

# NETWORK LAYER

Lets Talk

# Two key network-layer functions

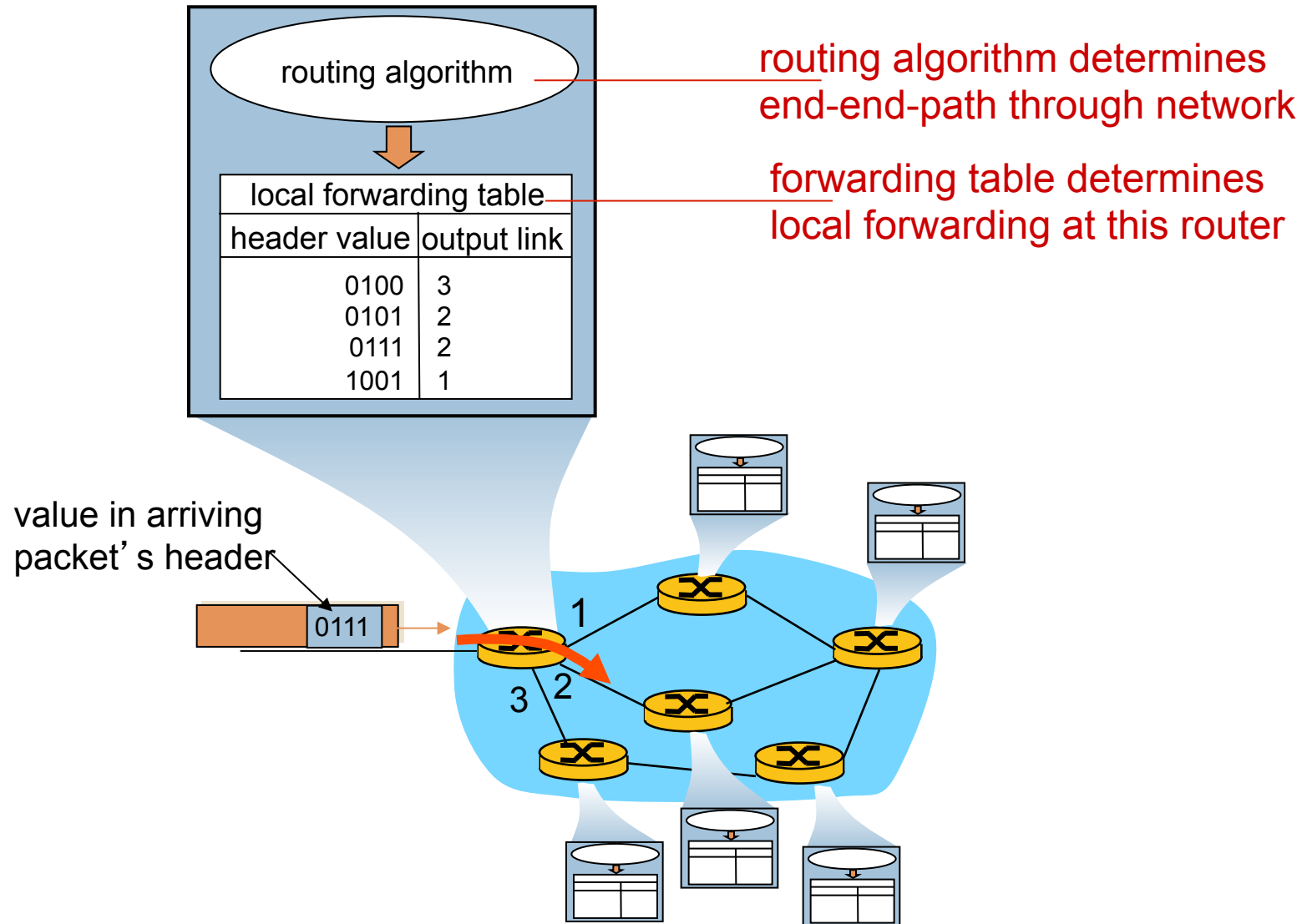
- *forwarding*: move packets from router's input to appropriate router output
- *routing*: determine route taken by packets from source to dest.

- *routing algorithms*

*analogy:*

- ❖ *routing*: process of planning trip from source to dest
- ❖ *forwarding*: process of getting through single interchange

# Interplay between routing and forwarding



# Connection setup

- before datagrams flow, two end hosts *and* intervening routers establish virtual connection
  - ▣ routers get involved
- network vs transport layer connection service:
  - ▣ *network*: between two hosts (may also involve intervening routers in case of VCs)
  - ▣ *transport*: between two processes

# Virtual circuits

“source-to-dest path behaves much like telephone circuit”

- ▣ performance-wise
- ▣ network actions along source-to-dest path

- call setup, teardown for each call *before* data can flow
- each packet carries VC identifier (not destination host address)
- every router on source-dest path maintains “state” for each passing connection
- link, router resources (bandwidth, buffers) may be *allocated* to VC (dedicated resources = predictable service)

