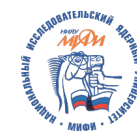
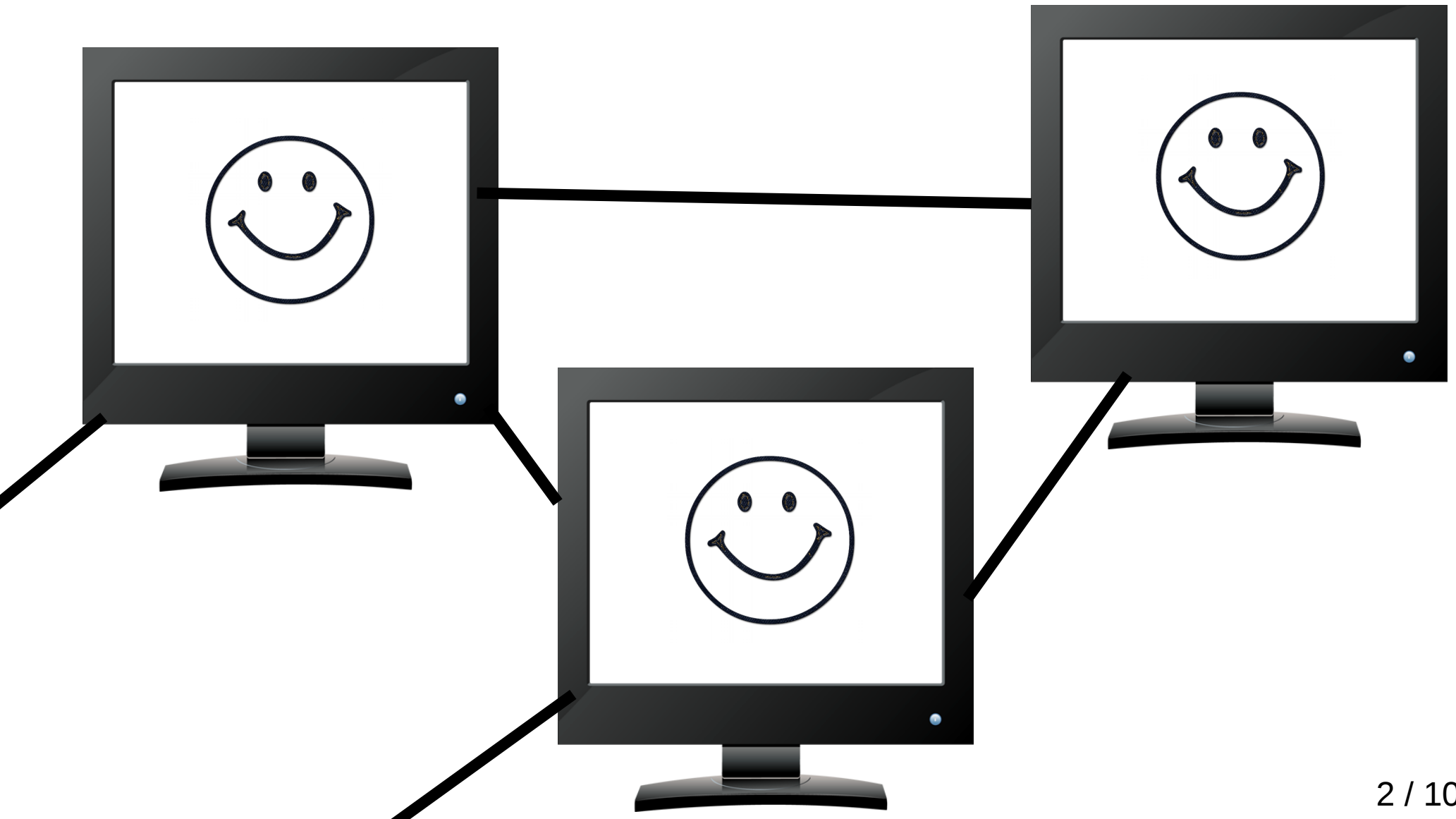


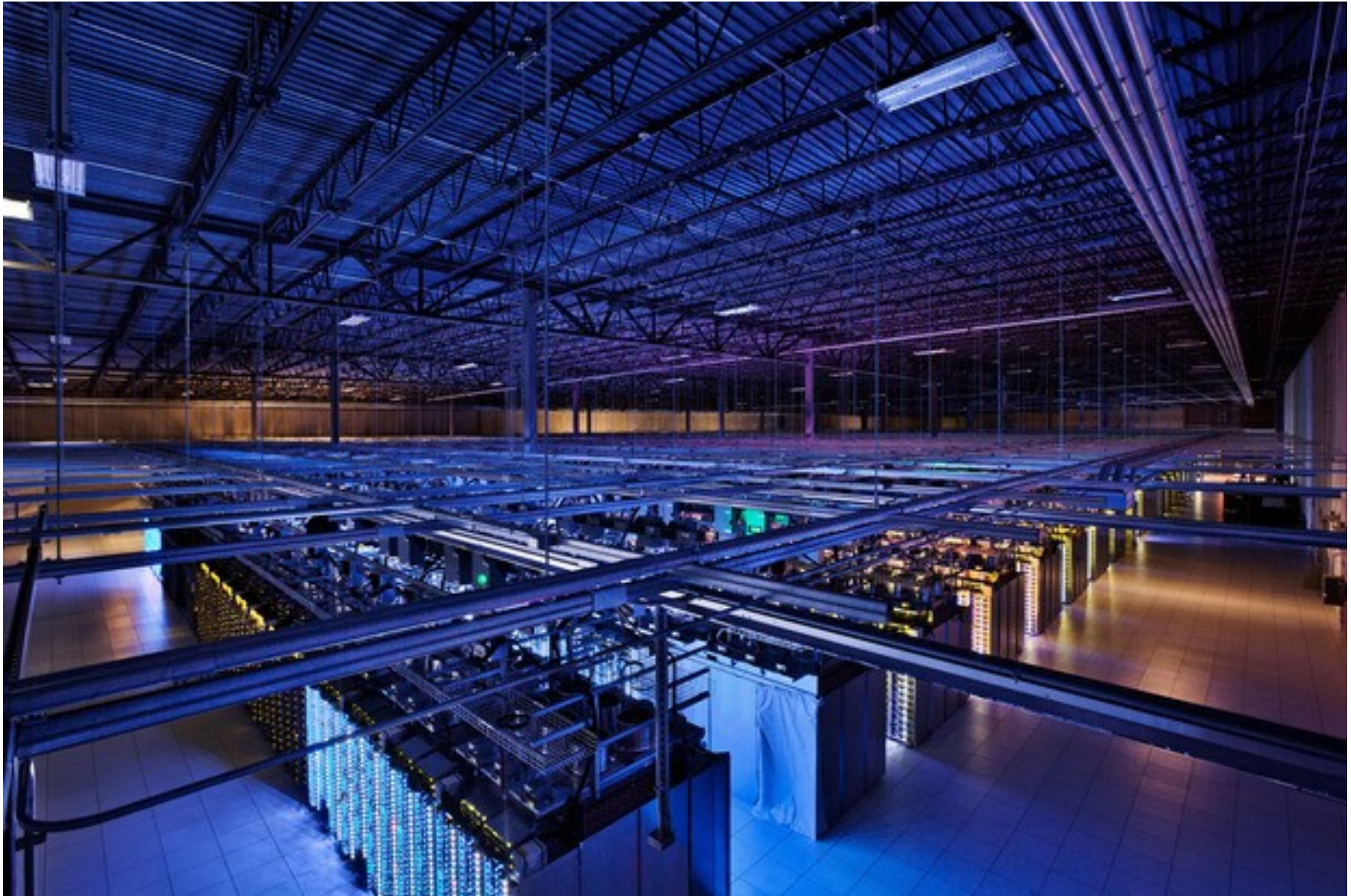
НРС



Remember?



Remember?

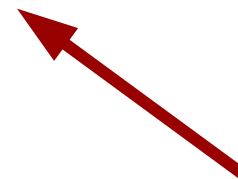


Основные характеристики вычислительного кластера

- Пиковая (теоретическая) производительность (Tflops)
- Linpack производительность (Tflops)
- VogoMIPS производительность?
- Пропускная способность Interconnect (Gbps)
- Задержки Interconnect (ns)

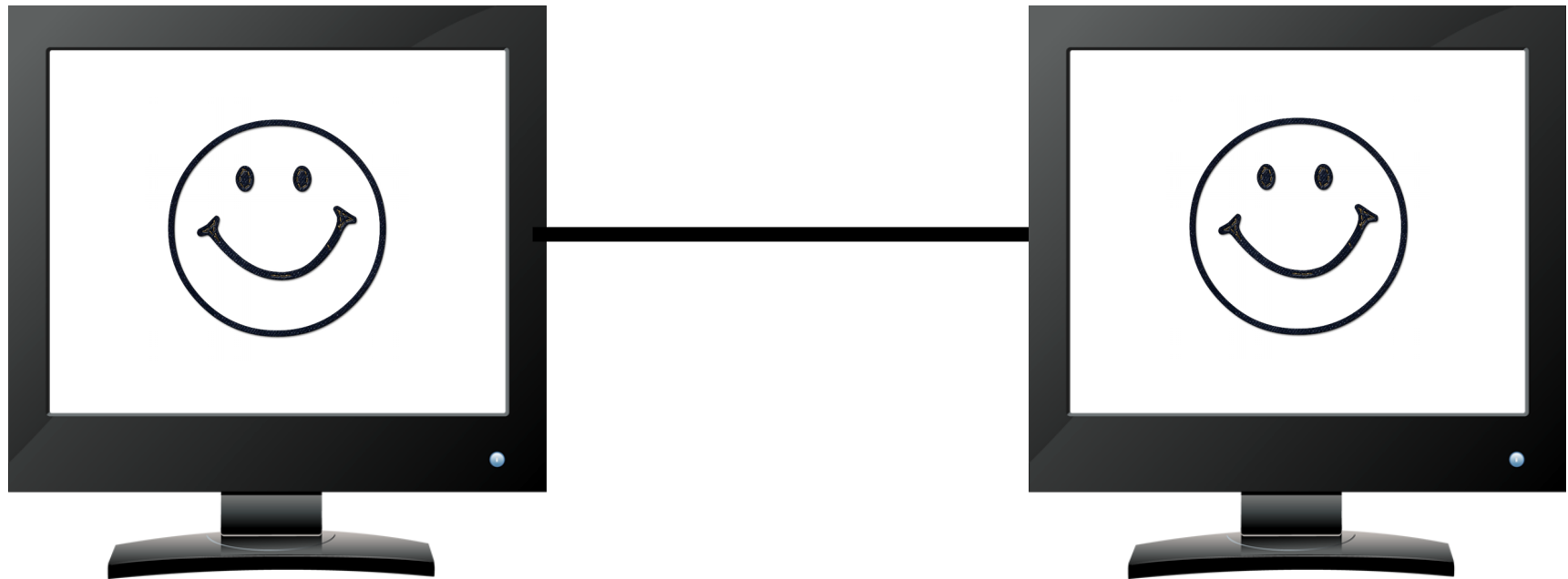
Основные характеристики вычислительного кластера

- Пиковая (теоретическая) производительность (Tflops)
- Linpack производительность (Tflops)
- VogoMIPS производительность?
- Пропускная способность Interconnect (Gbps)
- Задержки Interconnect (ns)



сеть

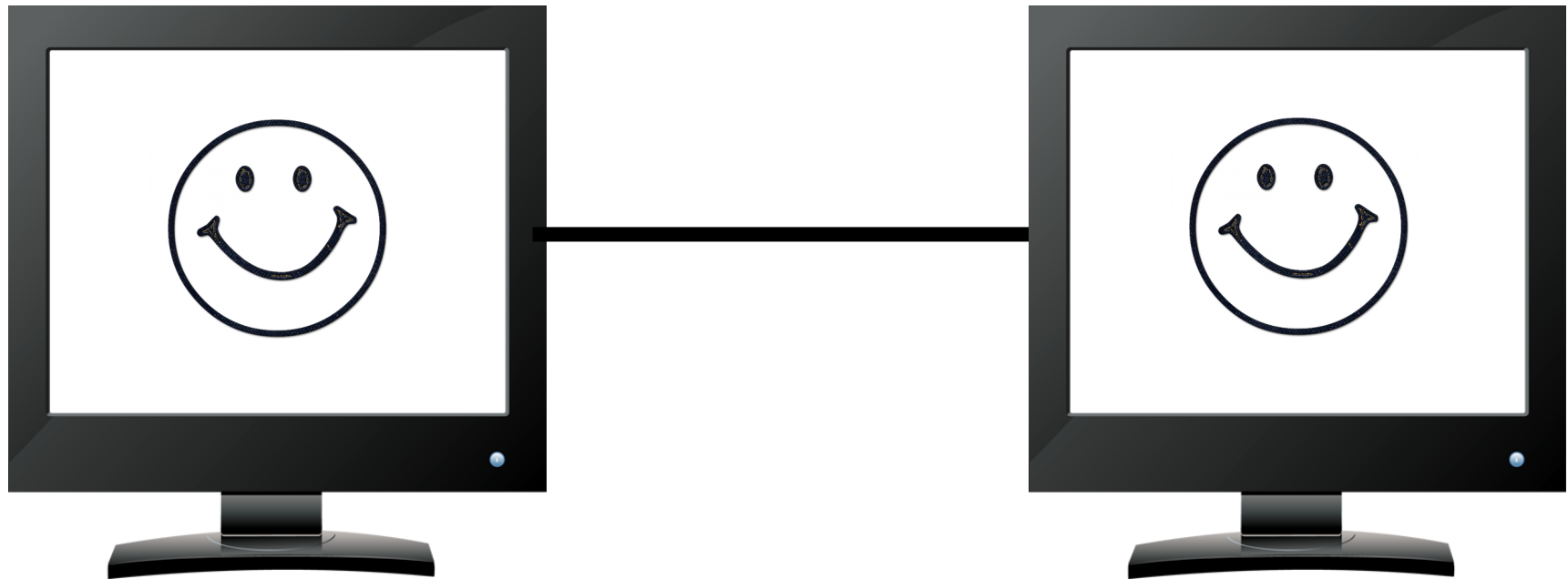
Что такое компьютерная сеть?



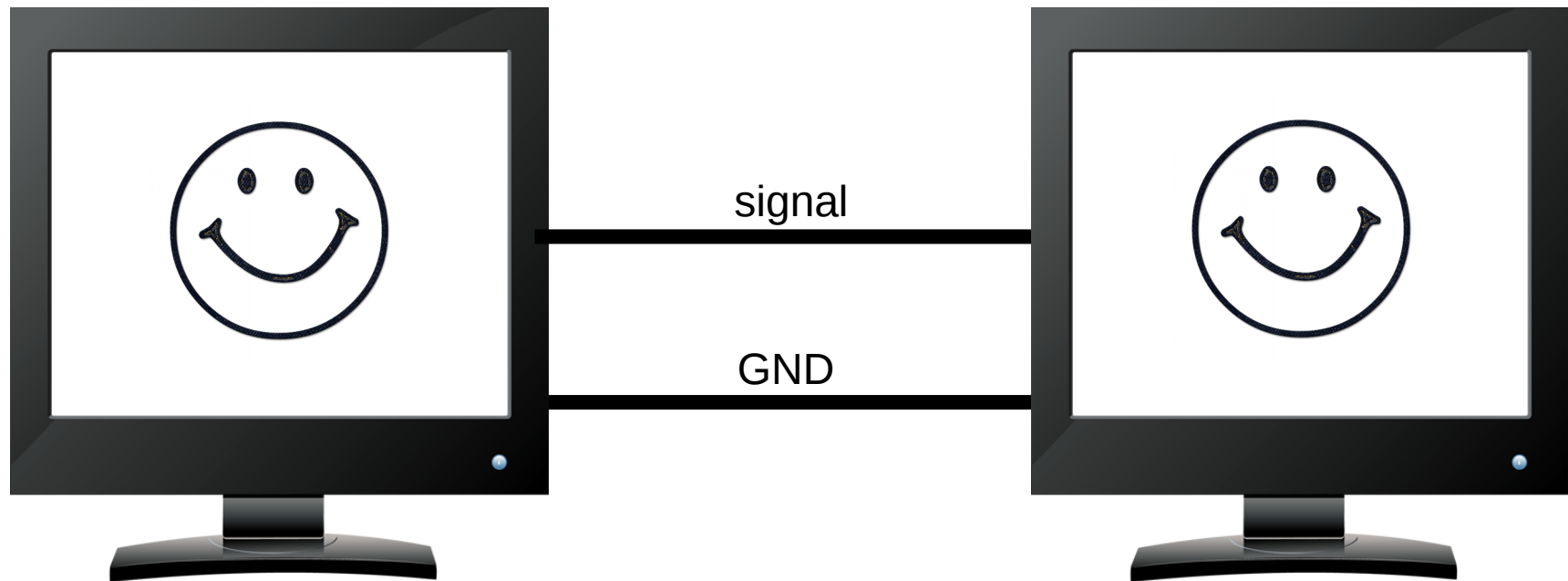
Что такое компьютерная сеть?

- Универсальность (любой компьютер)
- Масштабируемость
- Скорость
- Низкая стоимость

