

Distributed Systems – Efficient Communication Software

Today

✓ Distributed System Basics

✓ Lamport's clock

⇒ Efficient Communication Software

Therath + Levy { * Application interface to the Kernel
* Inside the Kernel

Wetterwald { * End-to-End QoS via Active Networks

Friday

Synthesizing Network Protocol Stacks

Latency vs. Throughput

Latency

Elapsed time

Throughput

Events per unit time

Bandwidth: throughput measure

Latency vs. Throughput

Latency

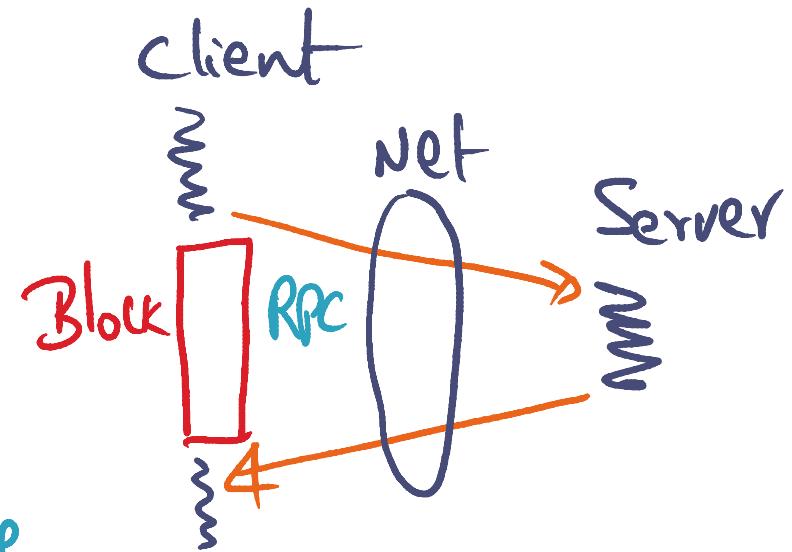
elapsed time

Throughput

Events per unit time

Bandwidth: throughput measure

RPC performance



Latency vs. Throughput

Latency

Elapsed time

Throughput

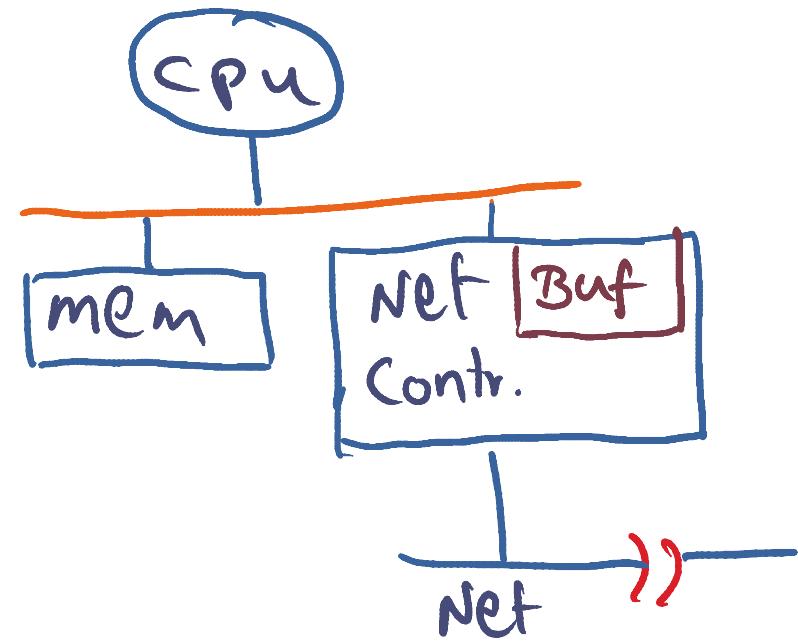
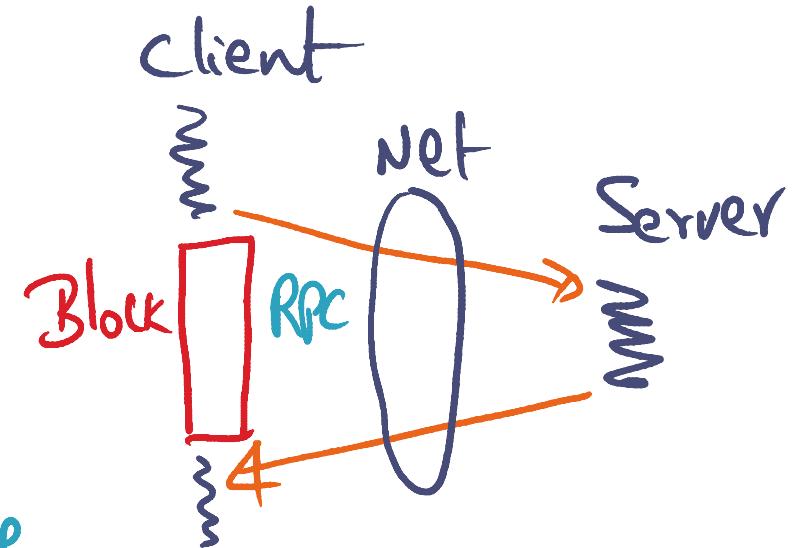
Events per unit time

Bandwidth: throughput measure

RPC performance

Hardware overhead

Software overhead



Latency vs. Throughput

Latency

Elapsed time

Throughput

Events per unit time

Bandwidth: throughput measure

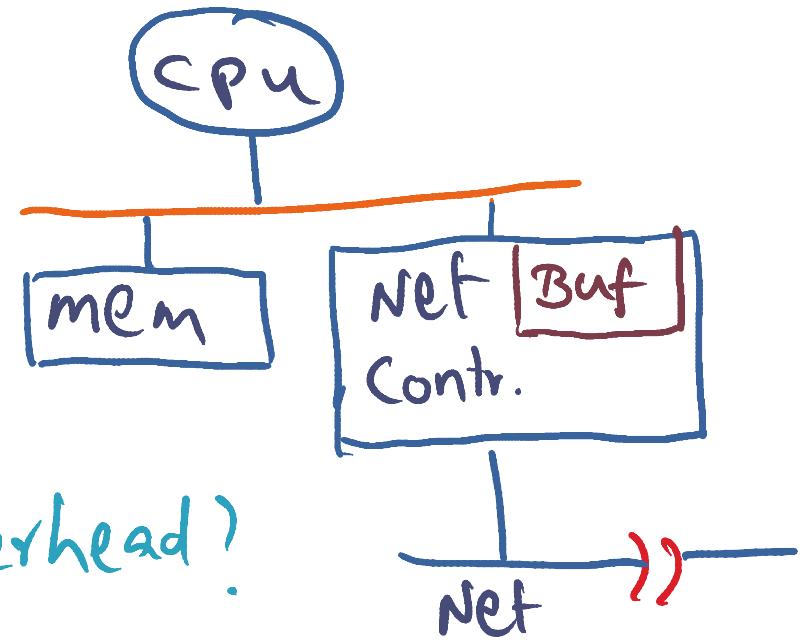
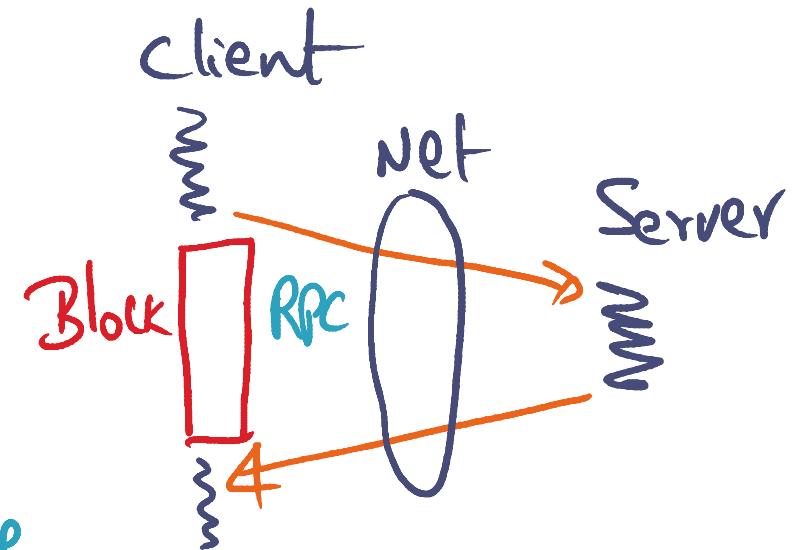
RPC performance

Hardware overhead

Software overhead

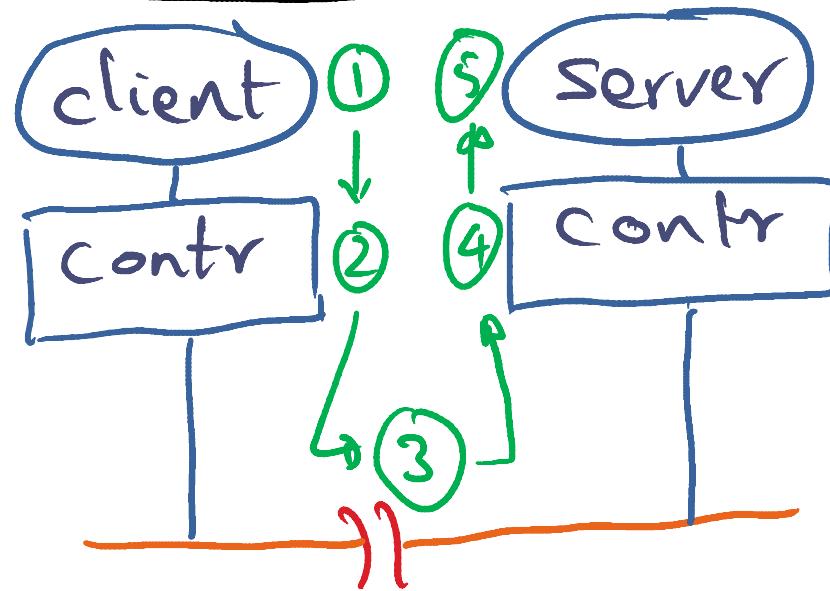
Focus

How to reduce software overhead?



