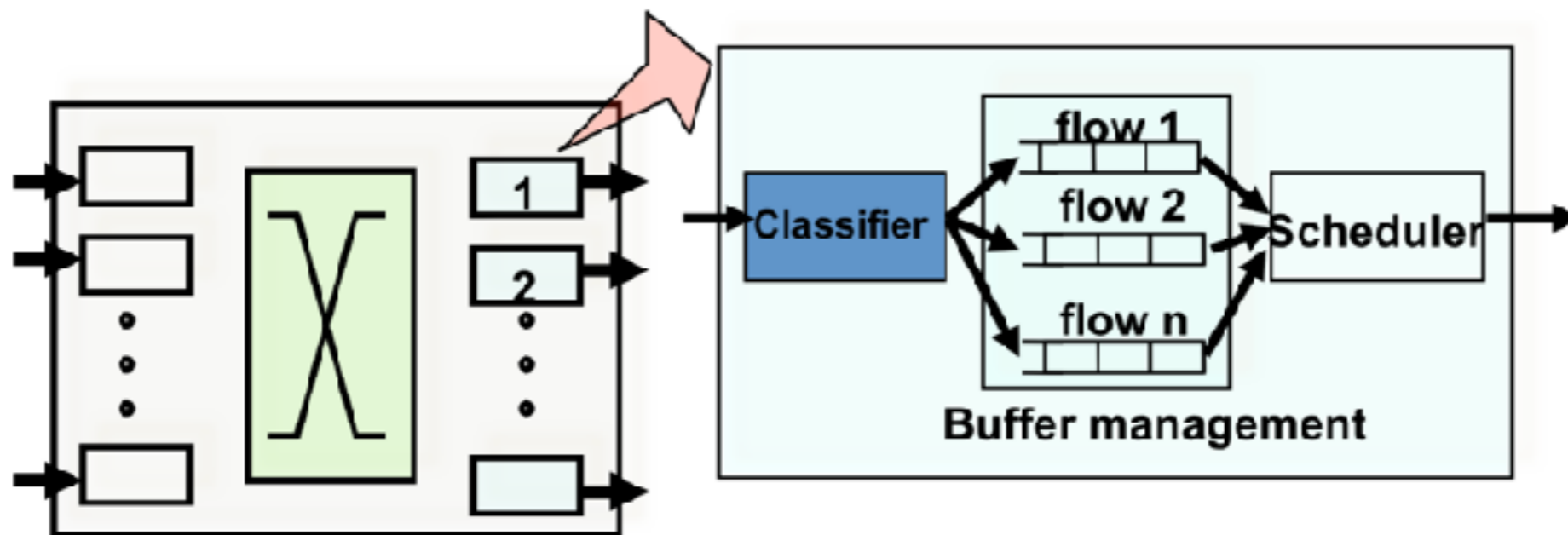
# Simplest FIFO Router

- No classification
- **Drop tail buffer management:** when buffer is full drop incoming packet
- **First In First Out (FIFO) Scheduling:** schedule packets in order of arrival