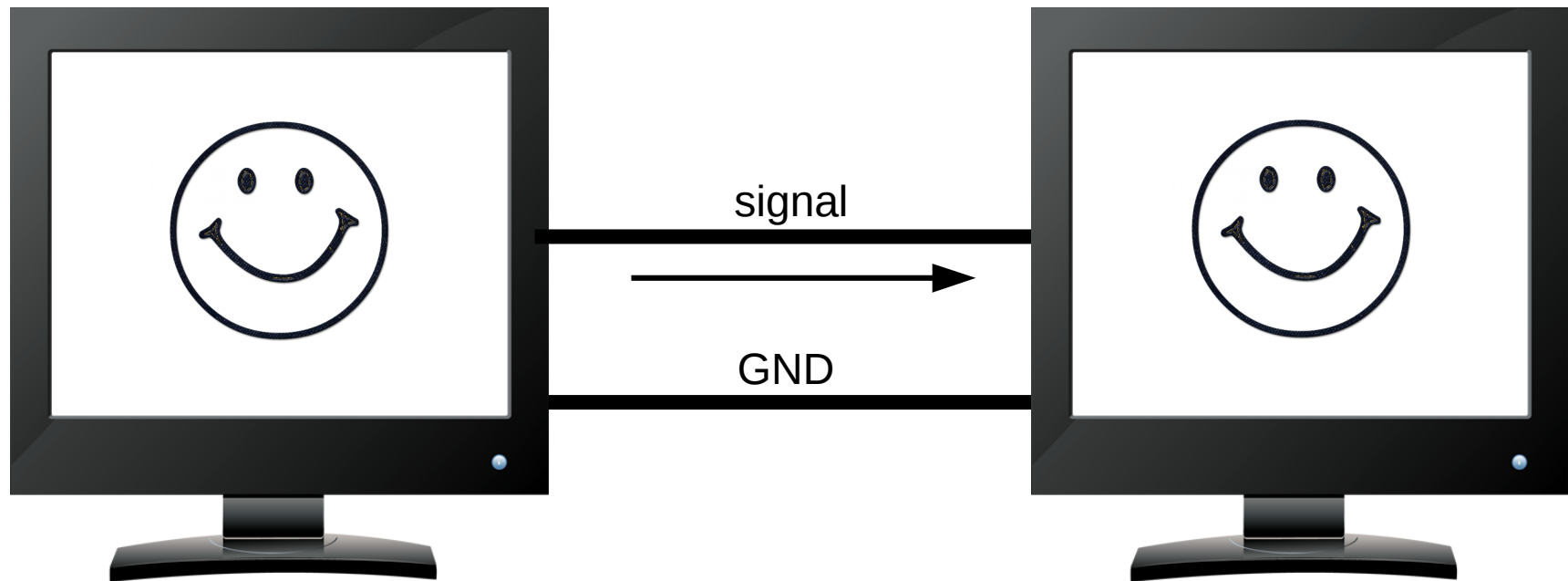
Ethernet



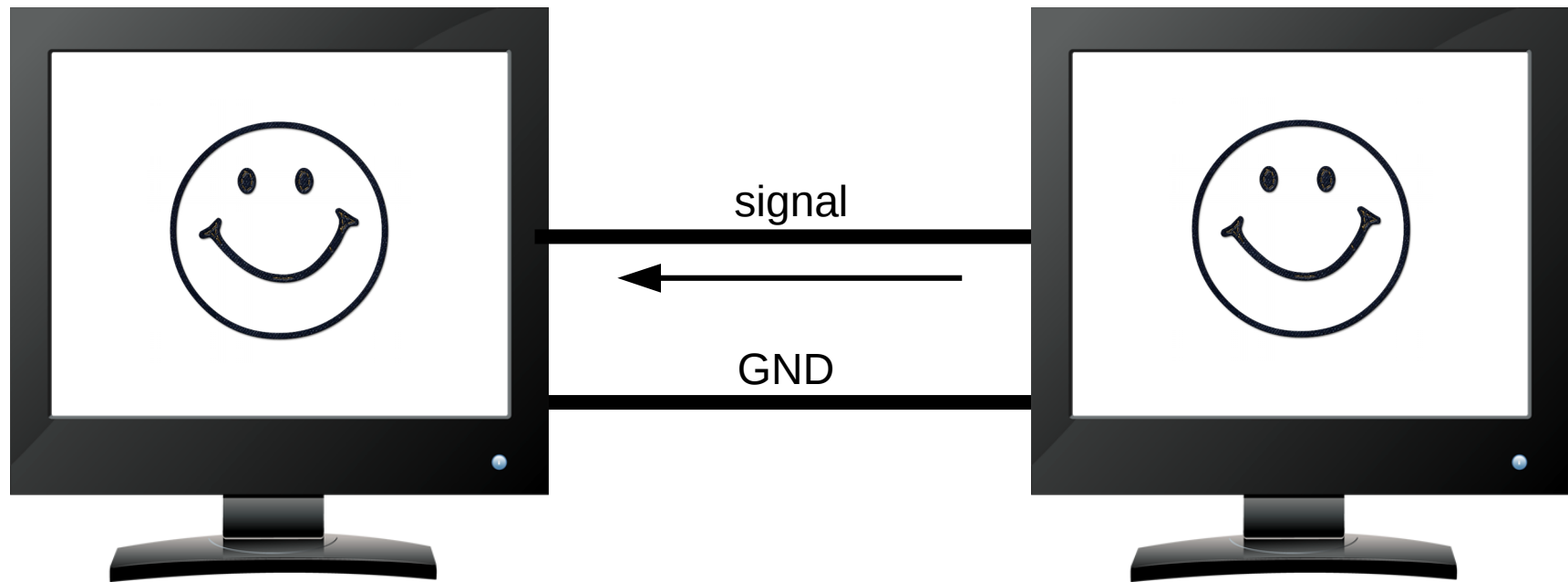
Ethernet



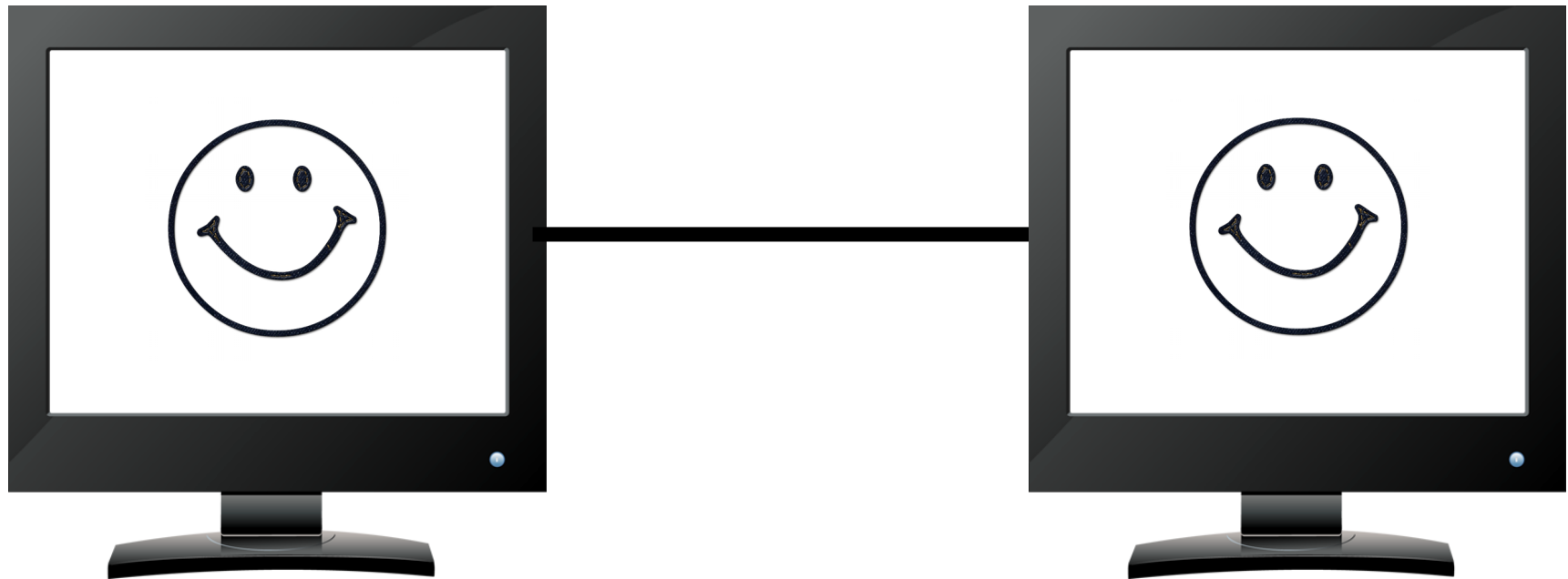
Ethernet



Ethernet



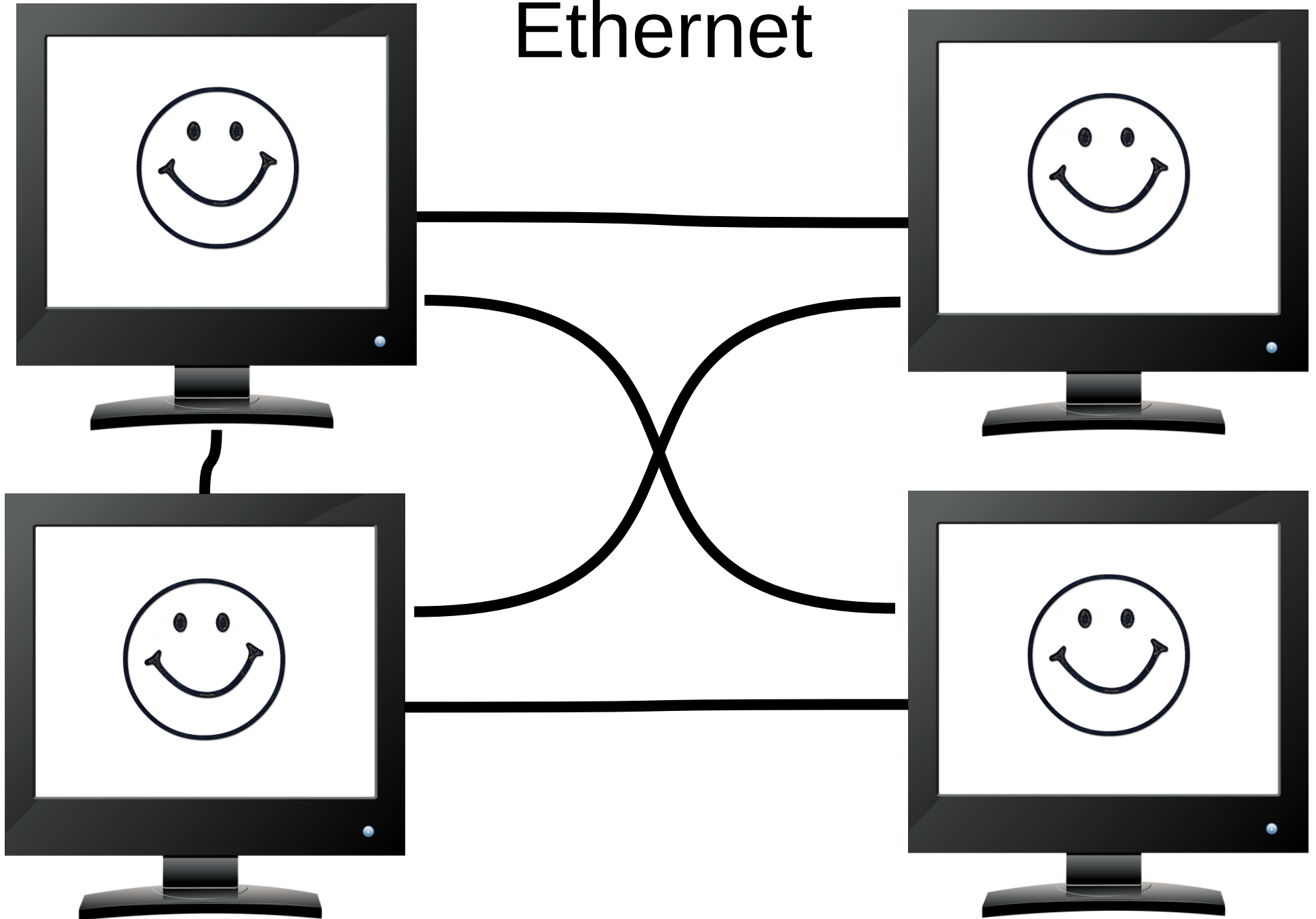
Ethernet



Ethernet



Ethernet



Ethernet

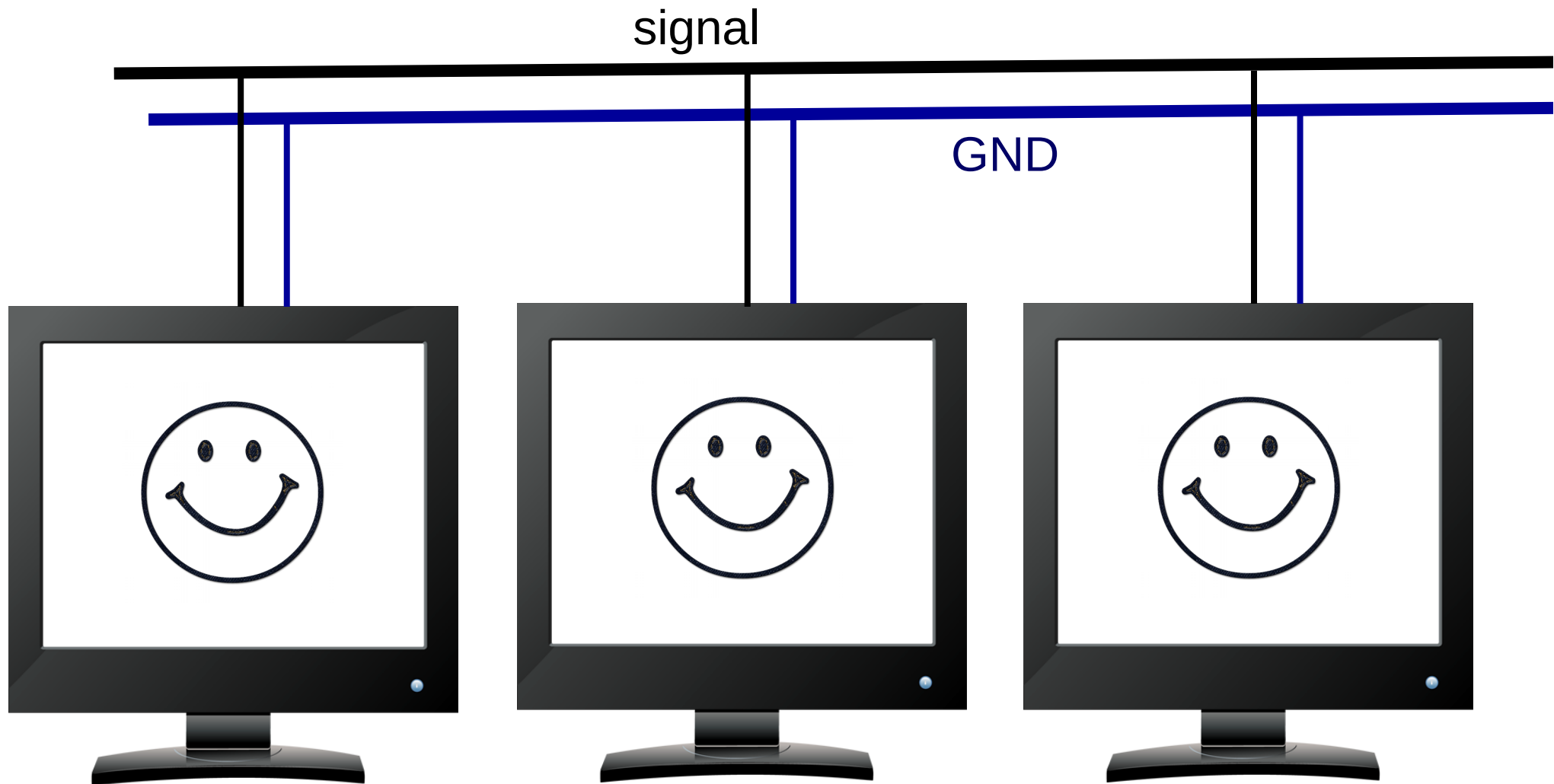


Ethernet

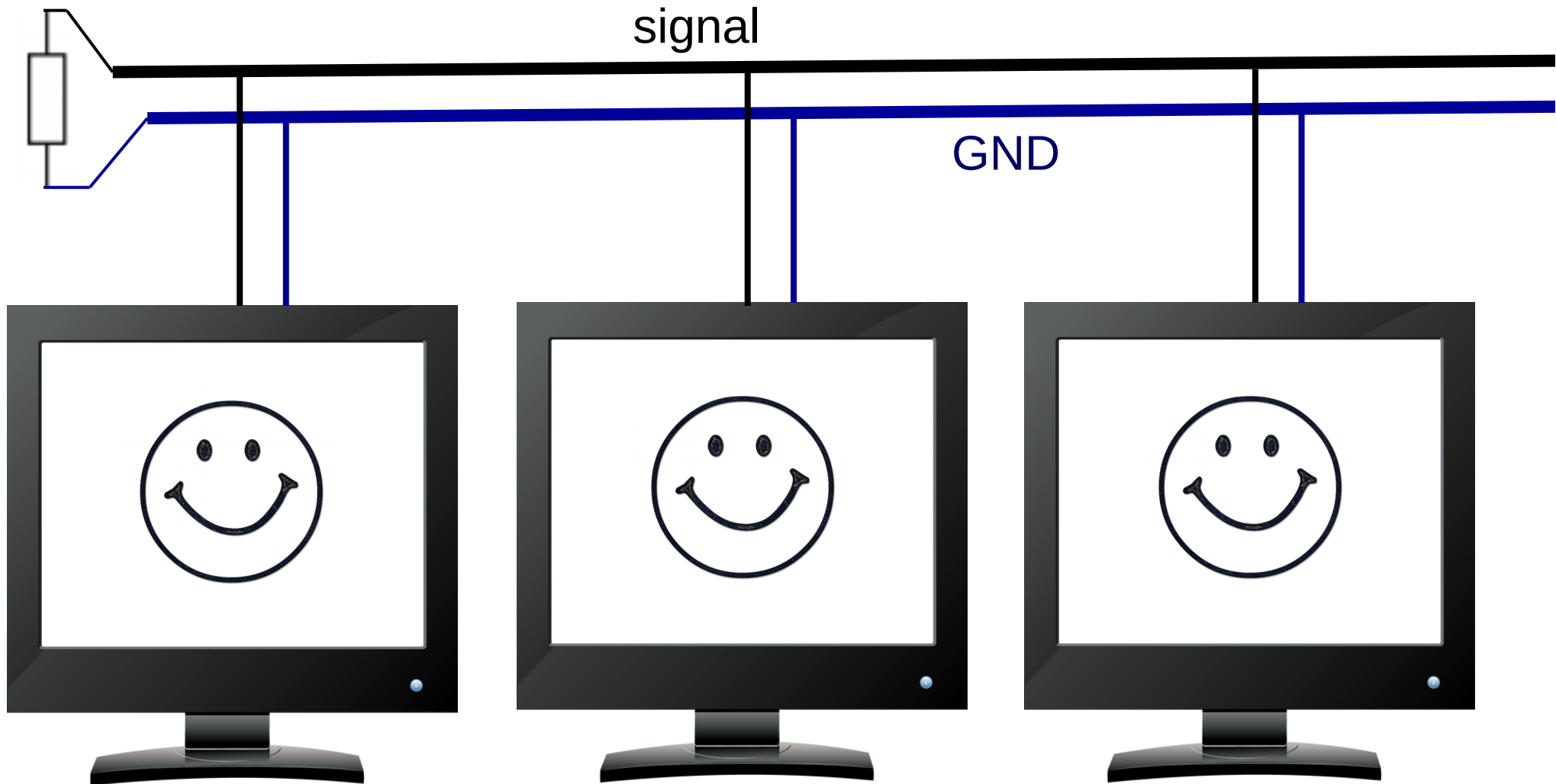
единая шина

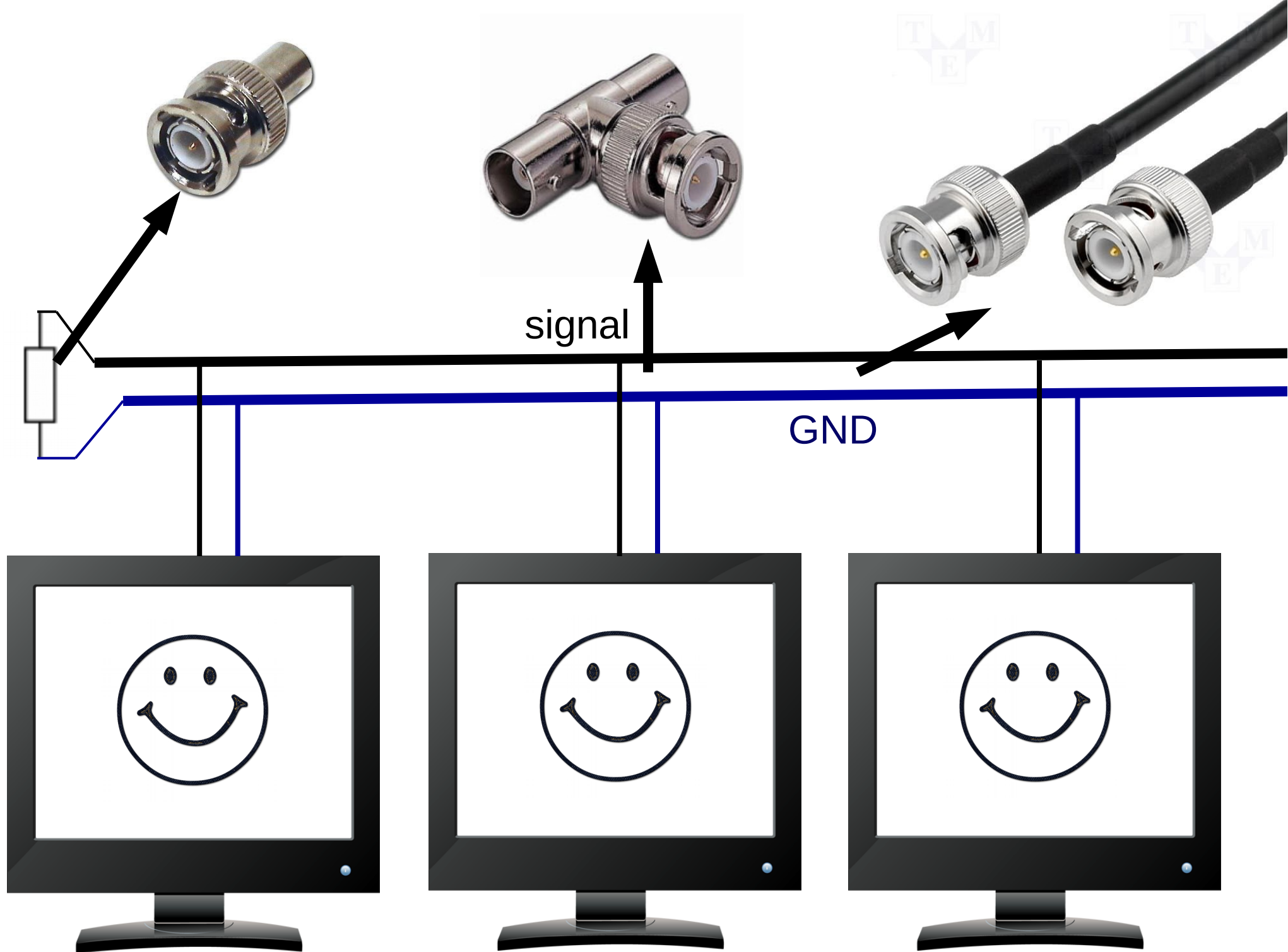


Ethernet

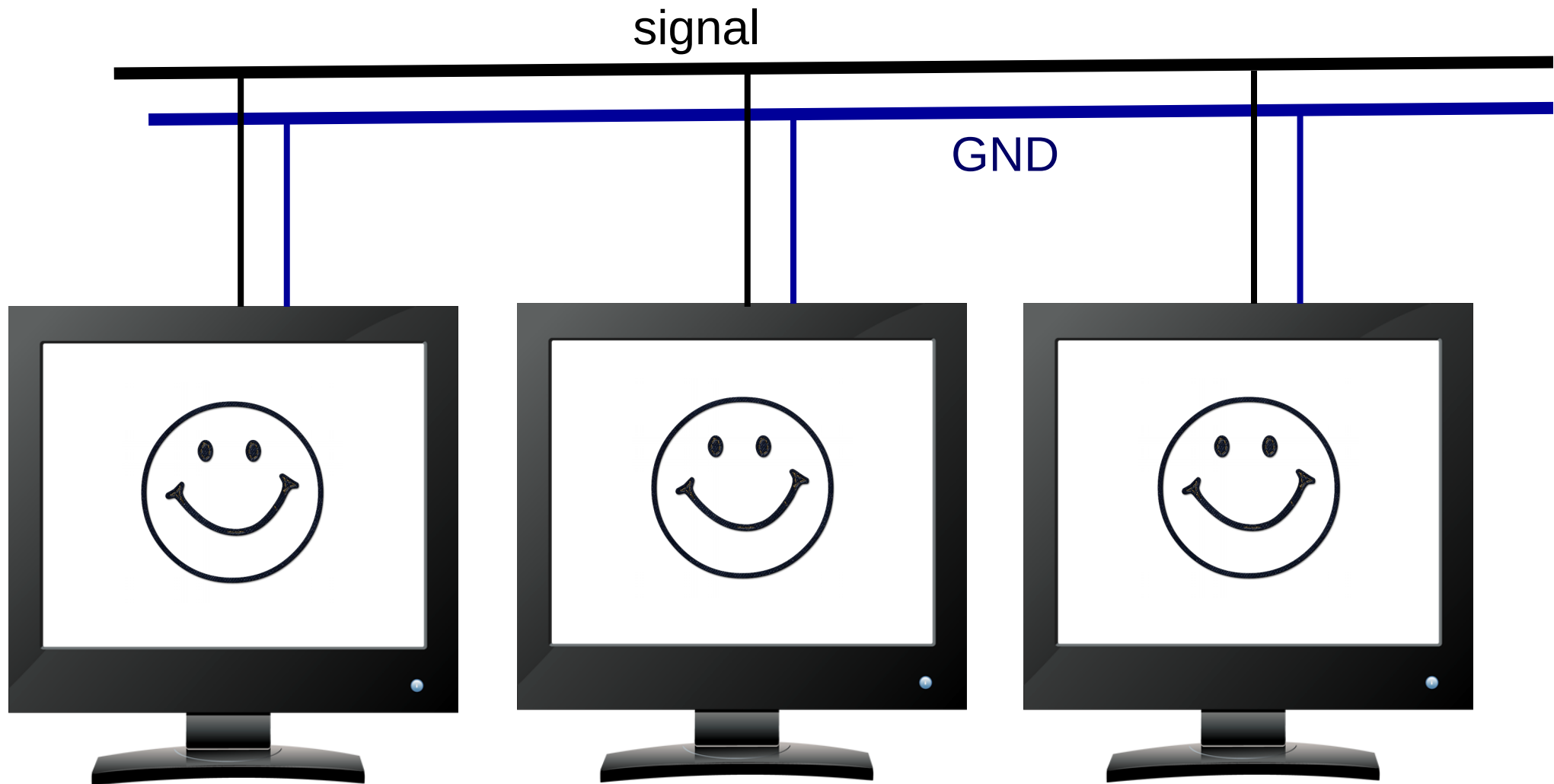


Ethernet

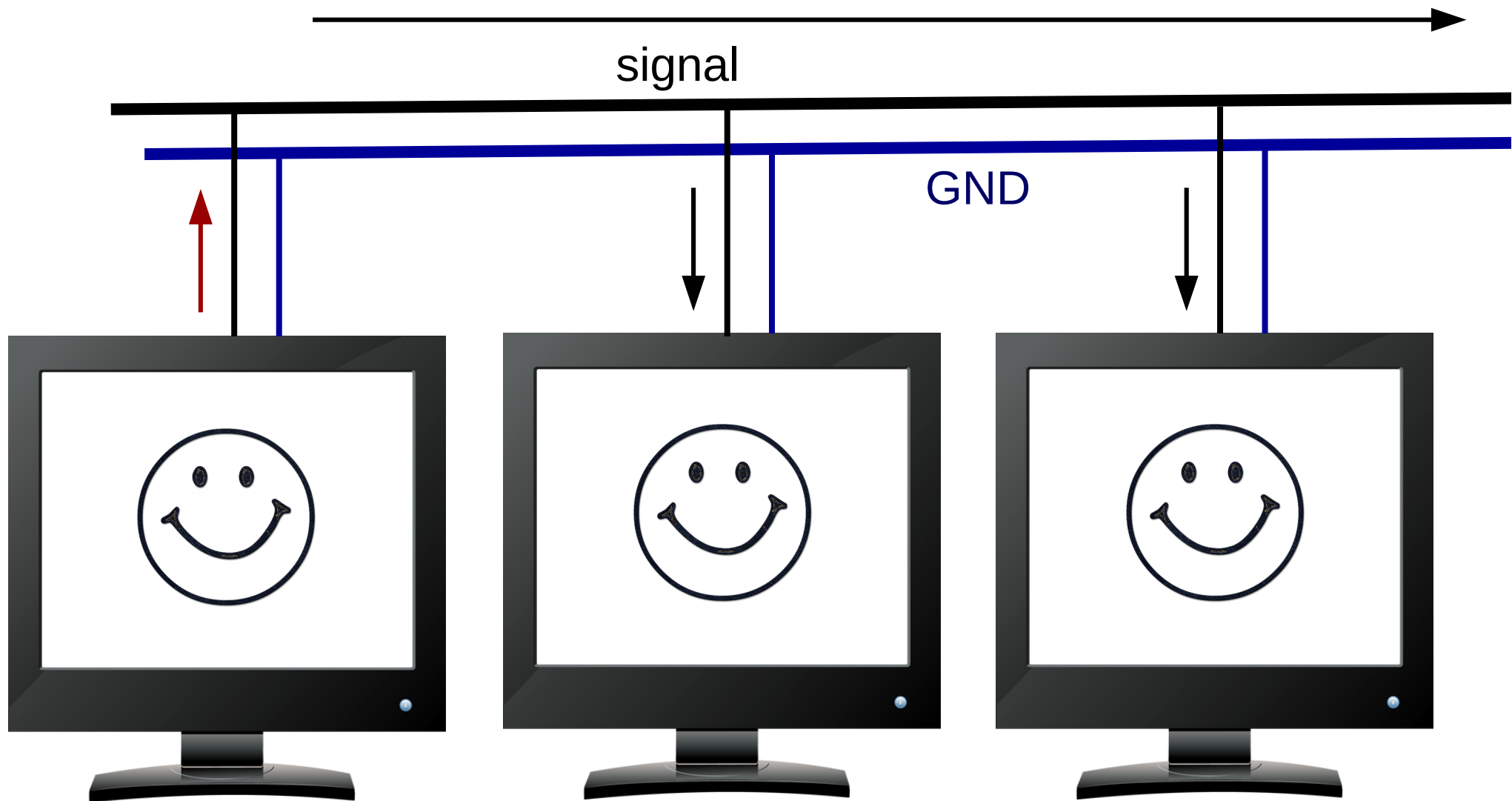




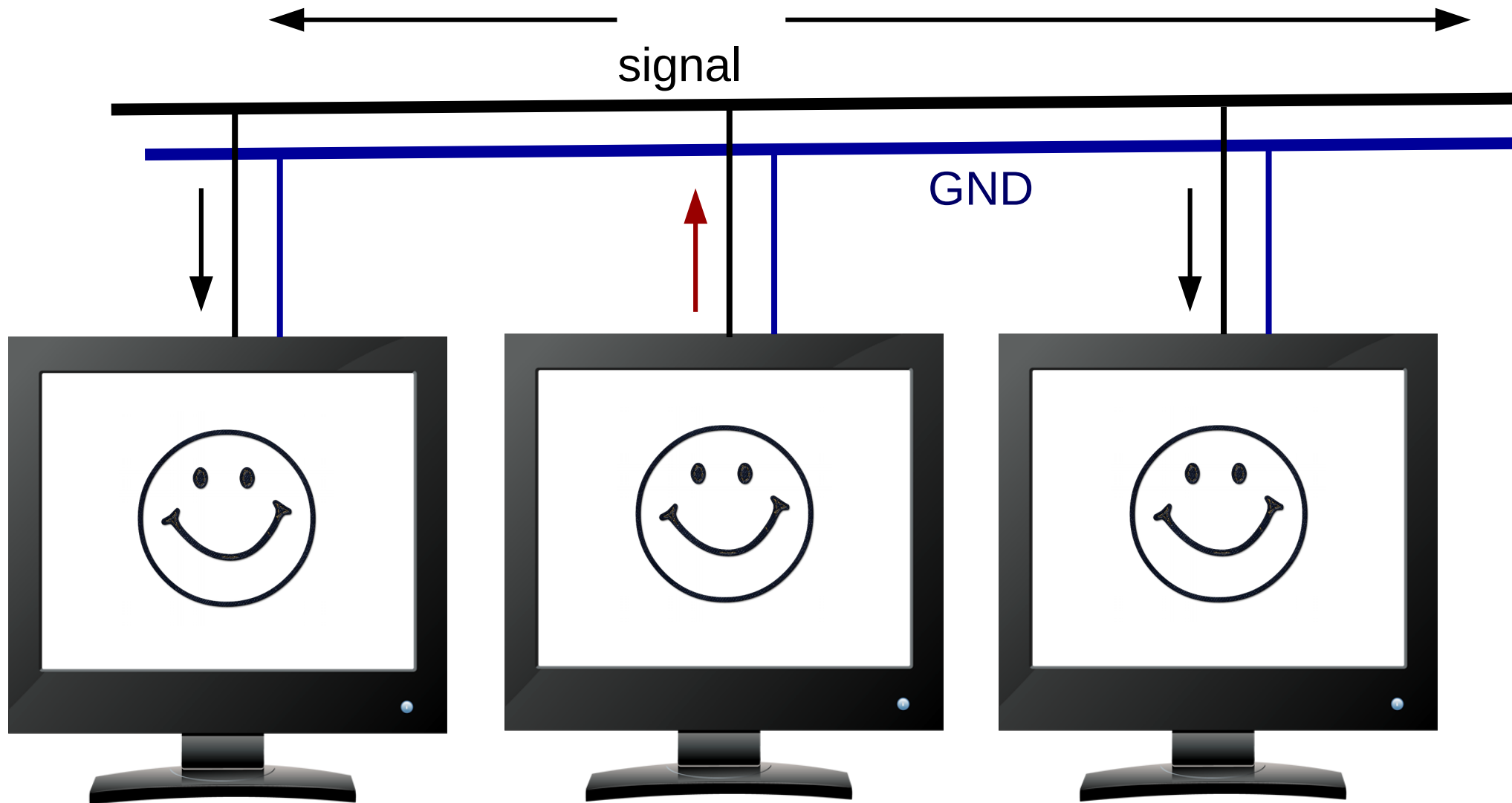
Ethernet



Ethernet

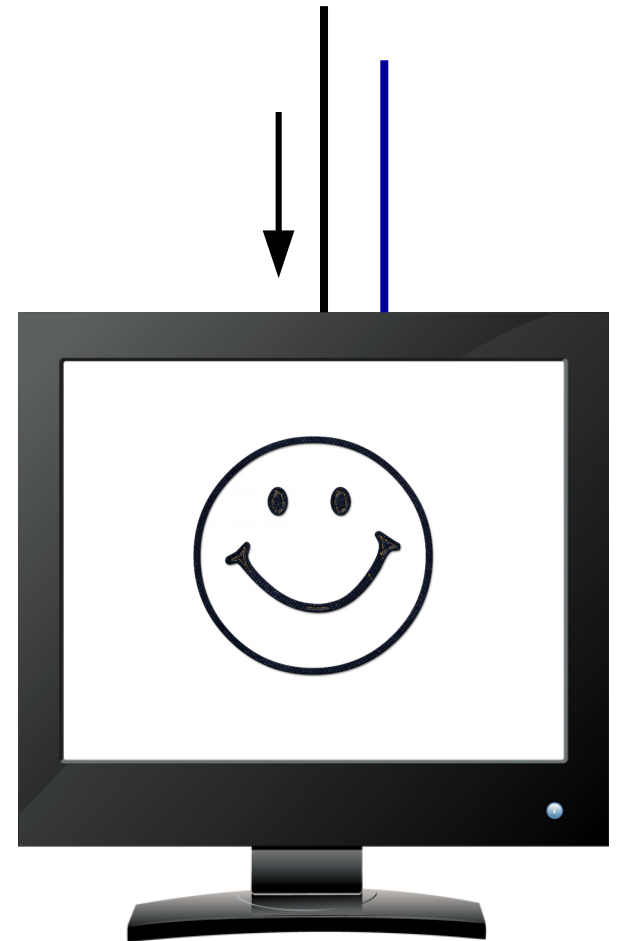


Ethernet



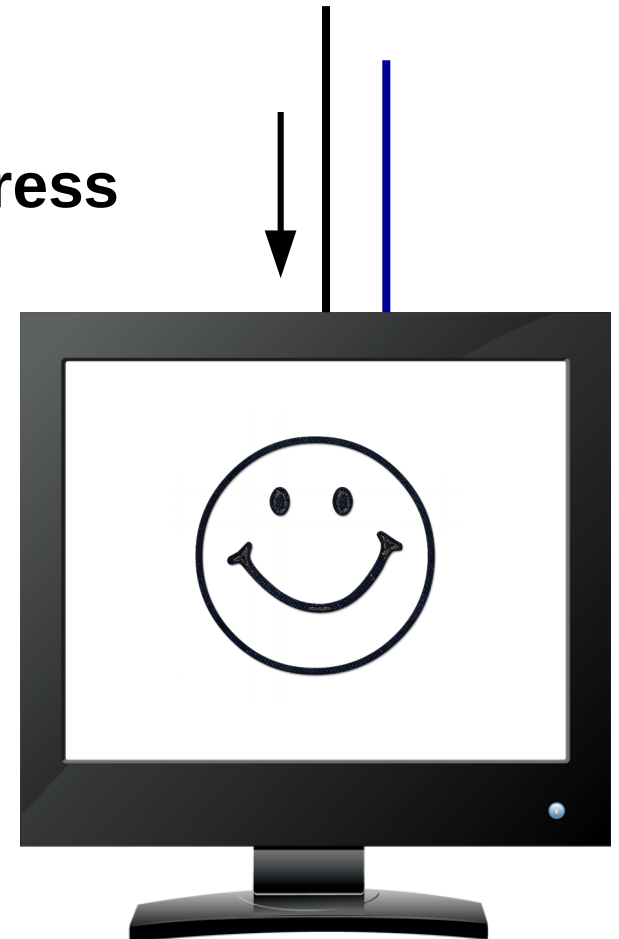
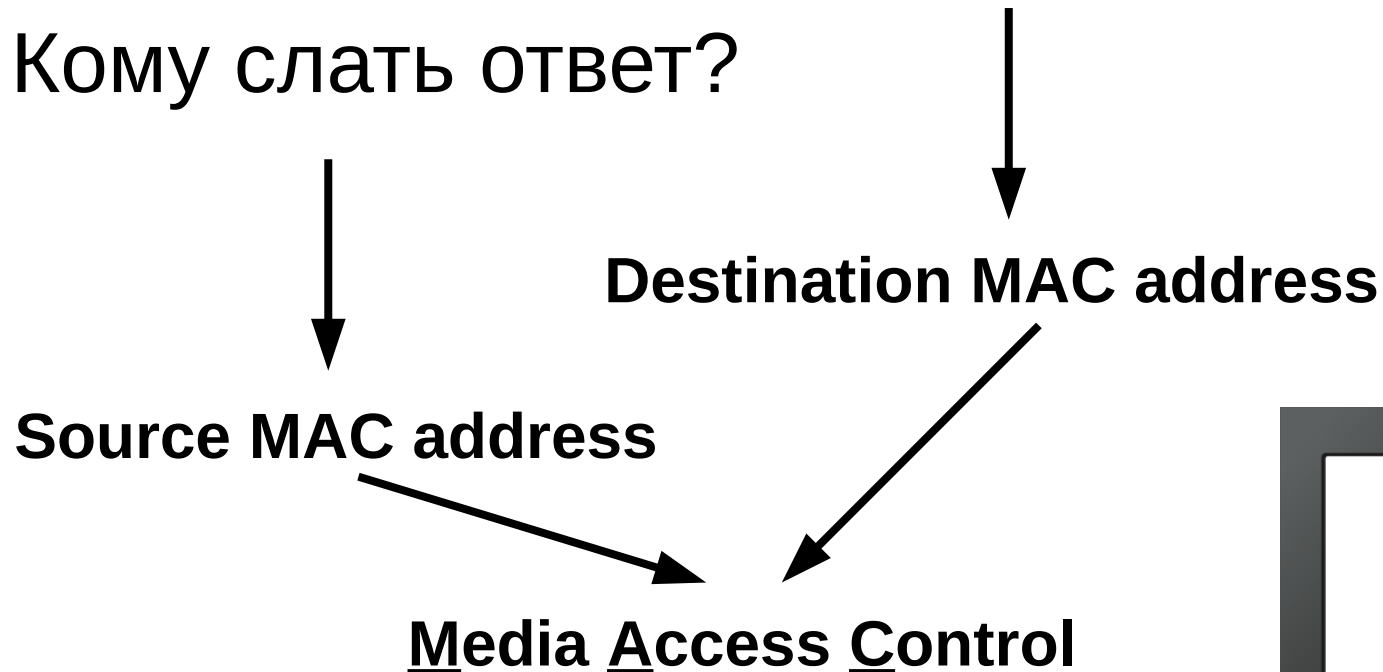
Ethernet

- Мне ли пришёл данный кадр?
- Кому слать ответ?



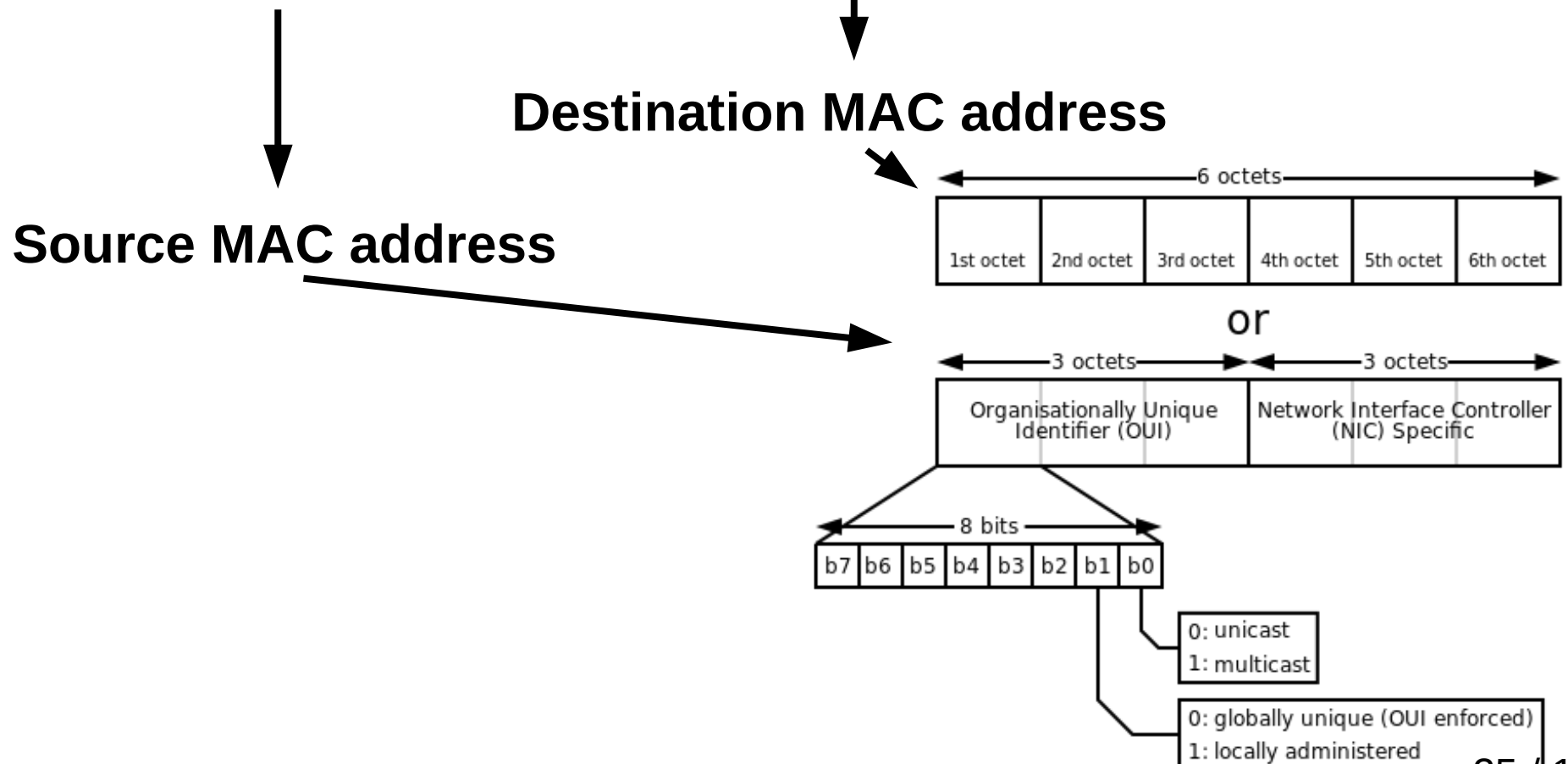
Ethernet

- Мне ли пришёл данный кадр?
- Кому слать ответ?

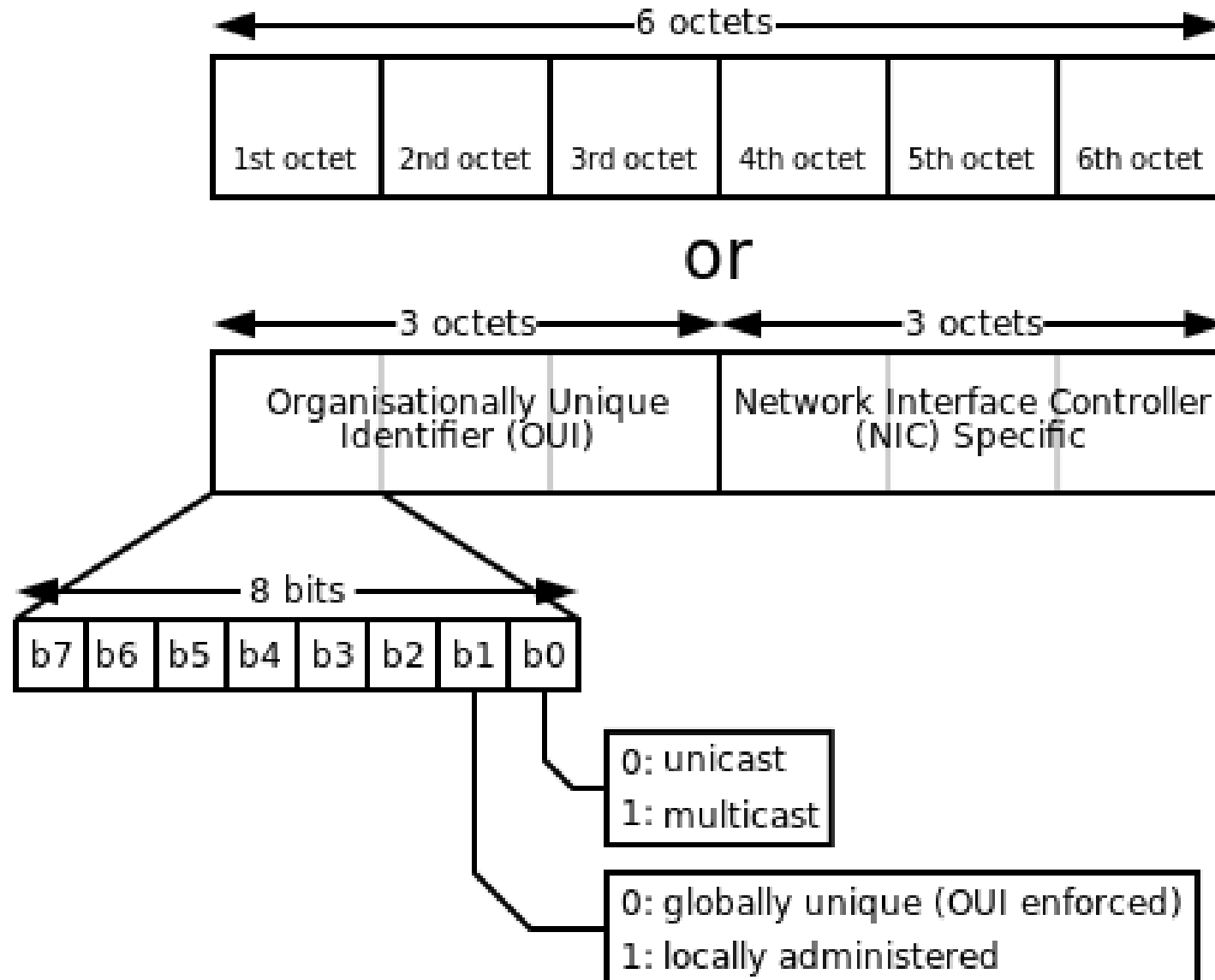


Ethernet

- Мне ли пришёл данный кадр?
- Кому слать ответ?

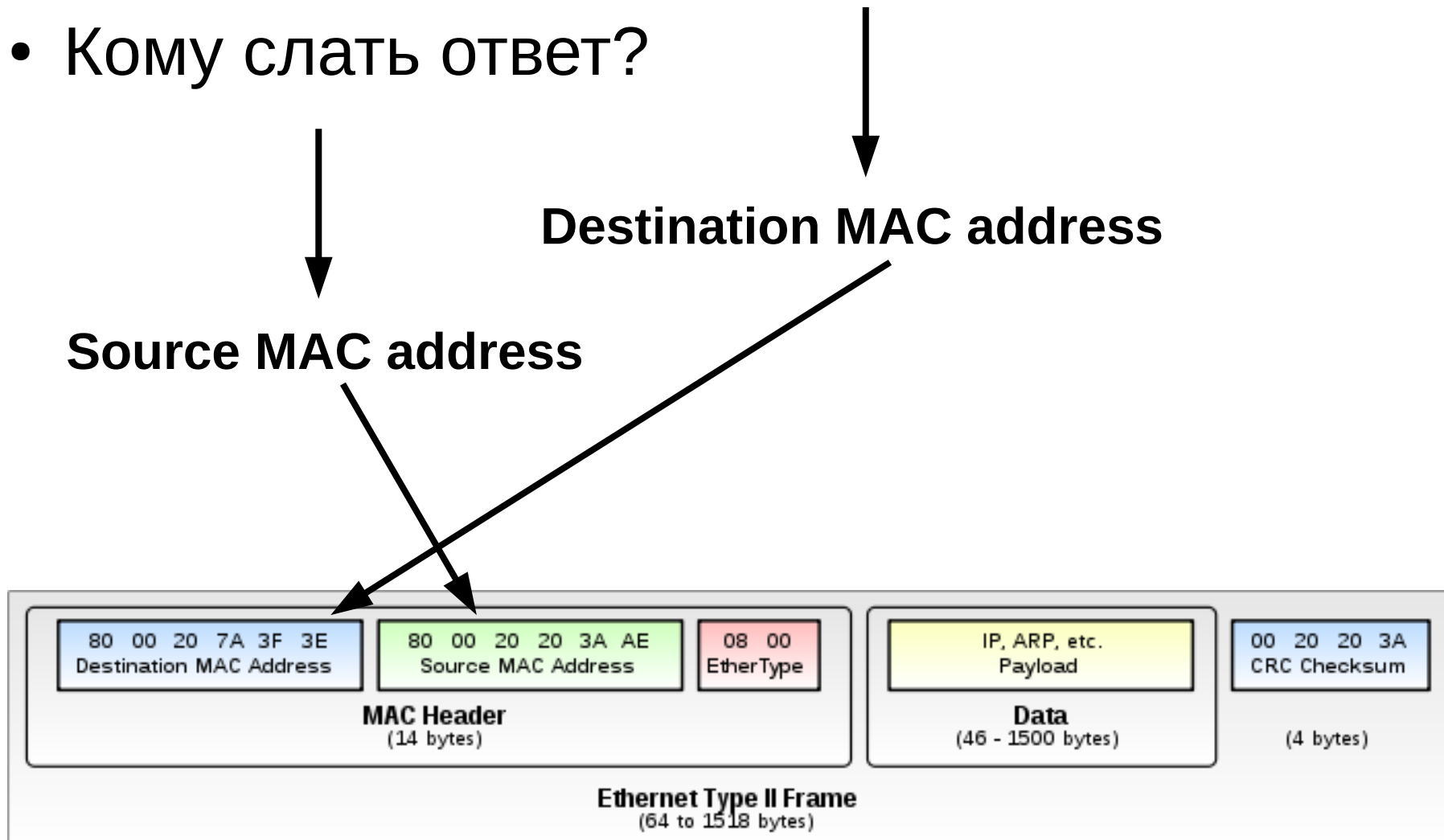


Ethernet

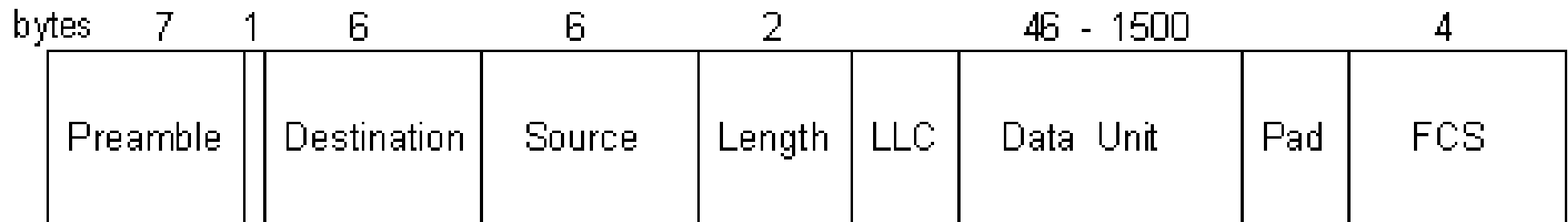


Ethernet

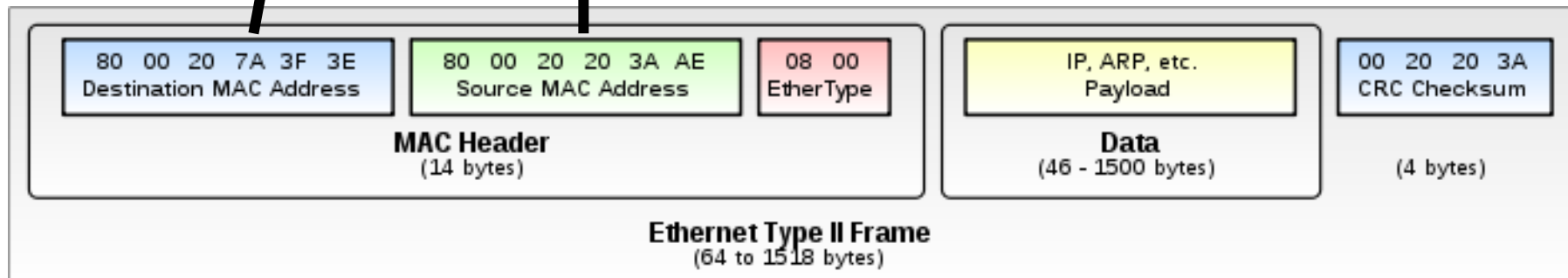
- Мне ли пришёл данный кадр?
- Кому слать ответ?