# VC implementation

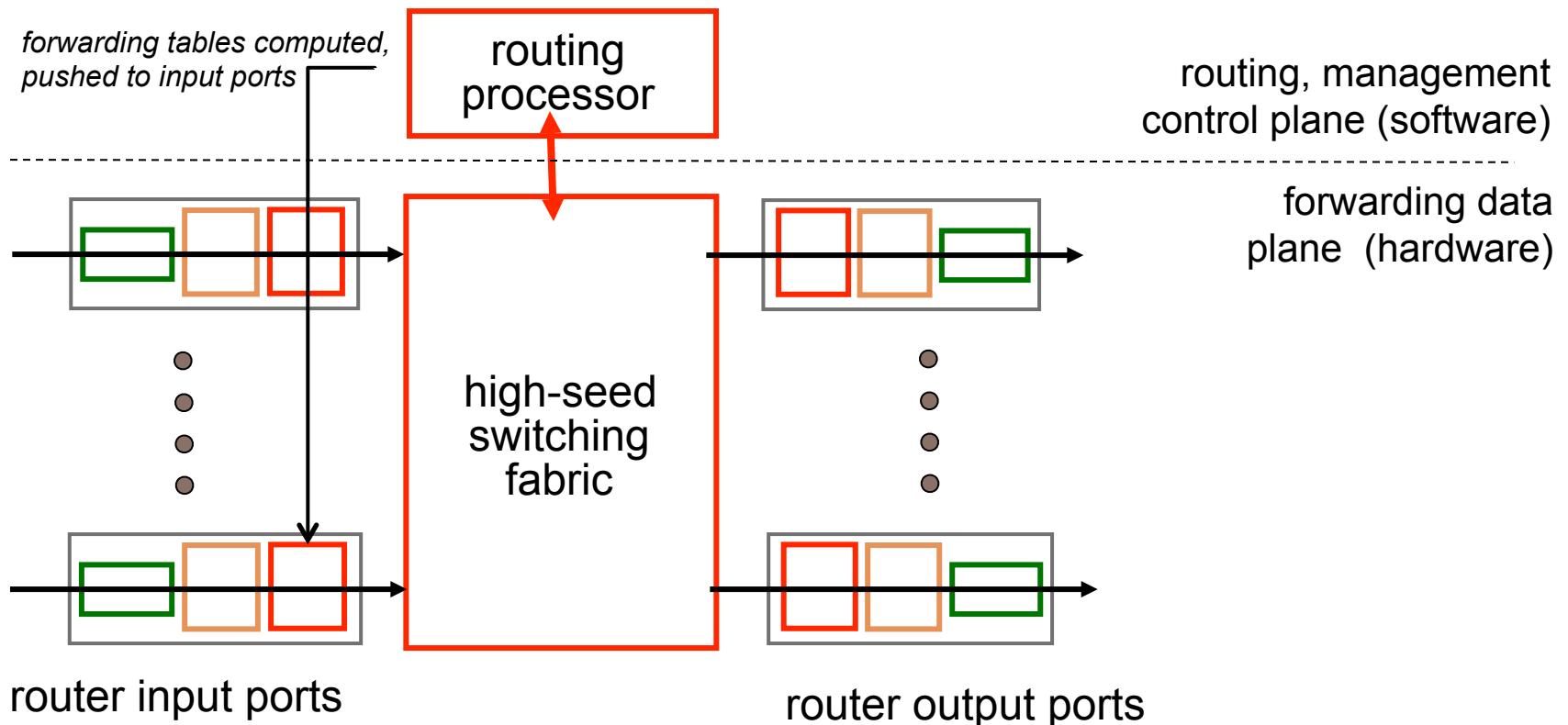
*a VC consists of:*

1. *path* from source to destination
  2. *VC numbers*, one number for each link along path
  3. *entries in forwarding tables* in routers along path
- packet belonging to VC carries VC number (rather than dest address)
  - VC number can be changed on each link.
    - ▣ new VC number comes from forwarding table

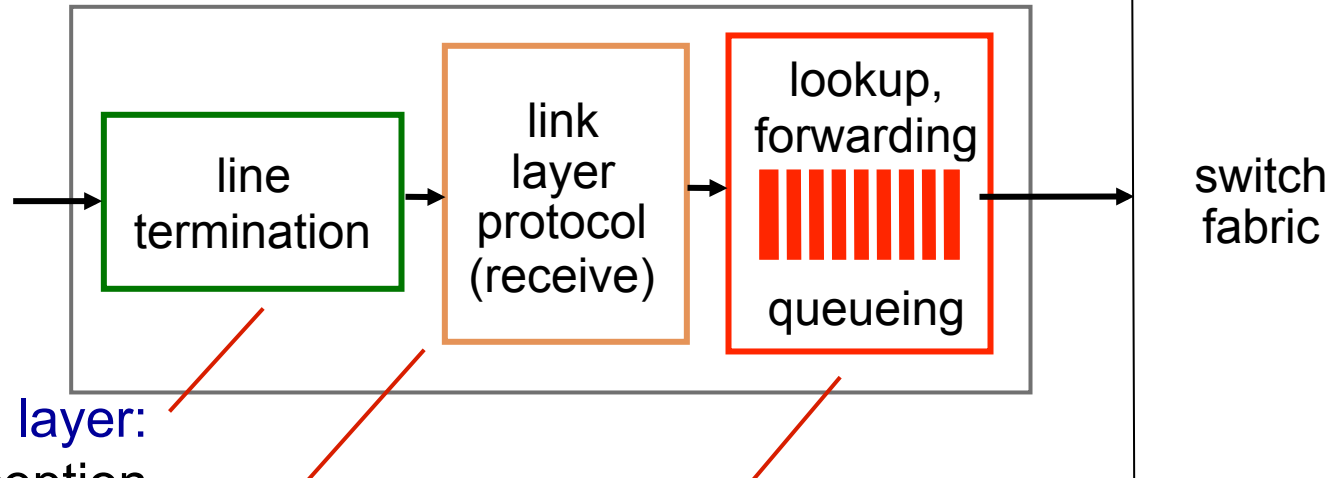
# Router Architecture Overview

two key router functions:

- run routing algorithms/protocol (RIP, OSPF, BGP)
- *forwarding* datagrams from incoming to outgoing link