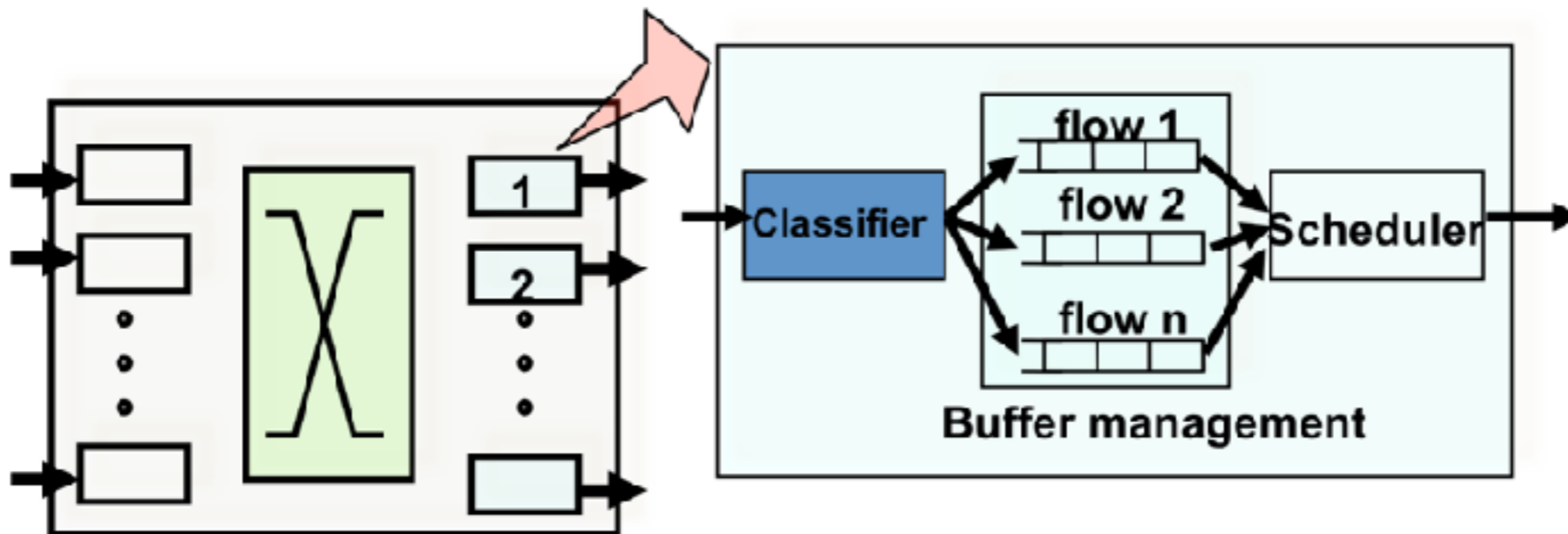
# Packet Classification

- Classify an IP packet based on the number of fields in the packet header
  - Source/destination IP address (32 bits)
  - Source/destination TCP port number (16 bits)
  - Type of Service (TOS) byte (8 bits)
  - Type of Protocol (8 bits)
- In general fields are specified by range
  - Classification requires a multi-dimensional range search