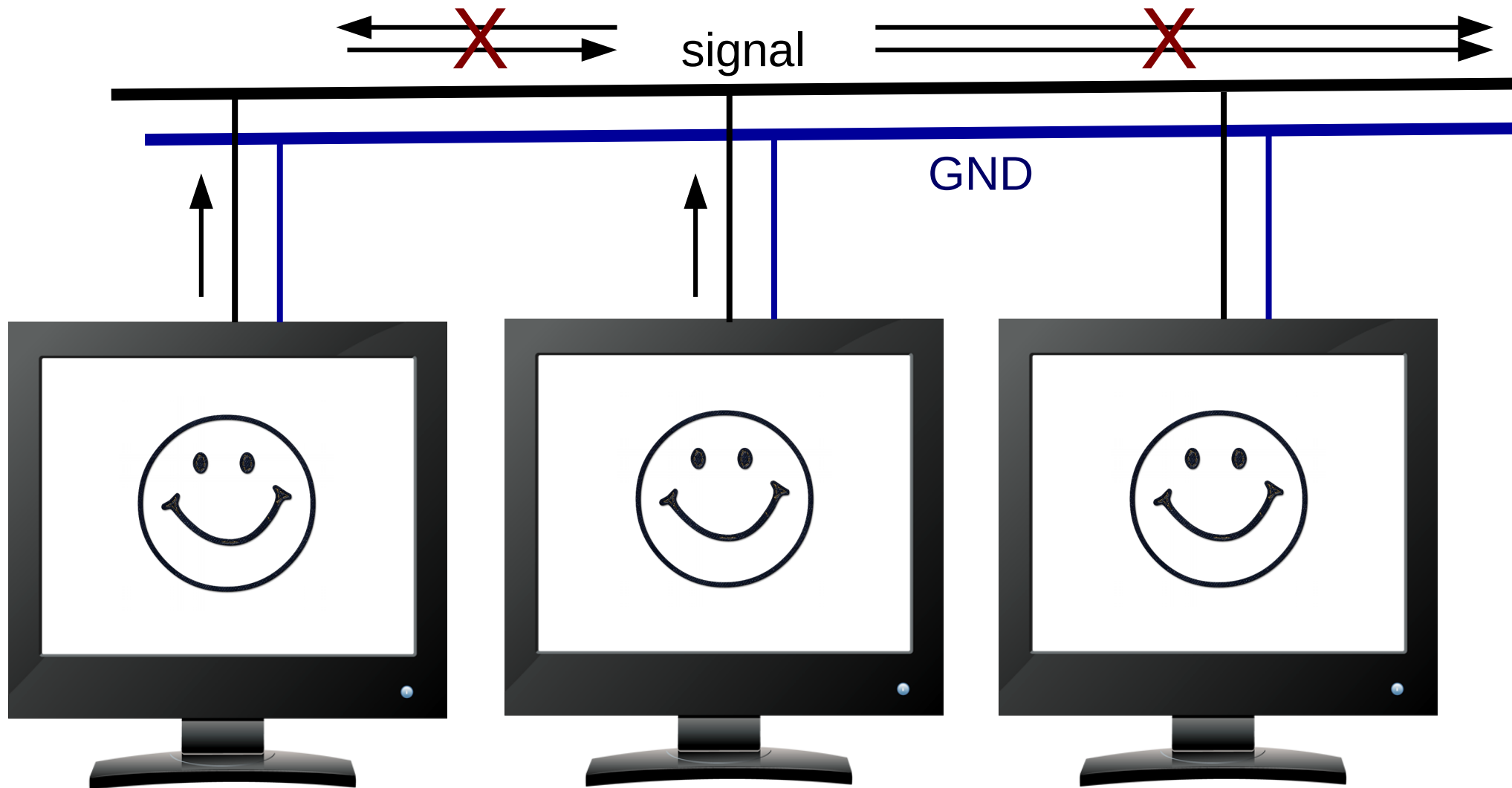
Ethernet



IEEE 802.3 Frame



Ethernet



Ethernet

Большая сеть



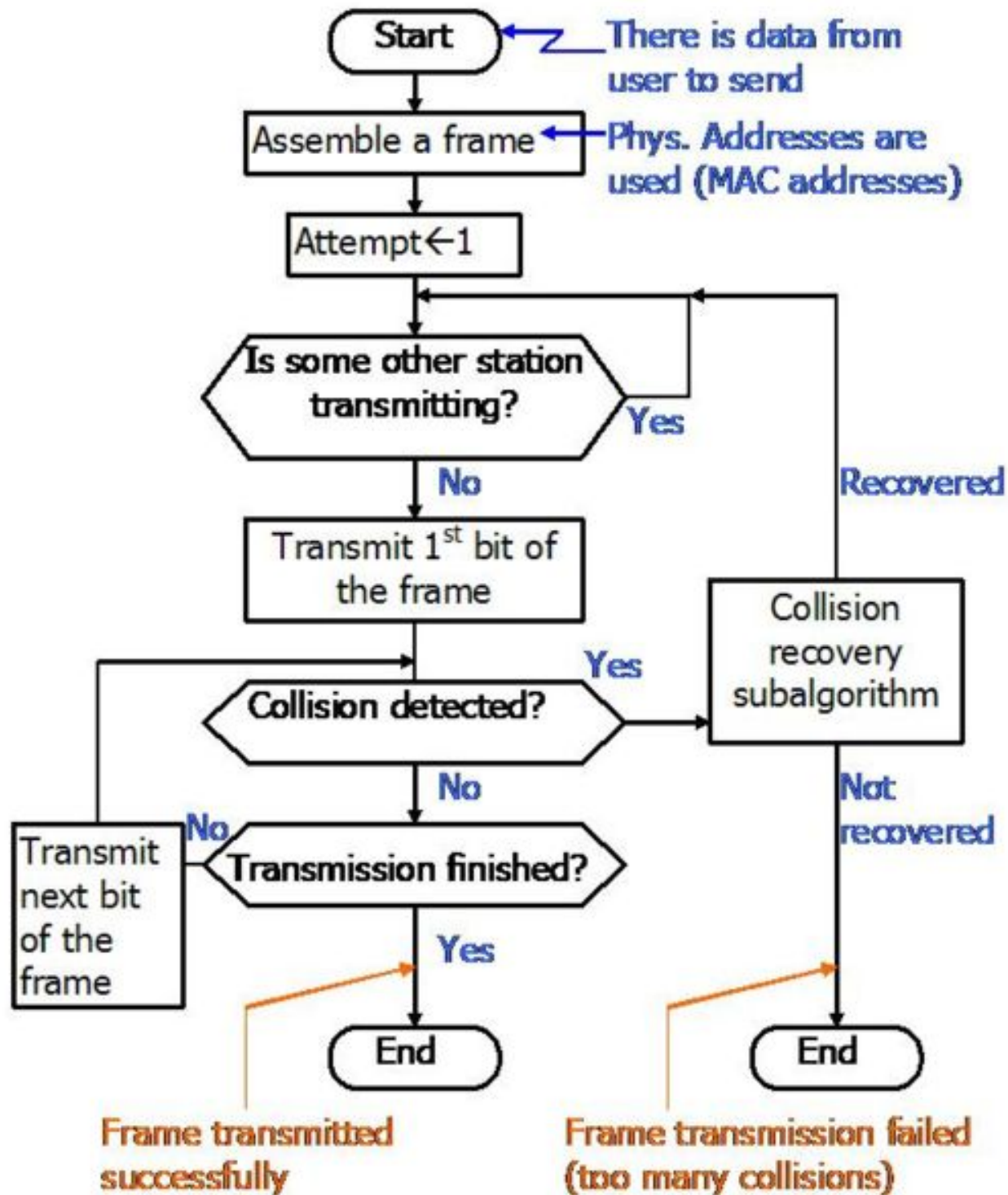
Ethernet



Ethernet

общая шина, КОЛЛИЗИИ

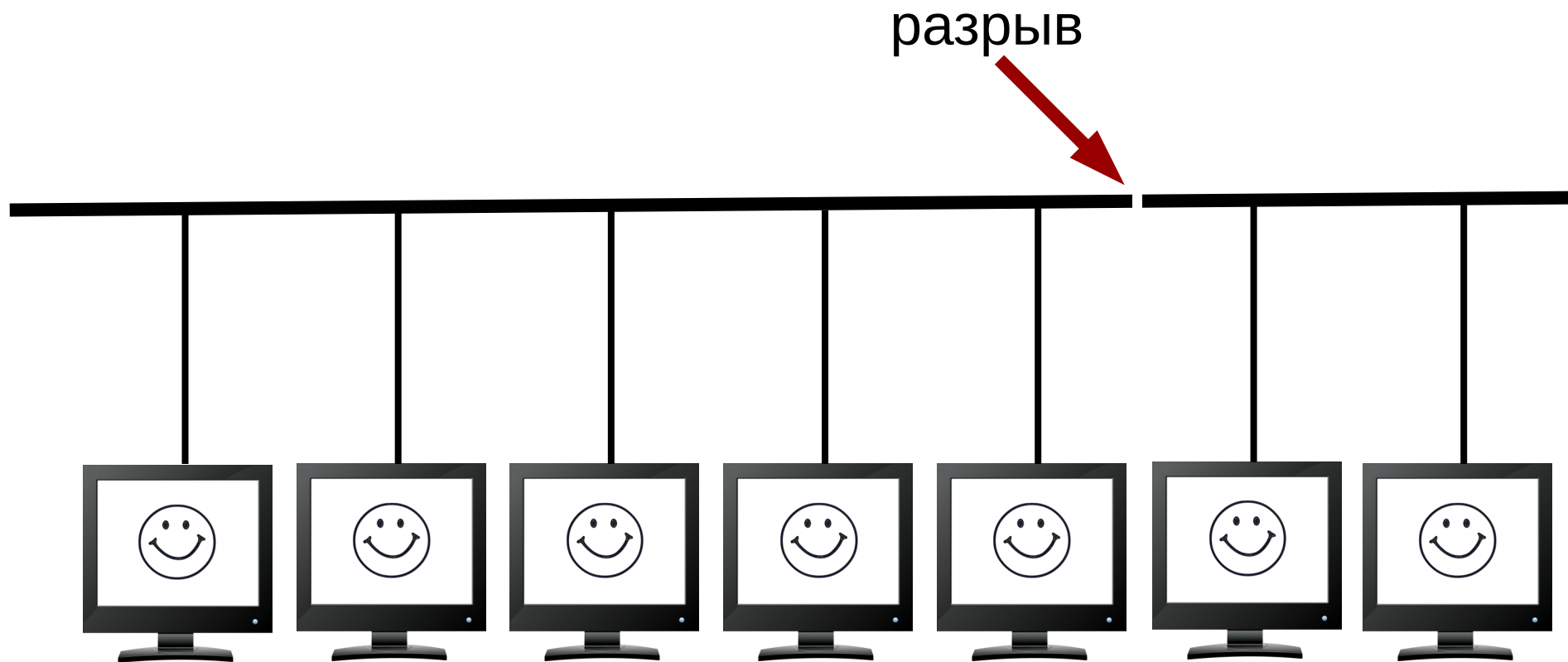




Ethernet

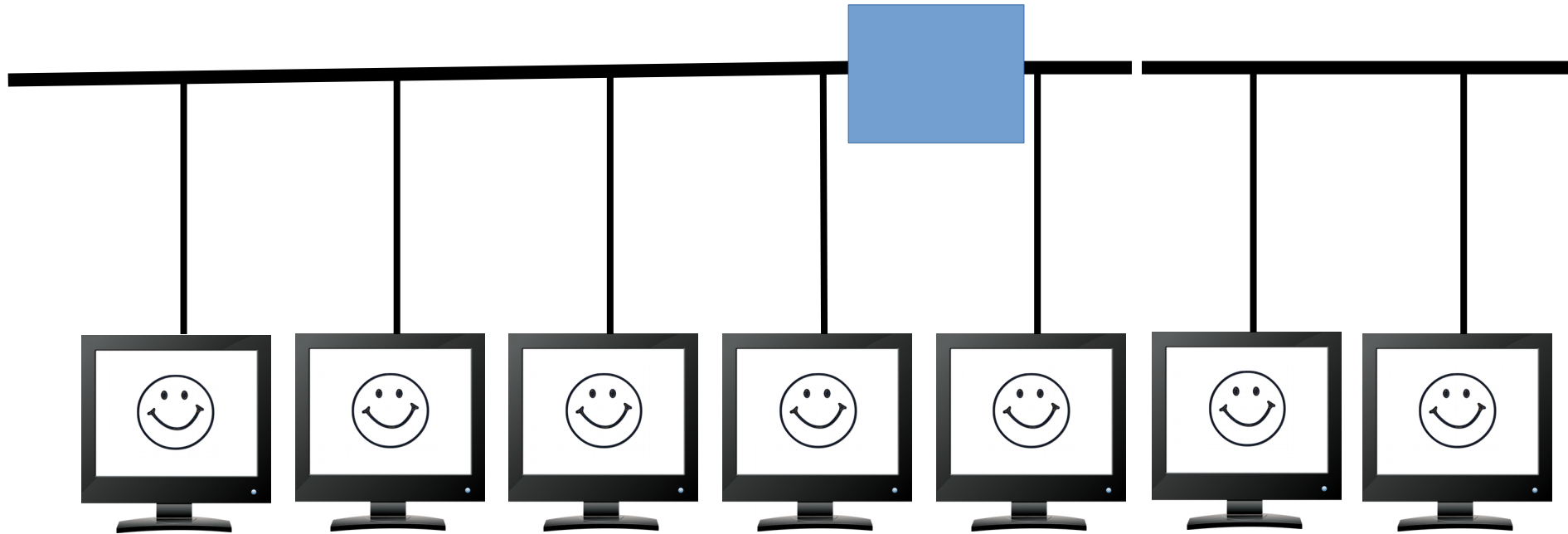


Ethernet



Ethernet

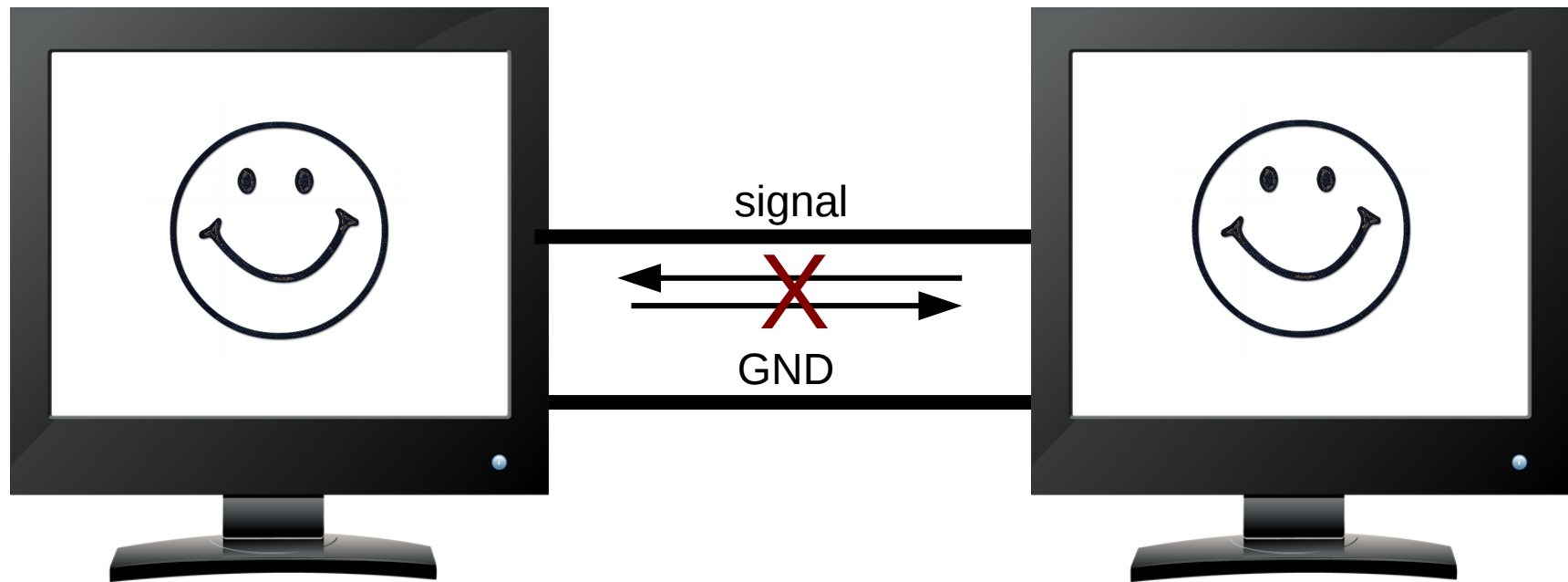
repeater



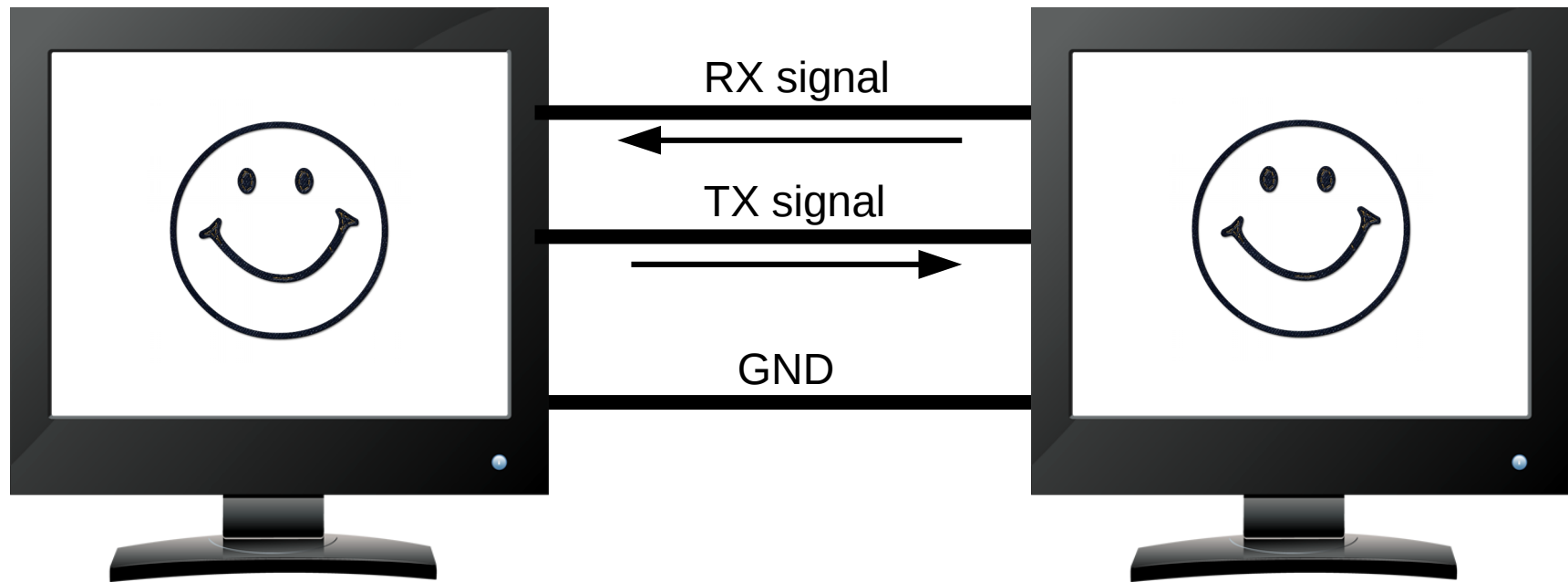
Ethernet

- Плохая масштабируемость
- Низкая надёжность
- Коллизии
- Низкая скорость

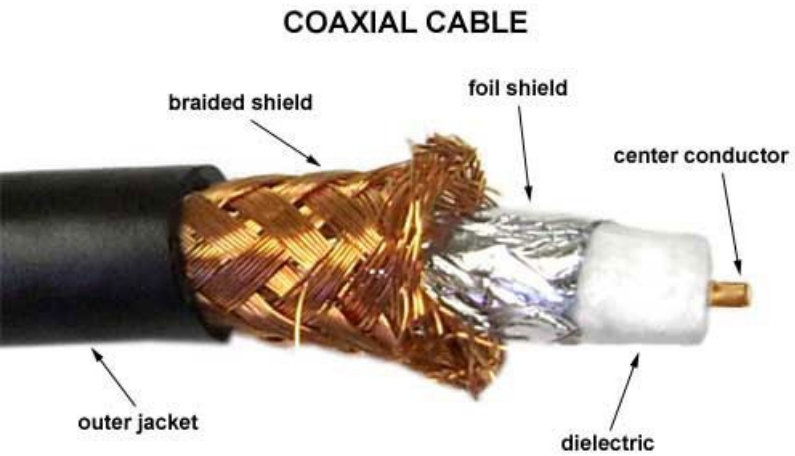
Ethernet



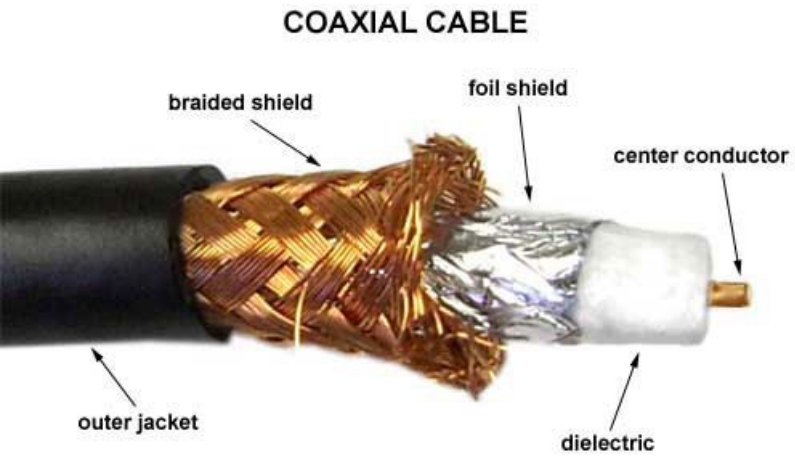
Ethernet



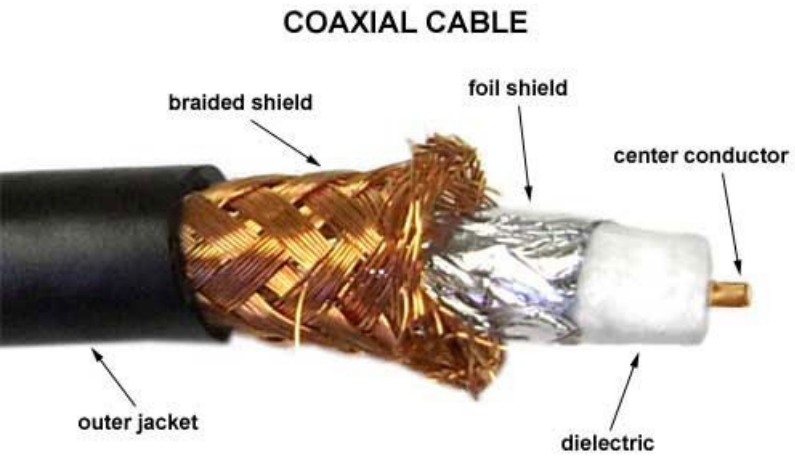
Ethernet



Ethernet



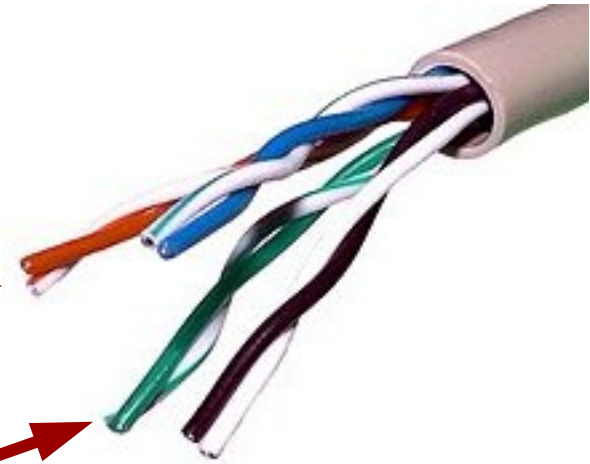
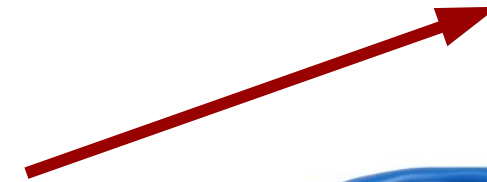
Ethernet



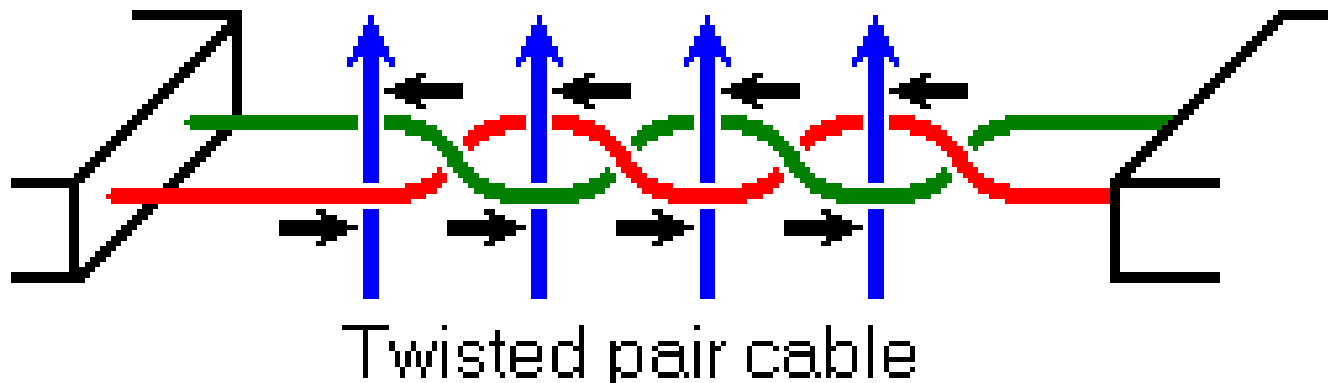
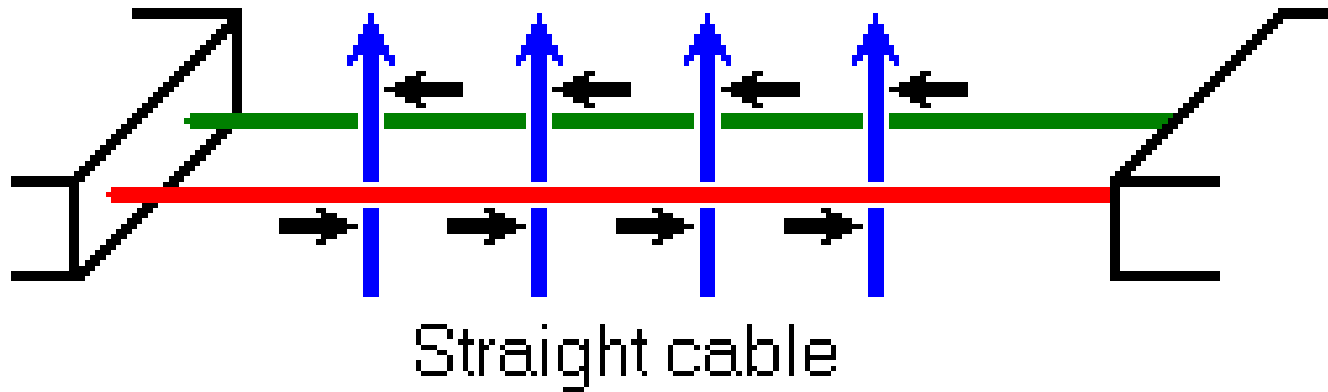
RX



TX

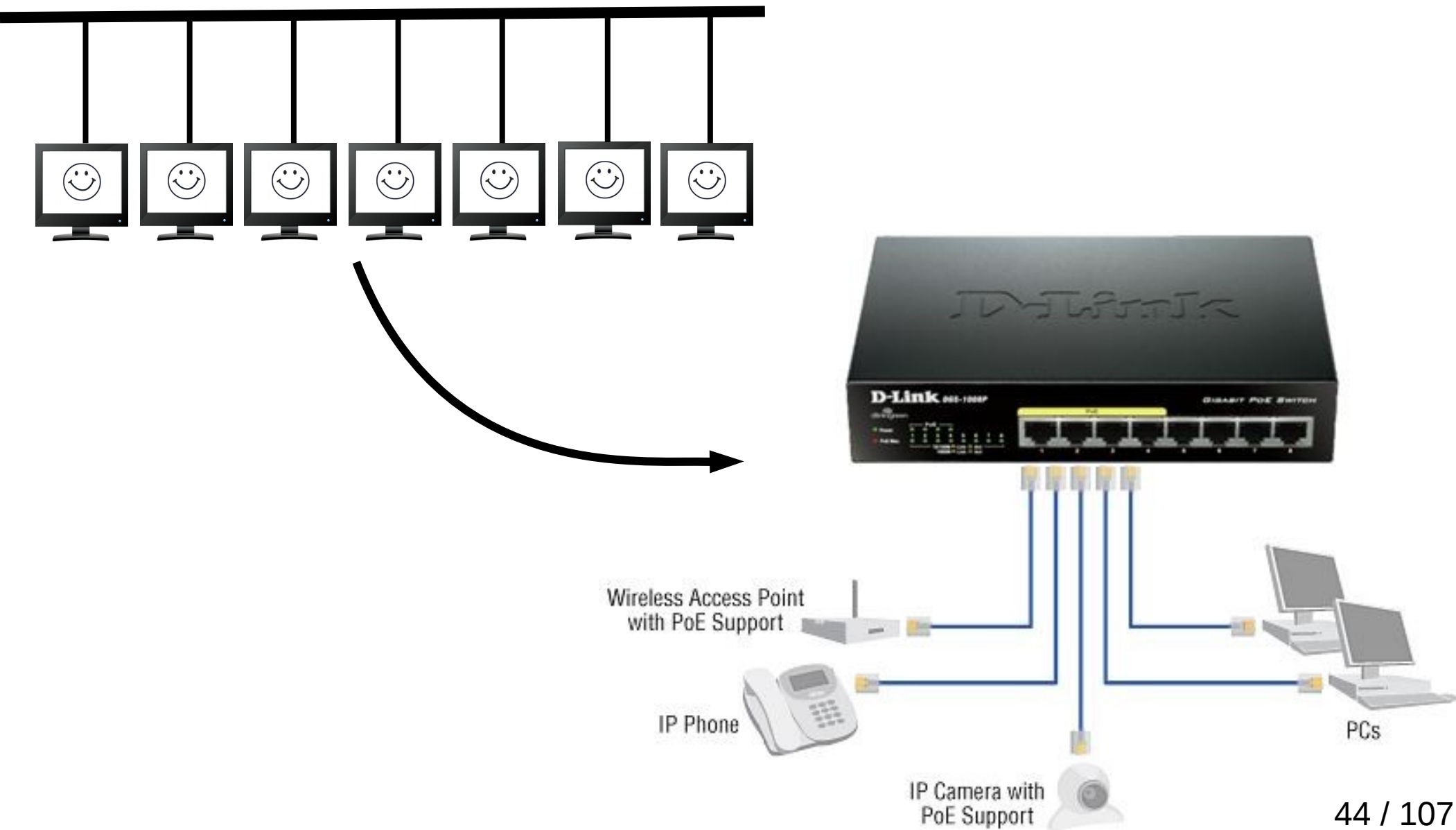


Ethernet

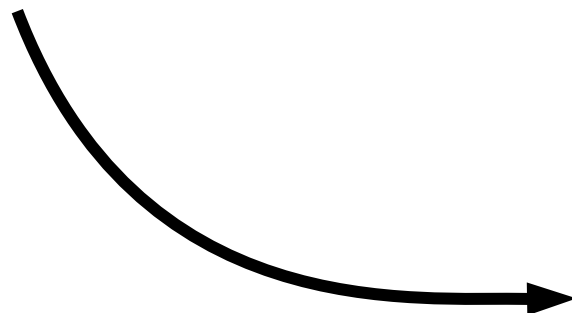


-  **Magnetic field**
-  **Induced noise current**

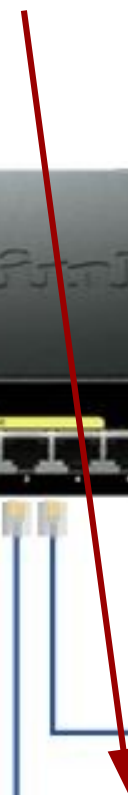
Ethernet



Ethernet



Разрыв, но сеть
работоспособна



Wireless Access Point
with PoE Support



IP Phone



IP Camera with
PoE Support



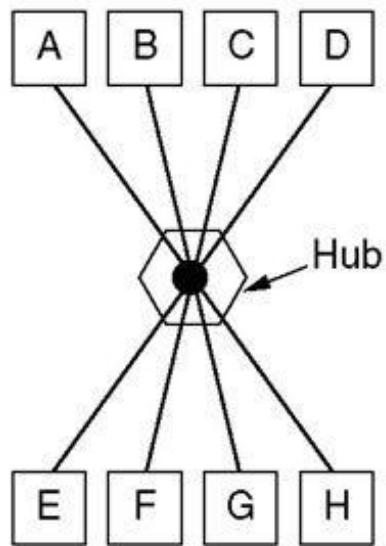
PCs



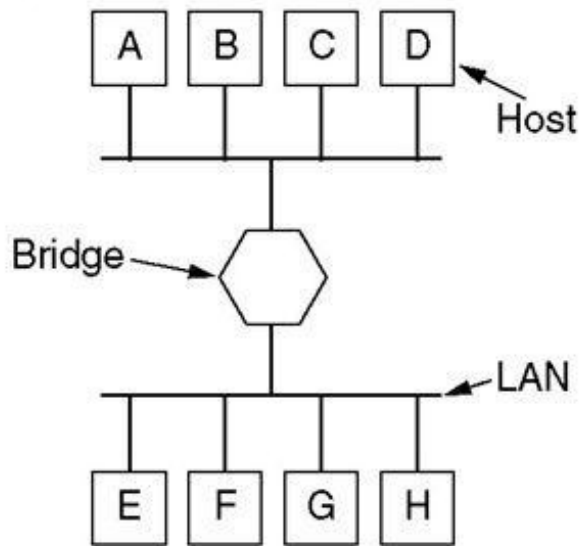
Ethernet



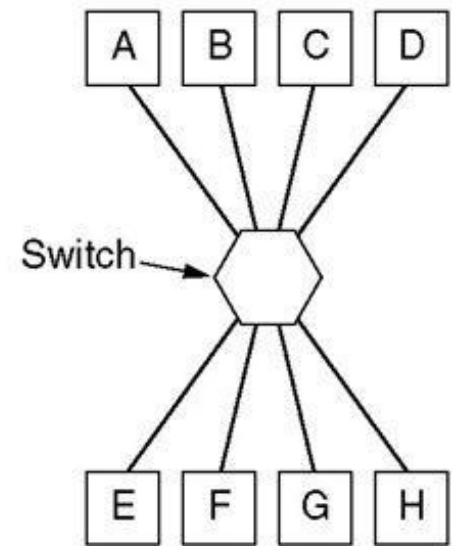
Ethernet



(a)



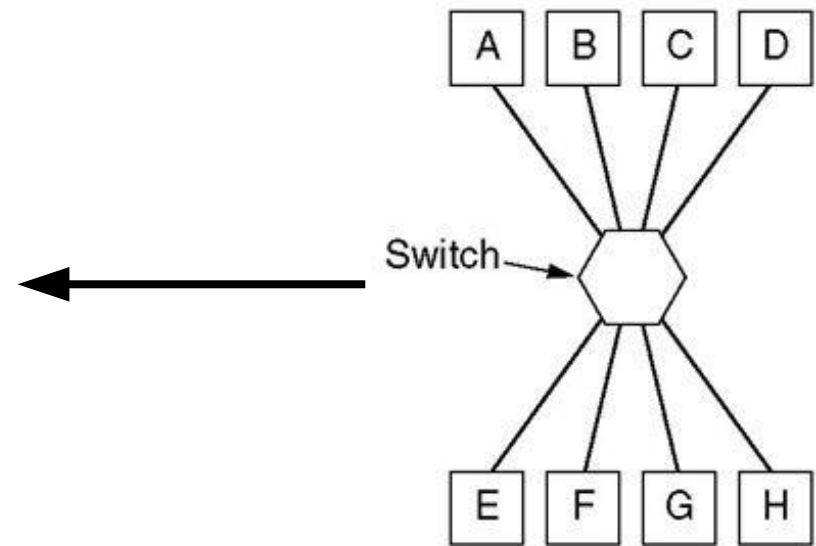
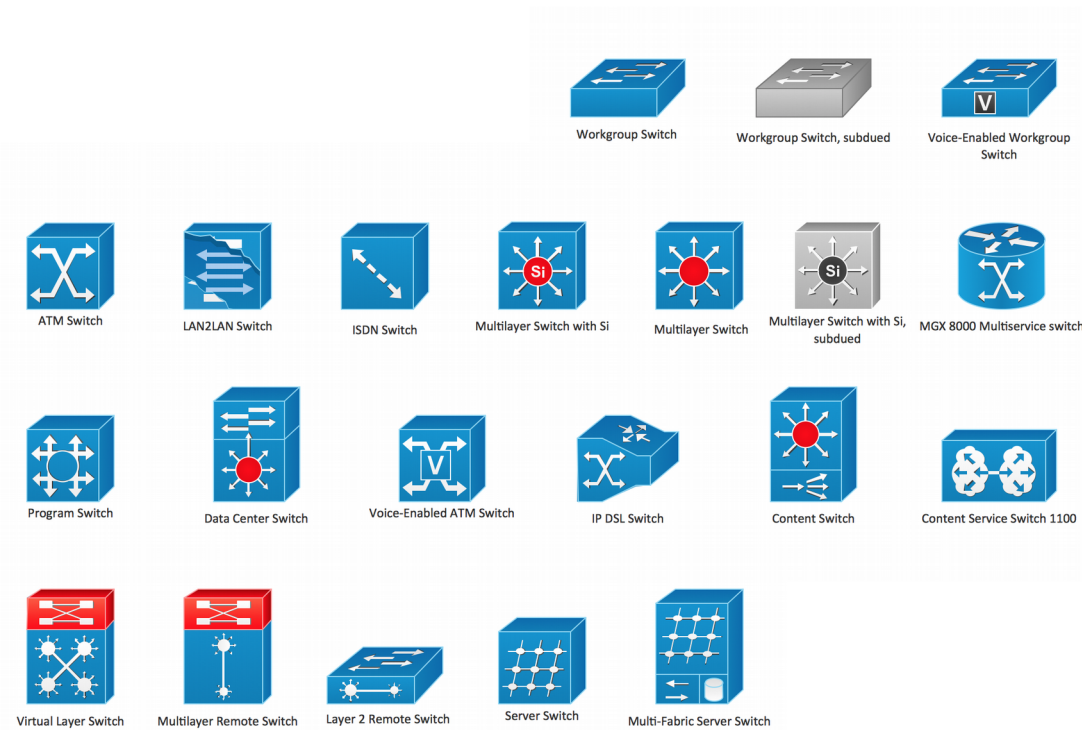
(b)