# Input Port Functions



physical layer:  
bit-level reception

data link layer:  
e.g., Ethernet  
see chapter 5

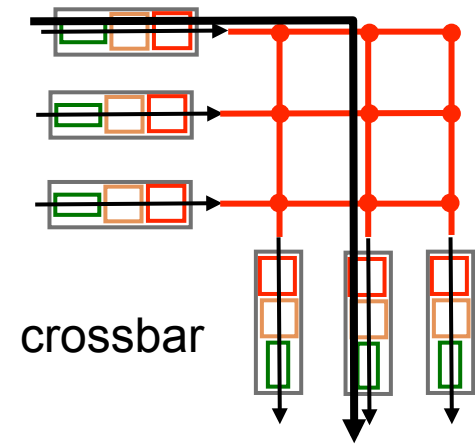
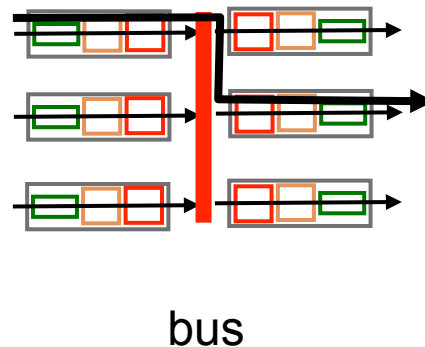
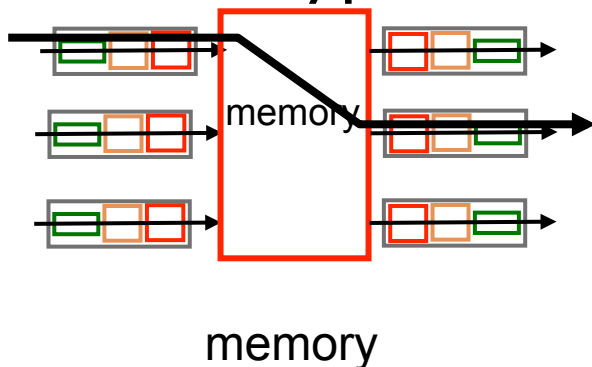
## decentralized switching:

- given datagram dest., lookup output port using forwarding table in input port memory (“*match plus action*”)
- goal: complete input port processing at ‘line speed’
- queuing: if datagrams arrive faster than forwarding rate into switch fabric



# Switching fabrics

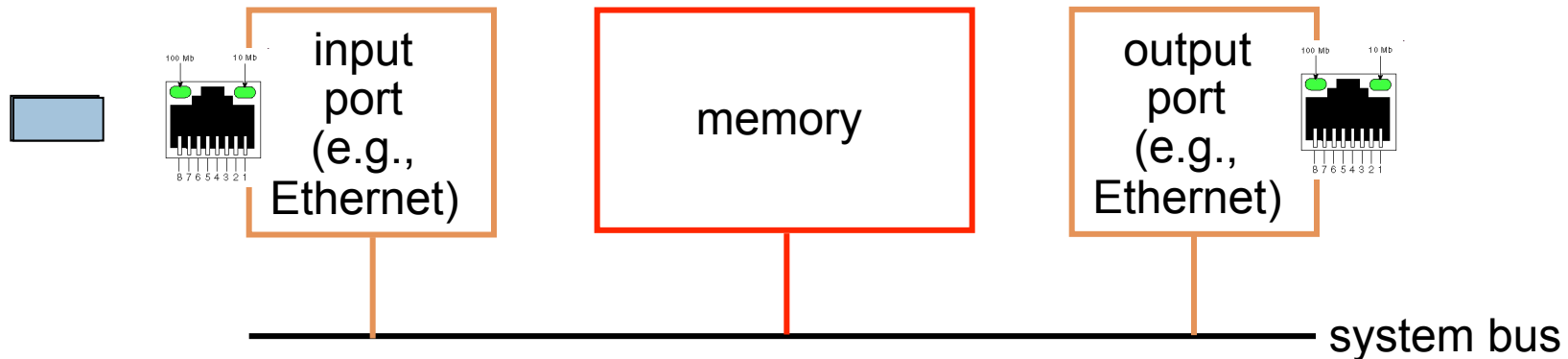
- transfer packet from input buffer to appropriate output buffer
- switching rate: rate at which packets can be transfer from inputs to outputs
  - ▣ often measured as multiple of input/output line rate
  - ▣ N inputs: switching rate N times line rate desirable
- three types of switching fabrics



# Switching via memory

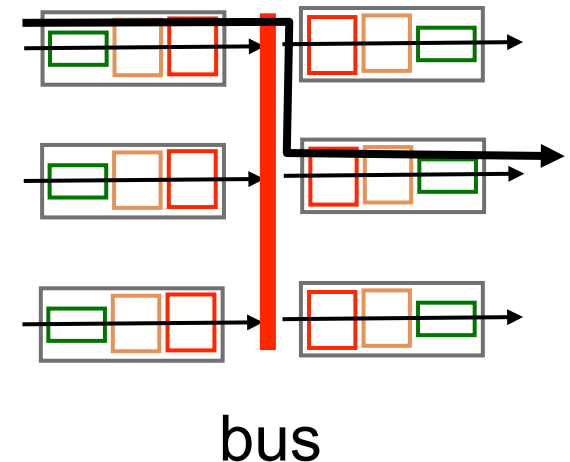
## *first generation routers:*

- traditional computers with switching under direct control of CPU
- packet copied to system's memory
- speed limited by memory bandwidth (2 bus crossings per datagram)