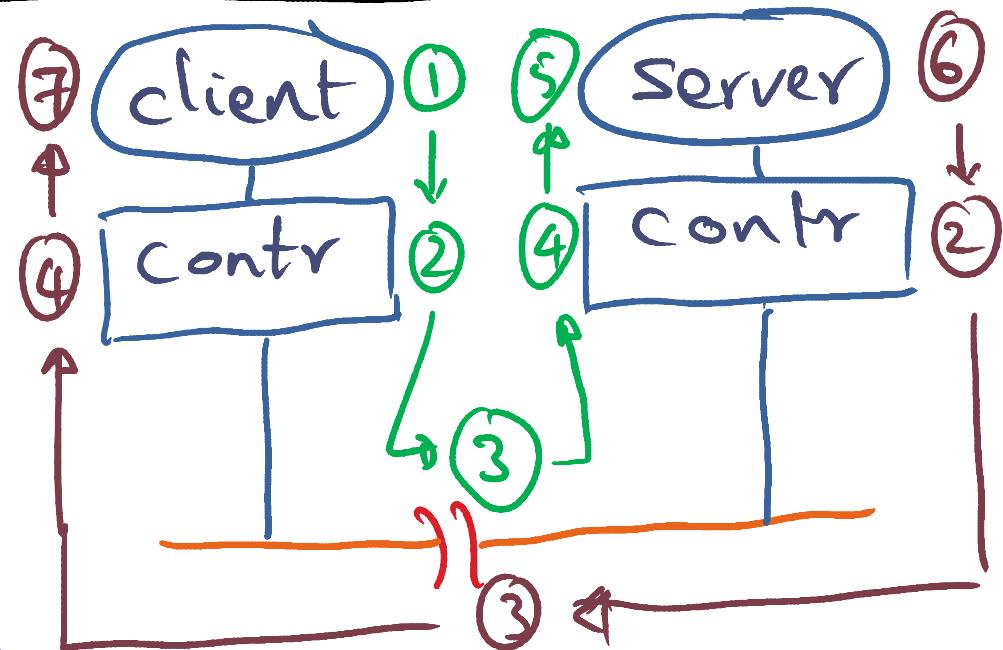
Components of RPC Latency

1. client call
2. Controller latency
3. Time on wire
4. Interrupt handling
5. server setup to execute call



Components of RPC Latency

1. client call
2. Controller latency
3. Time on wire
4. Interrupt handling
5. server setup to execute call
6. server execution + reply
7. client setup to receive results & restart