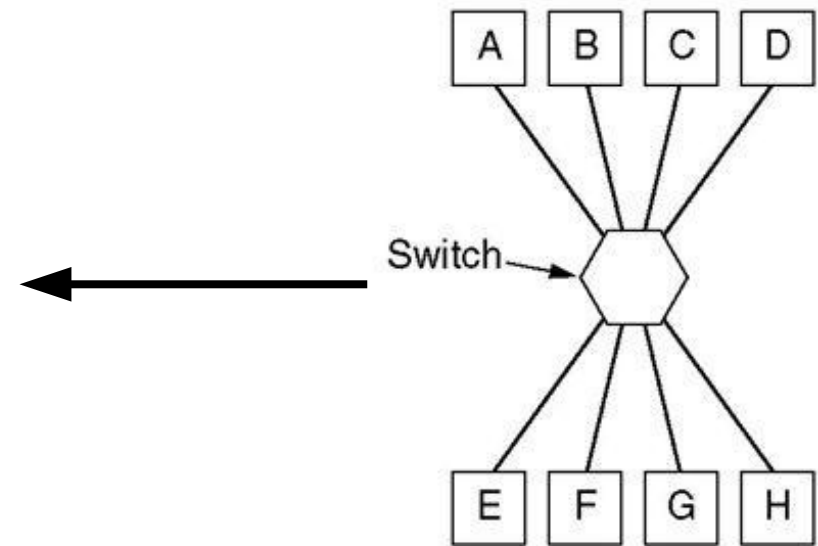
(c)

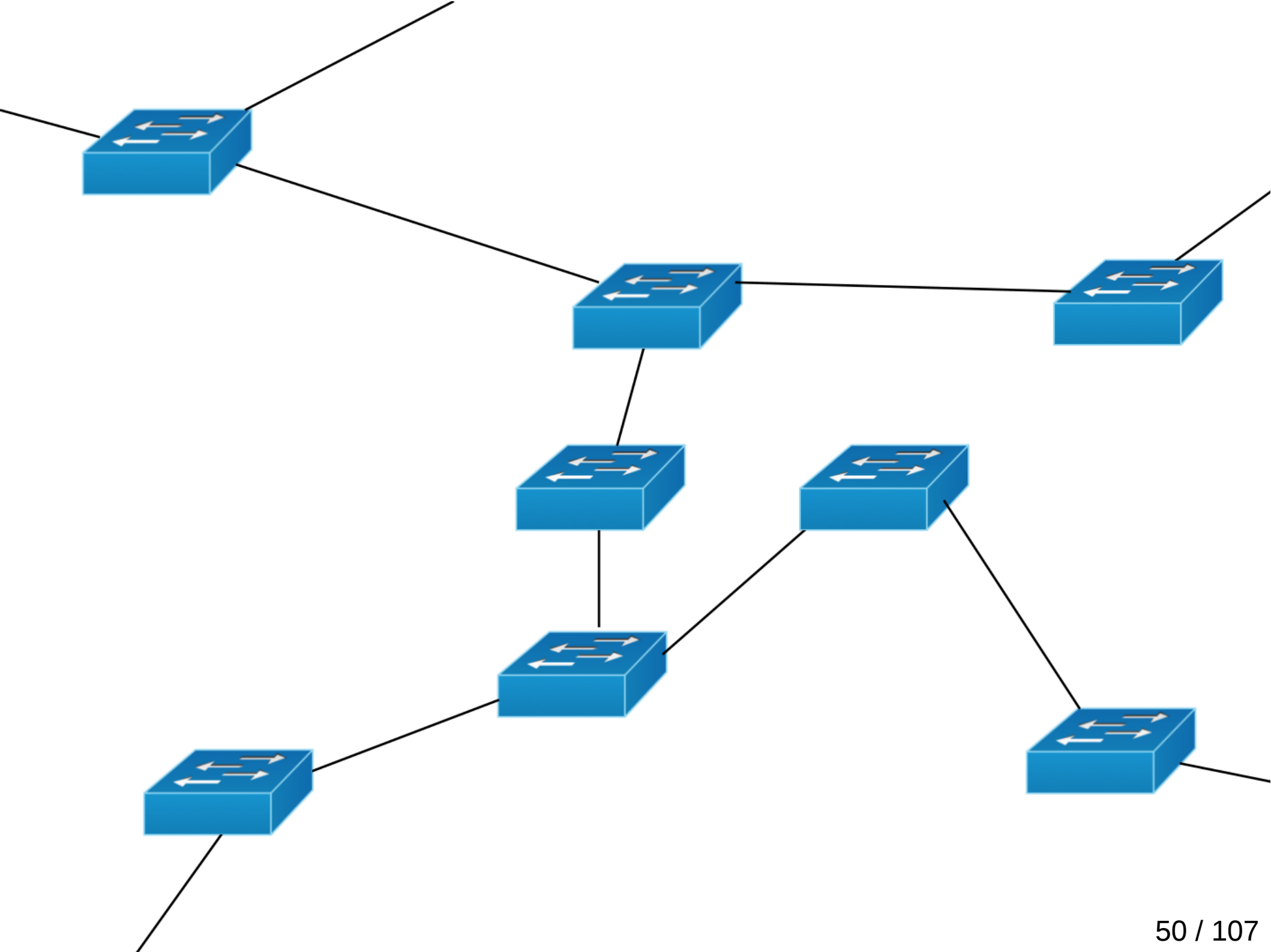
(a) Концентратор (b) Мост (c) Коммутатор

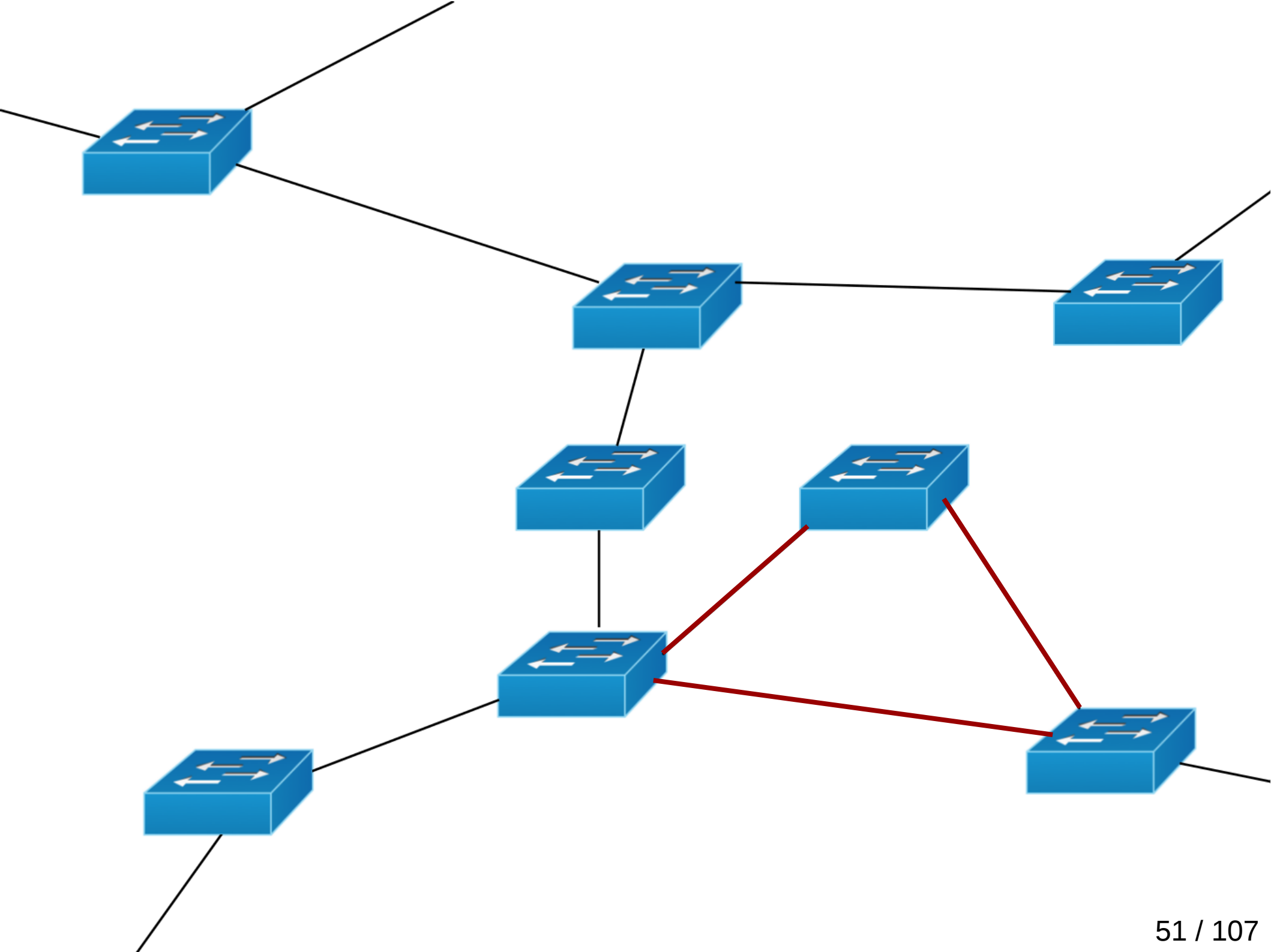
Ethernet

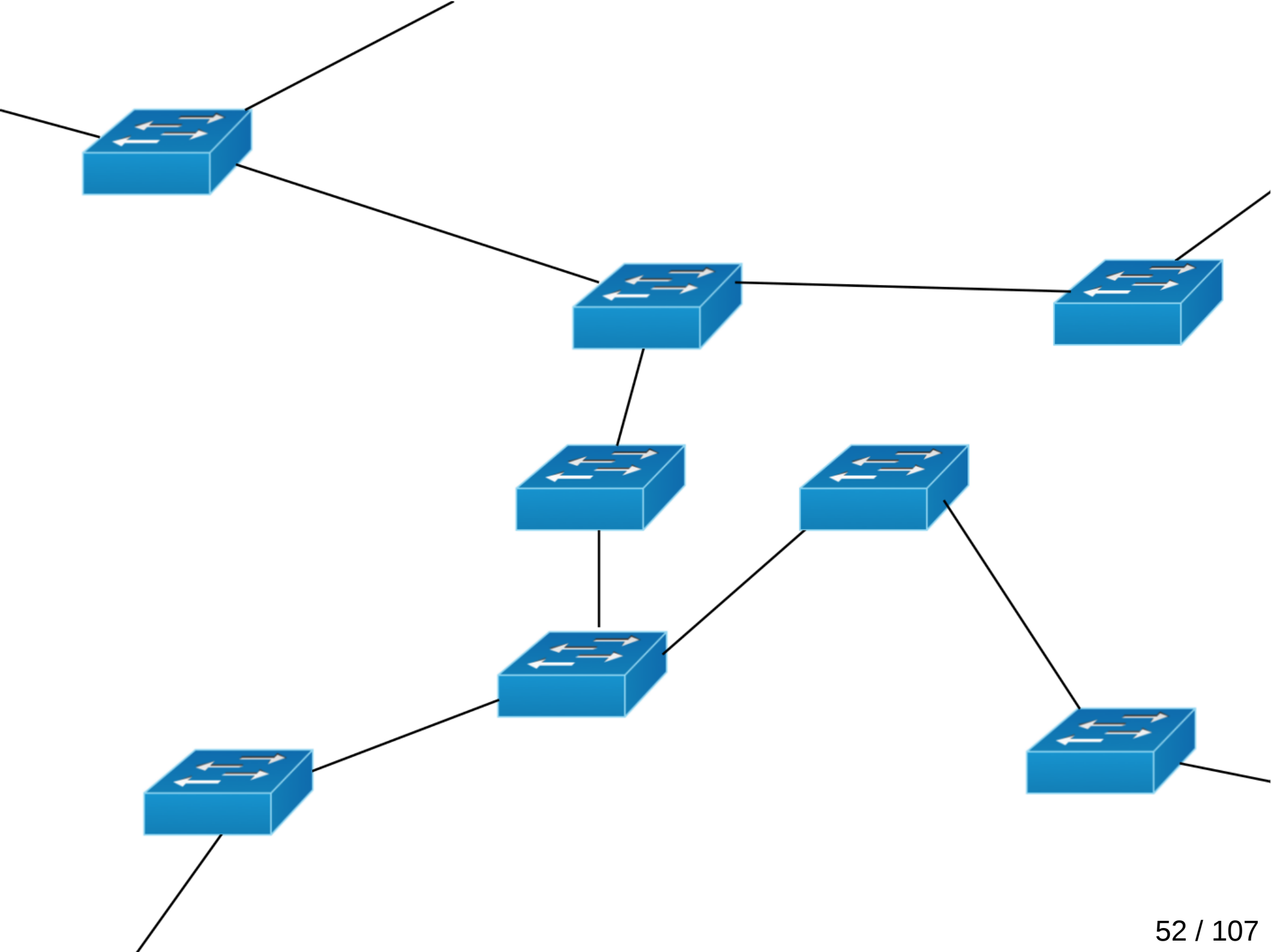


Ethernet





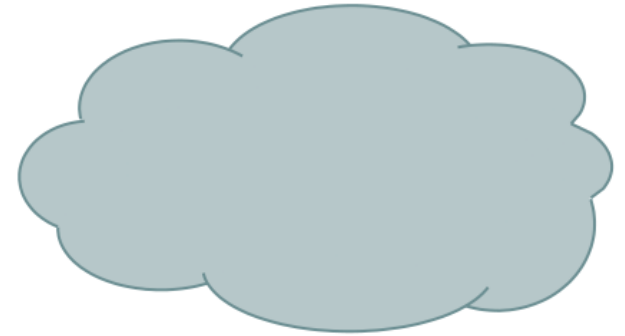
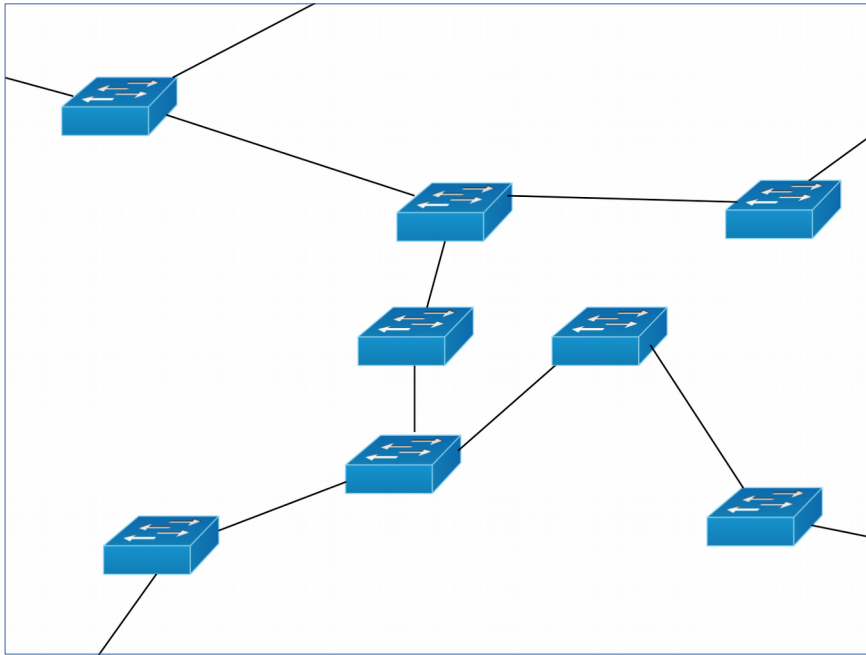