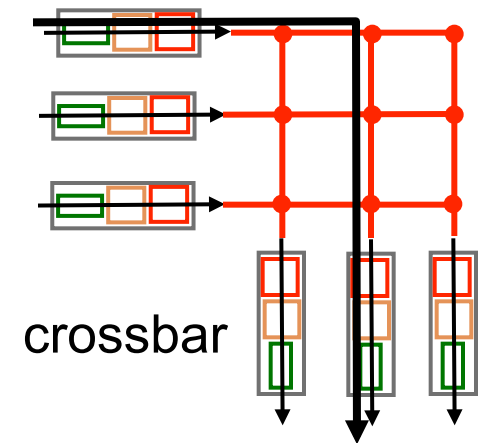
# Switching via a bus

- datagram from input port memory to output port memory via a shared bus
  - to output port memory via a shared bus
- *bus contention*: switching speed limited by bus bandwidth
- 32 Gbps bus, Cisco 5600: sufficient speed for access and enterprise routers

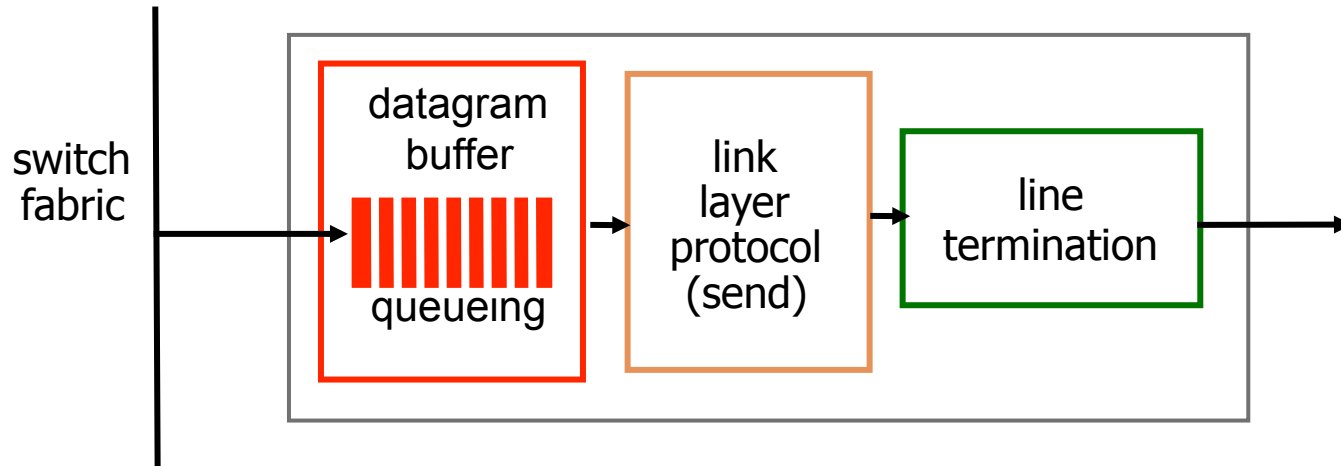


# Switching via interconnection network

- overcome bus bandwidth limitations
- banyan networks, crossbar, other interconnection nets initially developed to connect processors in multiprocessor
- advanced design: fragmenting datagram into fixed length cells, switch cells through the fabric.
- Cisco 12000: switches 60 Gbps through the interconnection network



# Output Ports

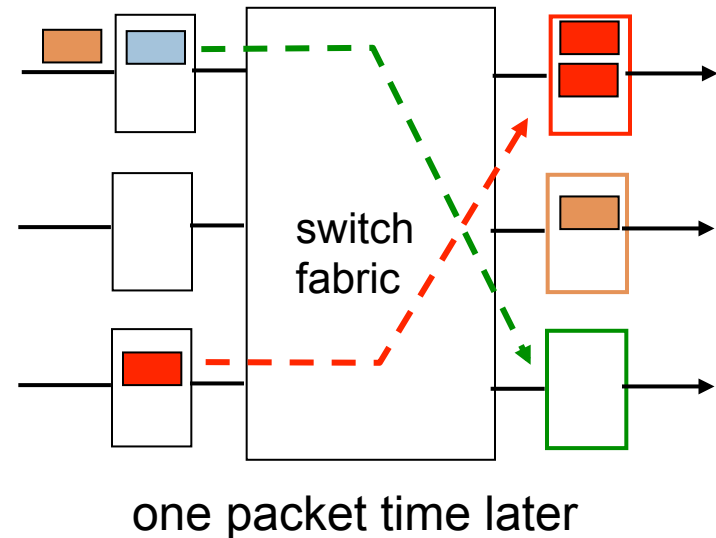
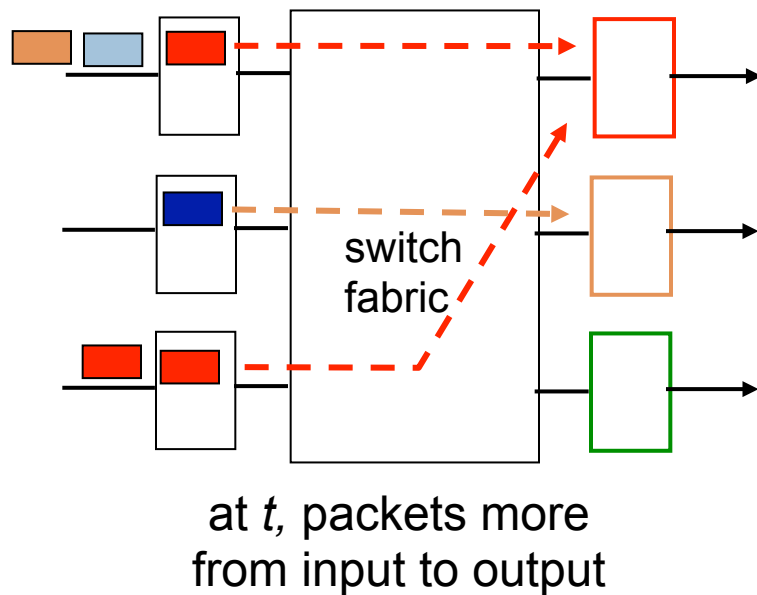