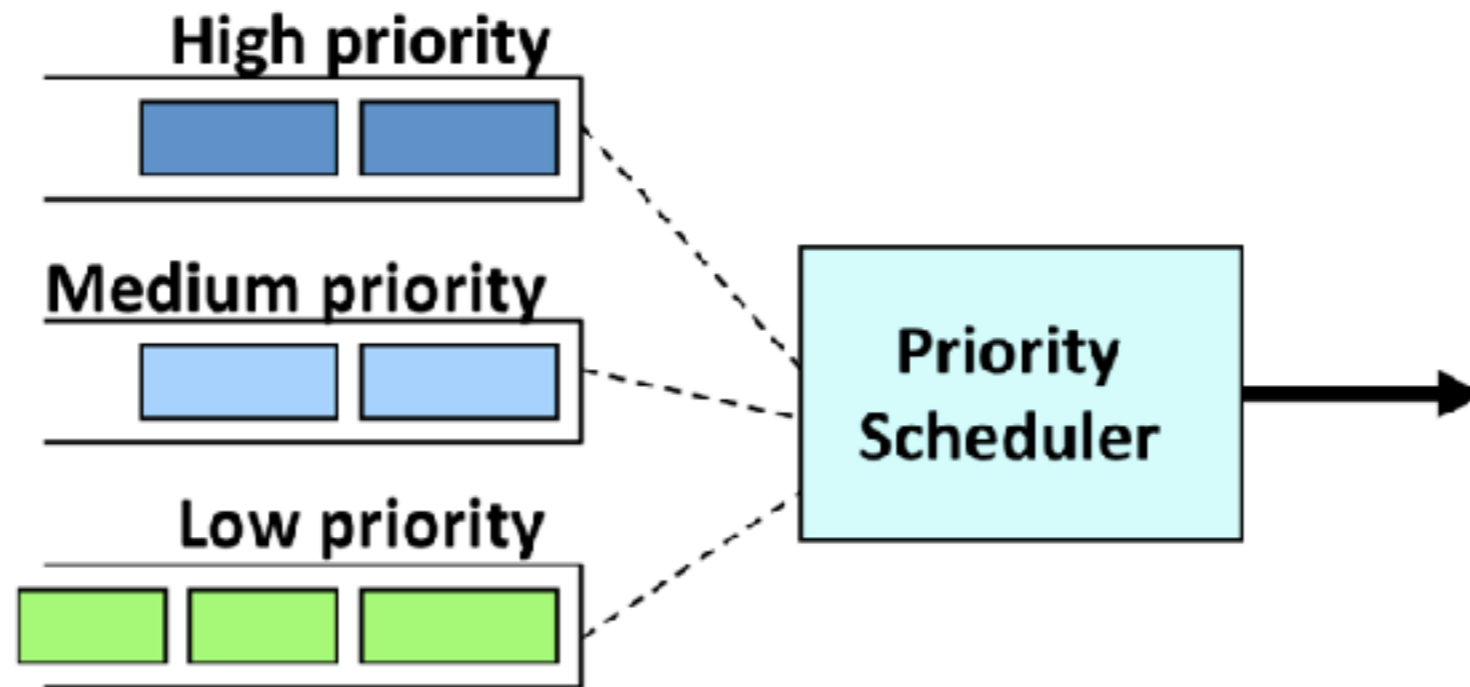
# Scheduler

- One queue per flow
- Scheduler decides from which queue to send a packet
- Goals of scheduling algorithm
  - Fast!
  - Depends on the policy being implemented (fairness, priority, etc.)