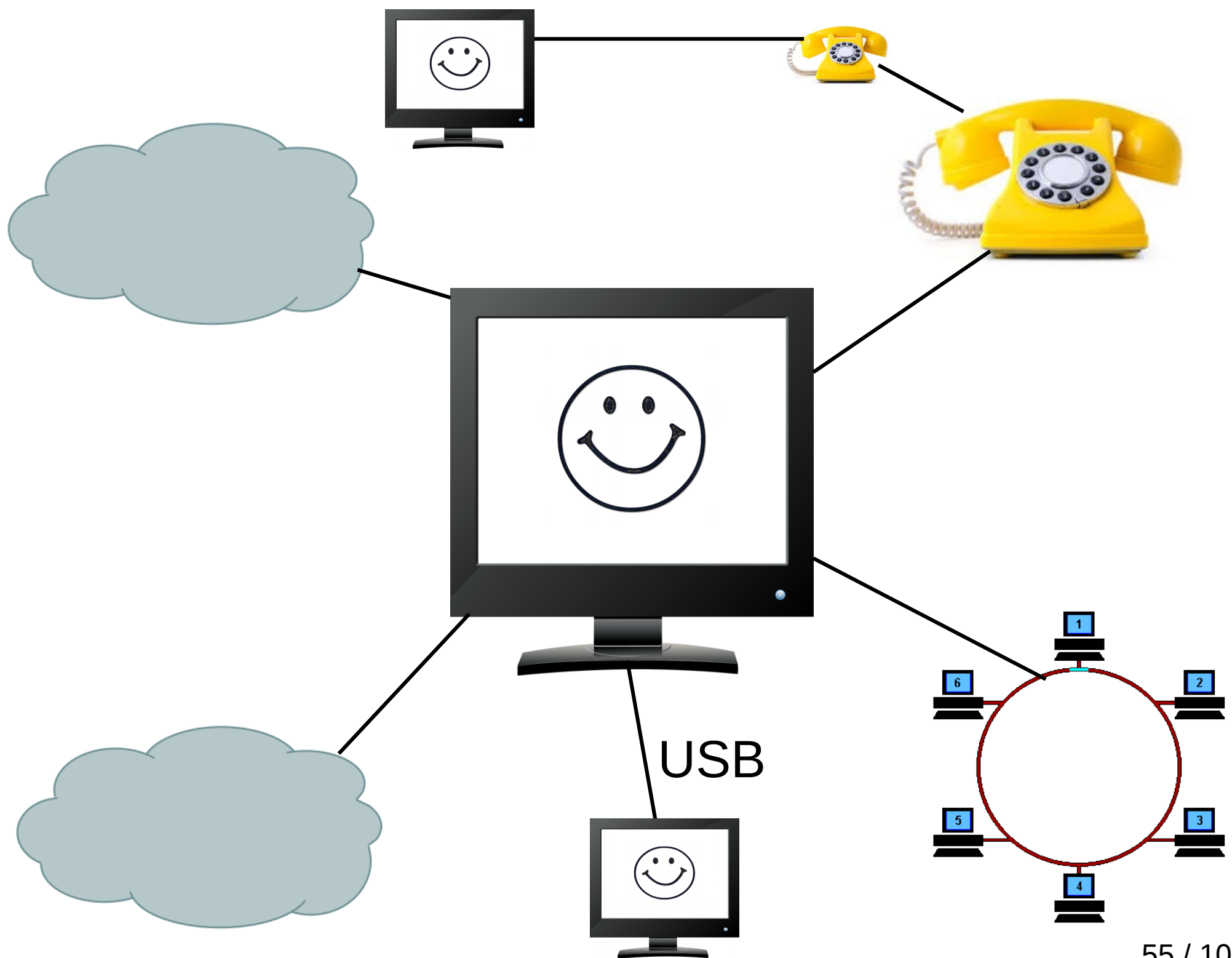


Ethernet

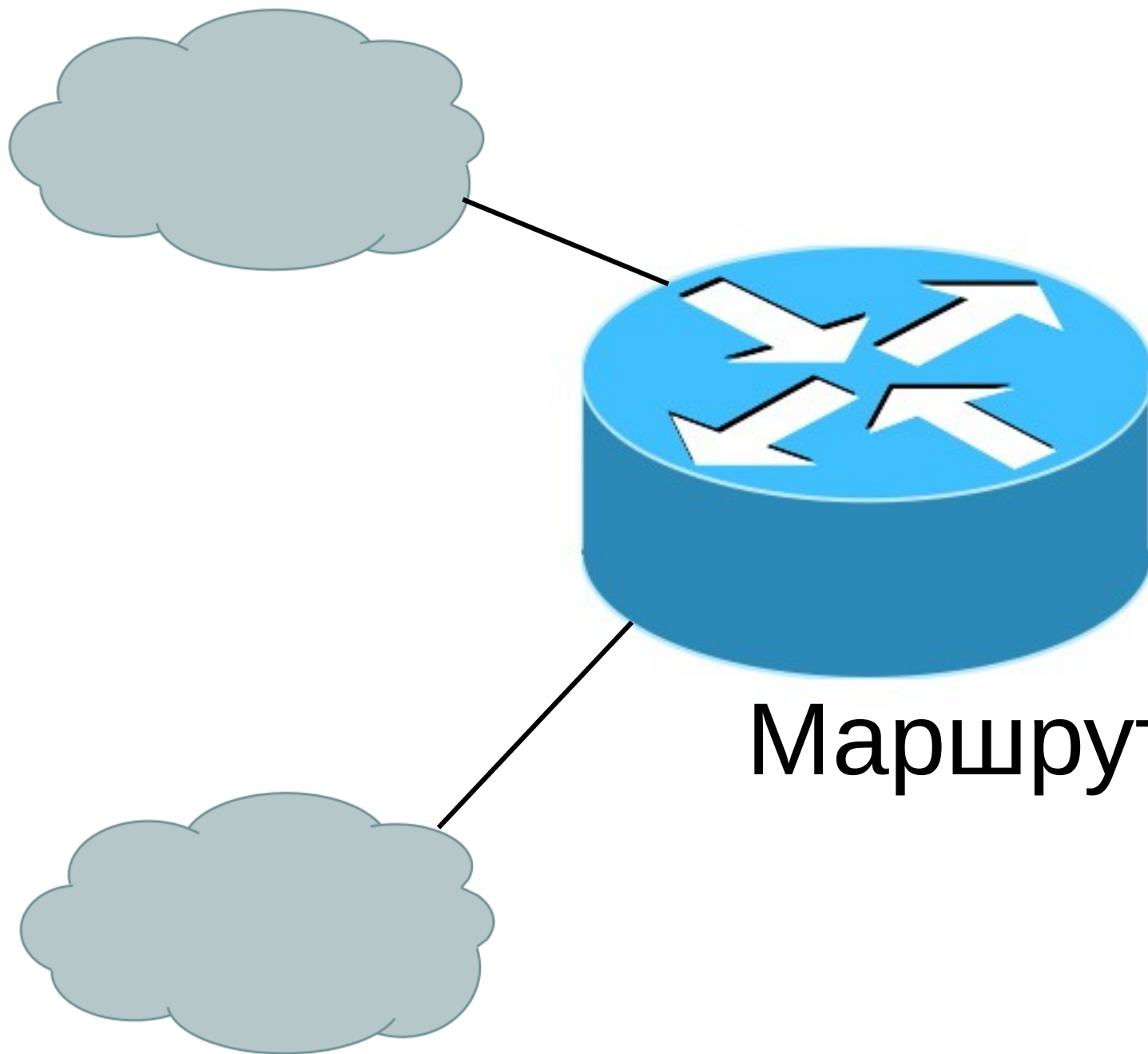
- Не универсально
- Плохо масштабируемо

Ethernet

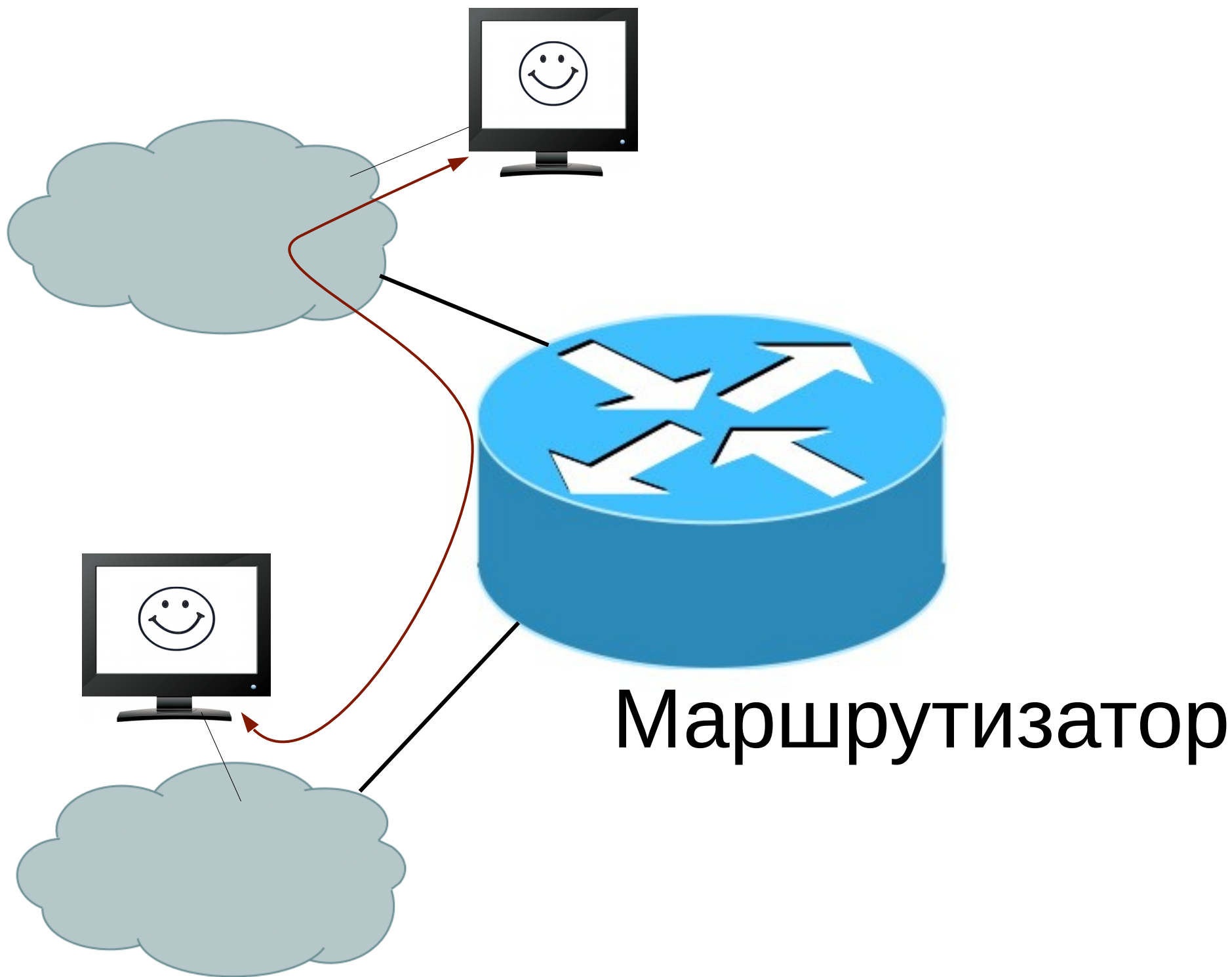


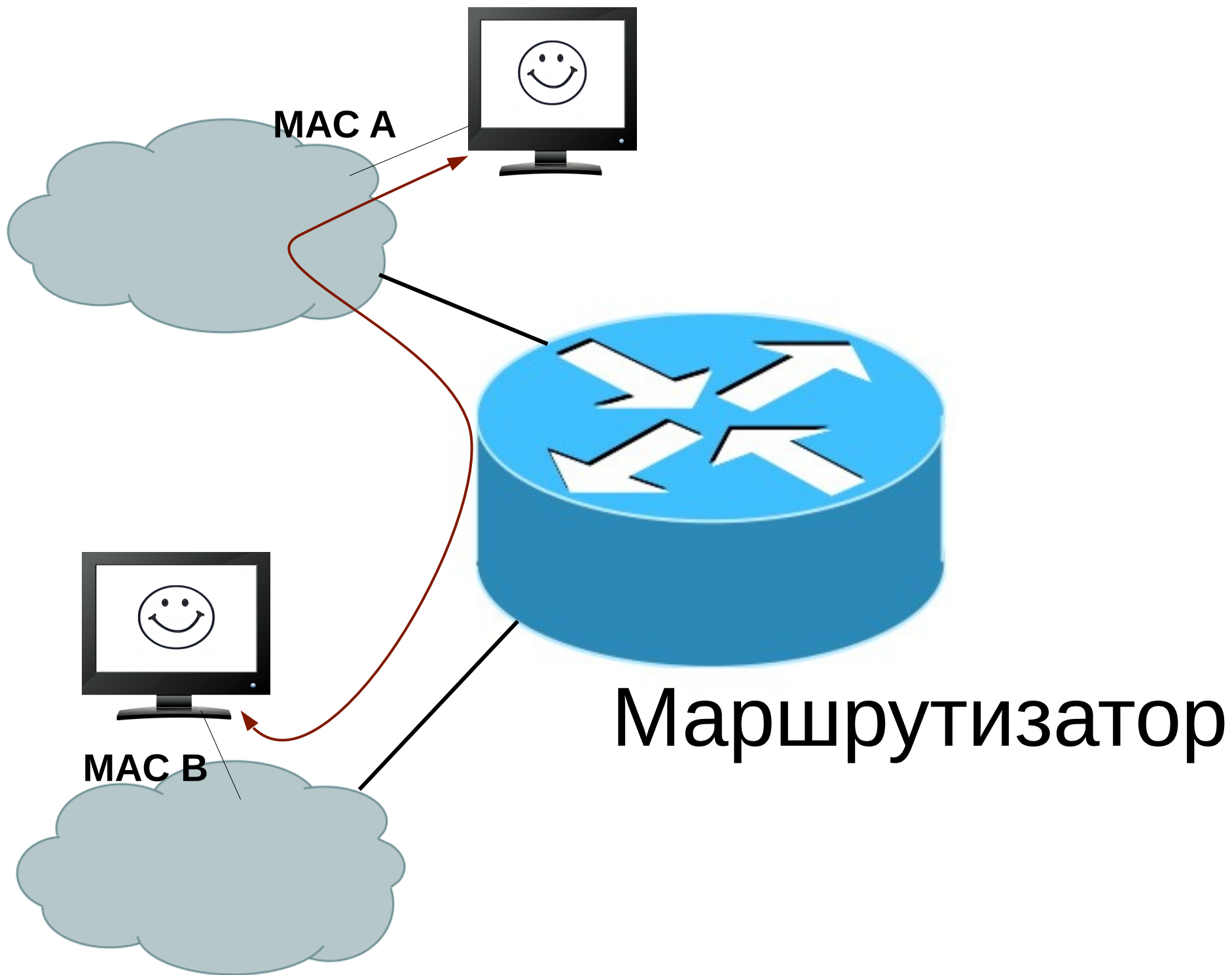


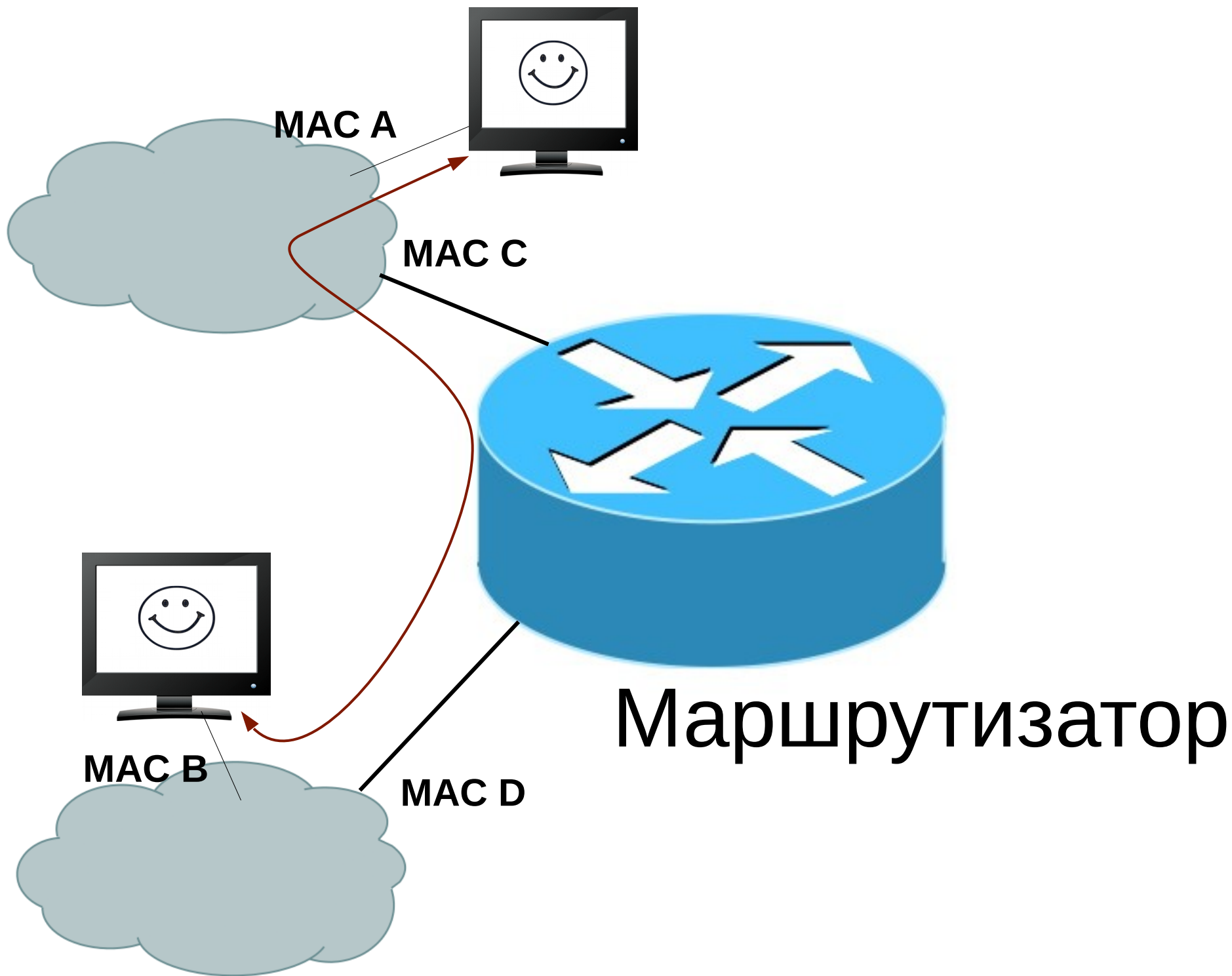


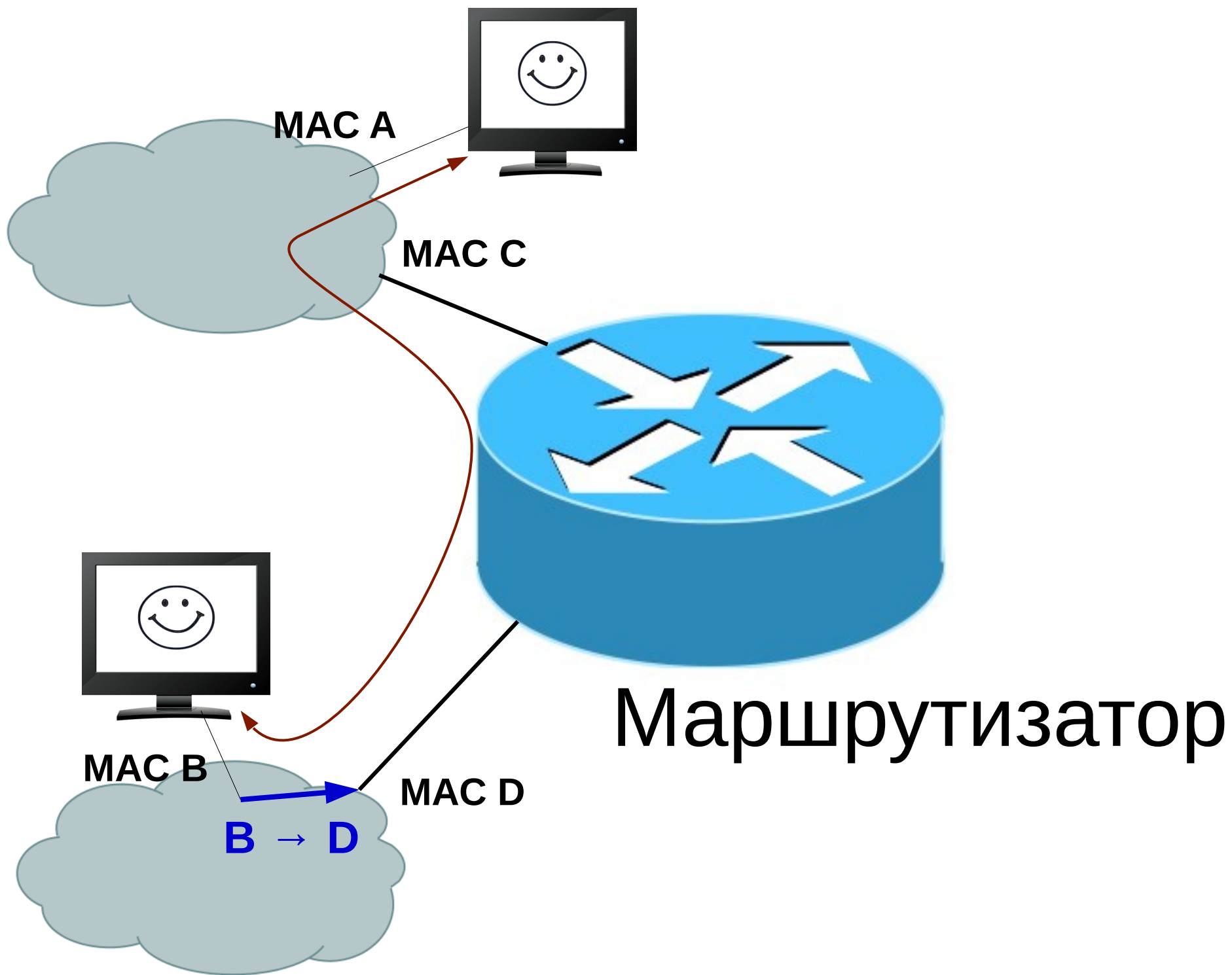


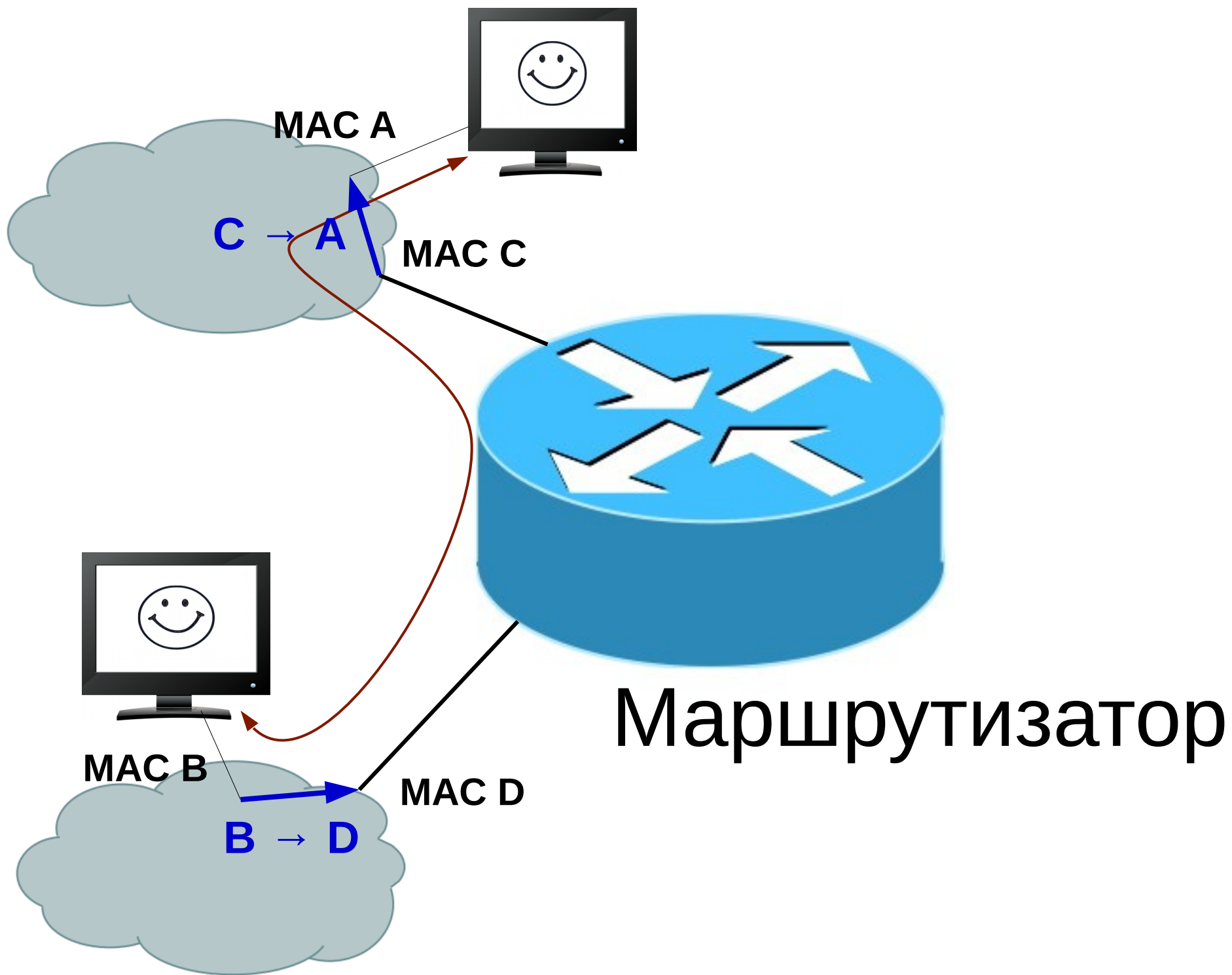
Маршрутизатор

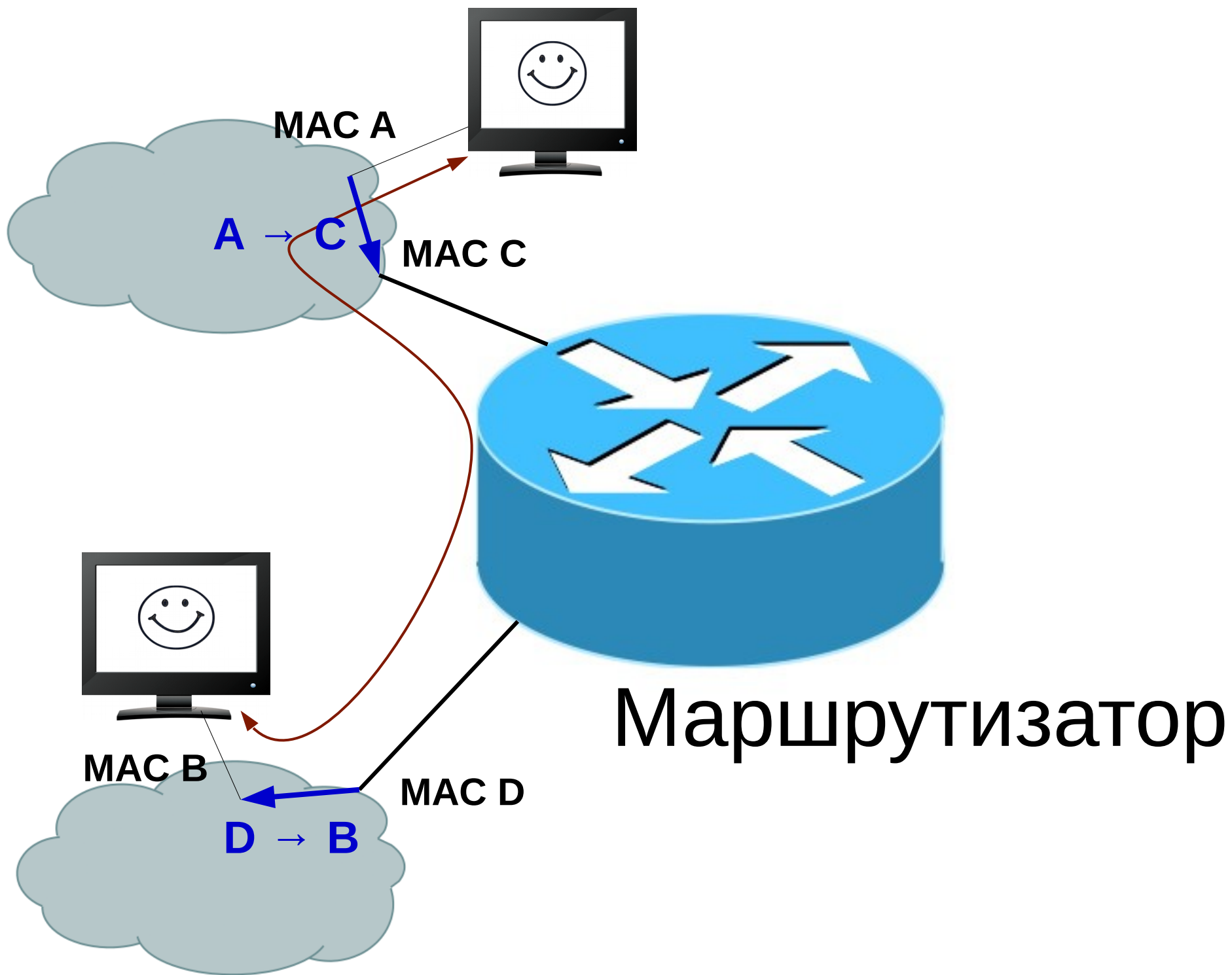




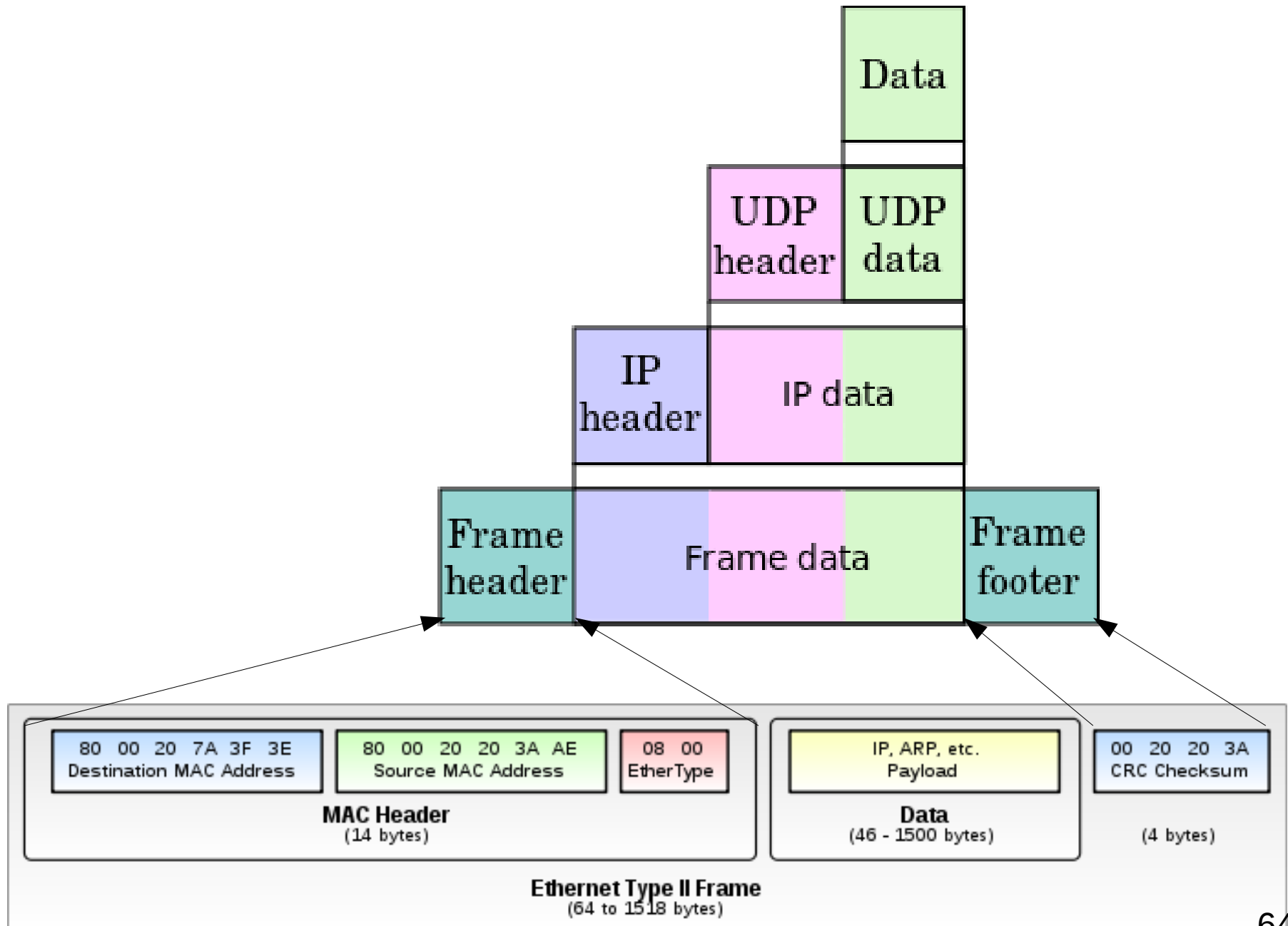




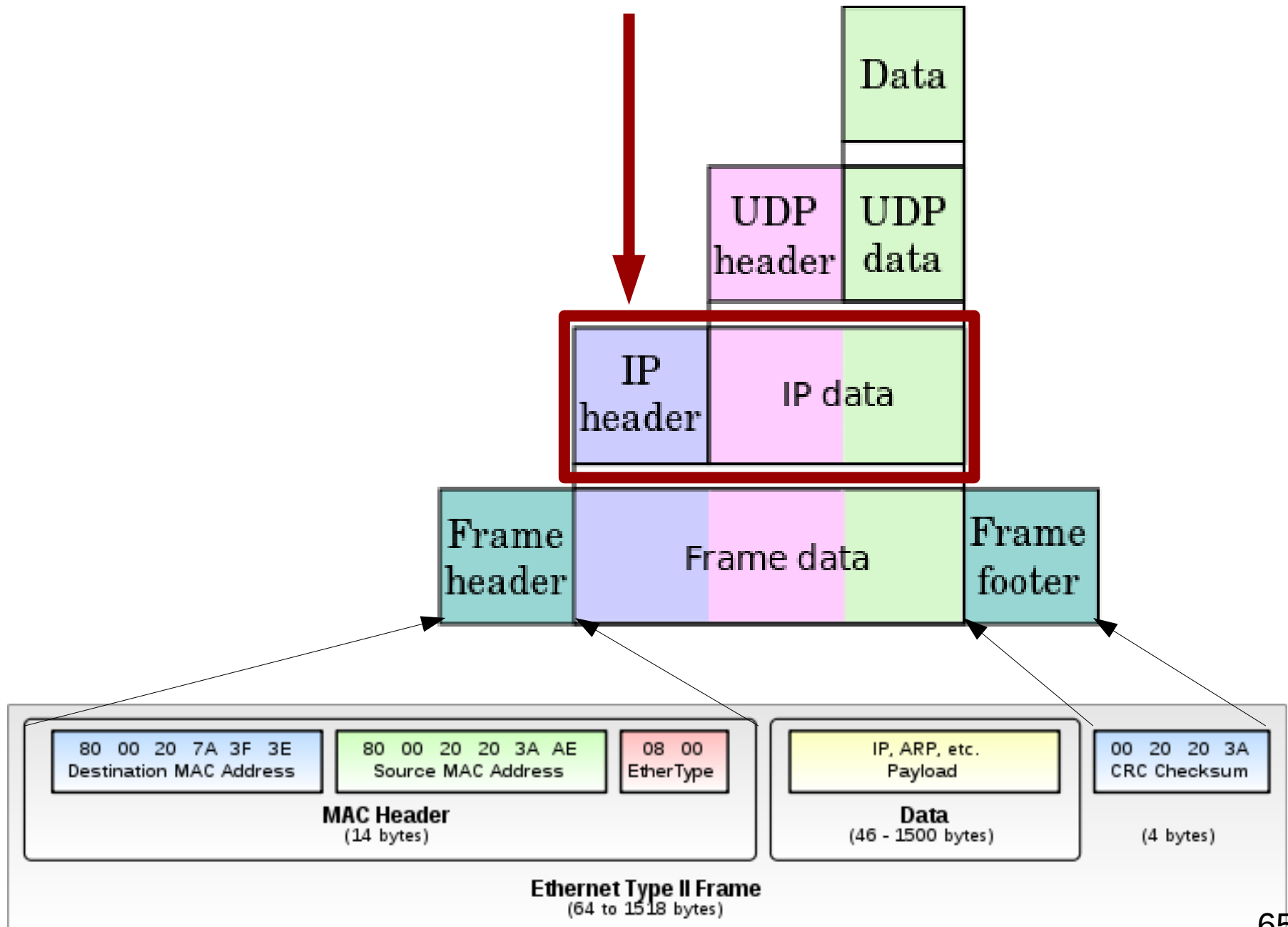




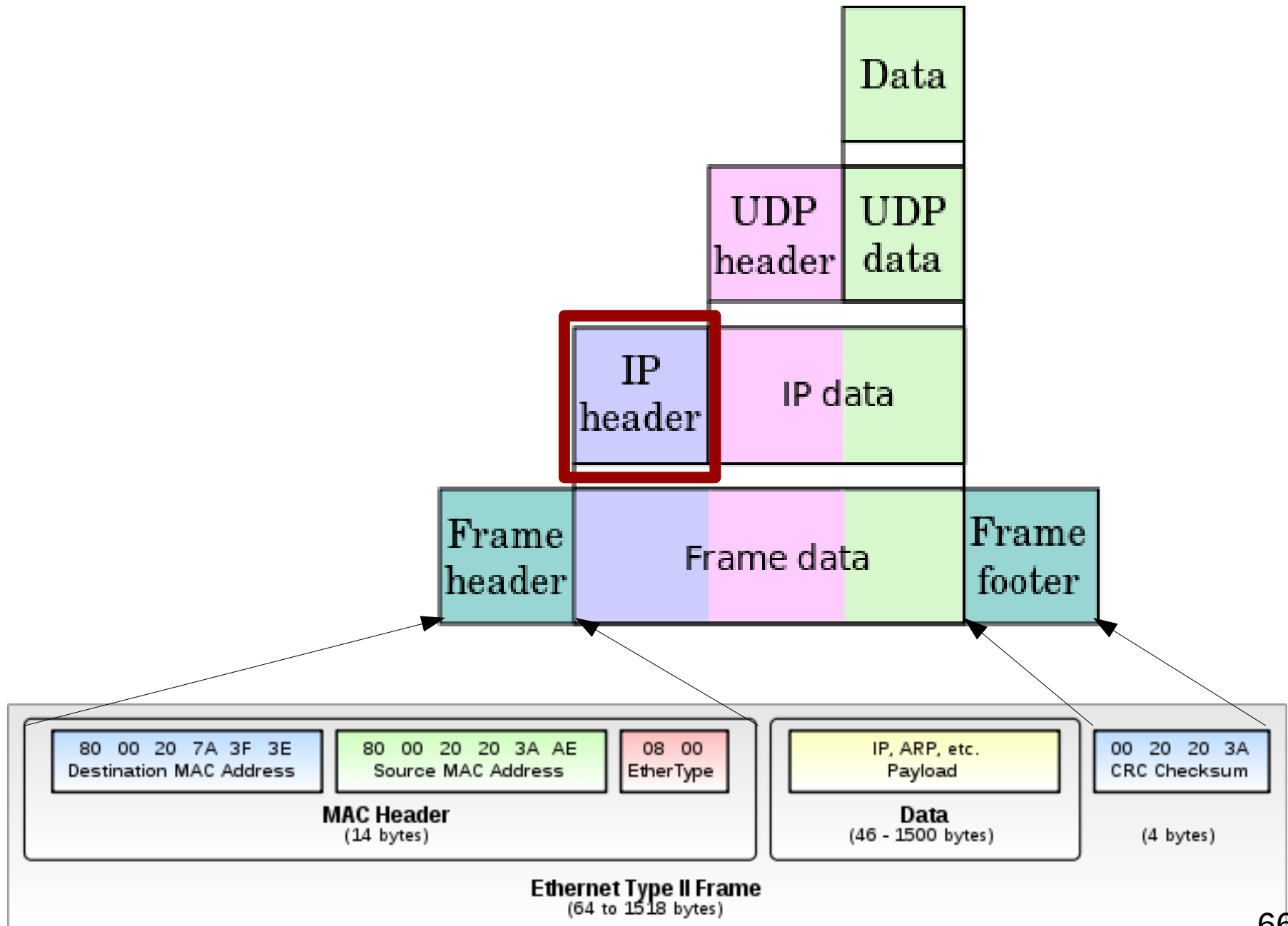
IP



IP



IP



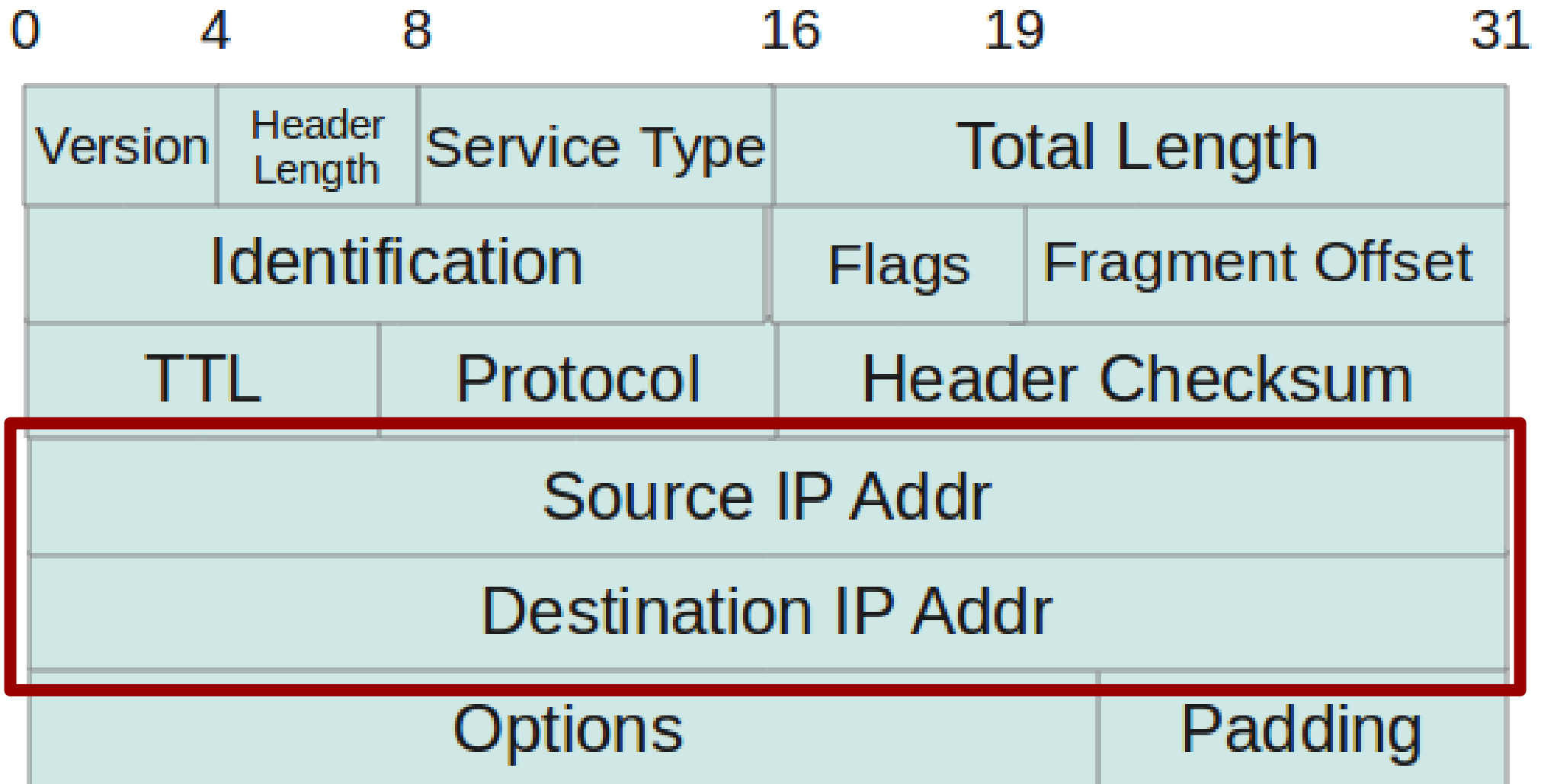


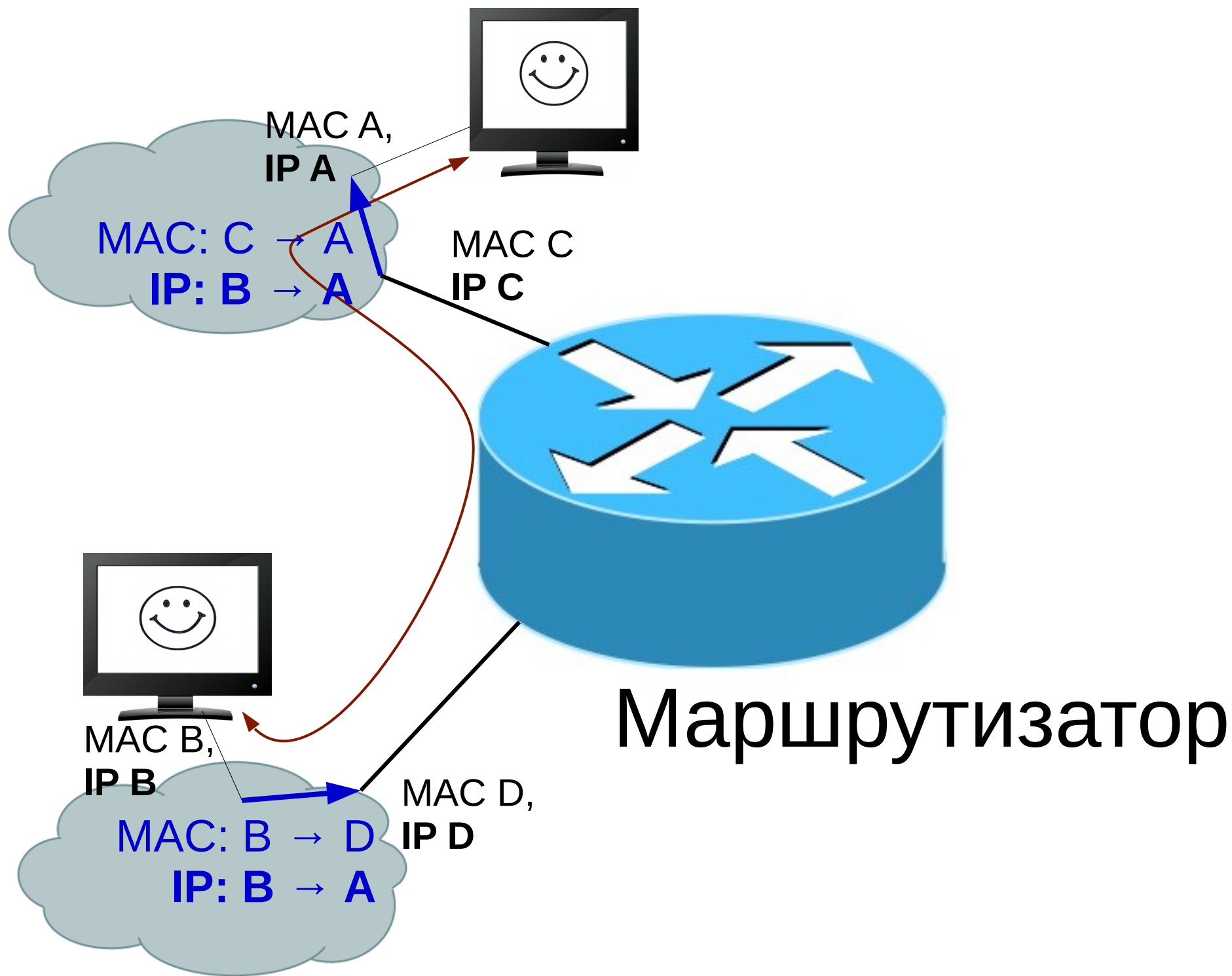
IP

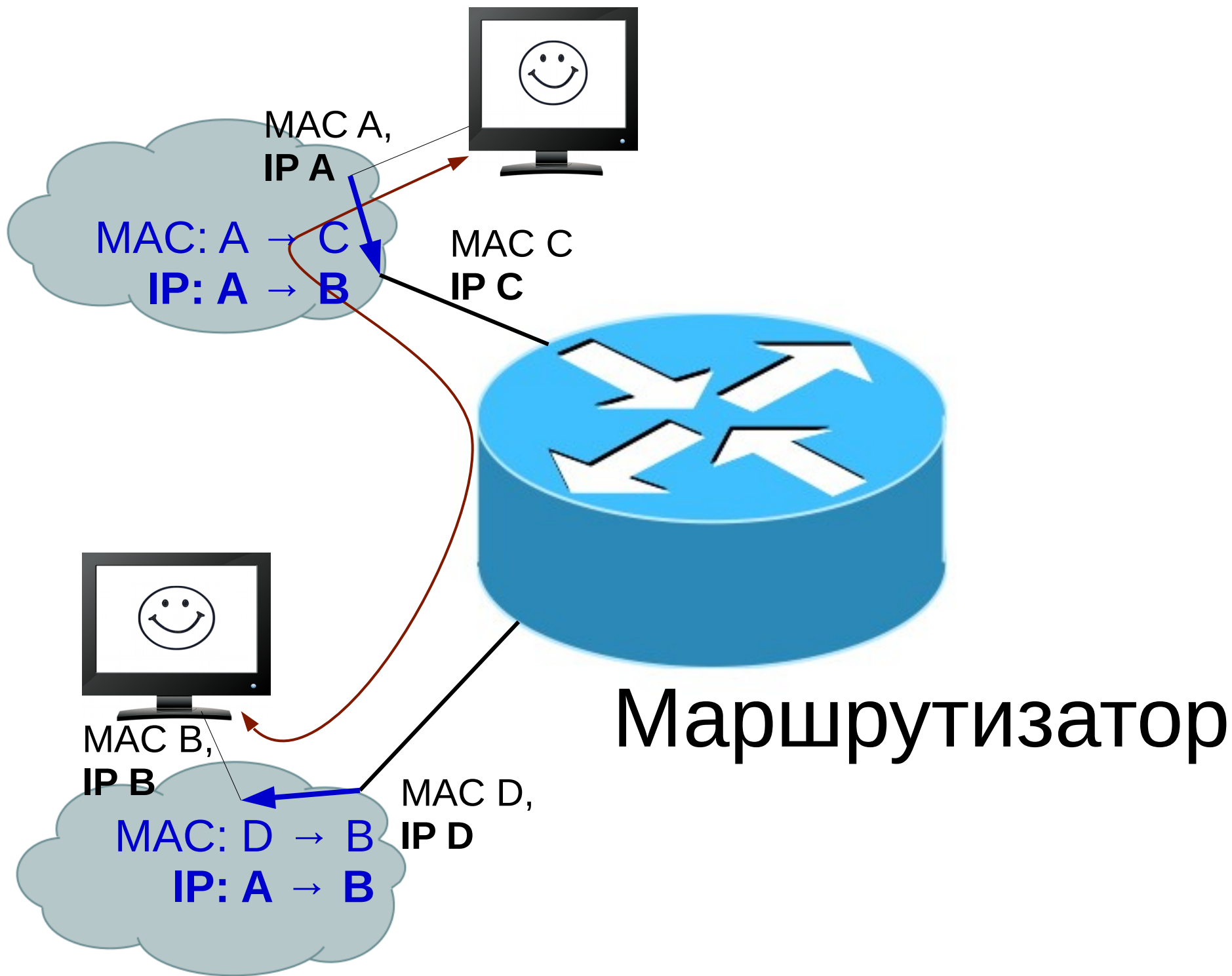
0 4 8 16 19 31

Version	Header Length	Service Type	Total Length	
Identification			Flags	Fragment Offset
TTL	Protocol		Header Checksum	
Source IP Addr				
Destination IP Addr				
Options				Padding