- *buffering* required when datagrams arrive from fabric faster than the transmission rate
- *scheduling discipline* chooses among queued datagrams for transmission

Datagram (packets) can be lost due to congestion, lack of buffers

Priority scheduling – who gets best performance, network neutrality

# Output Port Queuing



- buffering when arrival rate via switch exceeds output line speed
- *queueing (delay) and loss due to output port buffer overflow!*

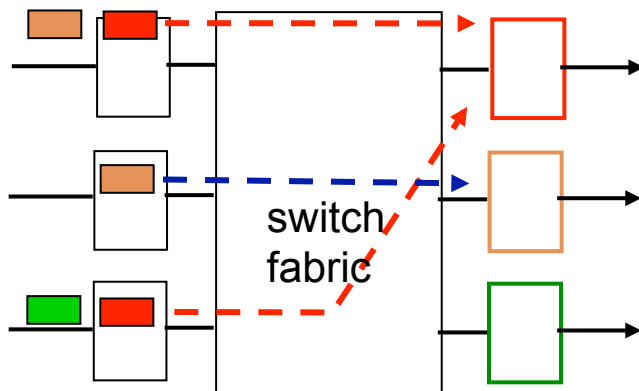
# How much buffering?

- RFC 3439 rule of thumb: average buffering equal to “typical” RTT (say 250 msec) times link capacity  $C$ 
  - e.g.,  $C = 10$  Gpbs link: 2.5 Gbit buffer
- recent recommendation: with  $N$  flows, buffering equal to

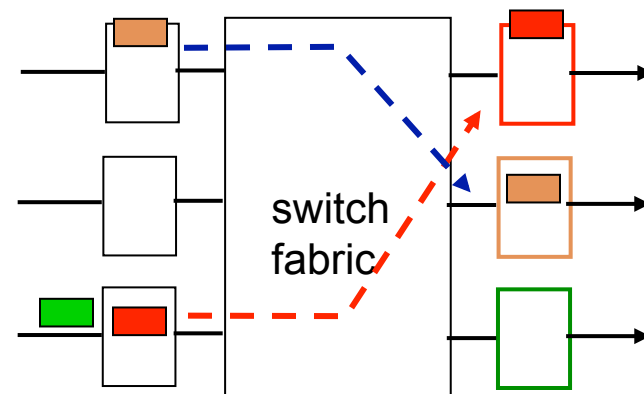
$$\frac{RTT \cdot C}{\sqrt{N}}$$

# Input port queuing

- fabric slower than input ports combined -> queueing may occur at input queues
  - ▣ *queueing delay and loss due to input buffer overflow!*
- **Head-of-the-Line (HOL) blocking:** queued datagram at front of queue prevents others in queue from moving forward



output port contention:  
only one red datagram can be  
transferred.  
*lower red packet is blocked*