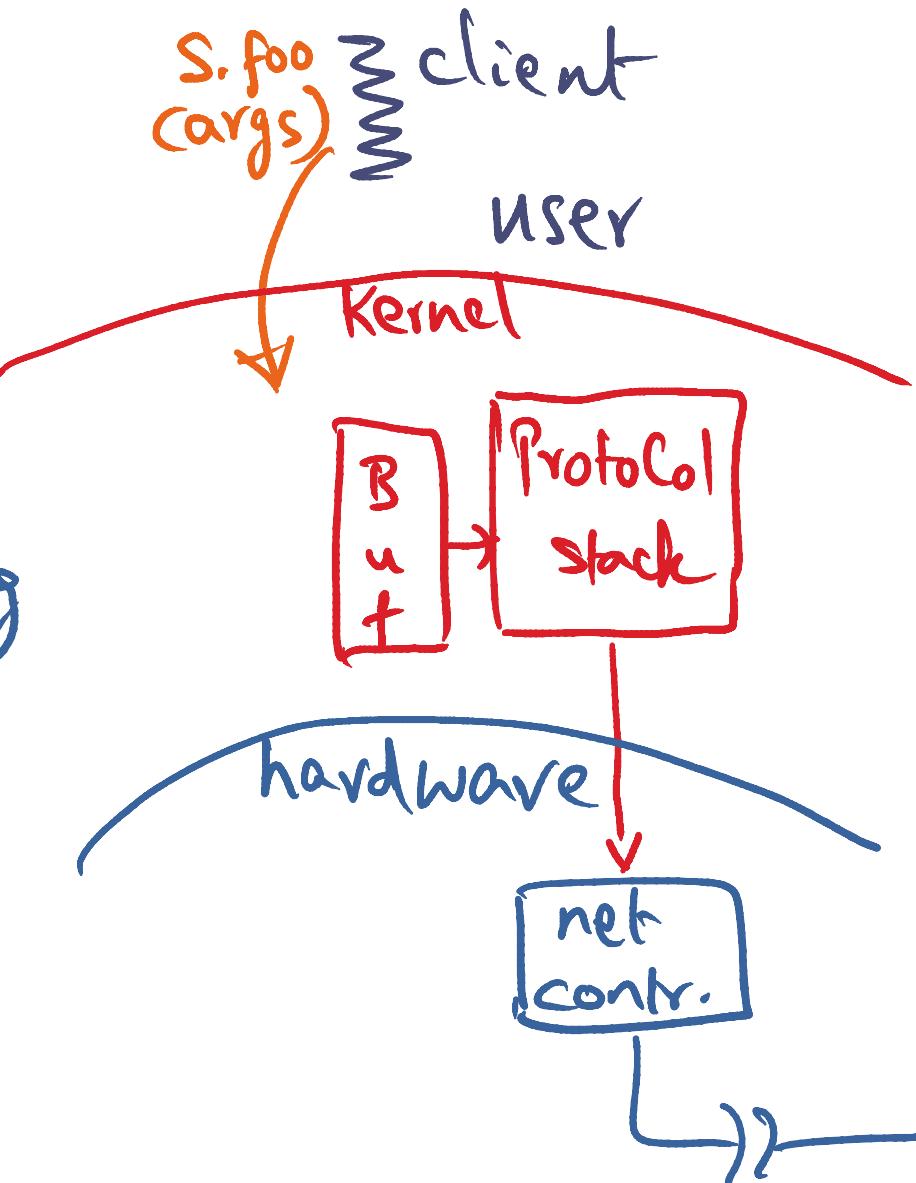
Sources of overhead in RPC

marshaling

Data copying

Control transfer

Protocol processing



Sources of overhead in RPC

marshaling

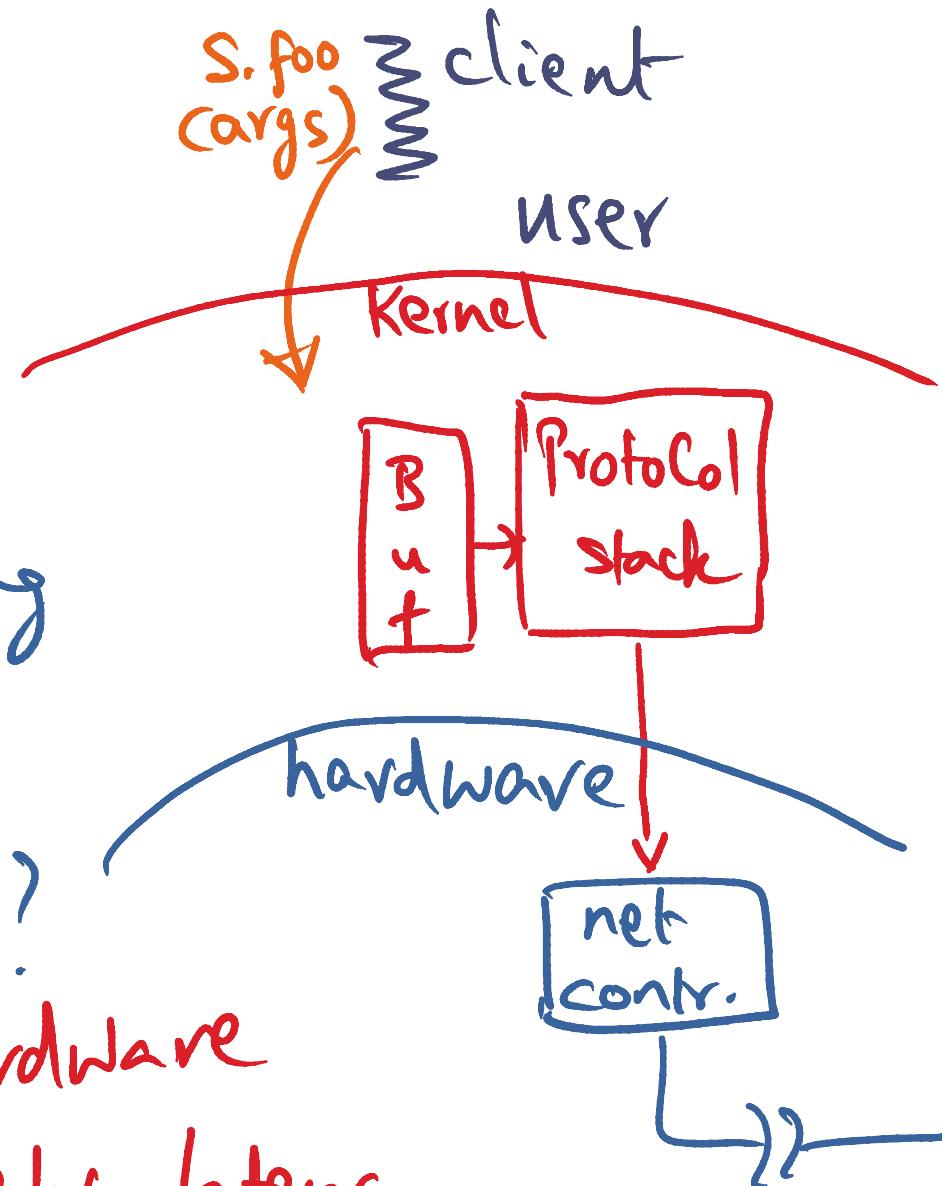
Data copying

Control transfer

Protocol processing

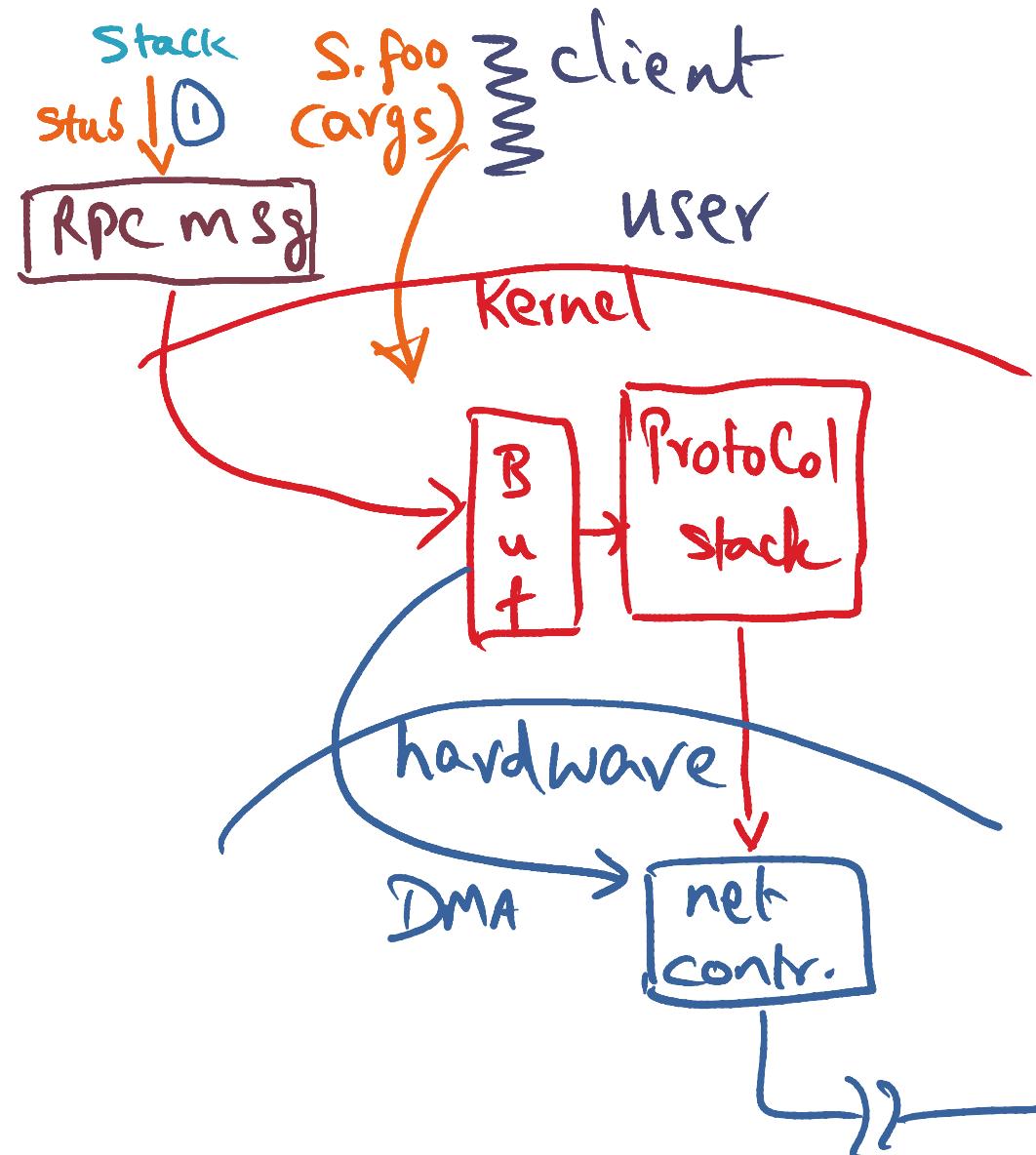
How to reduce
Kernel overhead?