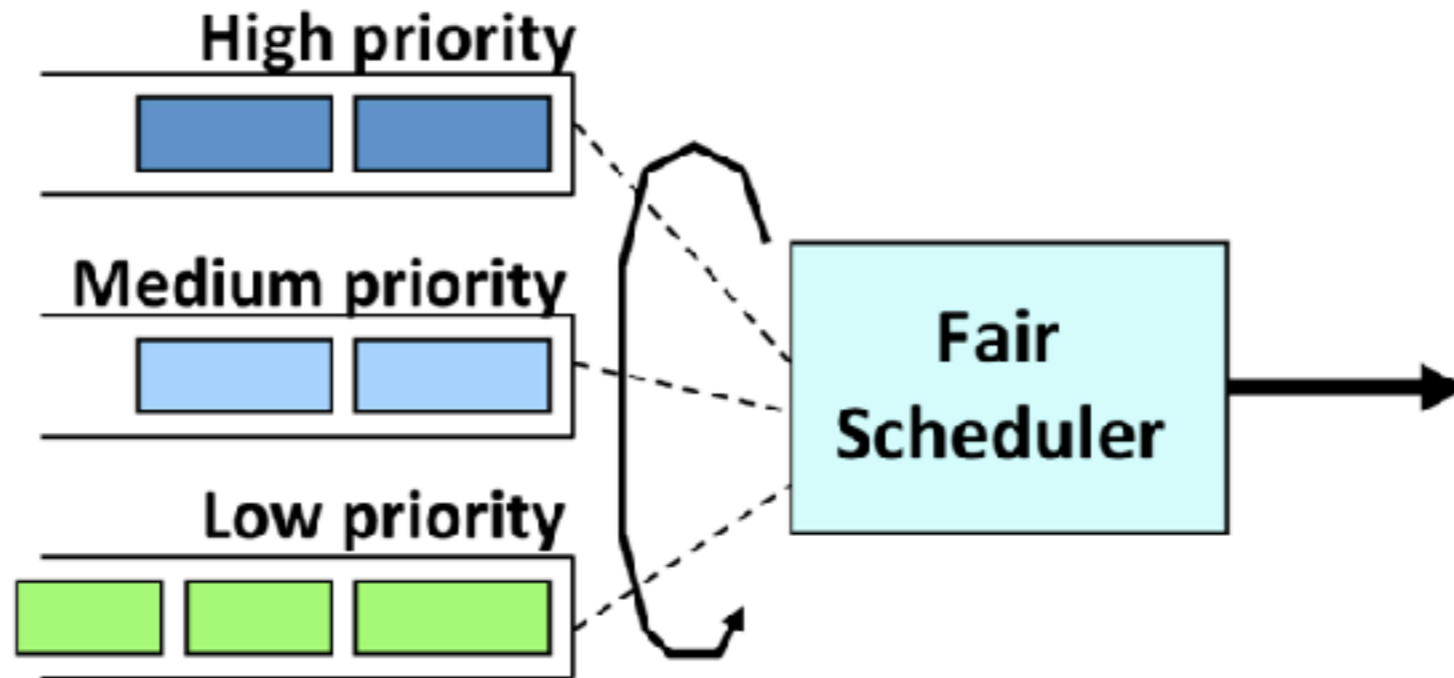
# Example: Priority Scheduler

- Packets in the highest priority queue are always served before the packets in the lower priority queues