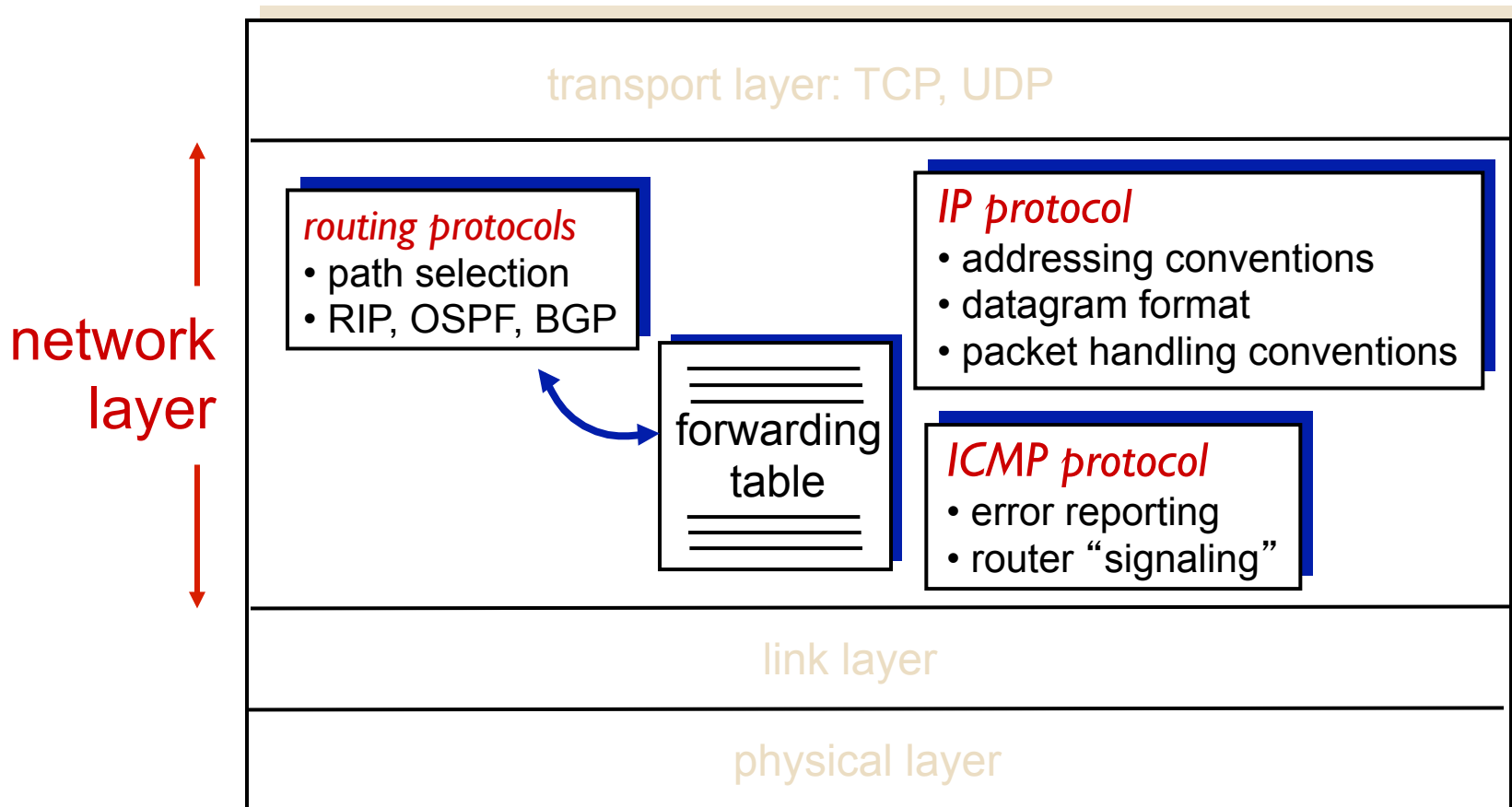


one packet time  
later: green packet  
experiences HOL  
blocking



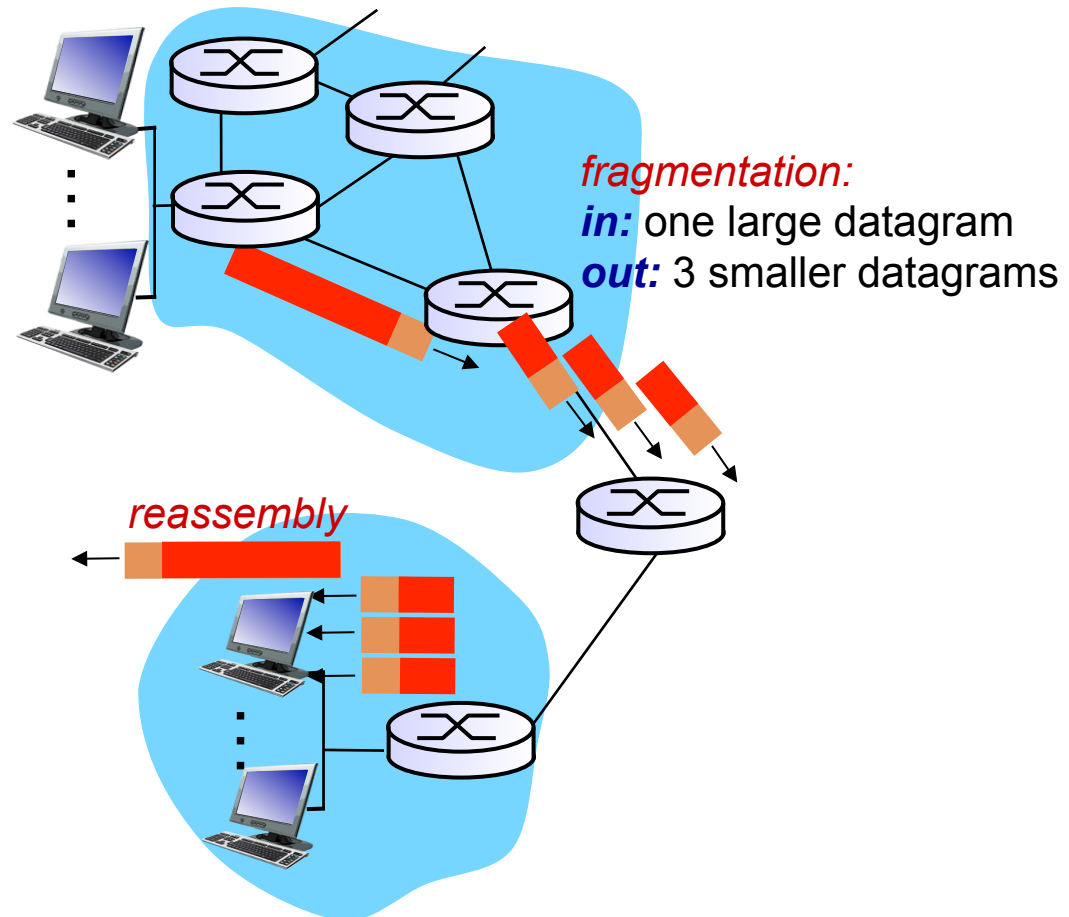
# The Internet Network Layer

host, router network layer functions:



# IP fragmentation, reassembly

- network links have MTU (max.transfer size) - largest possible link-level frame
  - ▣ different link types, different MTUs
- large IP datagram divided (“fragmented”) within net
  - ▣ one datagram becomes several datagrams
  - ▣ “reassembled” only at final destination
  - ▣ IP header bits used to identify, order related fragments



# IP fragmentation, reassembly

## *example:*

- ❖ 4000 byte datagram
- ❖ MTU = 1500 bytes

	length =4000	ID =x	fragflag =0	offset =0	
--	-----------------	----------	----------------	--------------	--