- Take what hardware
gives you to reduce latency



marshaling and data Copying

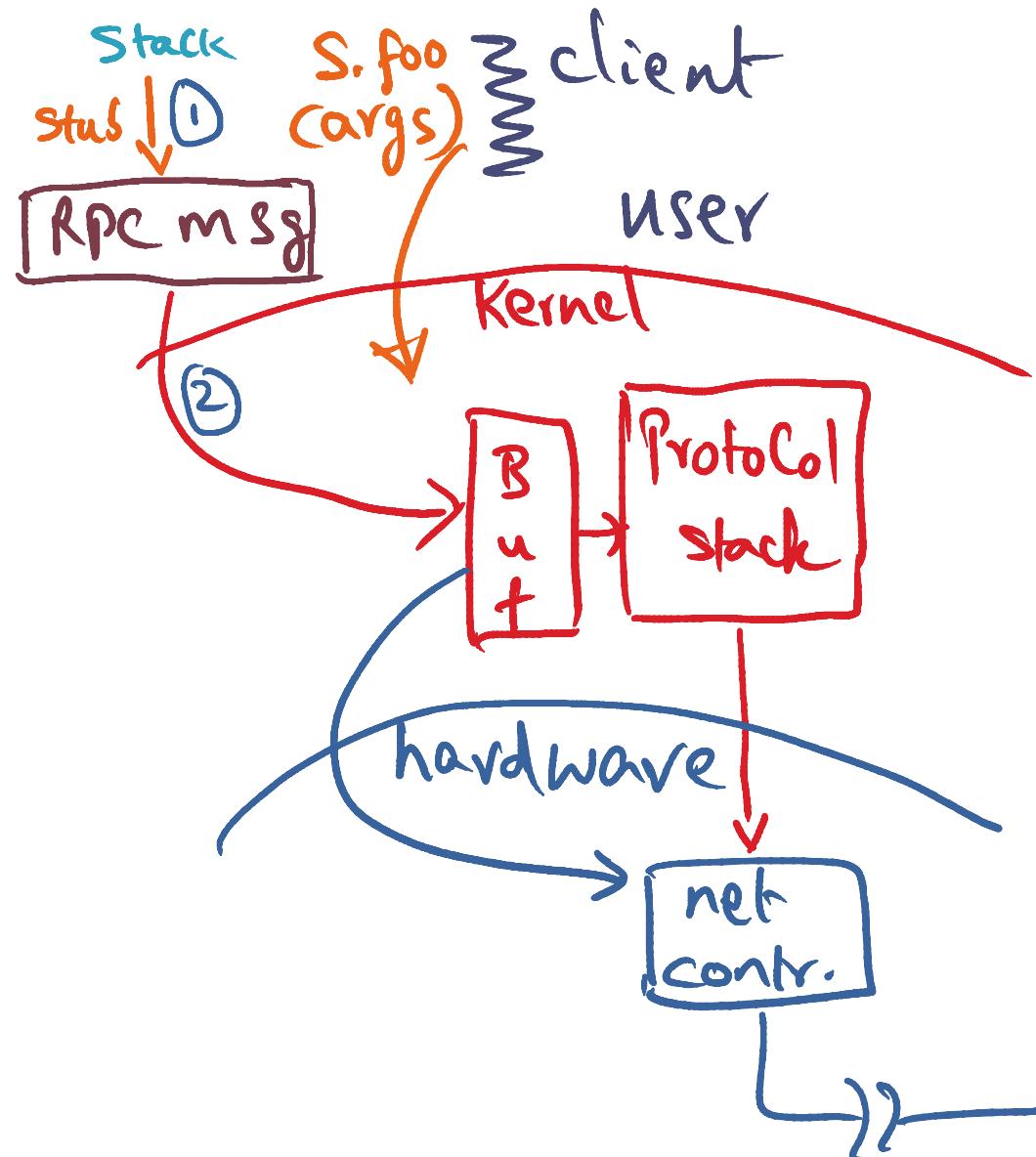
Three Copies
- client stubs



marshaling and data Copying

Three Copies

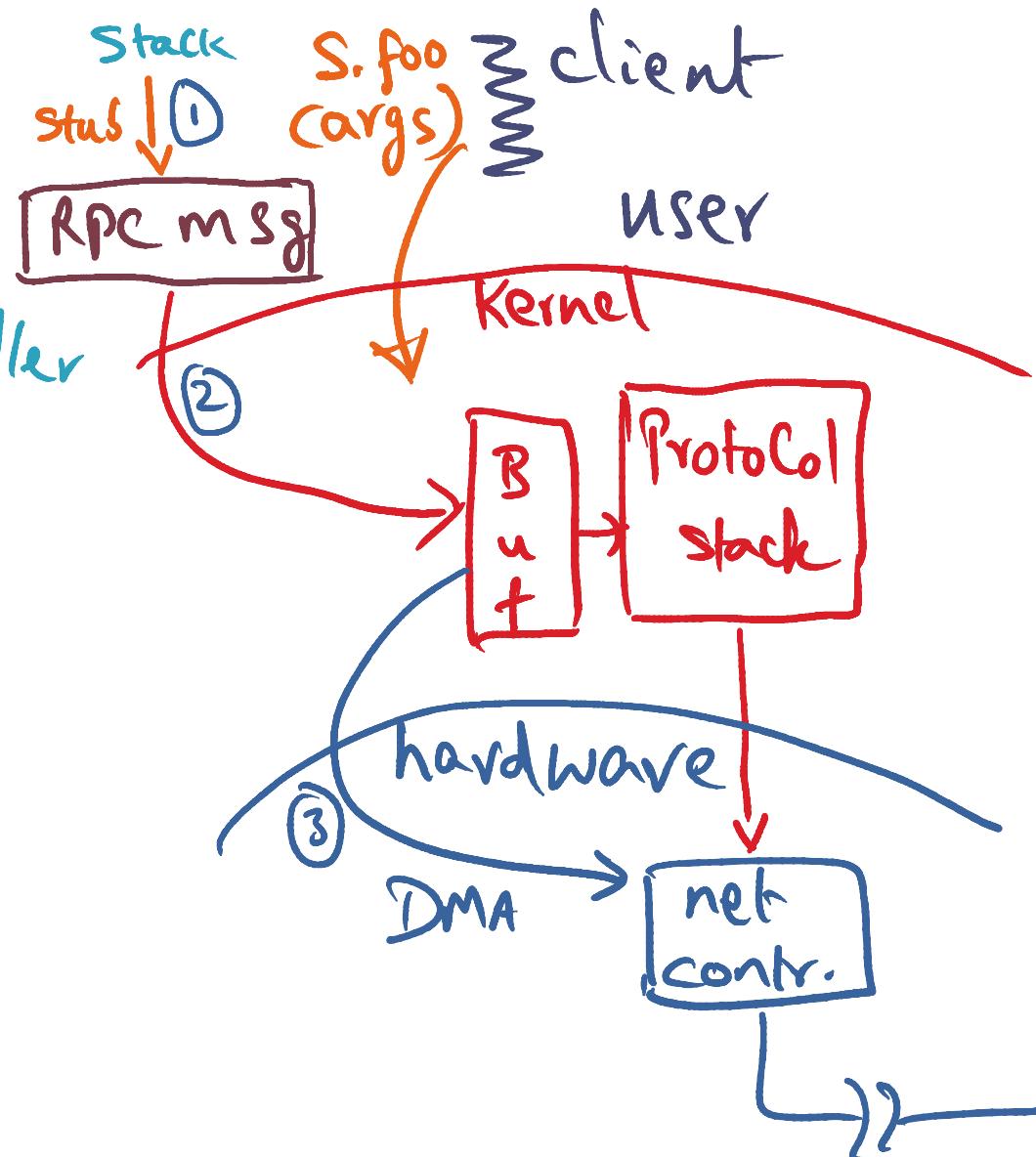
- client stubs
- Kernel buffer



marshaling and data Copying

Three Copies

- client stubs
- Kernel buffer
- DMA to controller

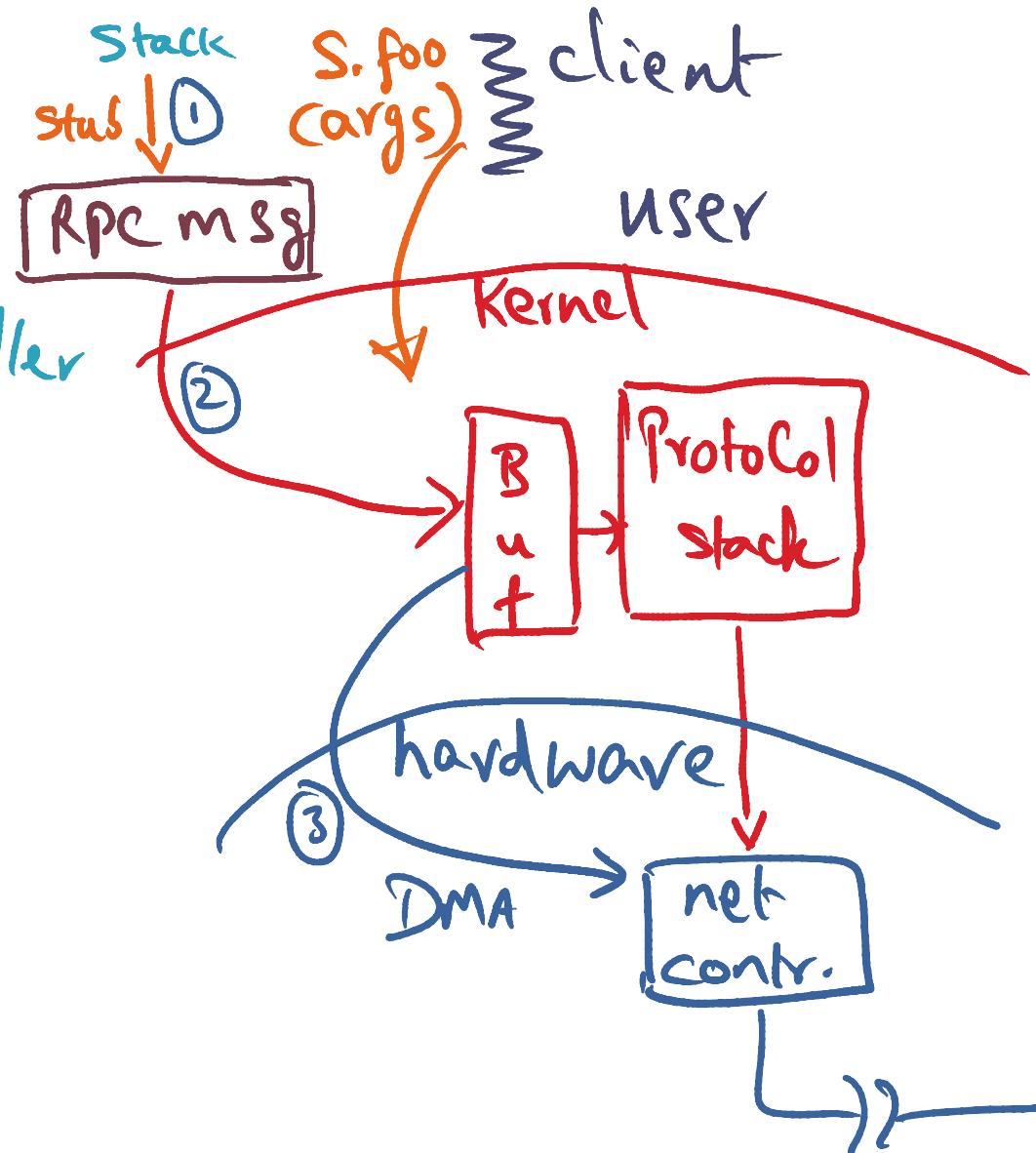


marshaling and data Copying

Three Copies

- client stubs
- Kernel buffer
- DMA to Controller

Reducing Copies?



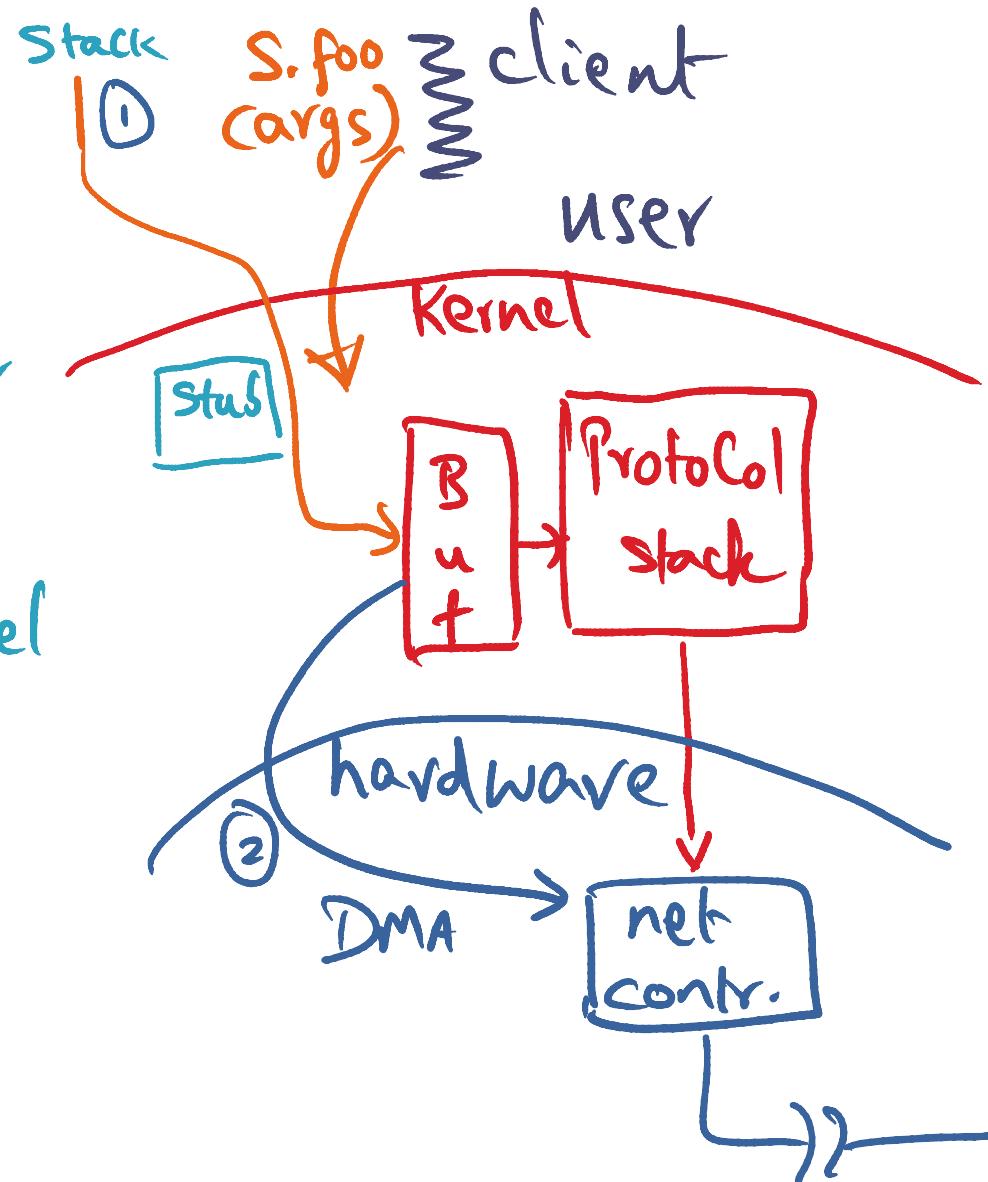
marshaling and data Copying

Three Copies

- client stubs
- Kernel buffer
- DMA to Controller

Reducing Copies?