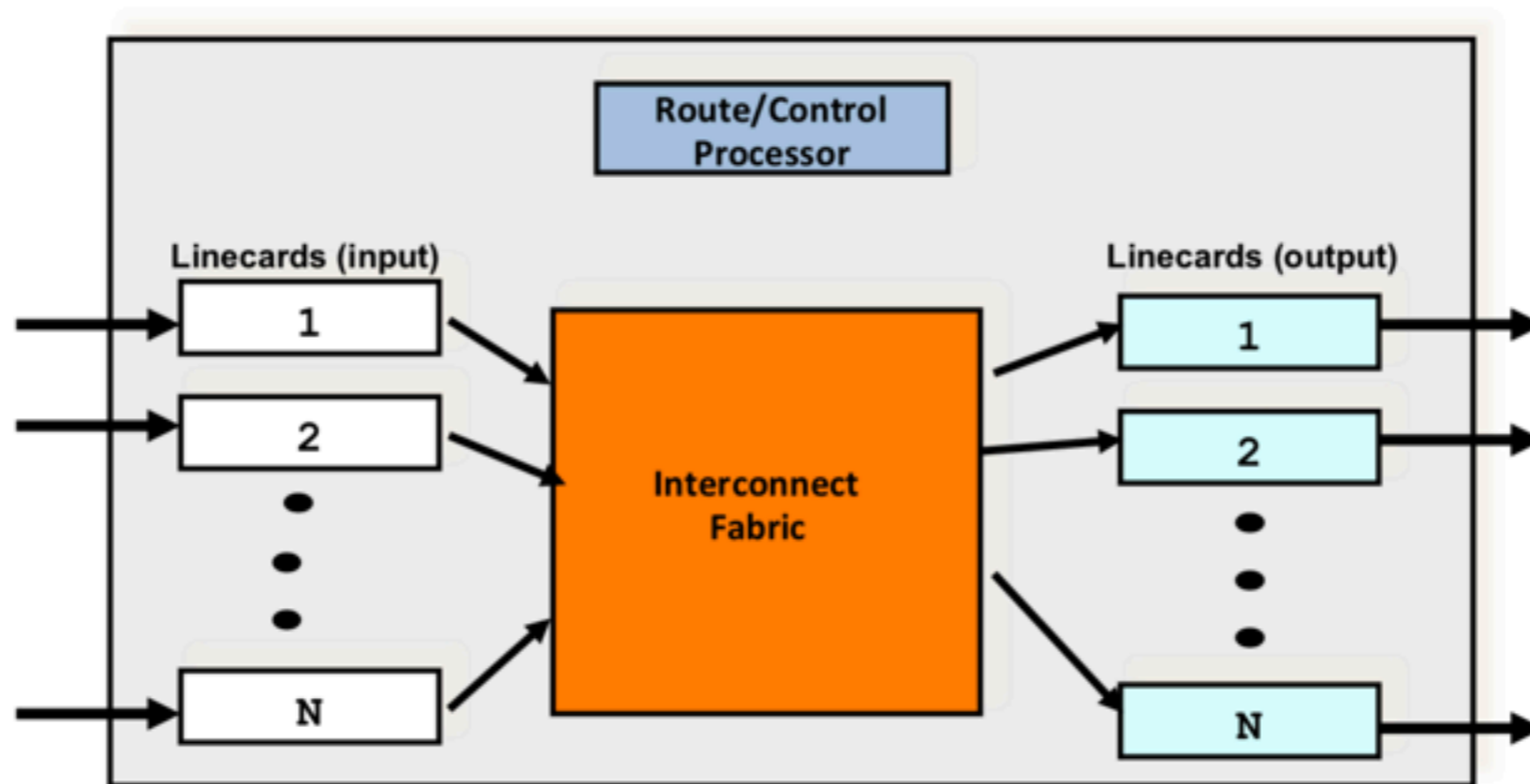
# Example: Round Robin Scheduler

- Packets are served from each queue in turn

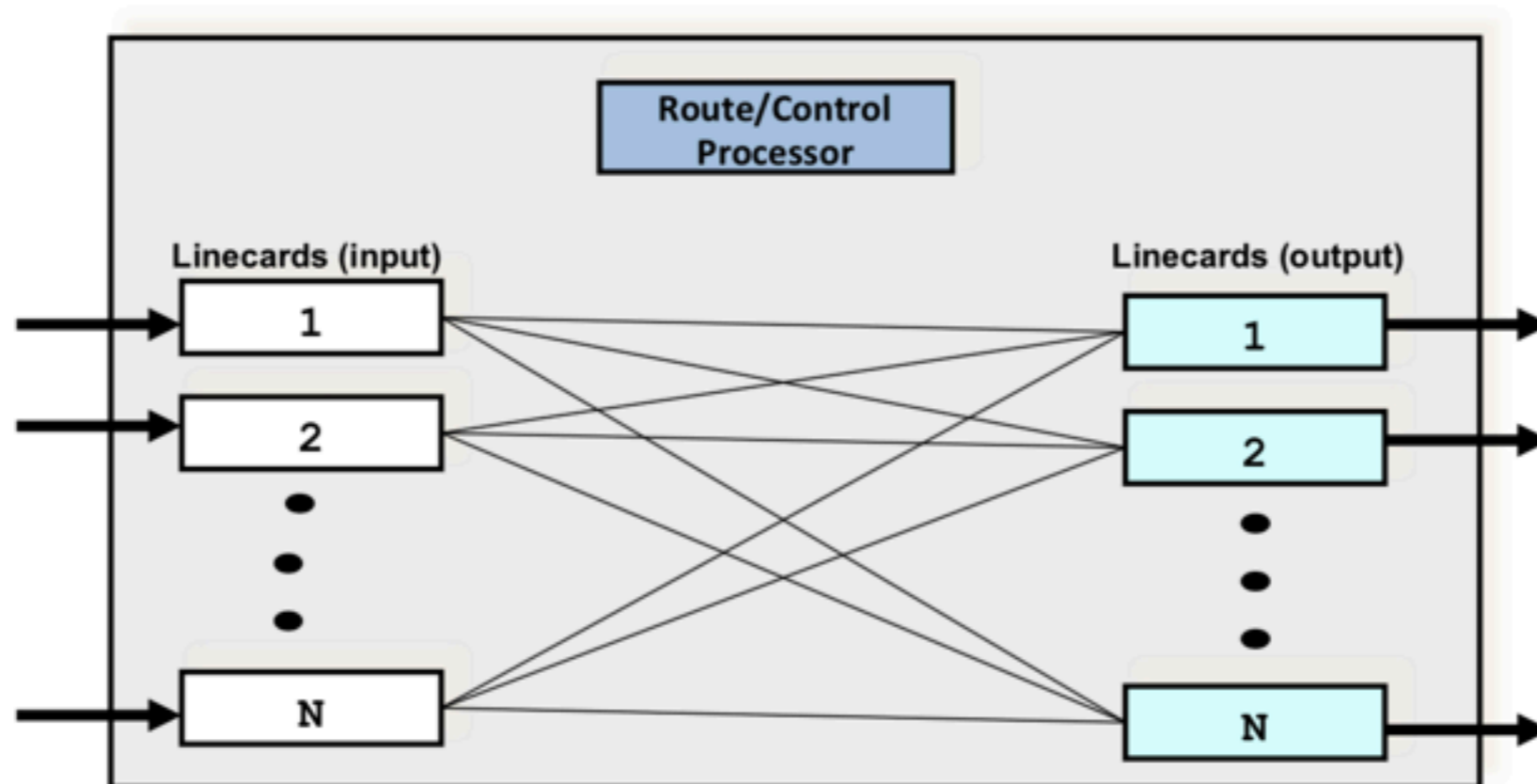