*one large datagram becomes  
several smaller datagrams*

1480 bytes in  
data field

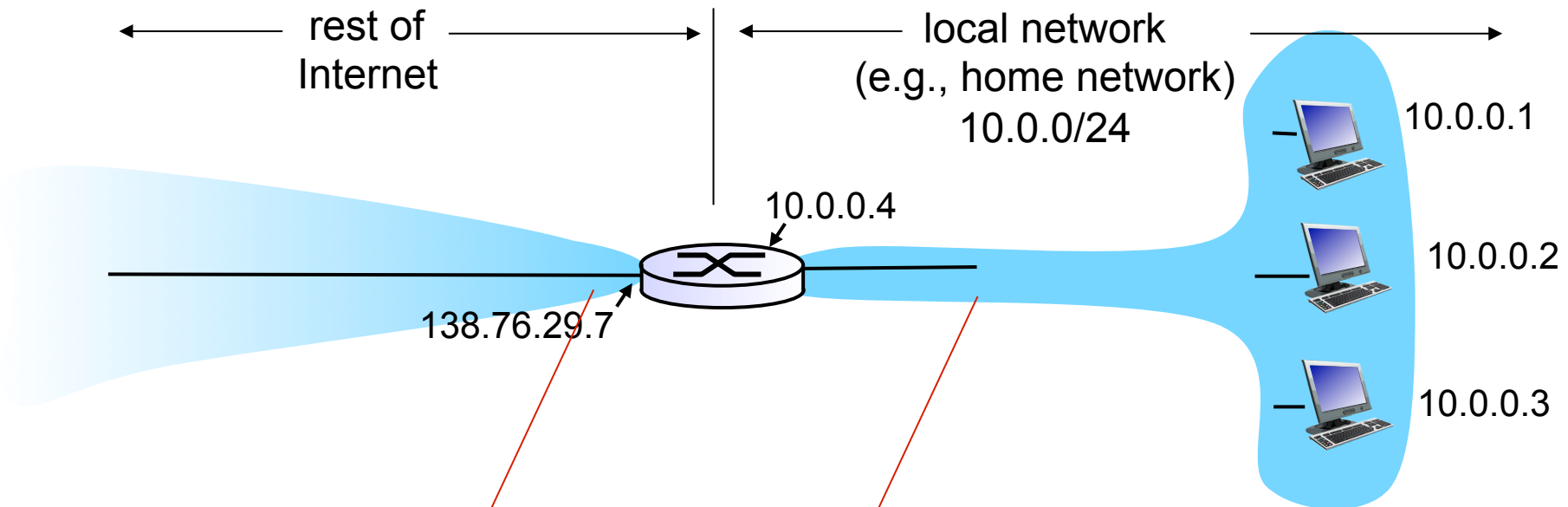
offset =  
 $1480/8$

	length =1500	ID =x	fragflag =1	offset =0	
--	-----------------	----------	----------------	--------------	--

	length =1500	ID =x	fragflag =1	offset =185	
--	-----------------	----------	----------------	----------------	--

	length =1040	ID =x	fragflag =0	offset =370	
--	-----------------	----------	----------------	----------------	--

# NAT: Network Address Translation



*all* datagrams *leaving* local network have *same* single source NAT IP address: 138.76.29.7, different source port numbers

datagrams with source or destination in this network have 10.0.0/24 address for source, destination (as usual)

# NAT: Motivation

local network uses just one IP address as far as outside world is concerned:

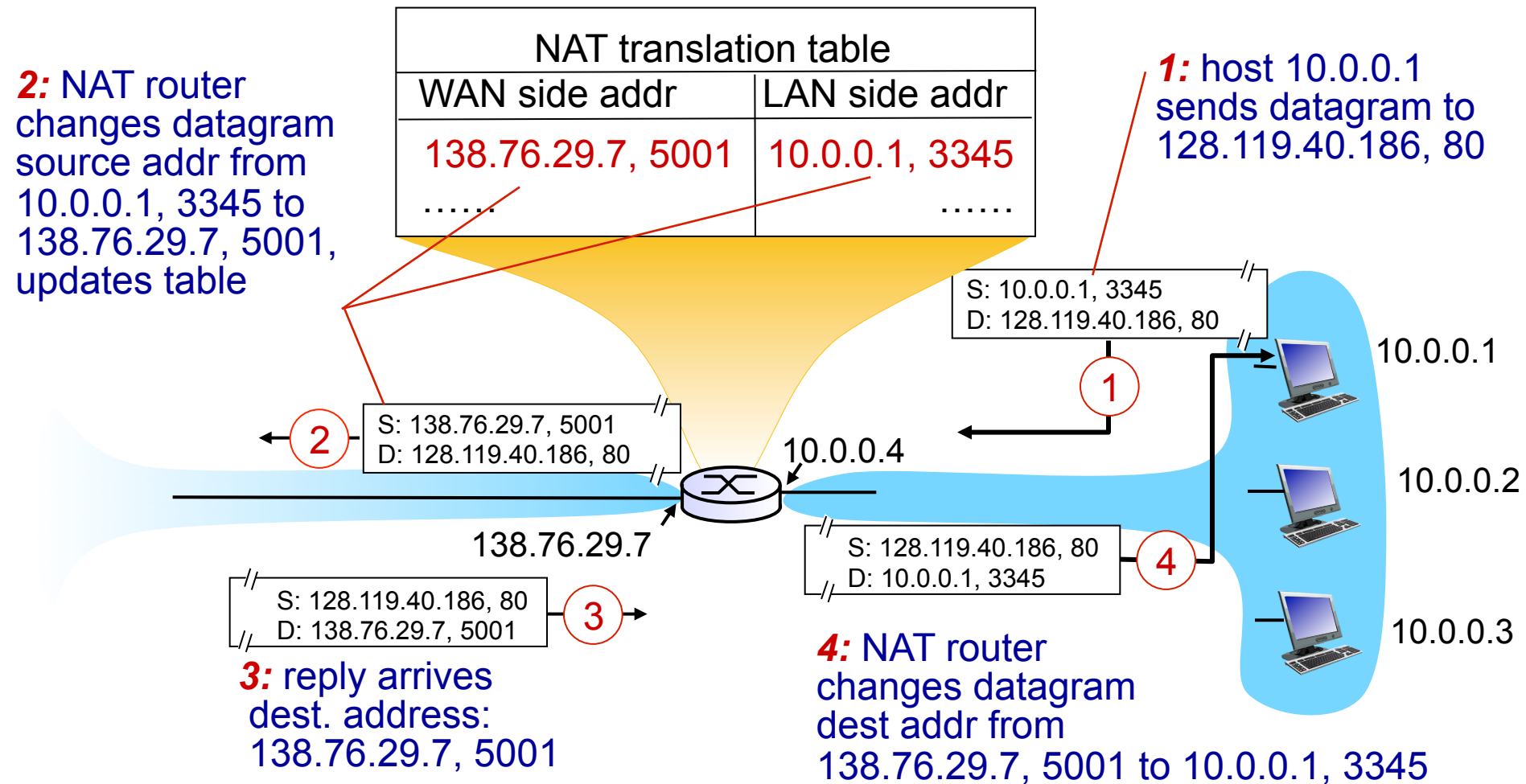
- ▣ range of addresses not needed from ISP: just one IP address for all devices
- ▣ can change addresses of devices in local network without notifying outside world
- ▣ can change ISP without changing addresses of devices in local network
- ▣ devices inside local net not explicitly addressable, visible by outside world (a security plus)