- marshal into kernel buffer directly



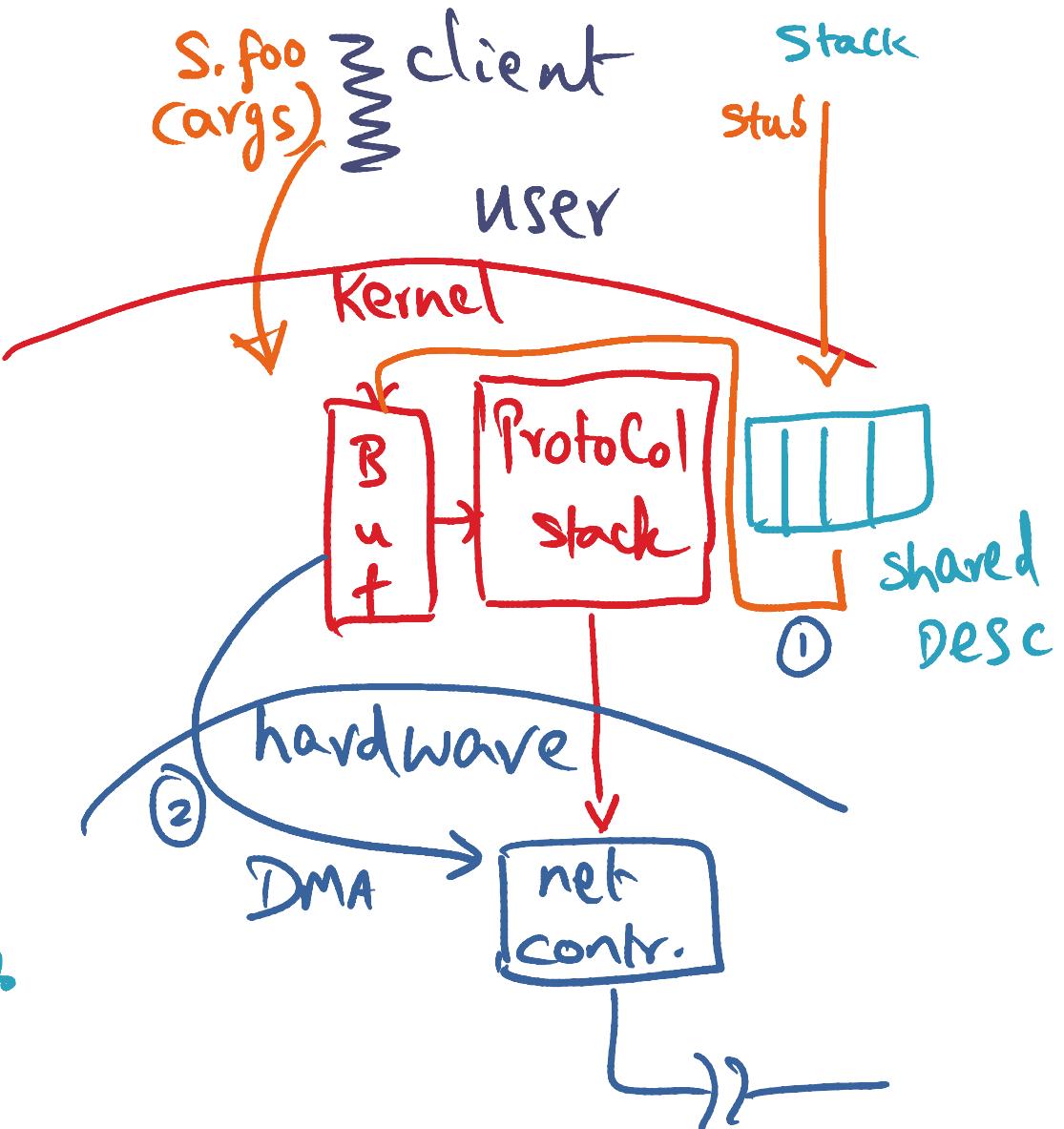
marshaling and data Copying

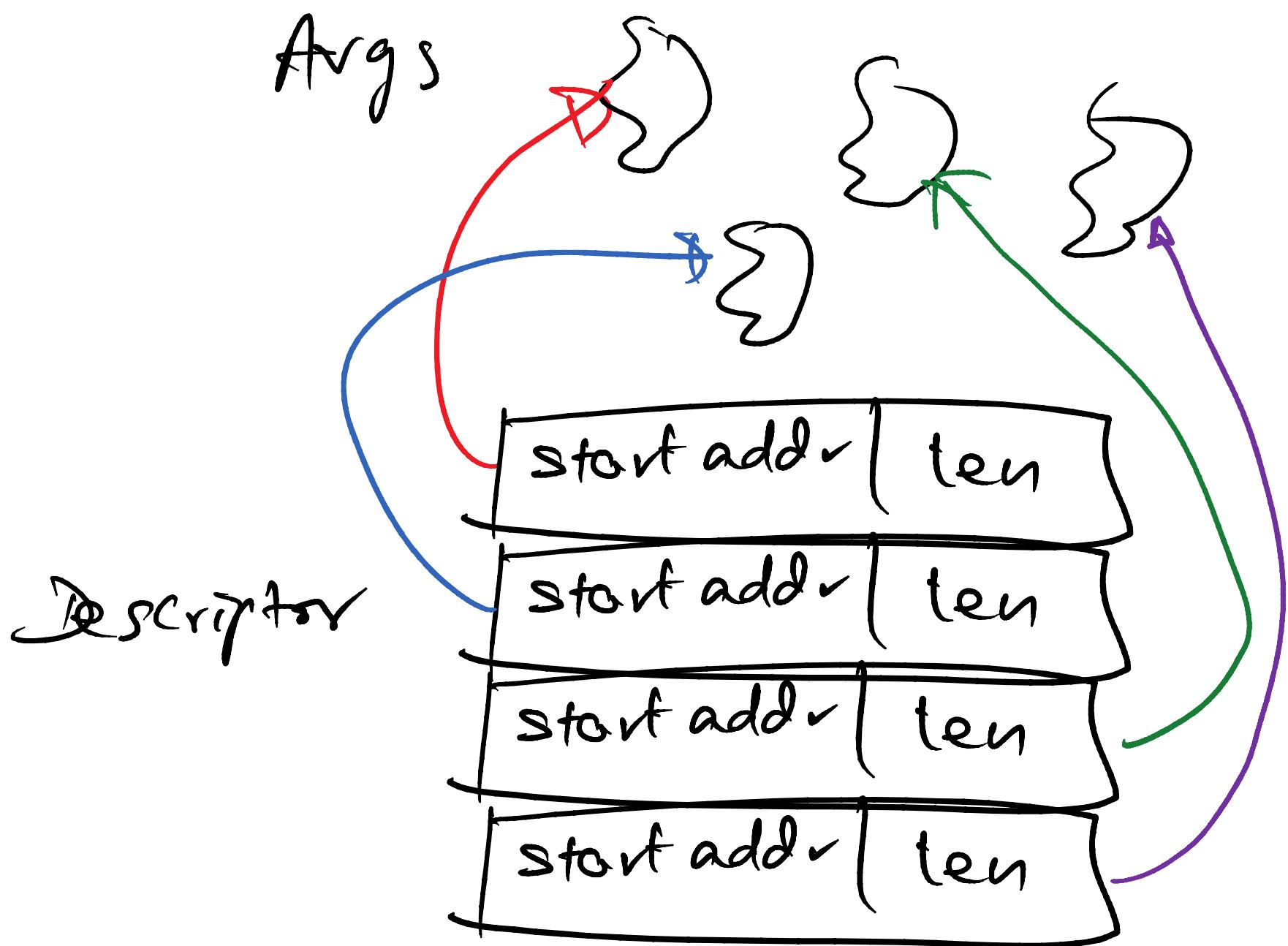
Three Copies

- client stub
- Kernel buffer
- DMA to Controller

Reducing Copies?

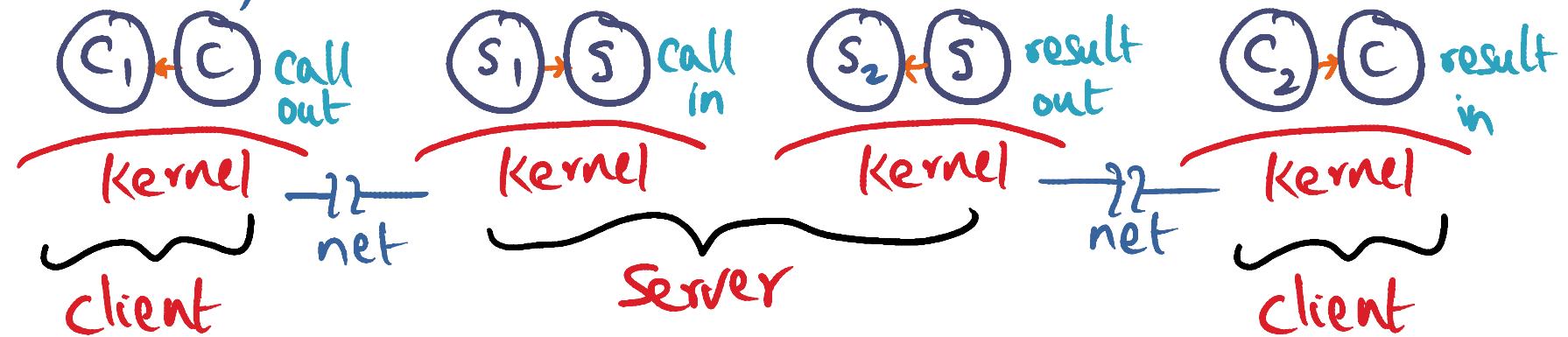
- marshal into kernel buffer directly
OR
- Shared descriptors between client stub and kernel





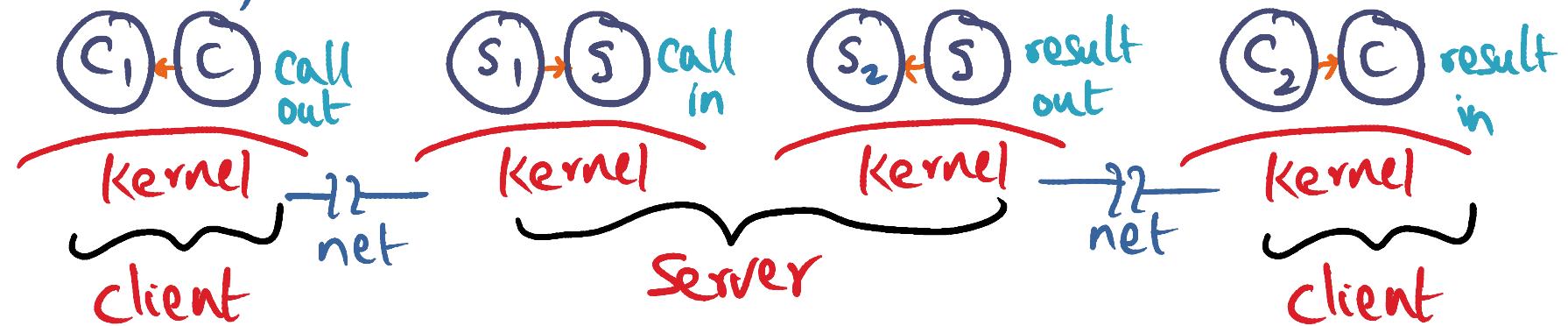
Control Transfer

Potentially 4

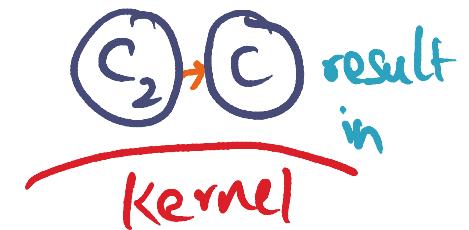
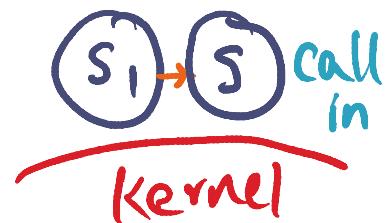