IP

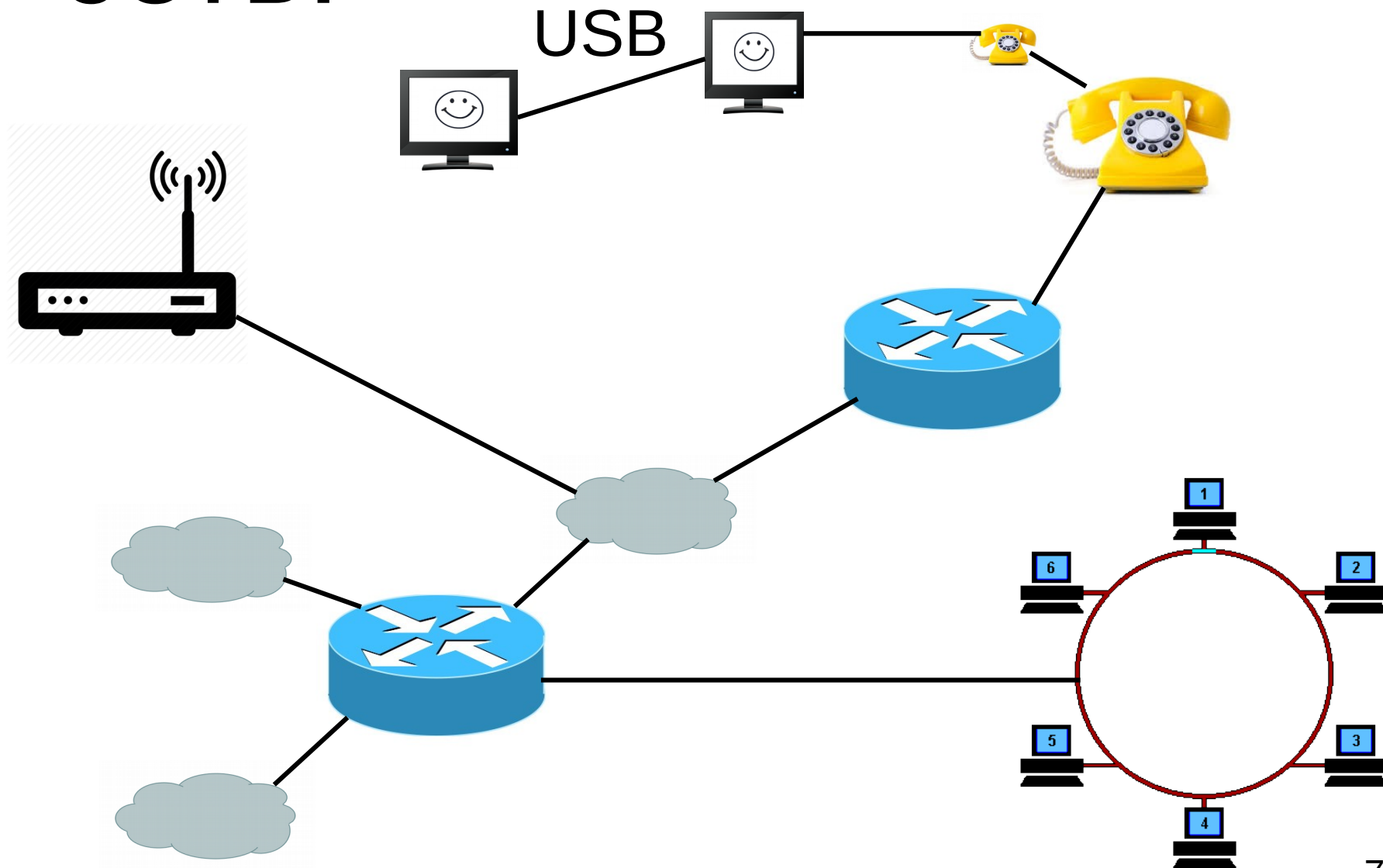




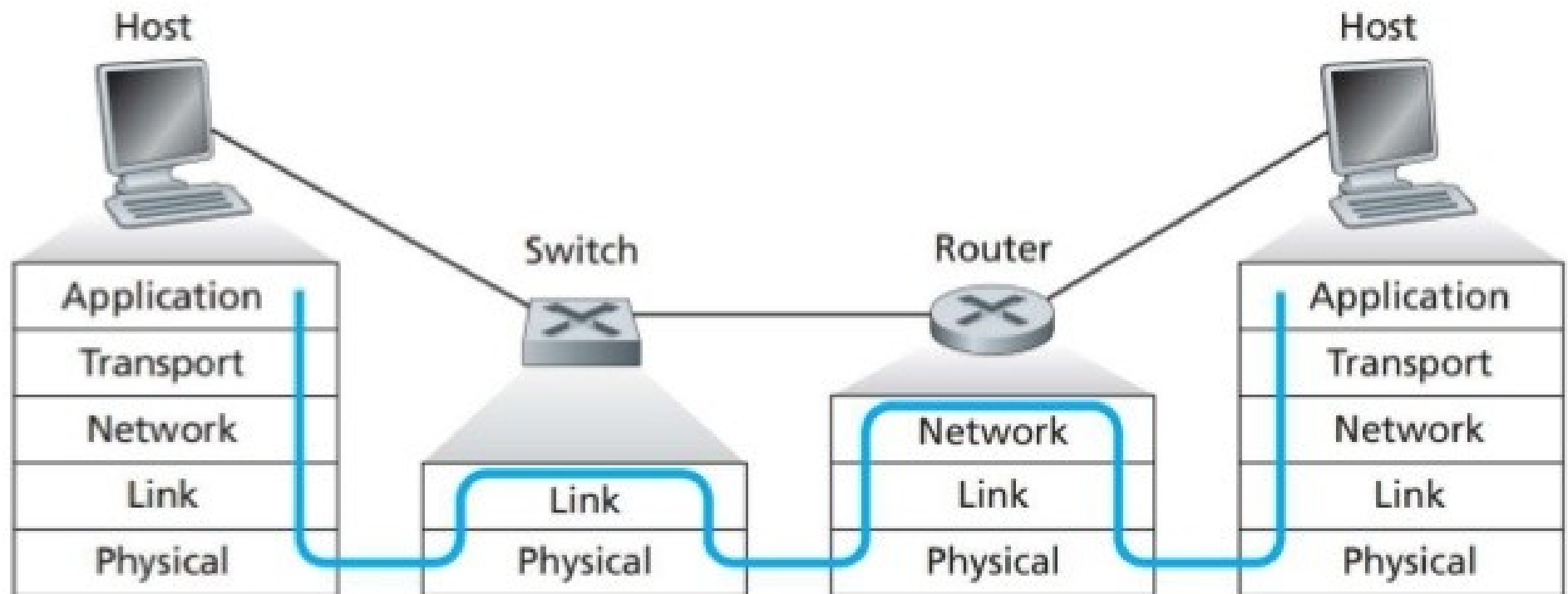


IP

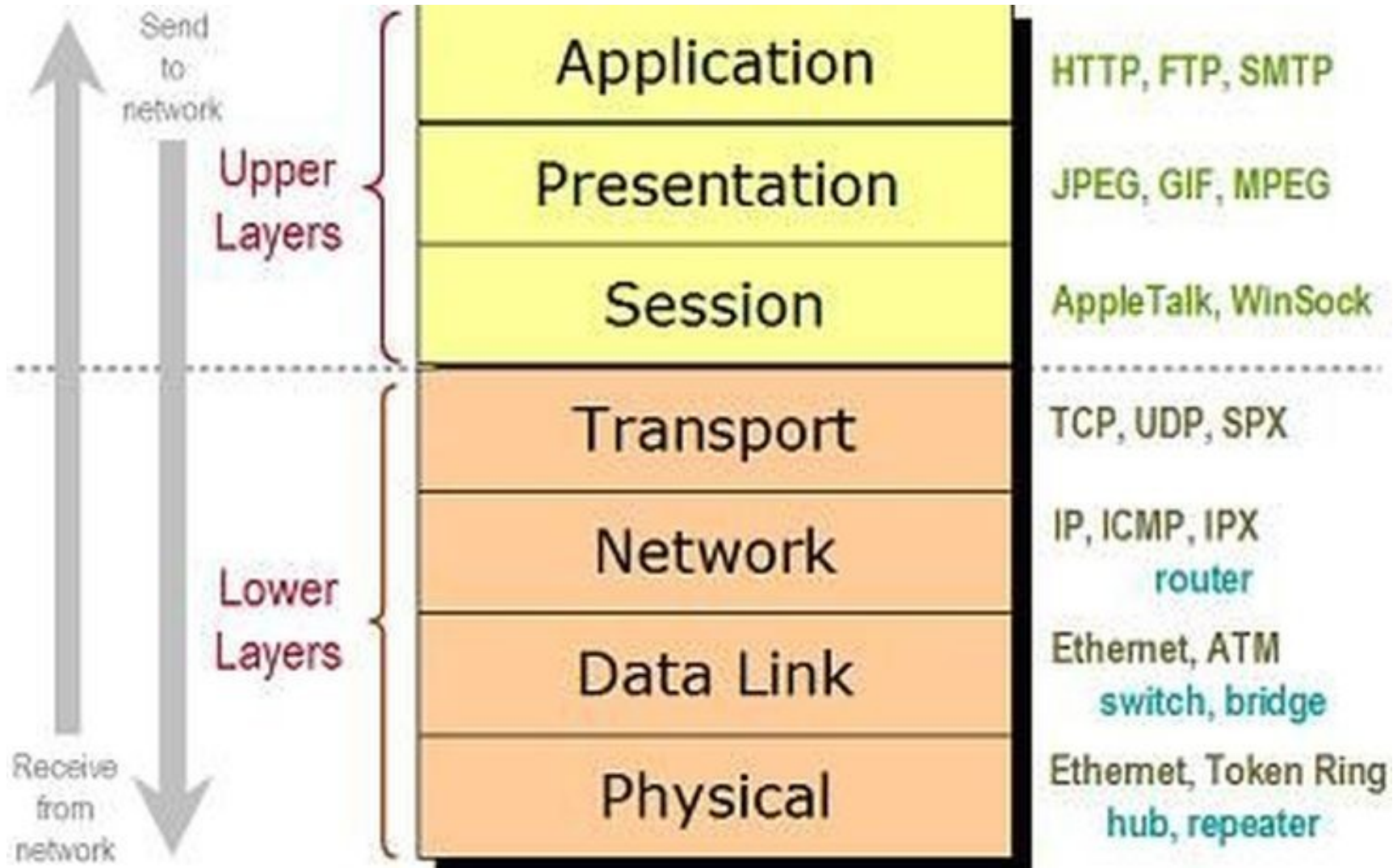
IP-сеть:



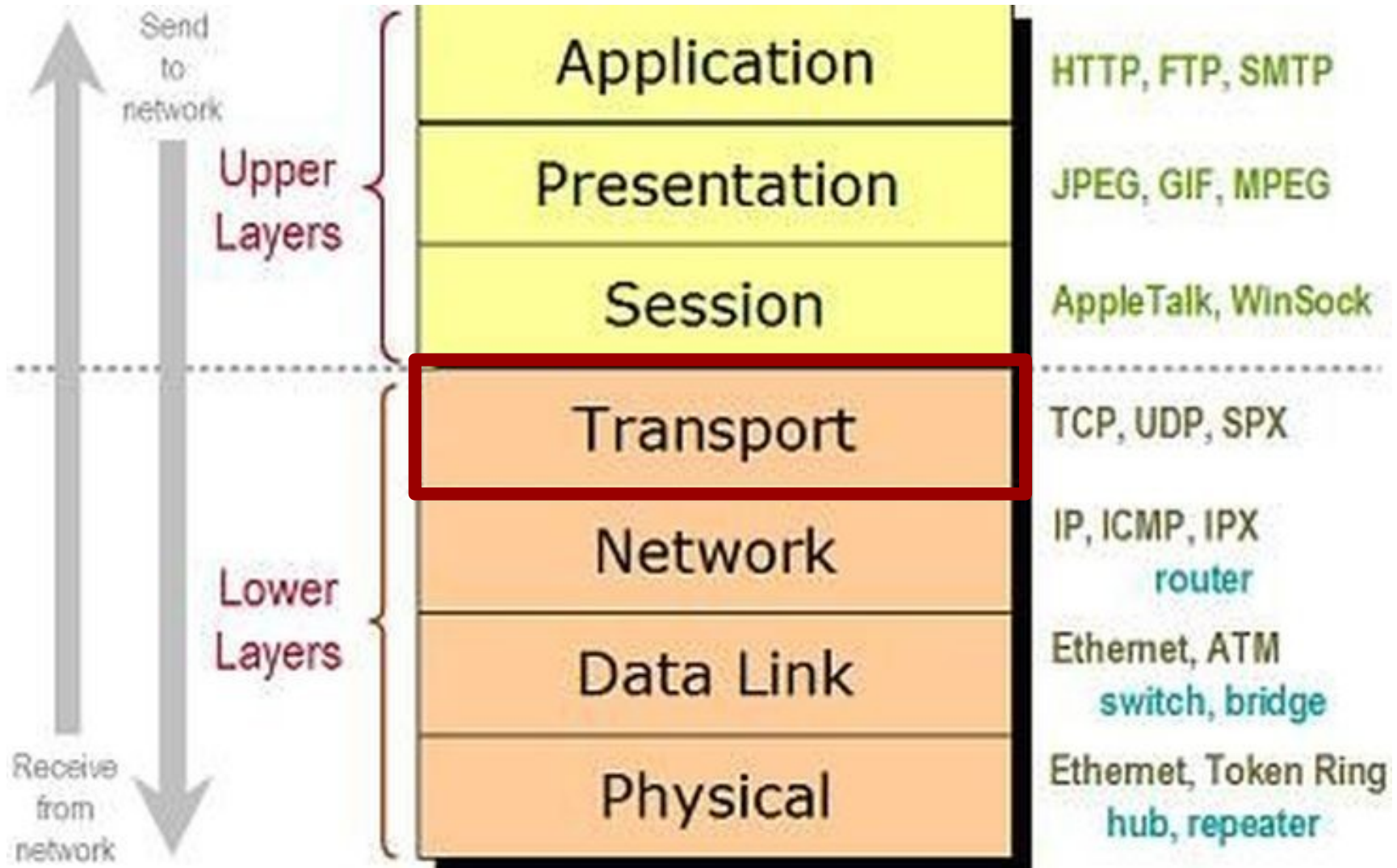
IP

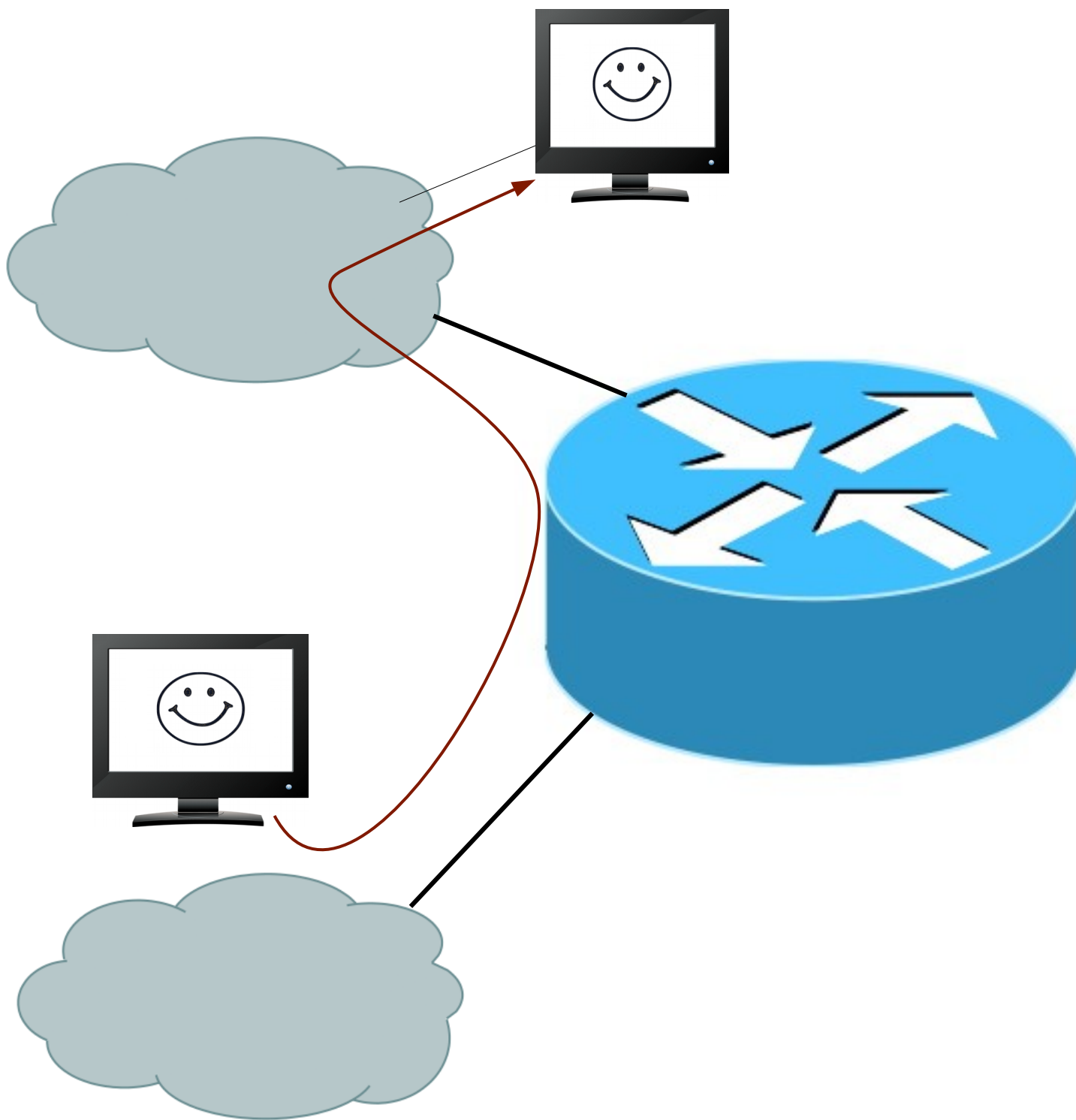


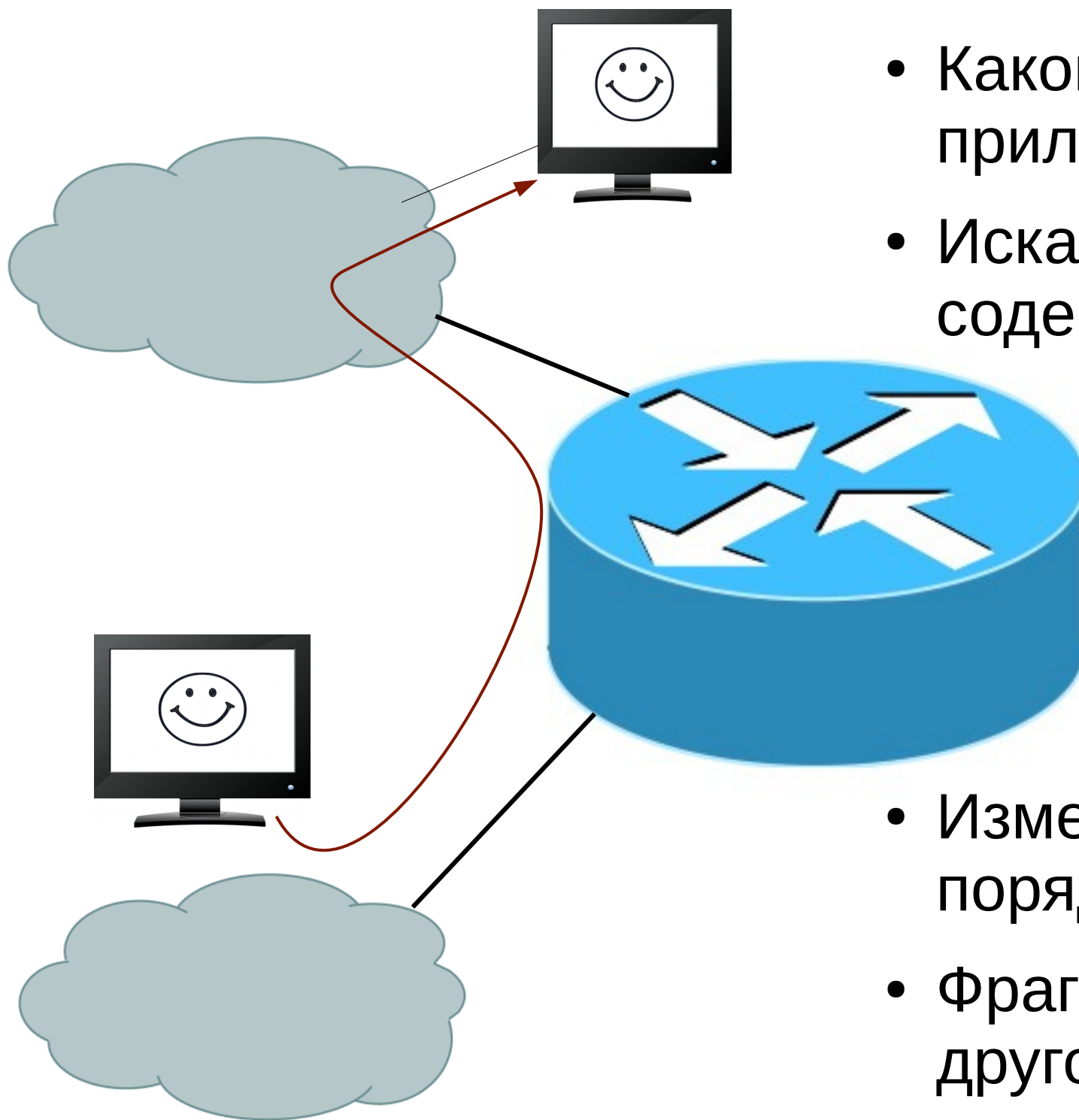
OSI



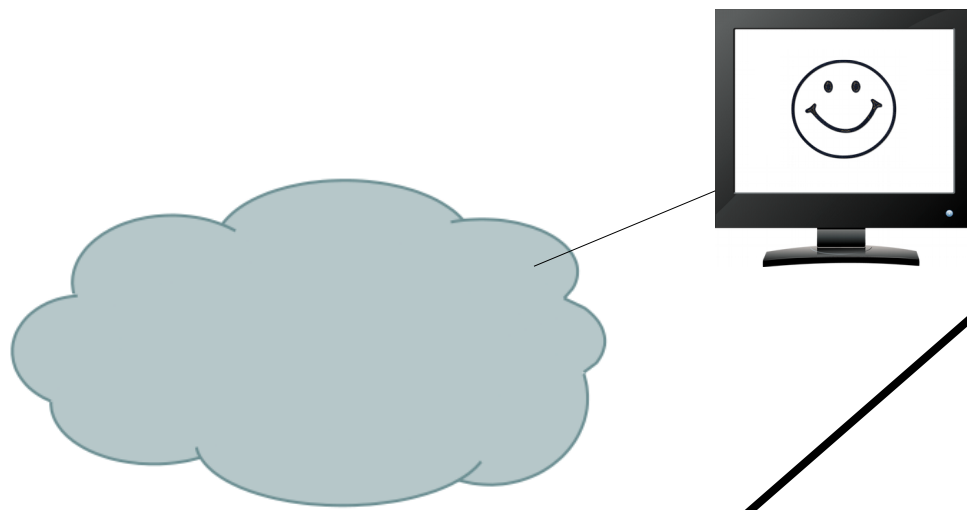
OSI







- Какому приложению?
- Искажение содержимого?
- Изменение порядка пакетов?
- Фрагментация и другое...



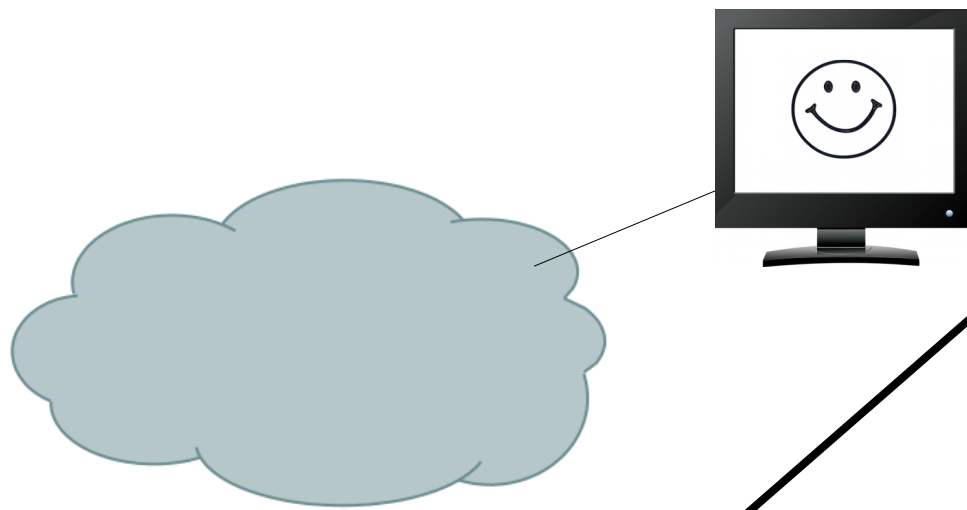
- Какому приложению?
- Искажение содержимого?

контрольная сумма

порт

порядковый номер

- Изменение порядка пакетов?
- Фрагментация и другое...



- Какому приложению?
- Искажение содержимого?

порт

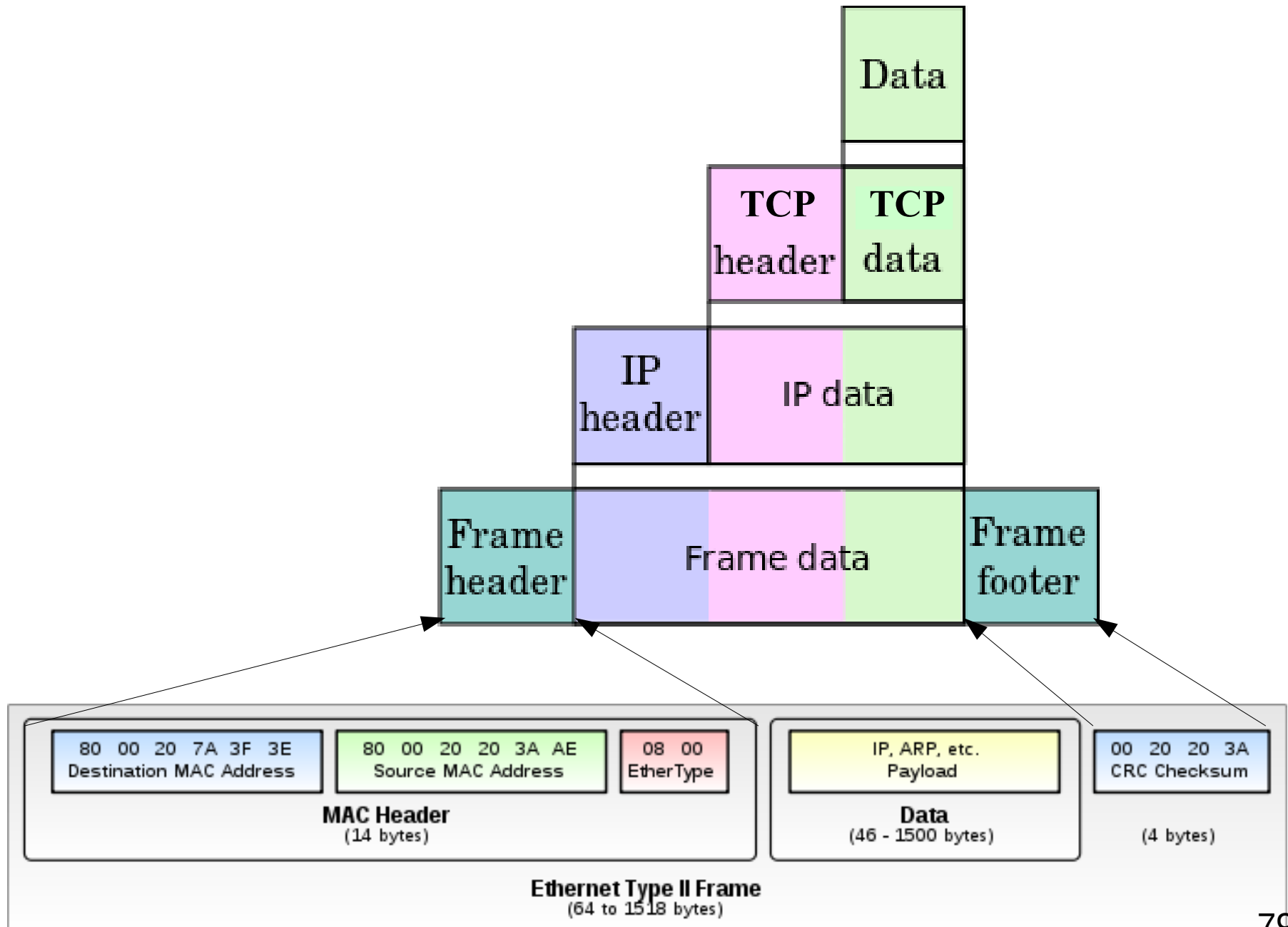
контрольная сумма

порядковый номер

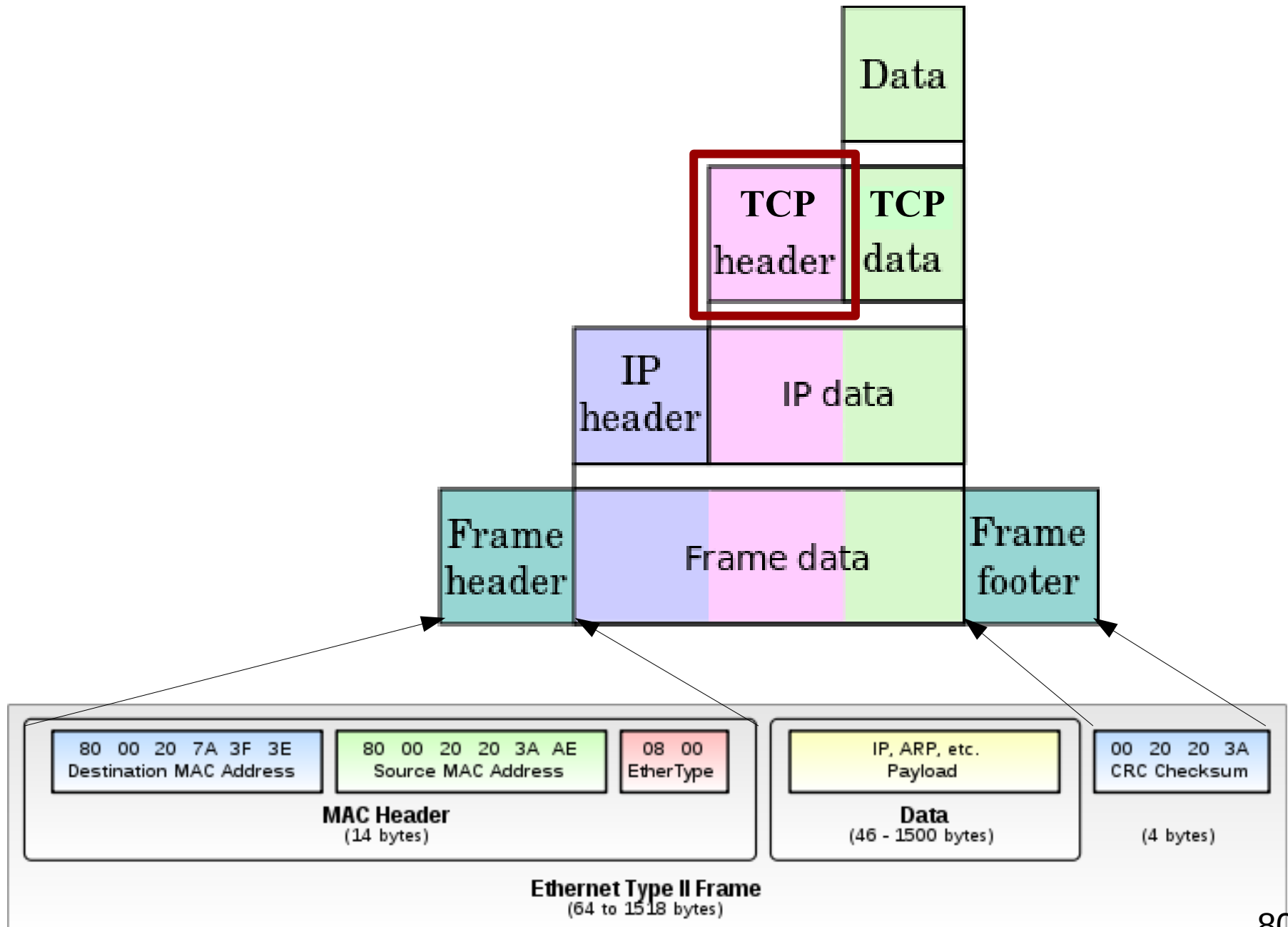
- Изменение порядка пакетов?
- Фрагментация и другое...

И прочее...