# Connecting Input to Output: Switch Fabric

- Priority Scheduler: packets are served from each queue in turn



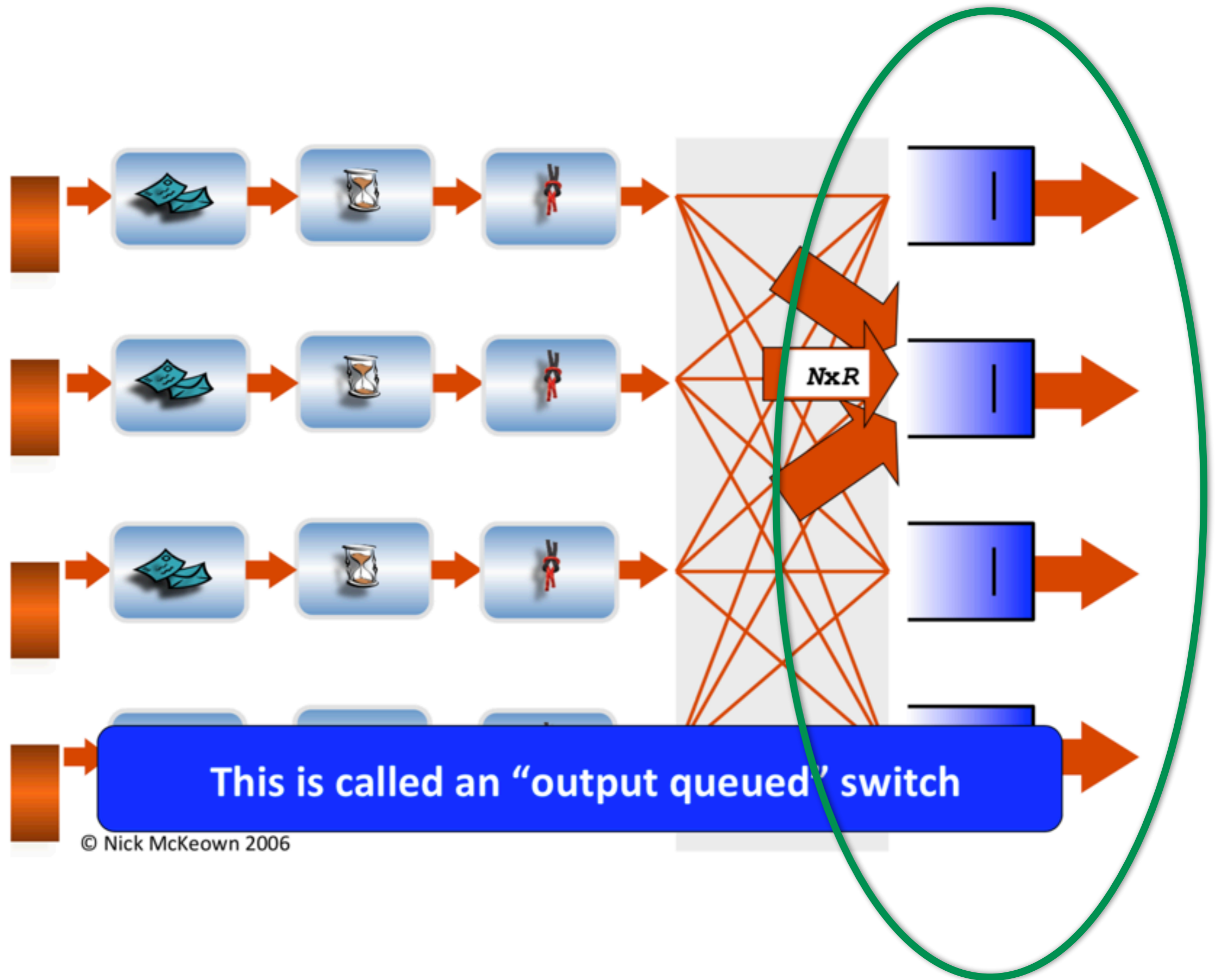
# Today's Switch Fabrics: Mini Network!