# NAT: Implementation

NAT router must:

- ▣ *outgoing datagrams: replace* (source IP address, port #) of every outgoing datagram to (NAT IP address, new port #) ... remote clients/servers will respond using (NAT IP address, new port #) as destination addr
- ▣ *remember (in NAT translation table)* every (source IP address, port #) to (NAT IP address, new port #) translation pair
- ▣ *incoming datagrams: replace* (NAT IP address, new port #) in dest fields of every incoming datagram with corresponding (source IP address, port #) stored in NAT table

# NAT: Example



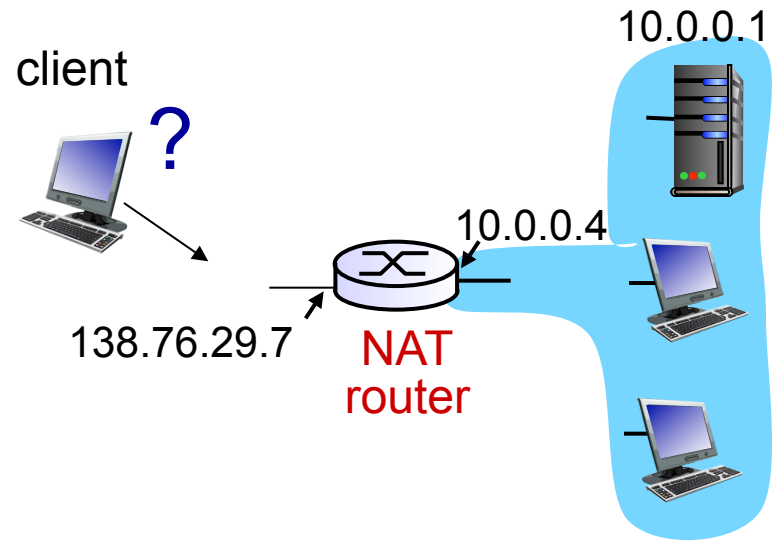
# NAT: Issues

- 16-bit port-number field:
  - ▣ 60,000 simultaneous connections with a single LAN-side address!
- NAT is controversial:
  - ▣ routers should only process up to layer 3
  - ▣ violates end-to-end argument
    - NAT possibility must be taken into account by app designers, e.g., P2P applications
  - ▣ address shortage should instead be solved by IPv6



# NAT: Traversal Problem – Solution1

- client wants to connect to server with address 10.0.0.1
  - ▣ server address 10.0.0.1 local to LAN (client can't use it as destination addr)
  - ▣ only one externally visible NATed address: 138.76.29.7
- **solution 1:** statically configure NAT to forward incoming connection requests at given port to server
  - ▣ e.g., (138.76.29.7, port 2500) always forwarded to 10.0.0.1 port 25000

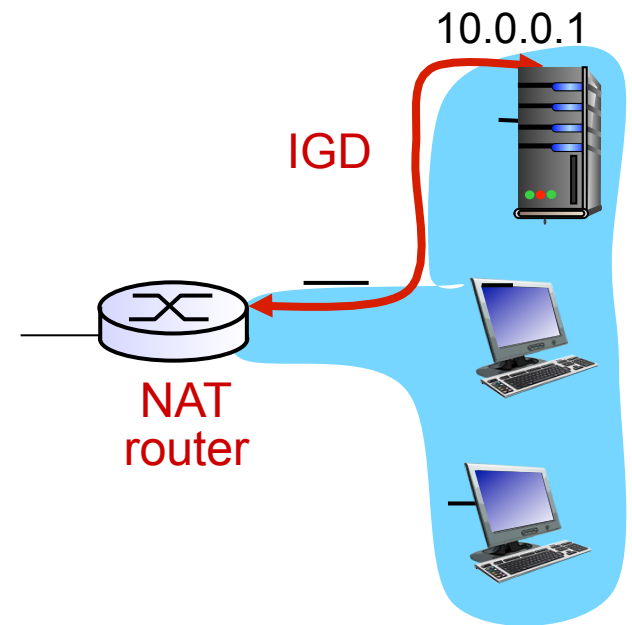


# NAT: Traversal Problem – Solution2

- *solution 2*: Universal Plug and Play (UPnP) Internet Gateway Device (IGD) Protocol. Allows NATed host to:

- ❖ learn public IP address (138.76.29.7)
- ❖ add/remove port mappings (with lease times)

i.e., automate static NAT port map configuration



# NAT: Traversal Problem – Solution 3

- **solution 3:** relaying (used in Skype)
  - ▣ NATed client establishes connection to relay
  - ▣ external client connects to relay
  - ▣ relay bridges packets between to connections