Control Transfer

Potentially 4

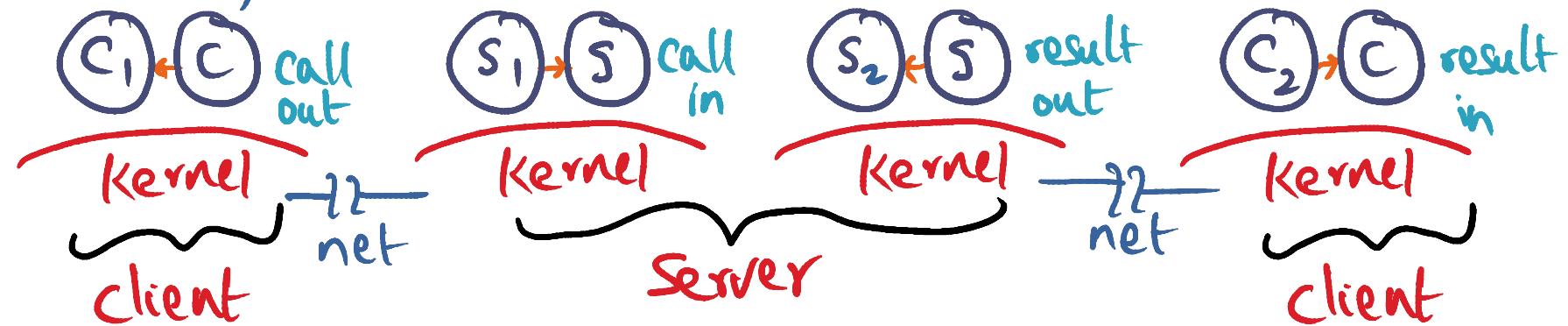


only two in critical path of latency



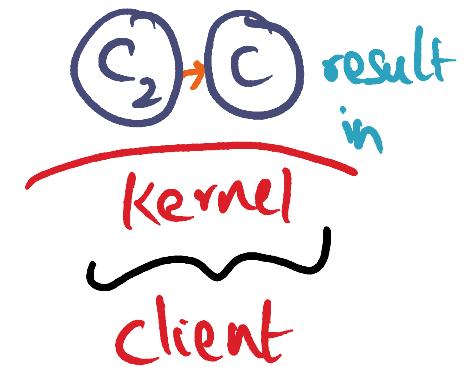
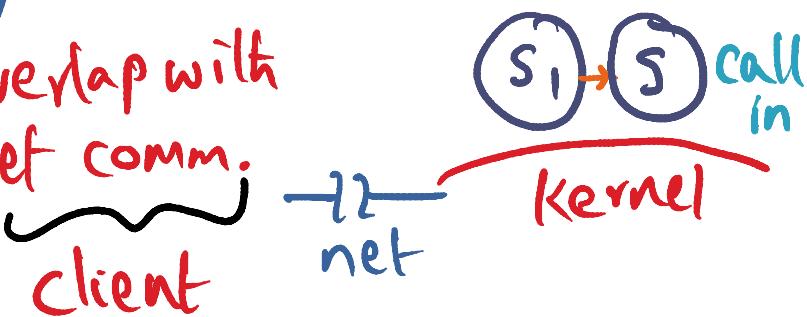
Control Transfer

Potentially 4



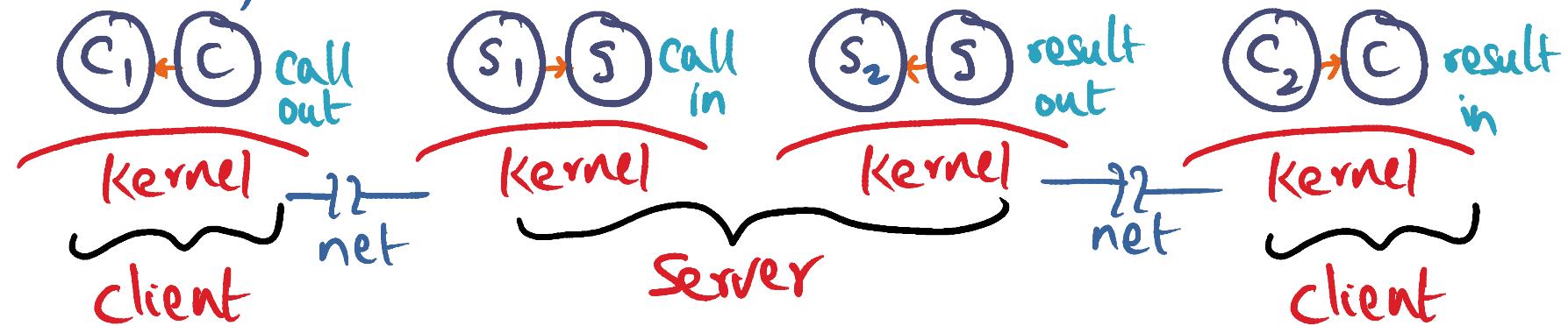
only two in critical path of latency

overlap with
net comm.



Control Transfer

Potentially 4



only two in critical path of latency

overlap with
net comm.

client

call in

net

Kernel

Server

overlap with
net comm.

Kernel

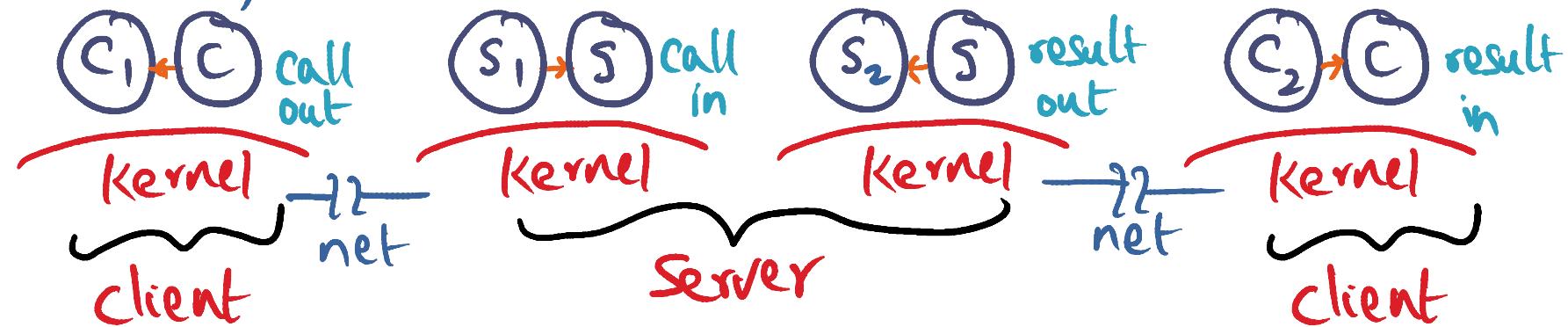
client

result in

result in

Control Transfer

Potentially 4



only two in critical path of latency \Rightarrow down to 2

overlap with
net comm.

client

call in

Kernel

Server

overlap with
net comm.

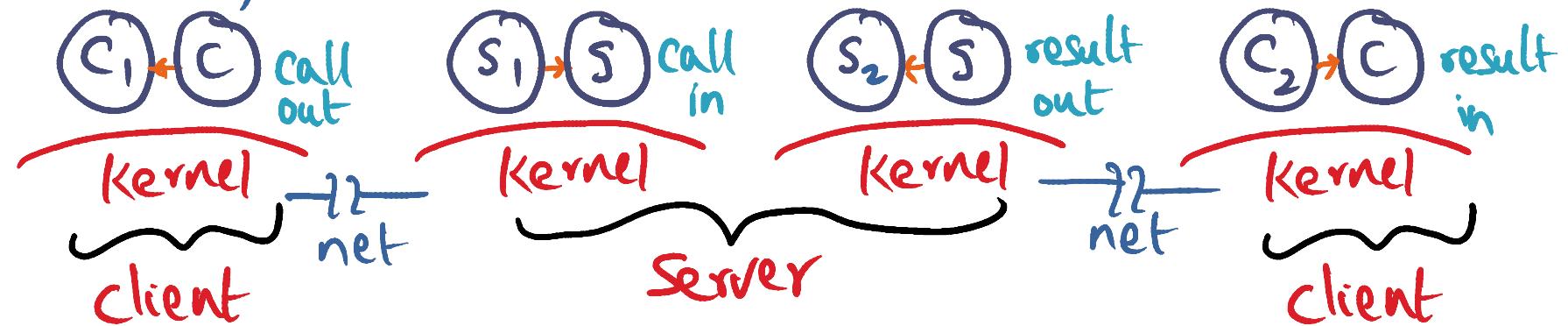
result in

Kernel

client

Control Transfer

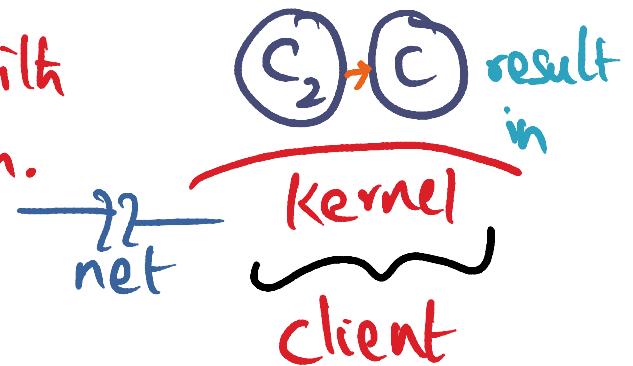
Potentially 4



only two in critical path of latency \Rightarrow down to 2

overlap with
net comm.