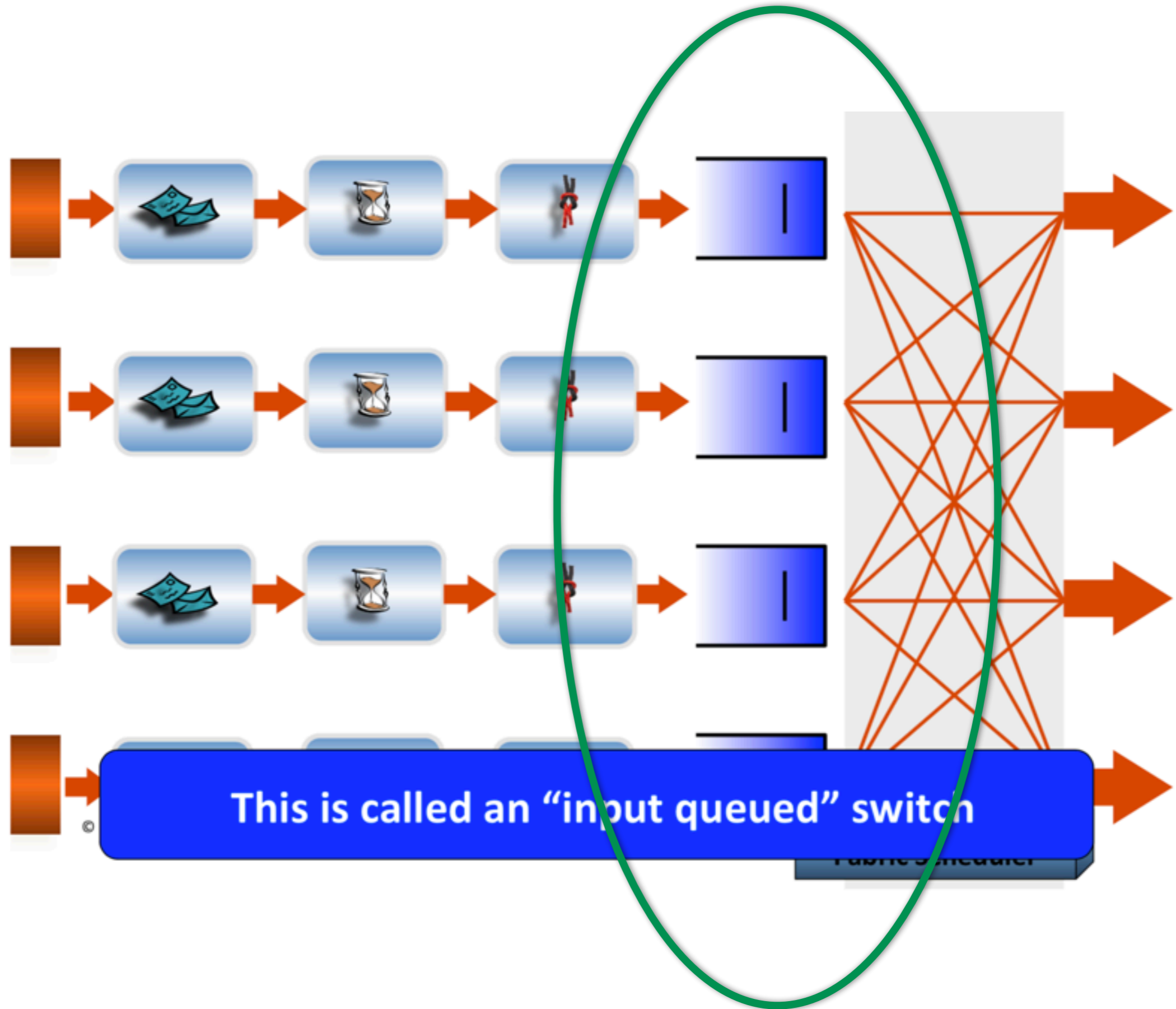
# What's Hard About the Switch Fabric?