TCP

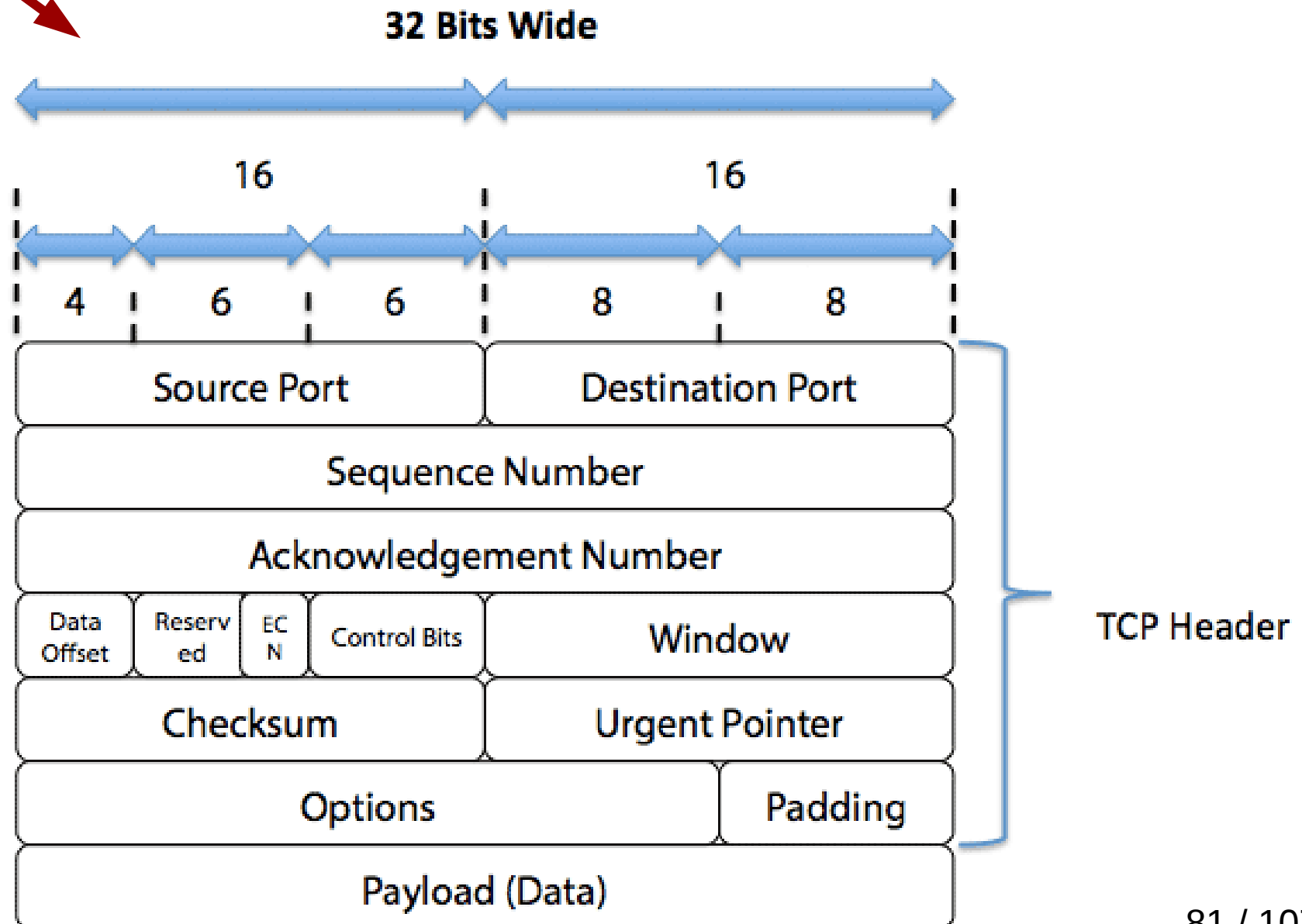


TCP



TCP
header

TCP



TCP



TCP State

TCP Packet



TCP State

CLOSED

LISTEN



SEQ = 1000, CTL = SYN



SYN-SENT

SYN-RECEIVED



SEQ = 750, ACK = 1001, CTL = SYN | ACK



ESTABLISHED

SYN-RECEIVED



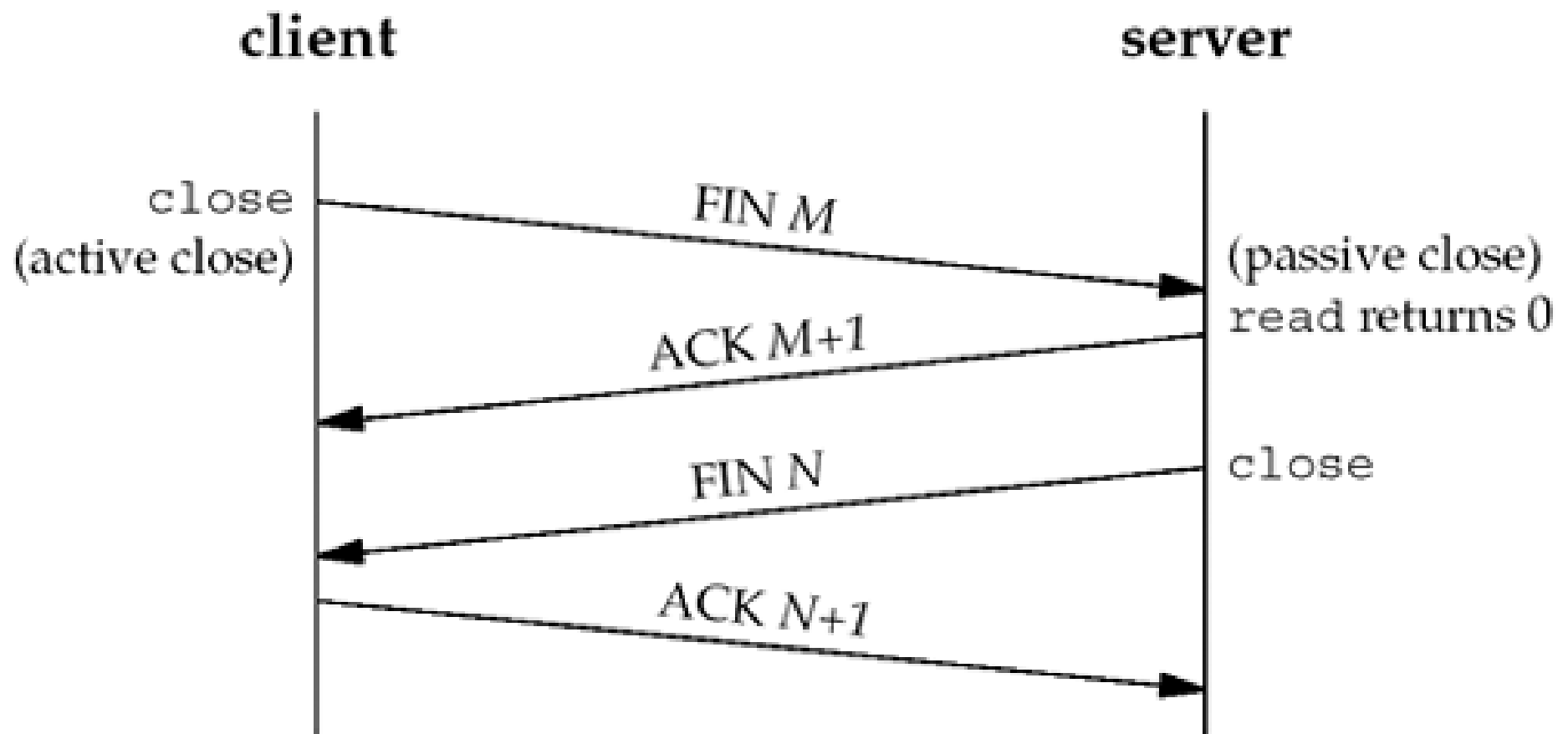
SEQ = 1000, ACK = 751, CTL = ACK



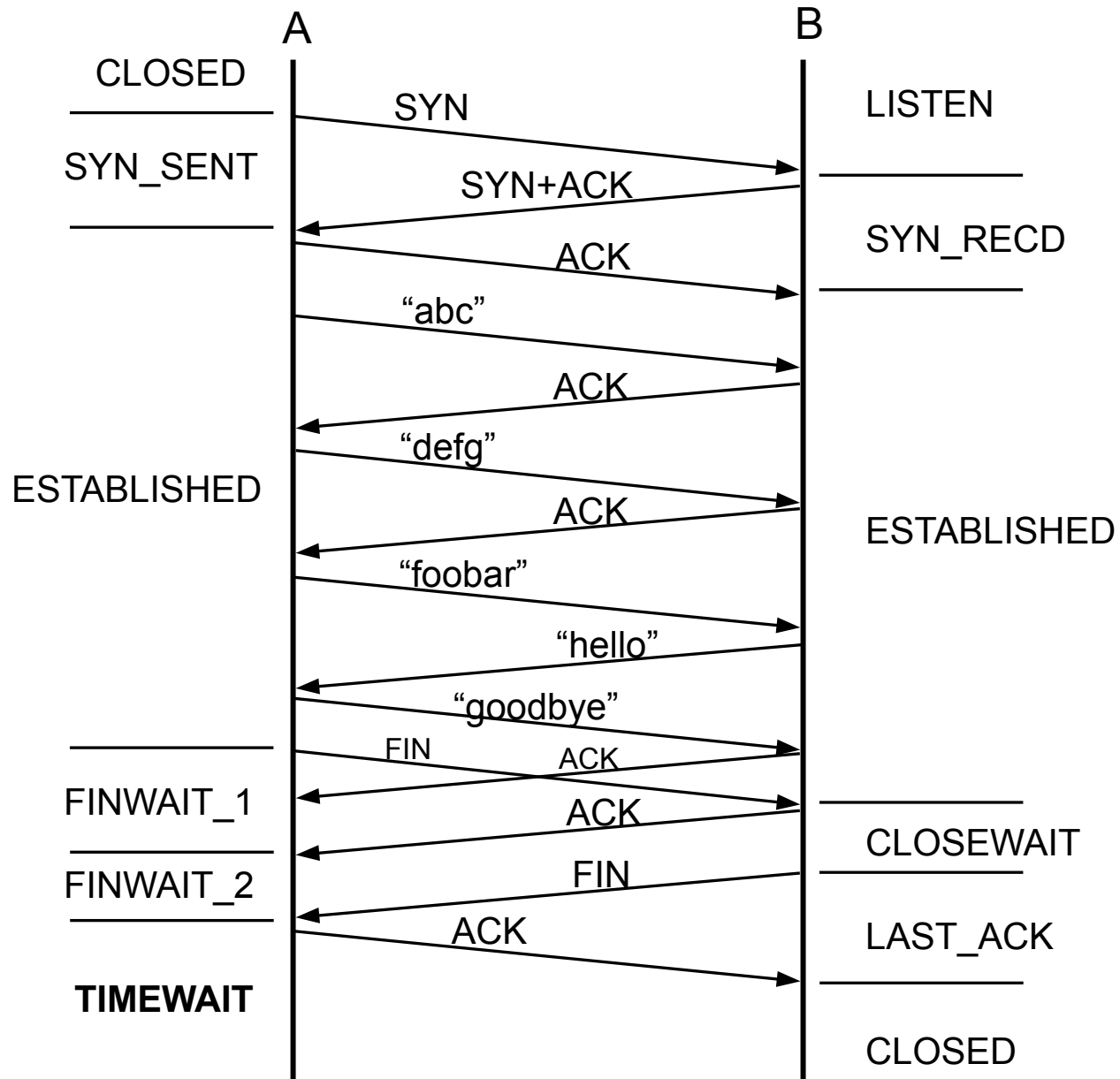
ESTABLISHED

ESTABLISHED

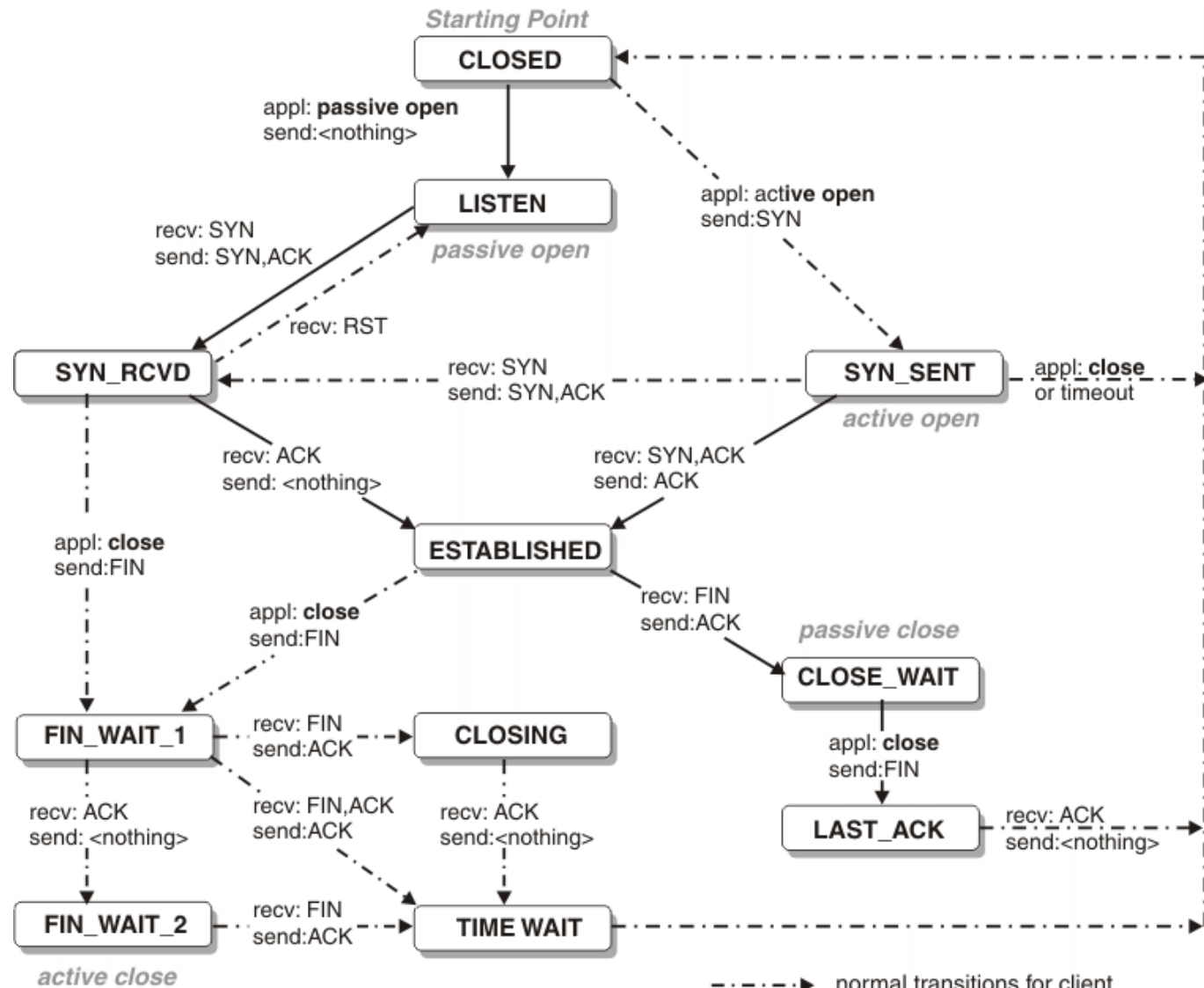
TCP



TCP



TCP



-----> normal transitions for client

————> normal transitions for server

appl: state transition taken when appl. issues operation

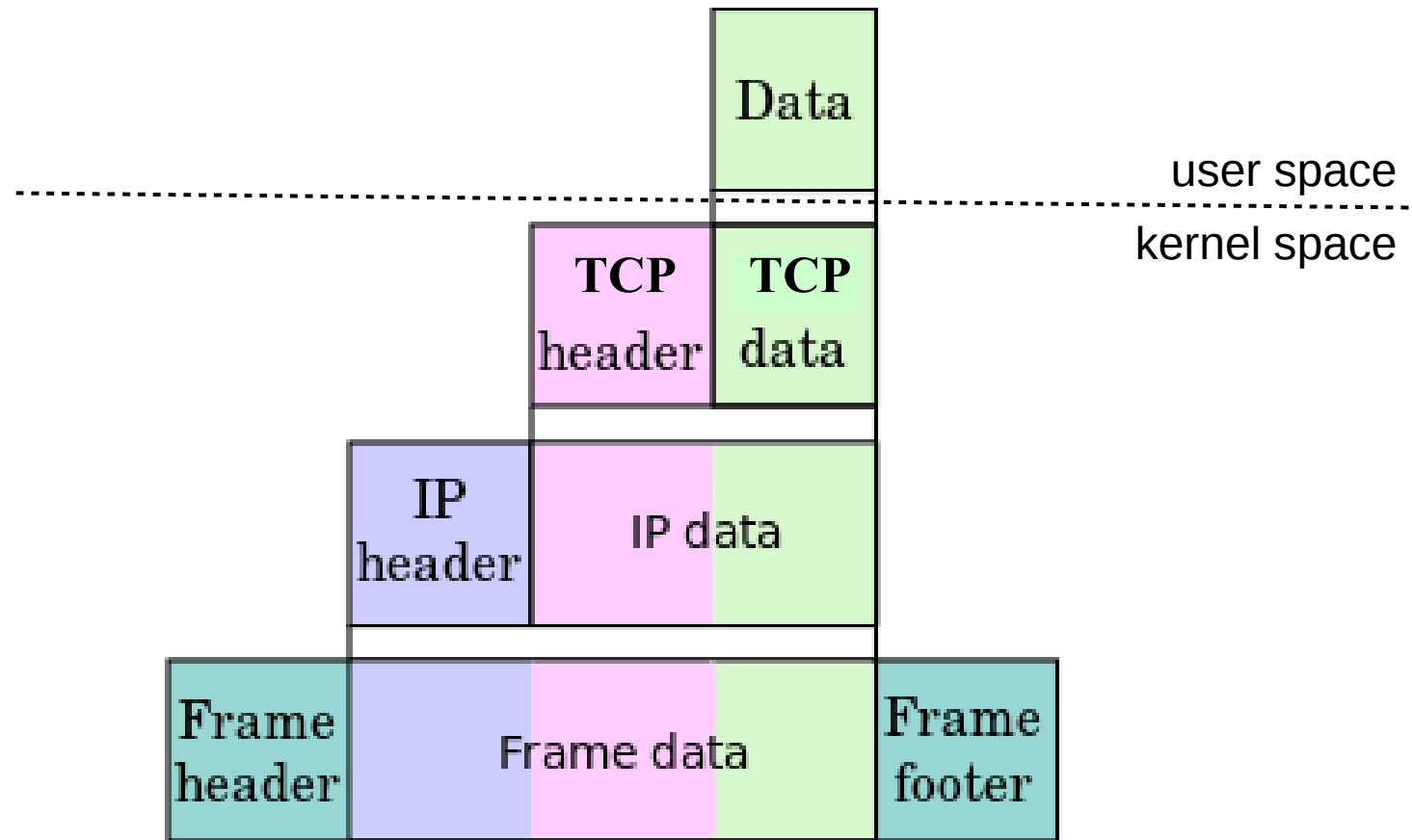
recv: state transition taken when segment is received

send: what is sent for this transition

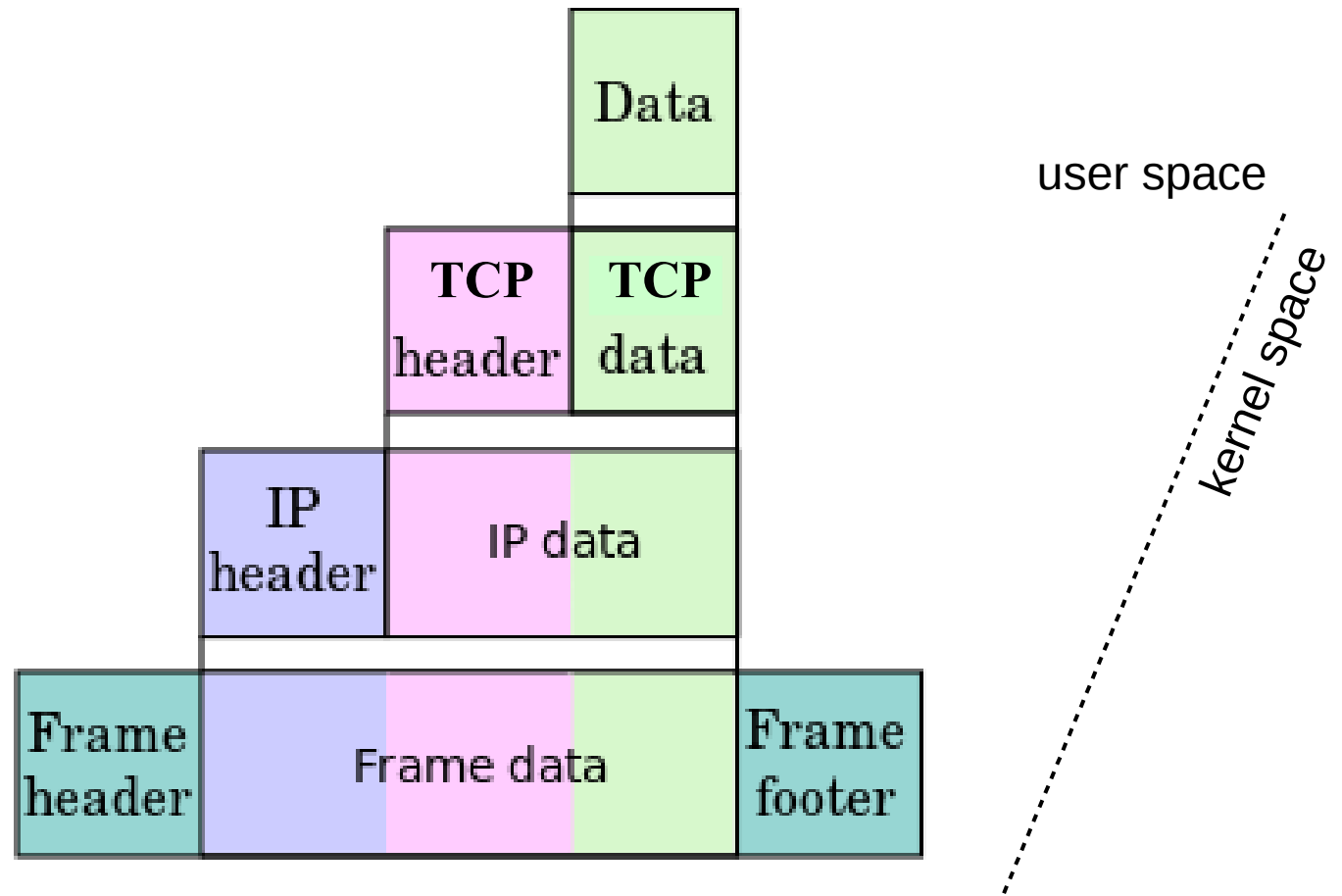
TCP vs UDP

TCP	UDP
Keeps track of lost packets. Makes sure that lost packets are re-sent	Doesn't keep track of lost packets
Adds sequence numbers to packets and reorders any packets that arrive in the wrong order	Doesn't care about packet arrival order
Slower, because of all added additional functionality	Faster, because it lacks any extra features
Requires more computer resources, because the OS needs to keep track of ongoing communication sessions and manage them on a much deeper level	Requires less computer resources
Examples of programs and services that use TCP: <ul style="list-style-type: none">- HTTP- HTTPS- FTP- Many computer games	Examples of programs and services that use UDP: <ul style="list-style-type: none">- DNS- IP telephony- DHCP- Many computer games

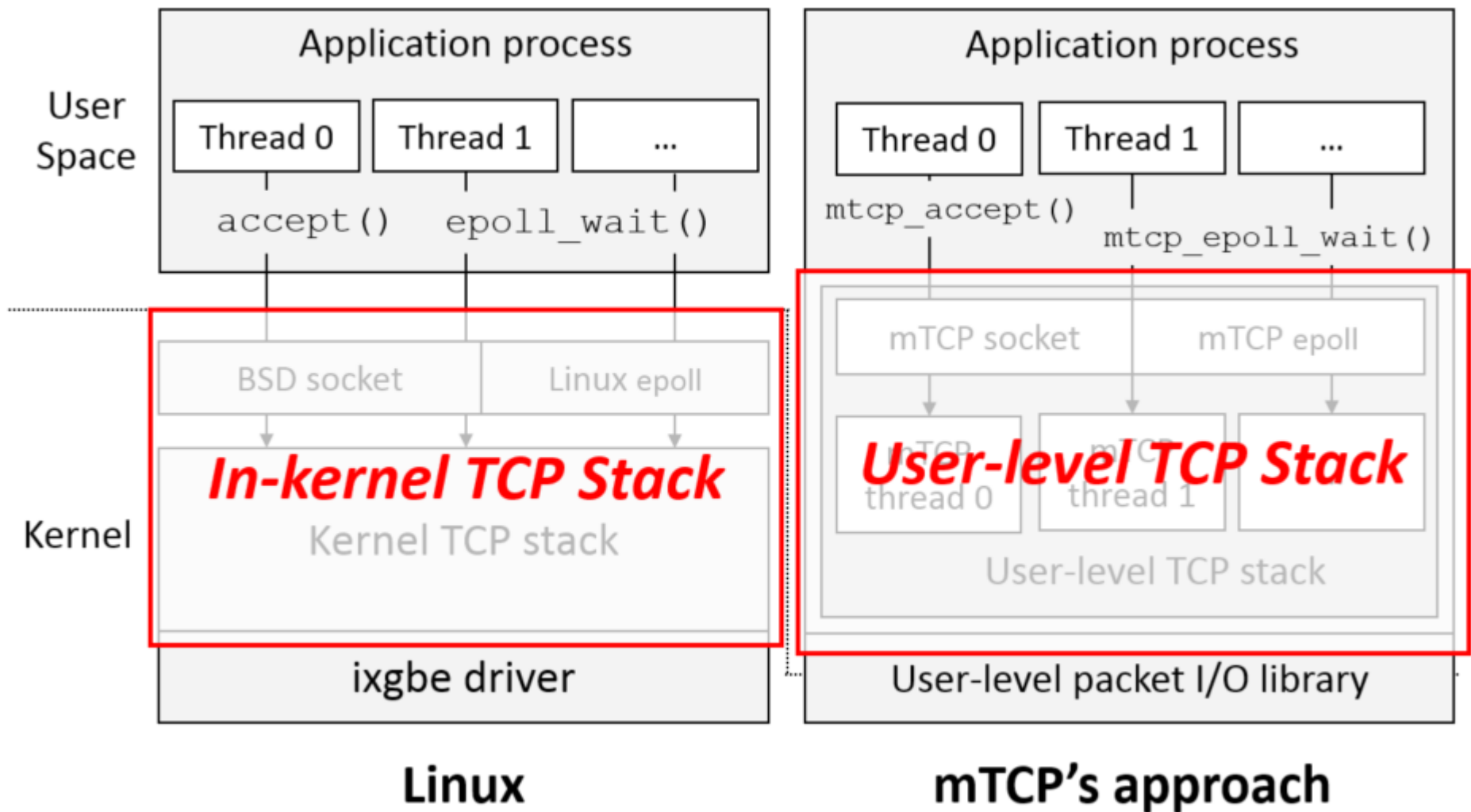
User TCP/IP



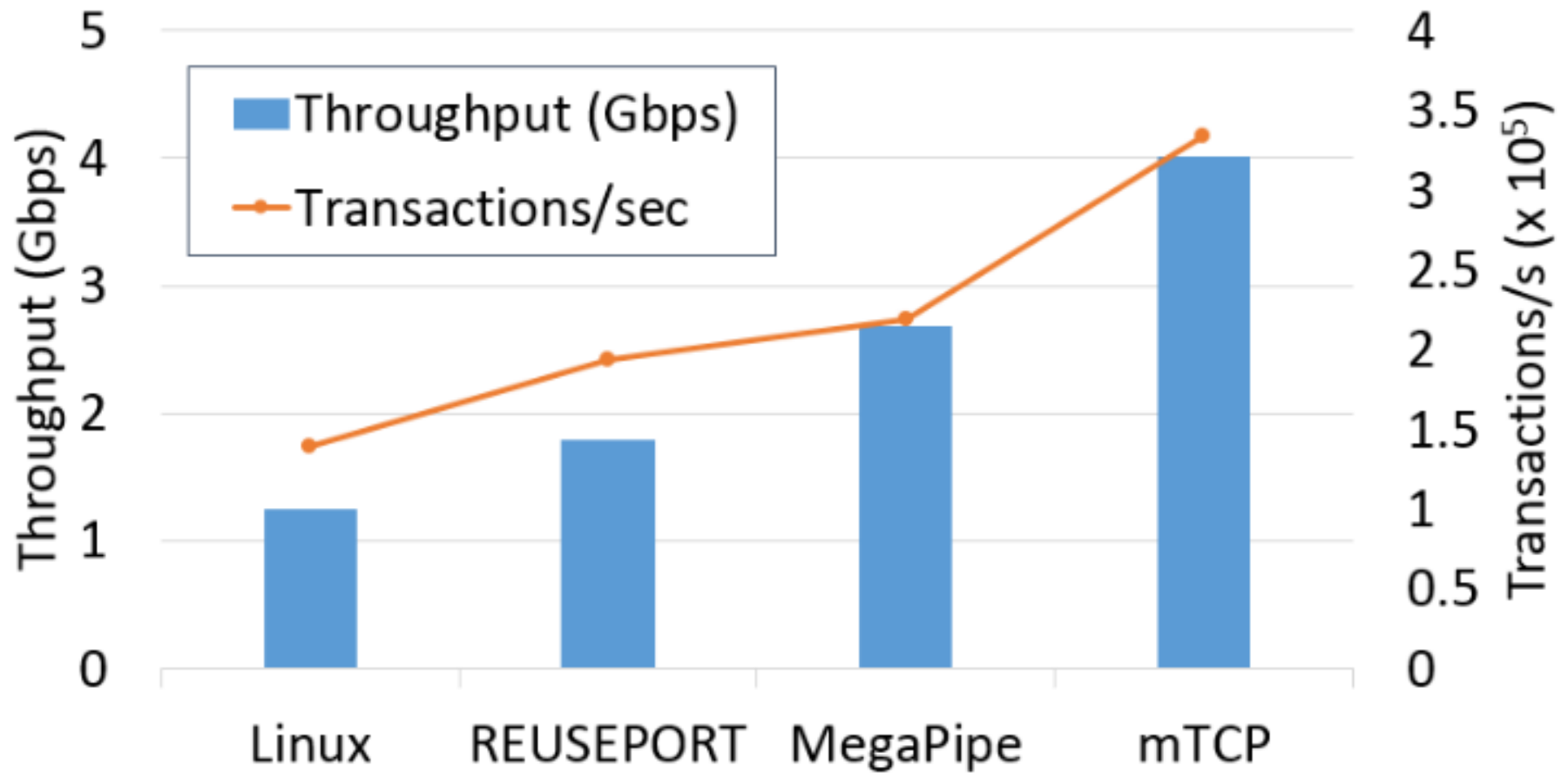
User TCP/IP



User TCP/IP

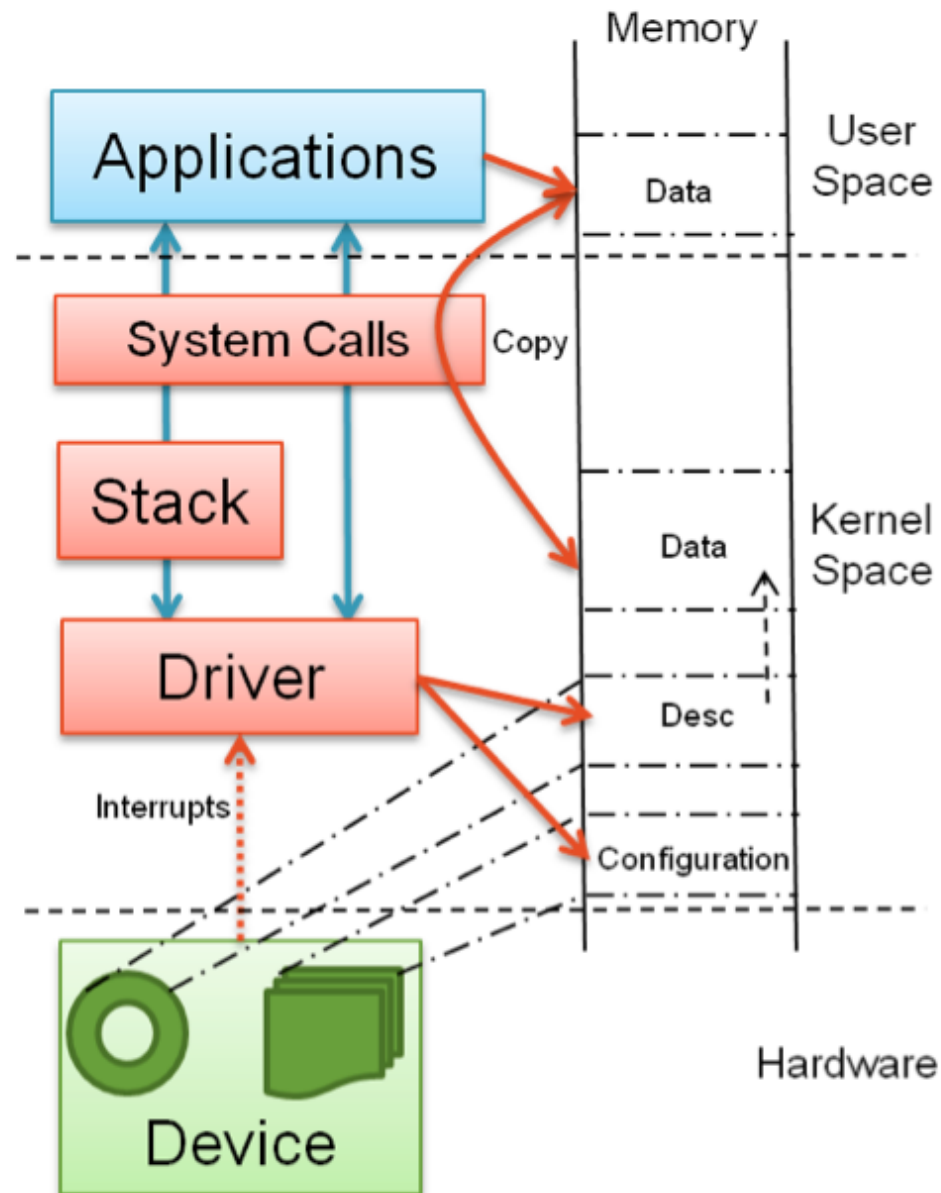


User TCP/IP



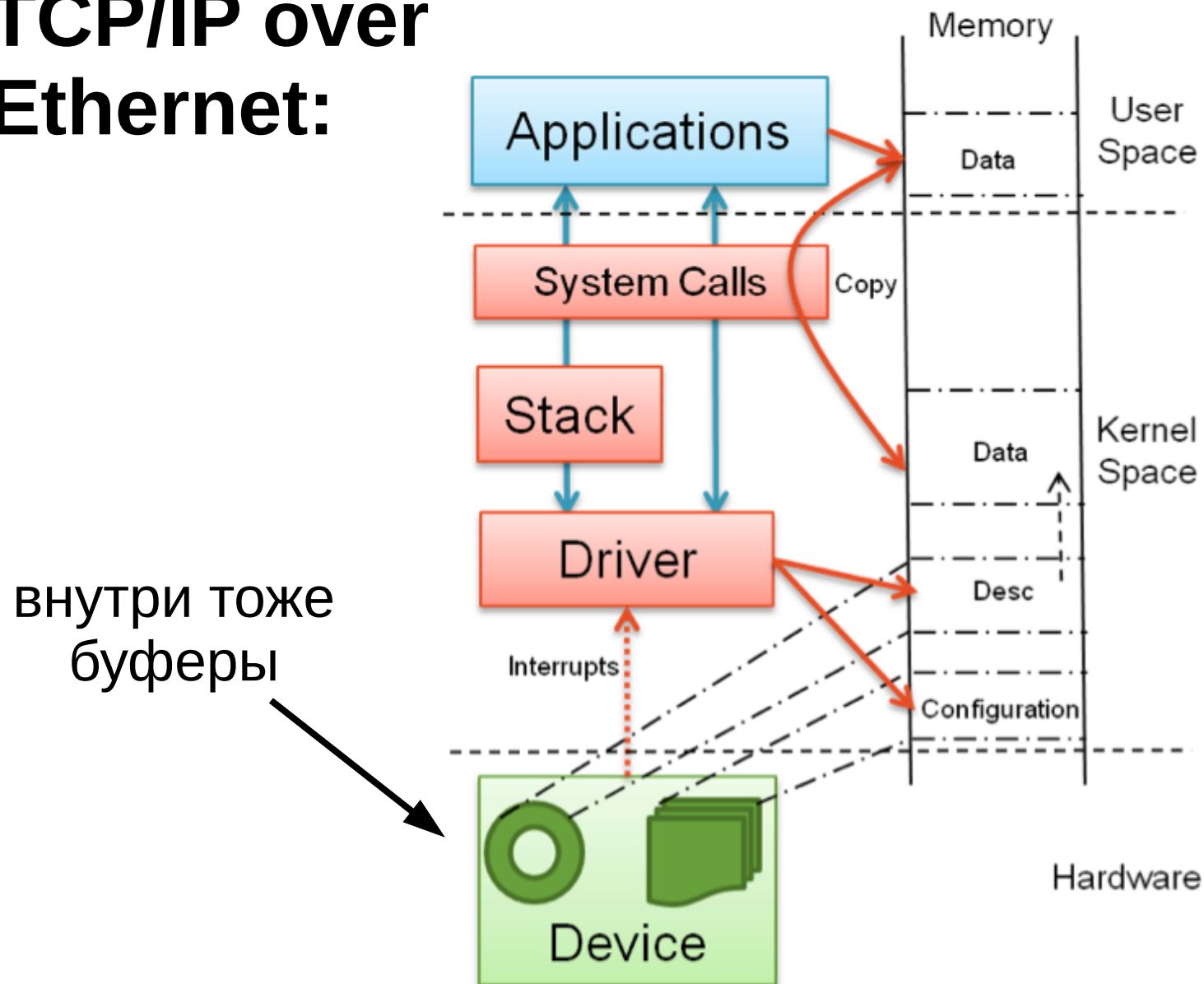
Infiniband & RDMA

TCP/IP over Ethernet:

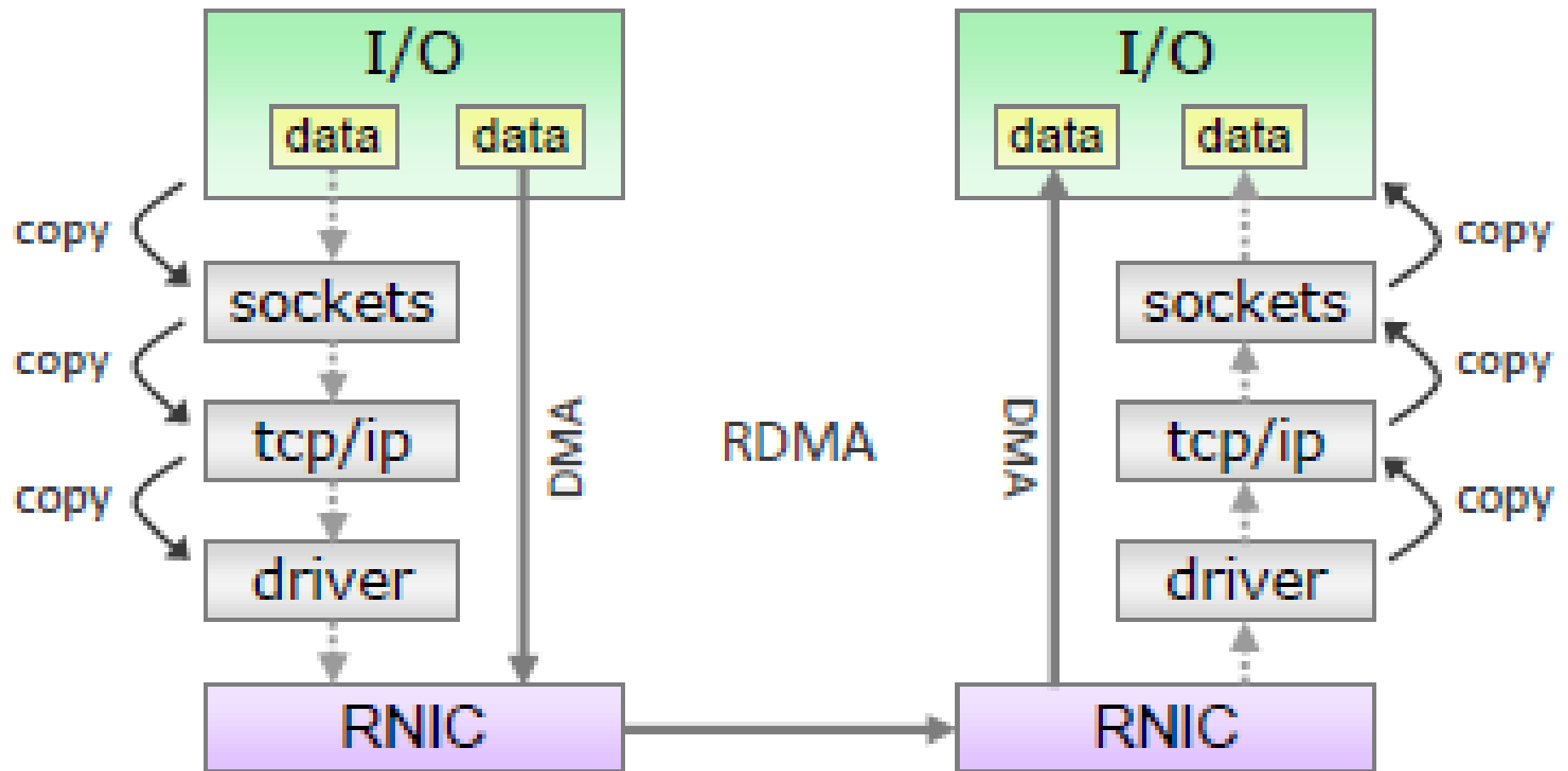


Infiniband & RDMA

TCP/IP over Ethernet:



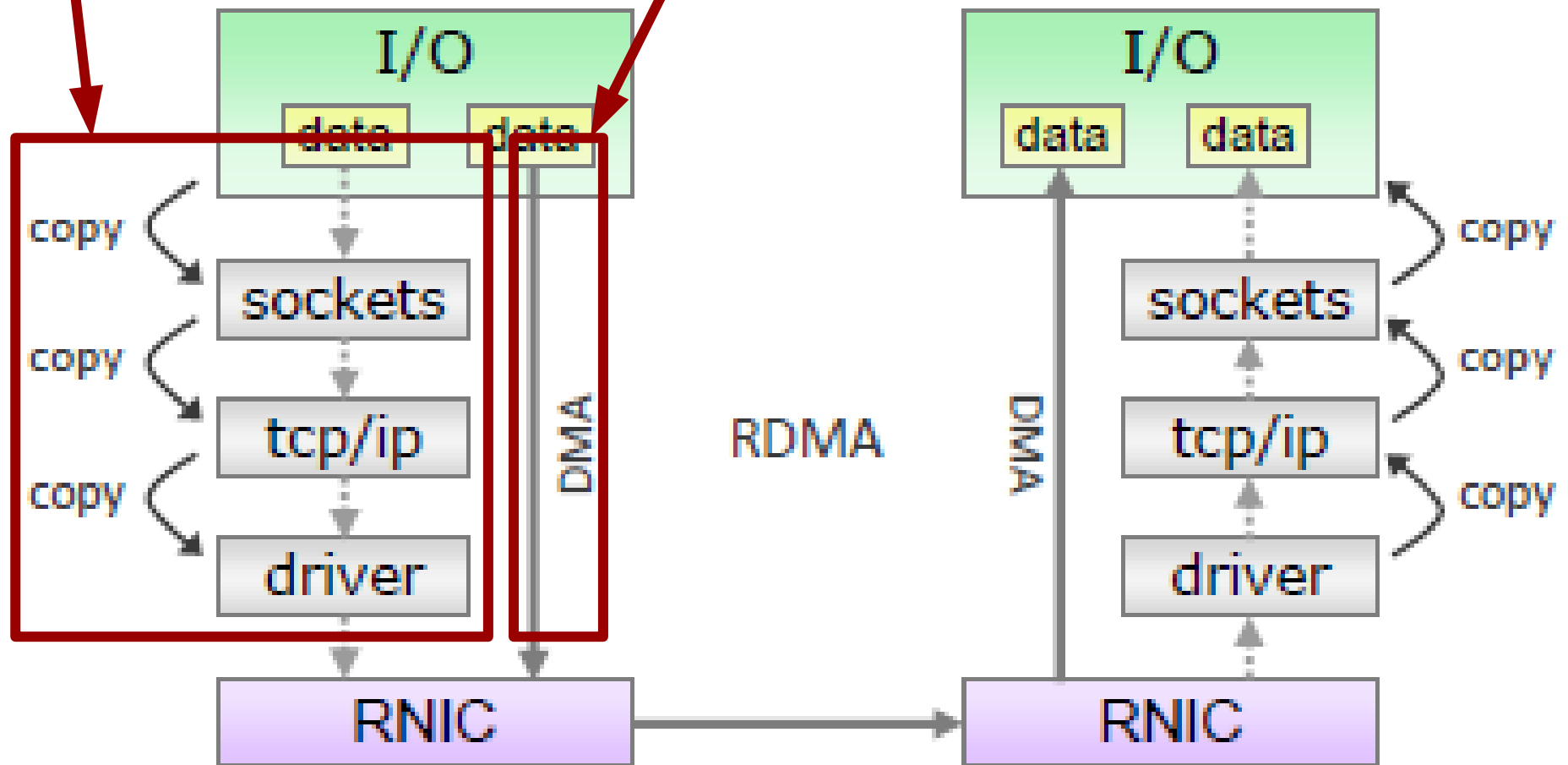
Infiniband & RDMA



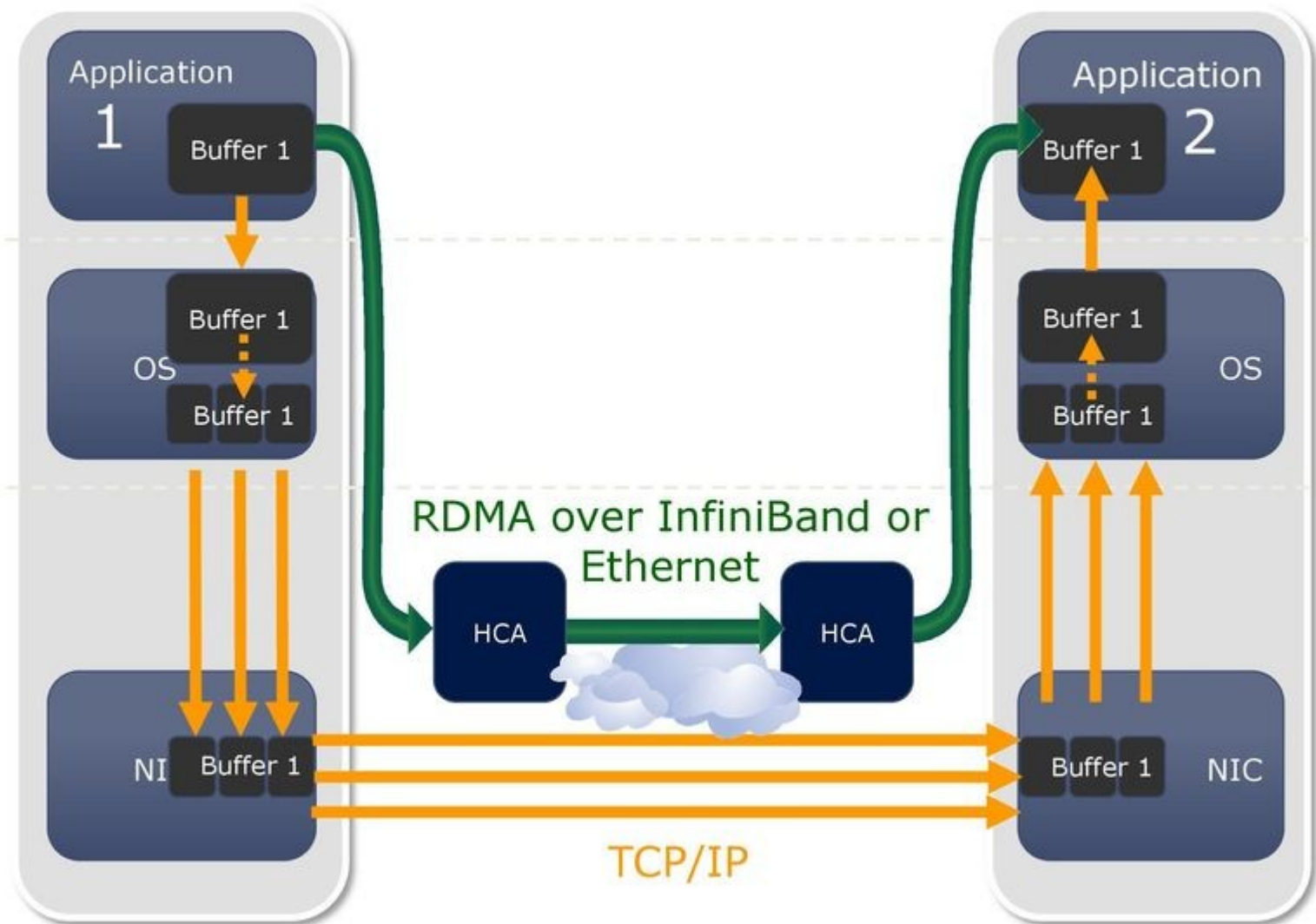
Infiniband & RDMA

TCP/IP

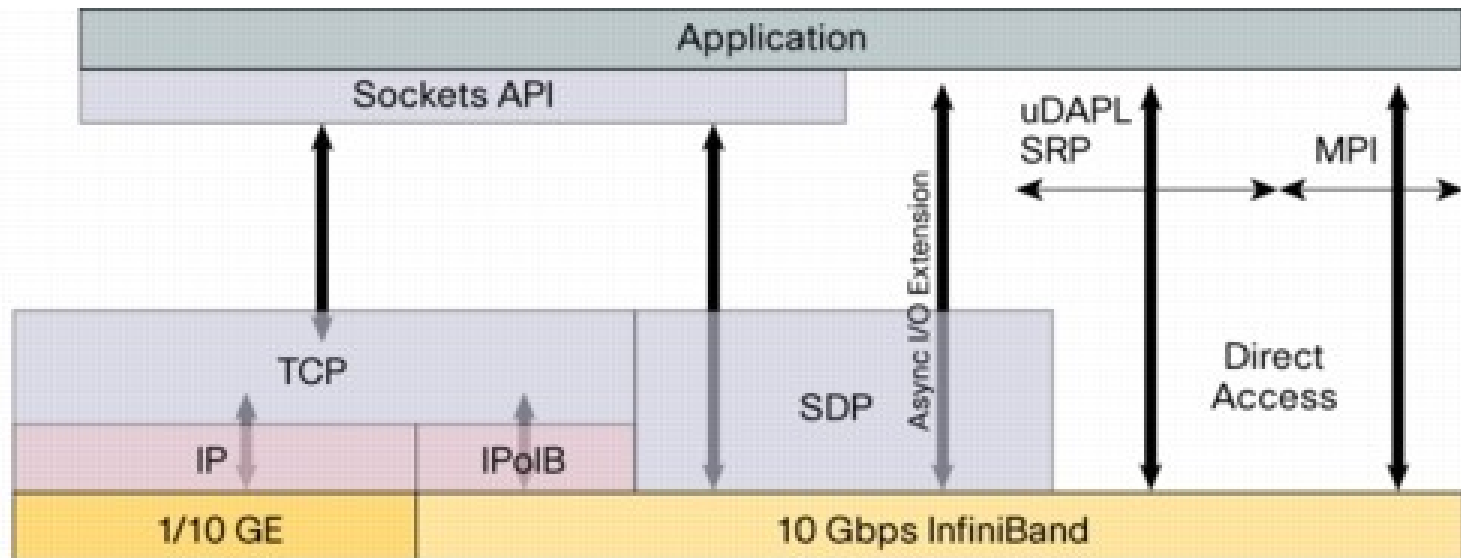
RDMA



Infiniband & RDMA



Infiniband & RDMA



Throughput	1 Gbps	3.9 Gbps	4.1 Gbps	4.5 Gbps	7.9 Gbps	8 Gbps	8 Gbs
Latency	40 usec	40 usec	20 usec	15 usec	15 usec	8 usec	4.5 usec
CPU Utilization	30%	80%	85%	25%	4%	1%	>1%

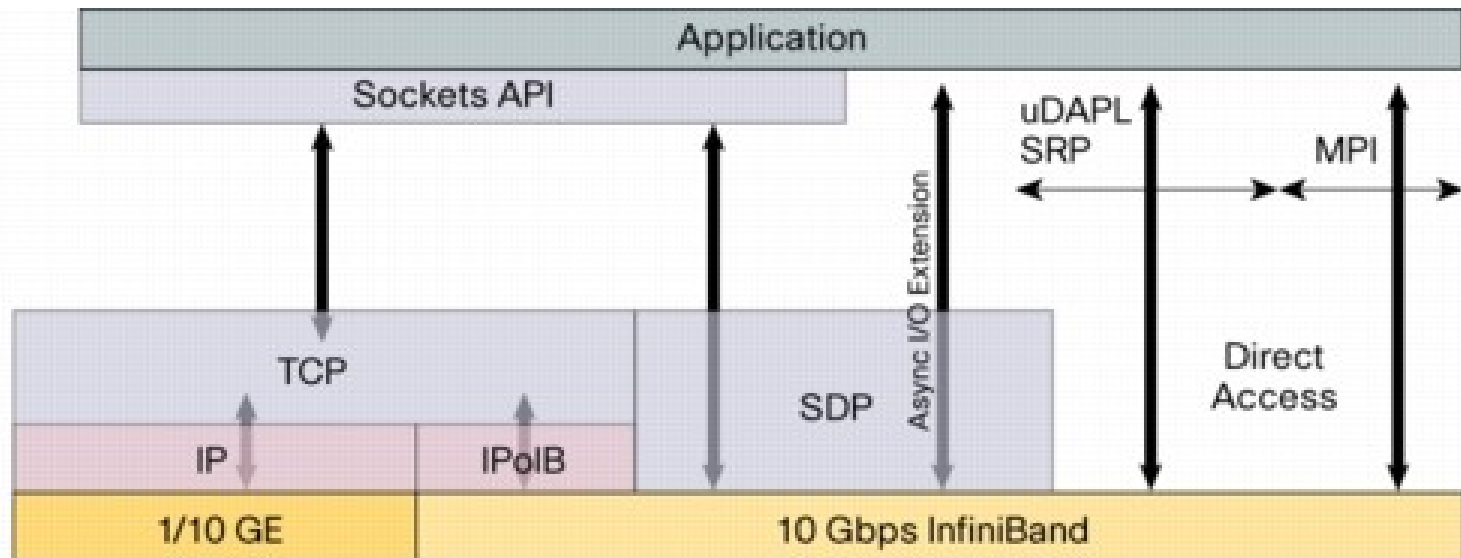
MPI: Message Passing Interface

SRP: SCSI Remote Protocol

uDAPL: User-level Direct Access Programming Language

* данные за 2006-ой год

Infiniband & RDMA



Throughput	1 Gbps	3.9 Gbps	4.1 Gbps	4.5 Gbps	7.9 Gbps	8 Gbps	8 Gbs
Latency	40 usec	40 usec	20 usec	15 usec	15 usec	8 usec	4.5 usec
CPU Utilization	30%	80%	85%	25%	4%	1%	> 1%

MPI: Message Passing Interface

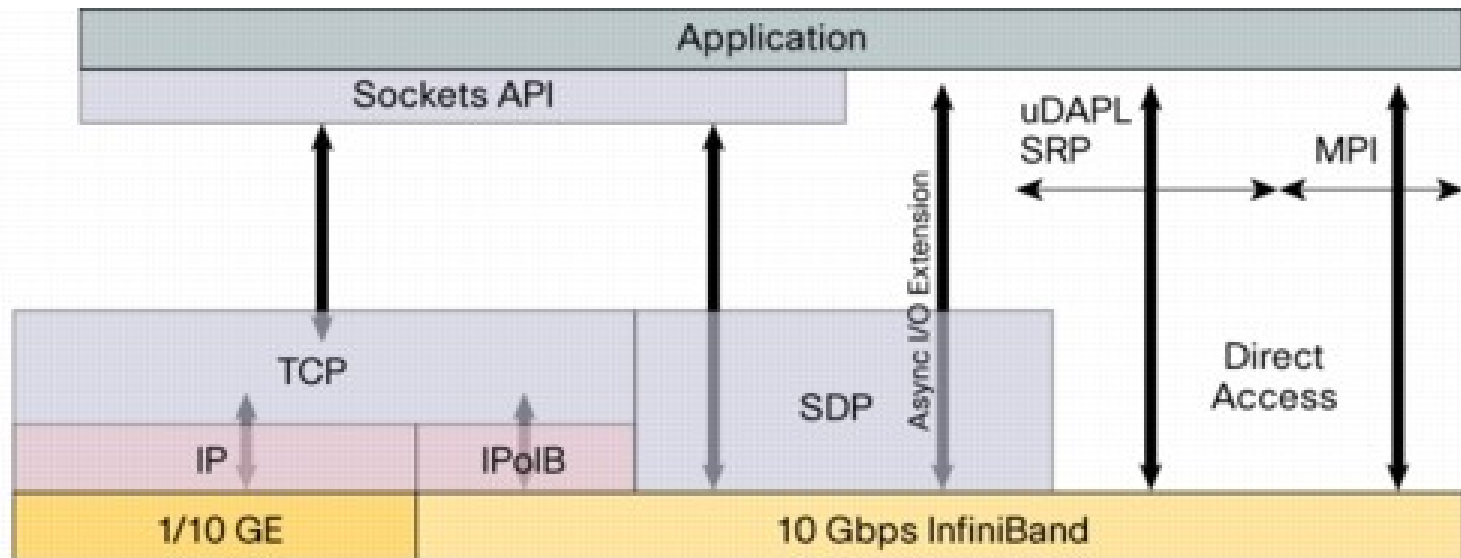
SRP: SCSI Remote Protocol

uDAPL: User-level Direct Access Programming Language

FDR: 700ns

EDR: 500ns

Infiniband & RDMA



Throughput	1 Gbps	3.9 Gbps	4.1 Gbps	4.5 Gbps	7.9 Gbps	8 Gbps	8 Gbs
Latency	40 usec	40 usec	20 usec	15 usec	15 usec	8 usec	4.5 usec
CPU Utilization	30%	80%	85%	25%	4%	1%	>1%

MPI: Message Passing Interface

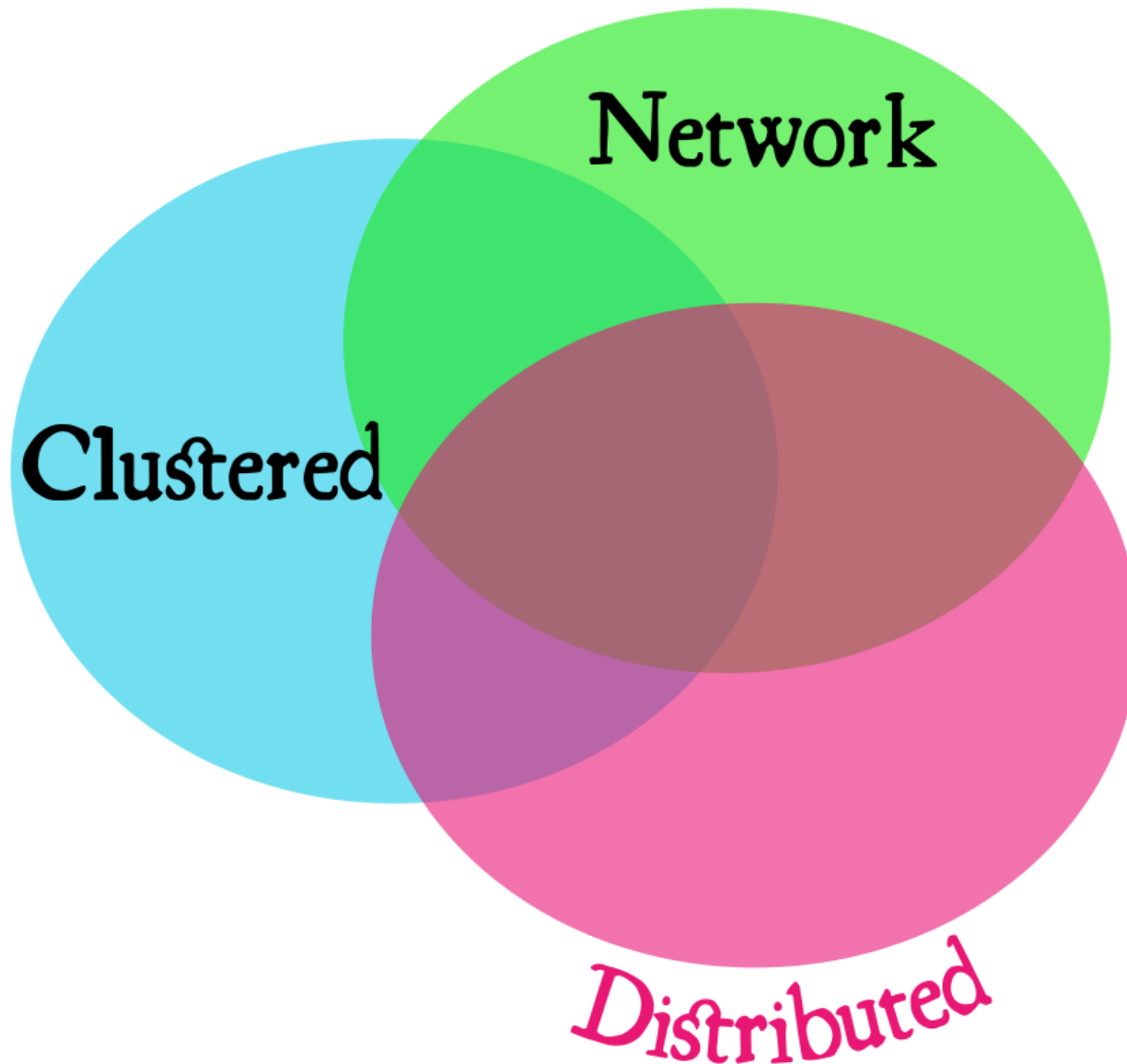
SRP: SCSI Remote Protocol

uDAPL: User-level Direct Access Programming Language

Latencies

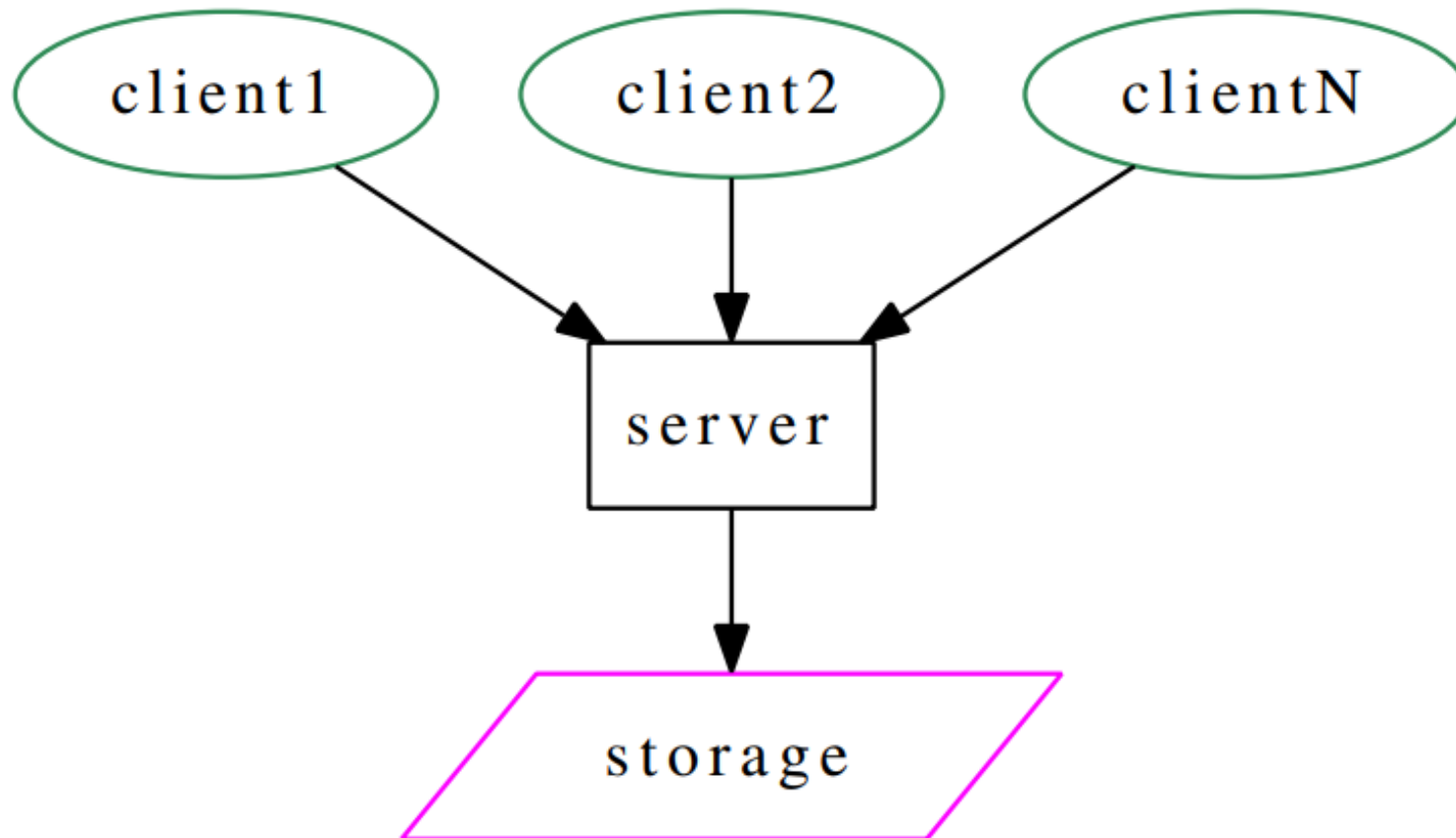
Event	Latency	Scaled
1 CPU cycle	0.3 ns	1 s
Level 1 cache access	0.9 ns	3 s
Level 2 cache access	2.8 ns	9 s
Level 3 cache access	12.9 ns	43 s
Main memory access (DRAM, from CPU)	120 ns	6 min
Solid-state disk I/O (flash memory)	50–150 µs	2–6 days
Rotational disk I/O	1–10 ms	1–12 months
Internet: San Francisco to New York	40 ms	4 years
Internet: San Francisco to United Kingdom	81 ms	8 years
Internet: San Francisco to Australia	183 ms	19 years
TCP packet retransmit	1–3 s	105–317 years
OS virtualization system reboot	4 s	423 years
SCSI command time-out	30 s	3 millennia
Hardware (HW) virtualization system reboot	40 s	4 millennia
Physical system reboot	5 m	32 millennia

File systems in HPC



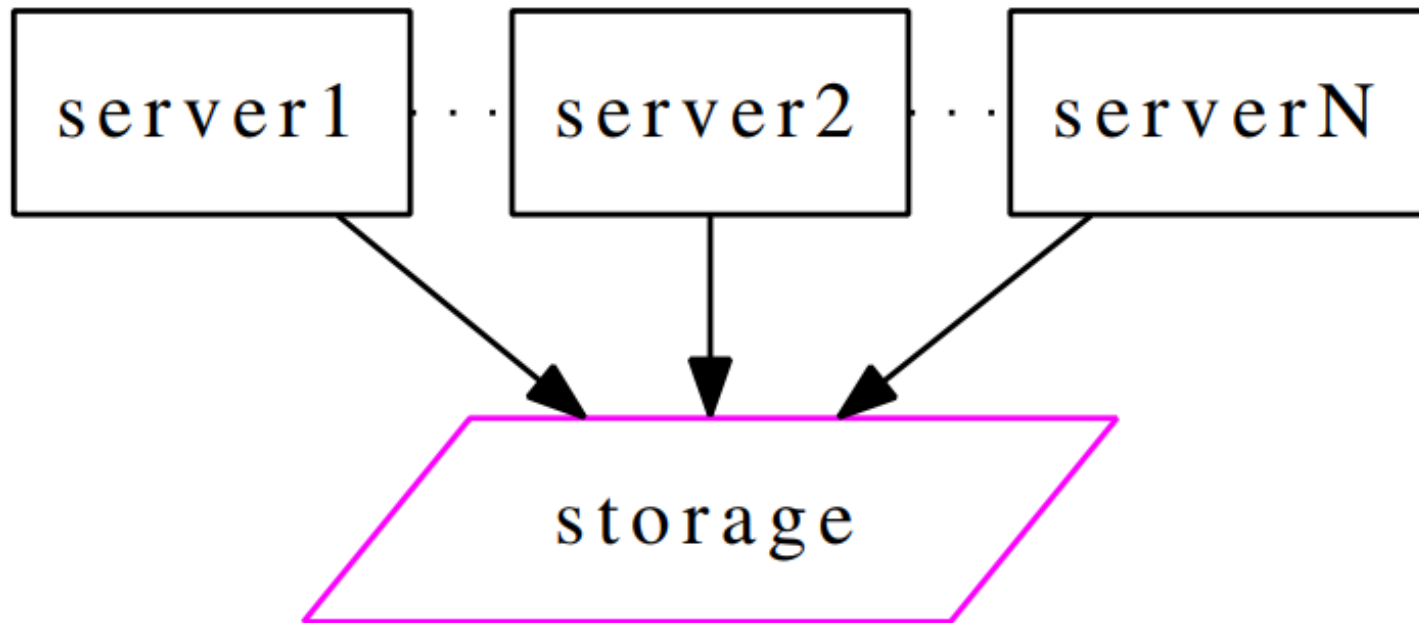
File systems in HPC

Network file system:



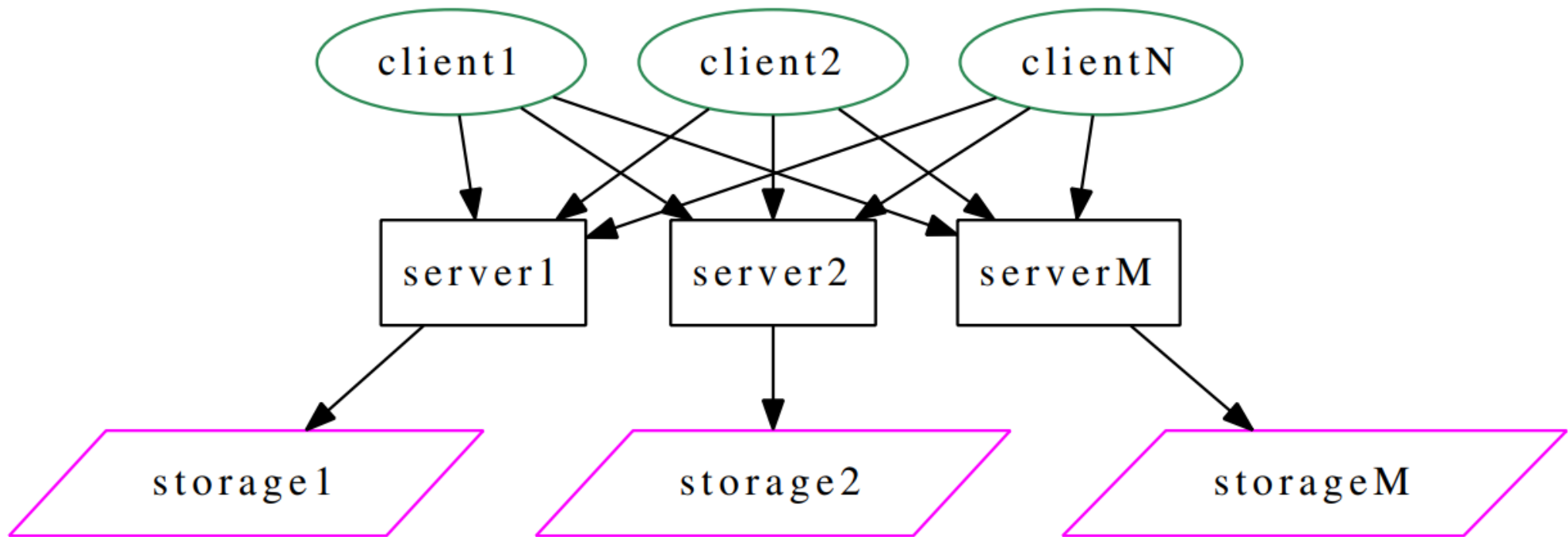
File systems in HPC

Clustered file system:

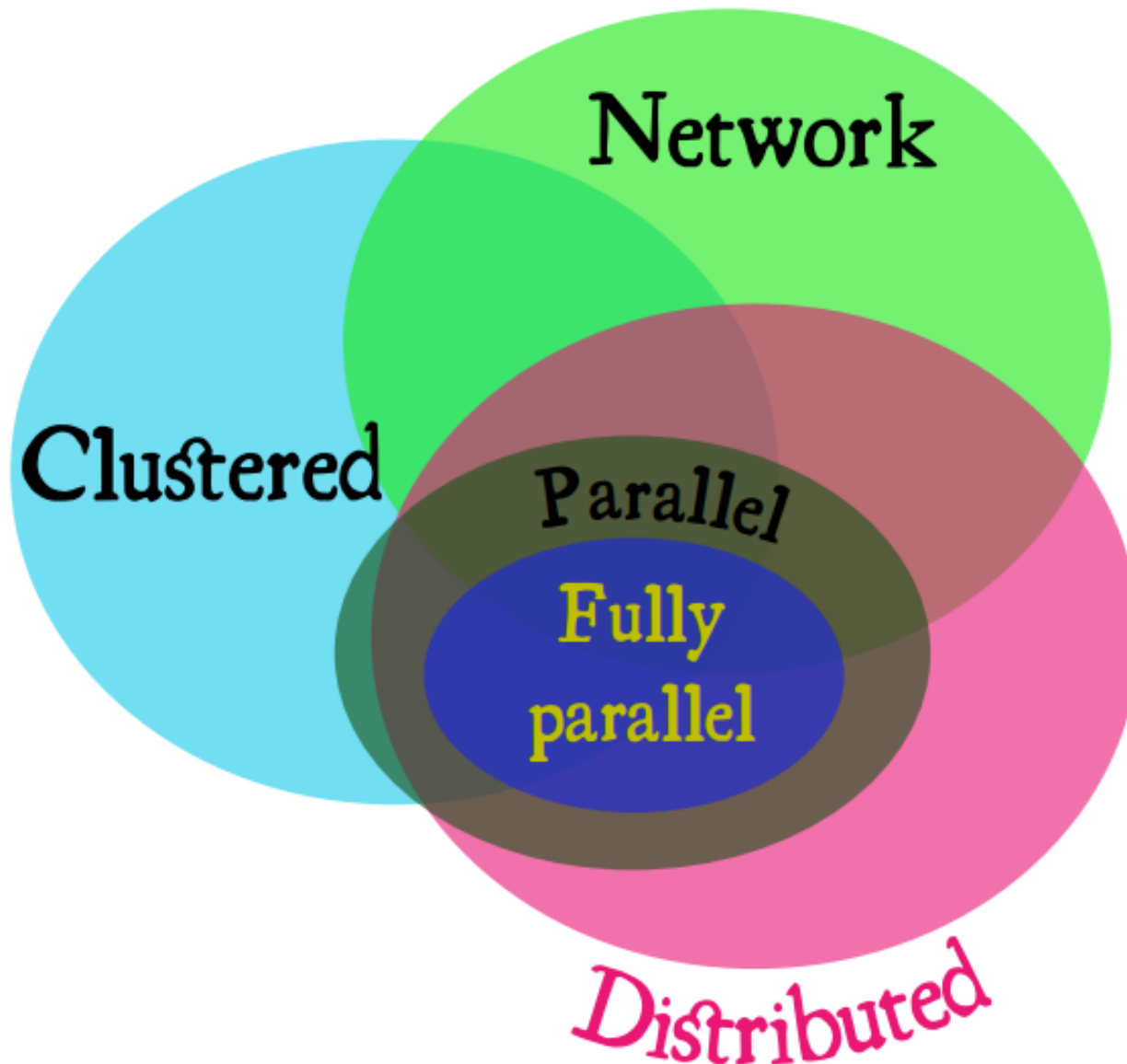


File systems in HPC

Distributed file system:



File systems in HPC



File systems in HPC

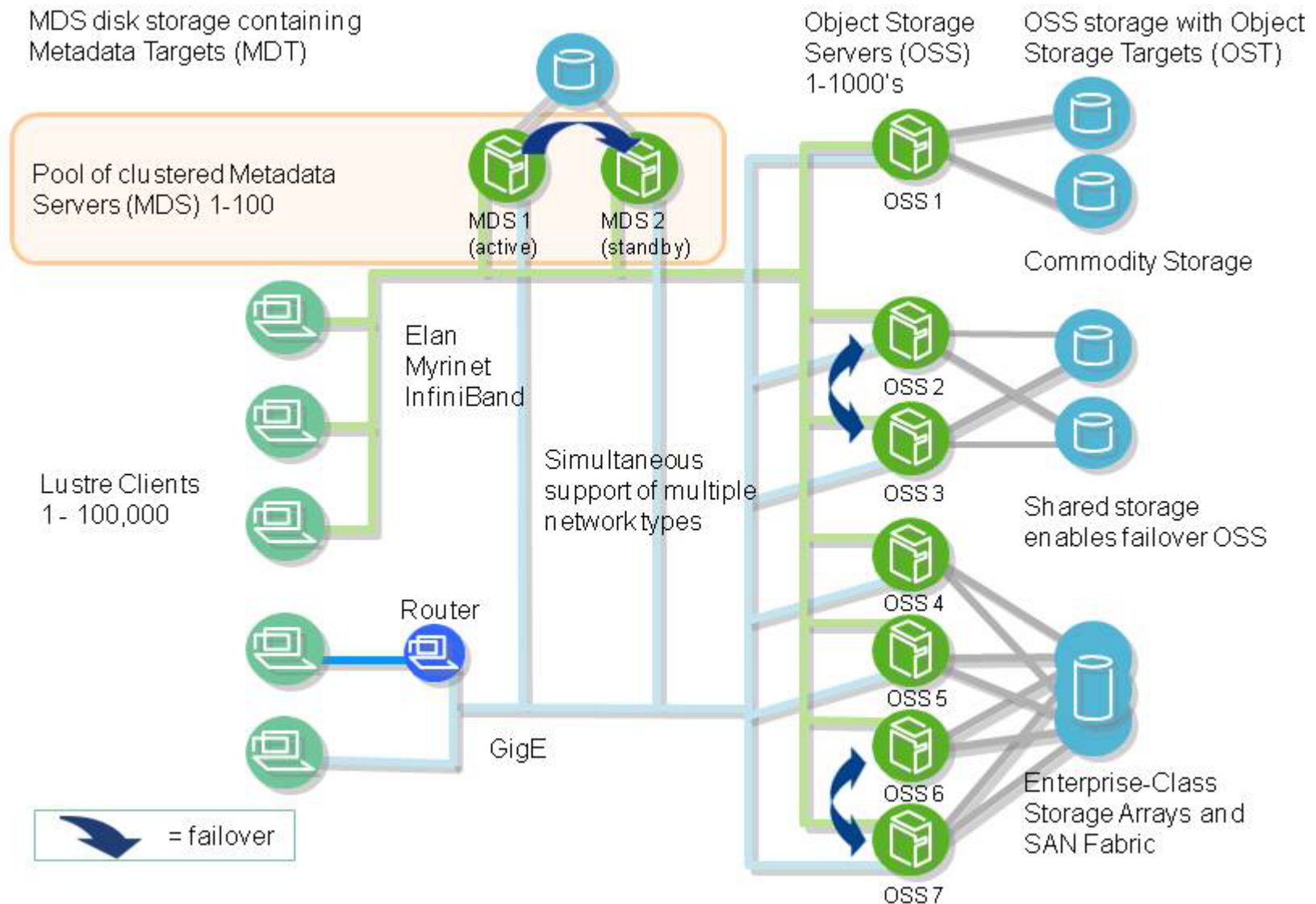
Parallel file system:

- Параллельный доступ ко всем серверам
- Параллельный доступ к данным
- Распределение нагрузки

Fully parallel file system:

+ Параллельный доступ к метаданным

File systems in HPC



Q&A