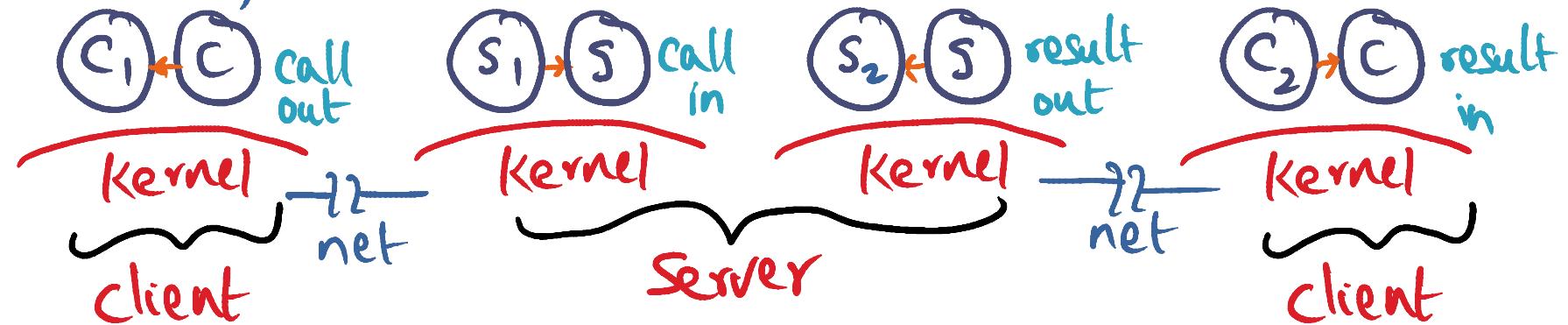
client



Can we reduce it to 1?

Control Transfer

Potentially 4



only two in critical path of latency \Rightarrow down to 2



Can we reduce it to 1?

- Yes if we **spin instead of switch** on client side

Protocol Processing

What transport for RPC?

LAN reliable \Rightarrow reduce latency

Protocol Processing

What transport for RPC?

LAN reliable \Rightarrow reduce latency

choices for reducing latency in transport

No low level ACKS

Hardware checksum for packet integrity

No client side buffering since client blocked

overlap server side buffering with

result transmission

Key Takeaways from Thekkath and Levy Paper

- Latency more important than bandwidth.
- High bandwidth networks do not imply low latency communication.
- It is possible to shave off software overheads and get close to hardware speeds by careful OS design of the RPC code path.

Read the paper by Thekkath and Levy fully to appreciate how they put these principles to practice and build a high-performance RPC system.

Today

✓ Distributed System Basics

✓ Lamport's clock

Efficient Communication Software

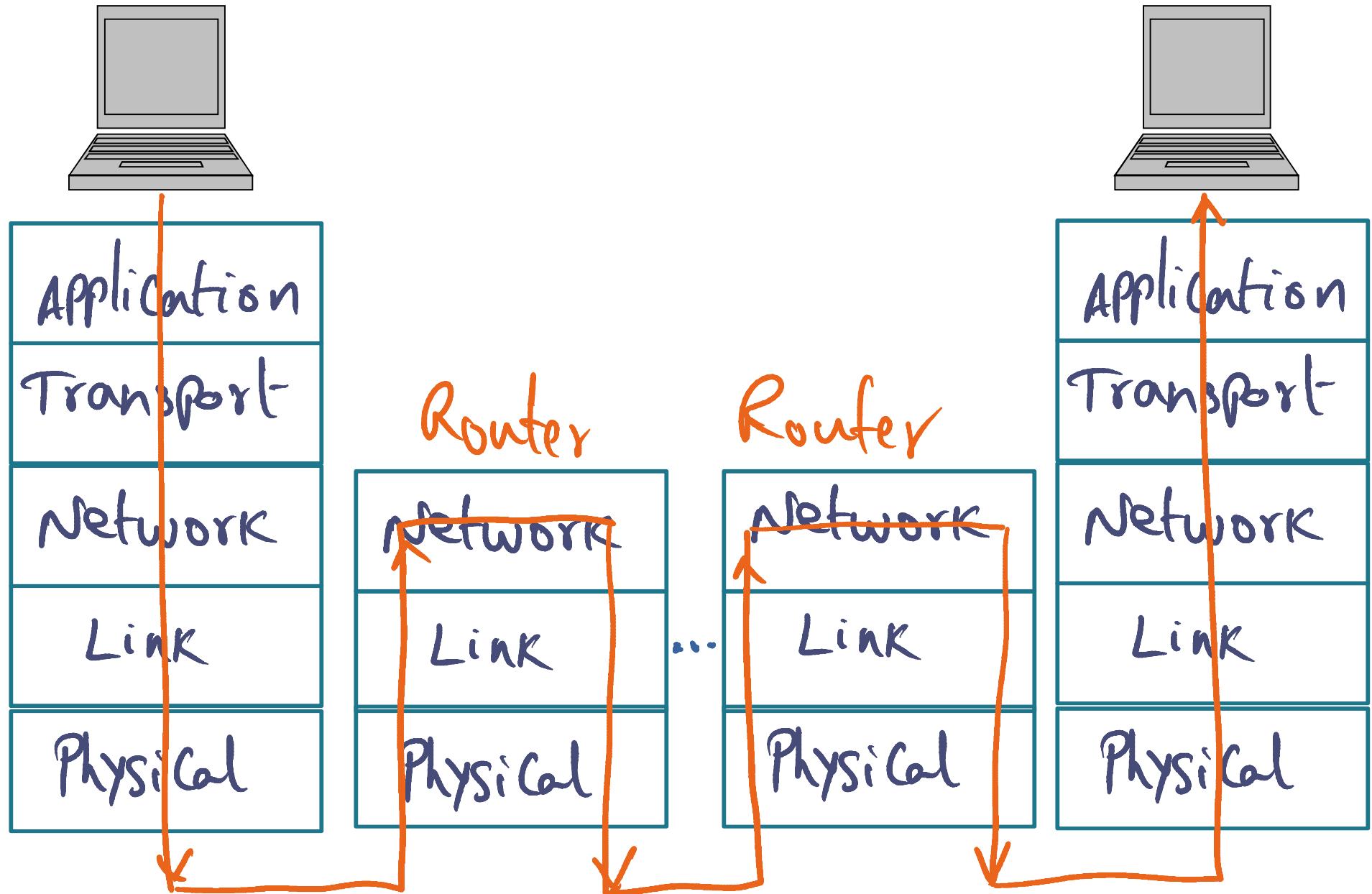
Therkauf + Levy { * Application interface to the Kernel
* Inside the Kernel