**Queueing!**

# Third Generation Router: Switched Interconnects



# Third Generation Router: Switched Interconnects



# Reality is More Complicated

- Commercial high-speed routers use
  - Combination of input and output queueing
  - Complex multi-stage “topologies”
  - Distributed multi-stage schedulers (for scalability)

# IP Routers Recap

- Core building block of Internet infrastructure
- Scalable Routing -> Longest Prefix Matching
- Need fast implementations for
  - Longest prefix matching
  - Switch fabric scheduling

# What do we know so far [1] ...

- **Network performance metrics**
  - Transmission delay, propagation delay, queueing delay, bandwidth
- **Sharing networks**
  - Circuit switching, packet switching, and associated tradeoffs
  - **Why** is Internet packet switched?
- **Architectural principles and design goals**
  - Layering principle, End-to-end principle, Fate sharing principle
  - Many important design goals from David Clark's paper
    - And many important missing goals
- **Addressing**
  - **Link layer MAC names**, and scalability challenges at the Internet
  - **Network layer IP addresses**: three requirements, aggregation, CIDR

# What do we know so far [2] ...

- **Link Layer**

- Sharing a Broadcast medium, associated challenges, CSMA/CD
- Link layer addressing: MAC names
- **Why Frames? Why Switched Ethernet?**
- The Spanning Tree Protocol (STP)

- **Network Layer**

- **Why Network Layer? Why not just use STP across the Internet?**
- **Routing Tables:** A collection of spanning trees, one per destination
- **Generating Valid Routing tables (within a domain):**
  - Global view (Link-State Protocol), and limitations
  - Local view (Distance-vector Protocol)
- **Generating Valid Routing tables (across domains):**
  - Border Gateway Protocol, Internet structure, routing policies

# Network Layer

- THE functionality: **delivering the data**
- **THE protocol: Internet Protocol (IP)**
- **Achieves its functionality (delivering the data), using three ideas:**
  - **Addressing** (IP addressing)
  - **Routing** (using a variety of protocols)
  - **Packet header as an interface** (Encapsulating data into packets)