wetterwald { * End-to-End QoS via Active Networks ↙

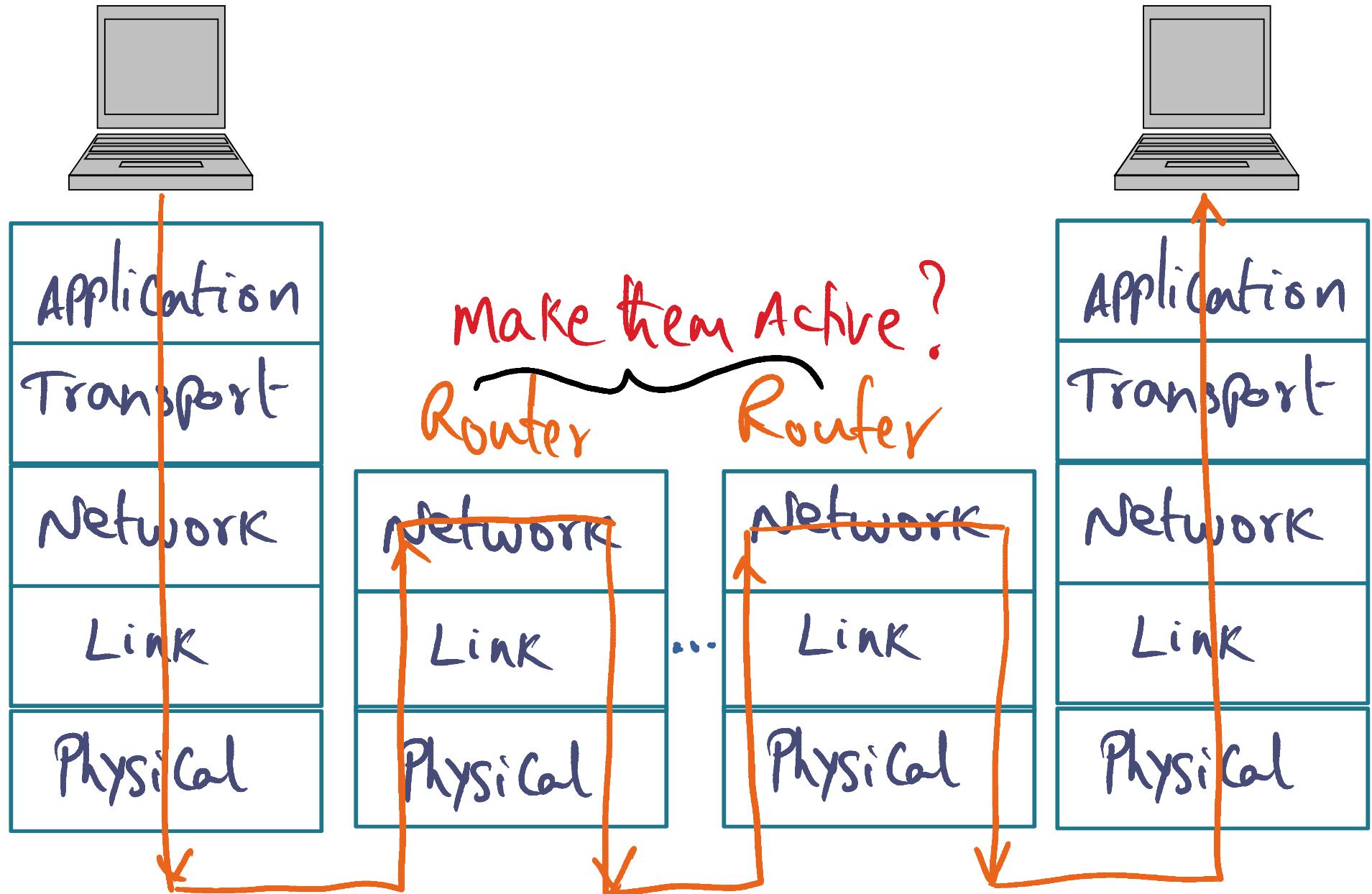
Friday

Synthesizing Network Protocol Stacks

Routing on the Internet

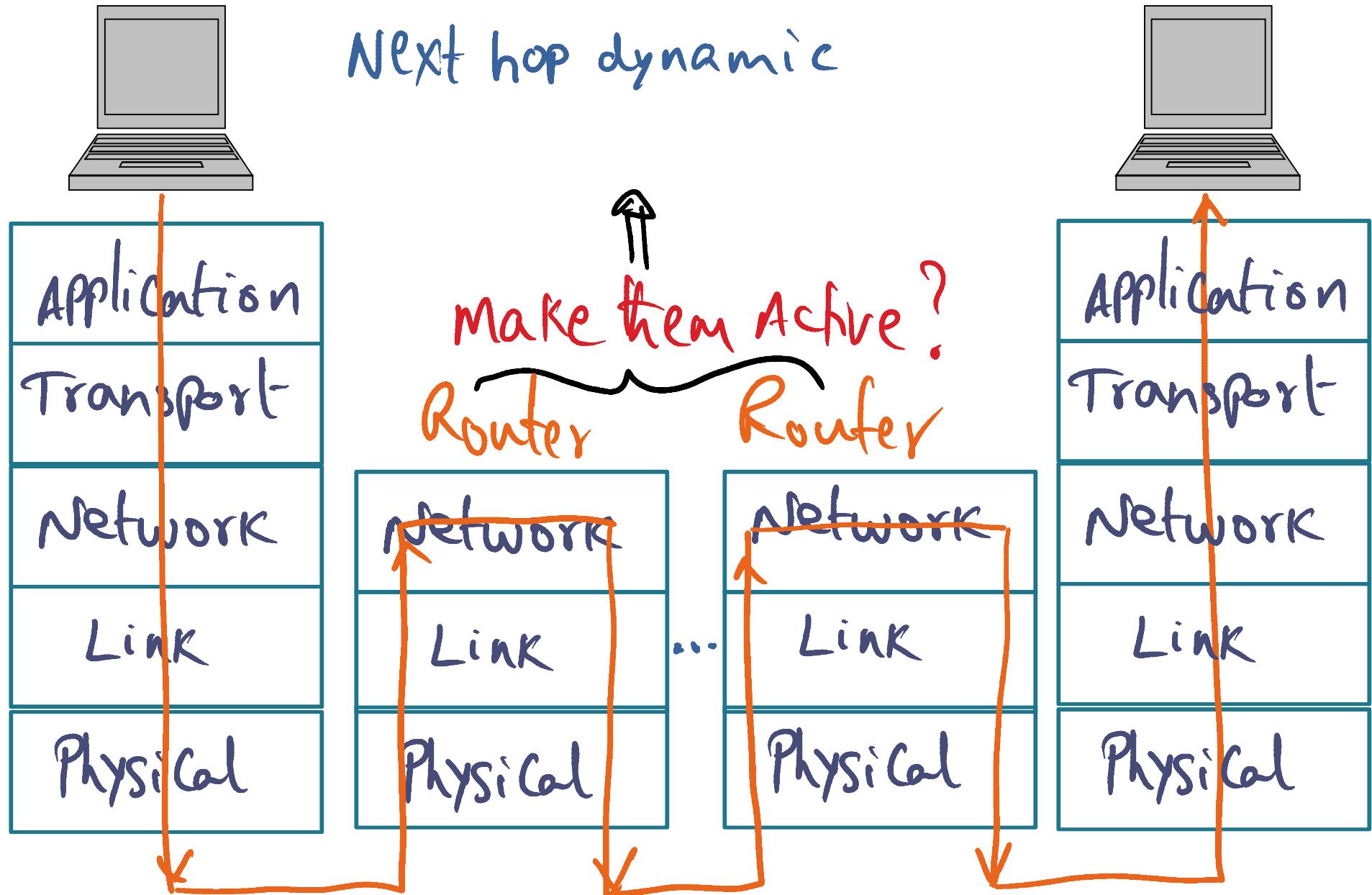


Routing on the Internet



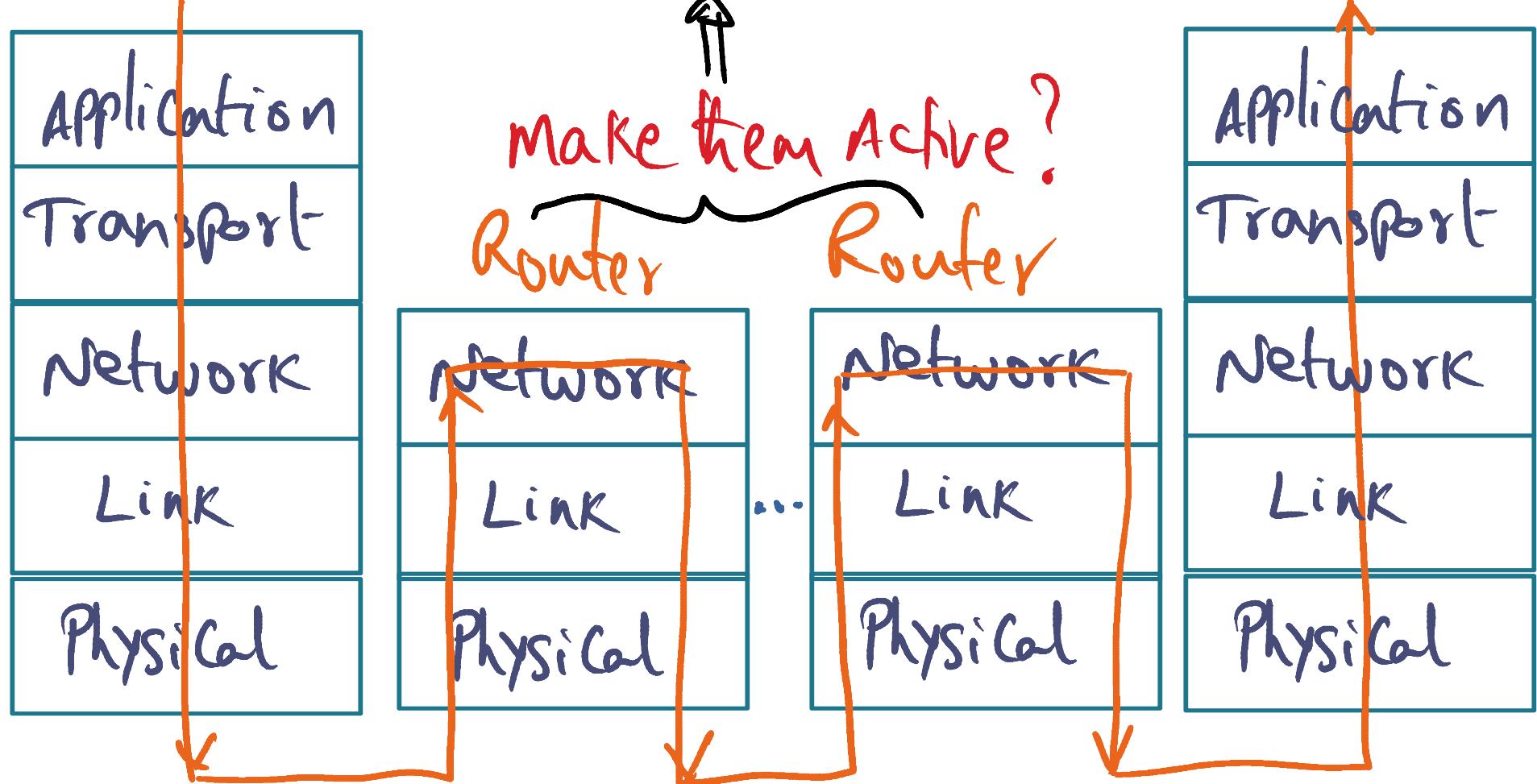
Routing on the Internet

Next hop dynamic



Routing on the Internet

Next hop dynamic
- How + who can write code?

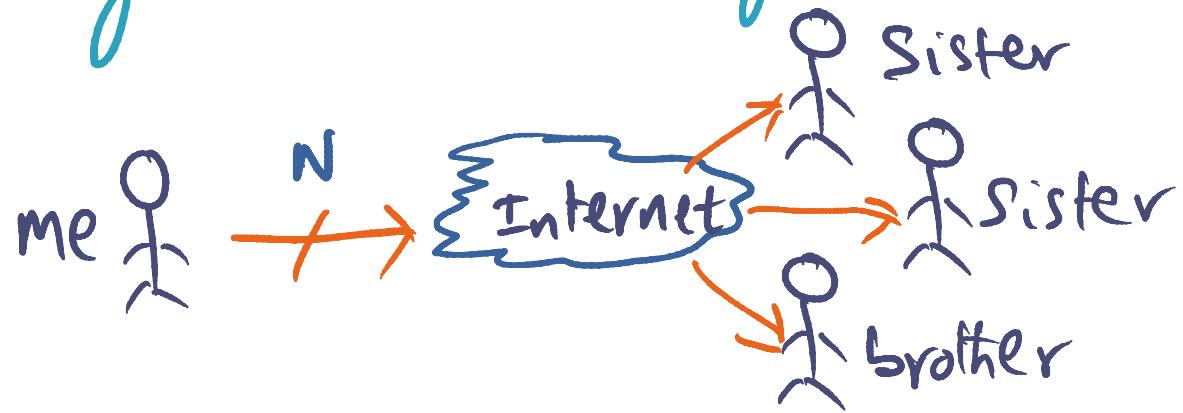


An Example

Me sending Diwali Greetings to my Siblings

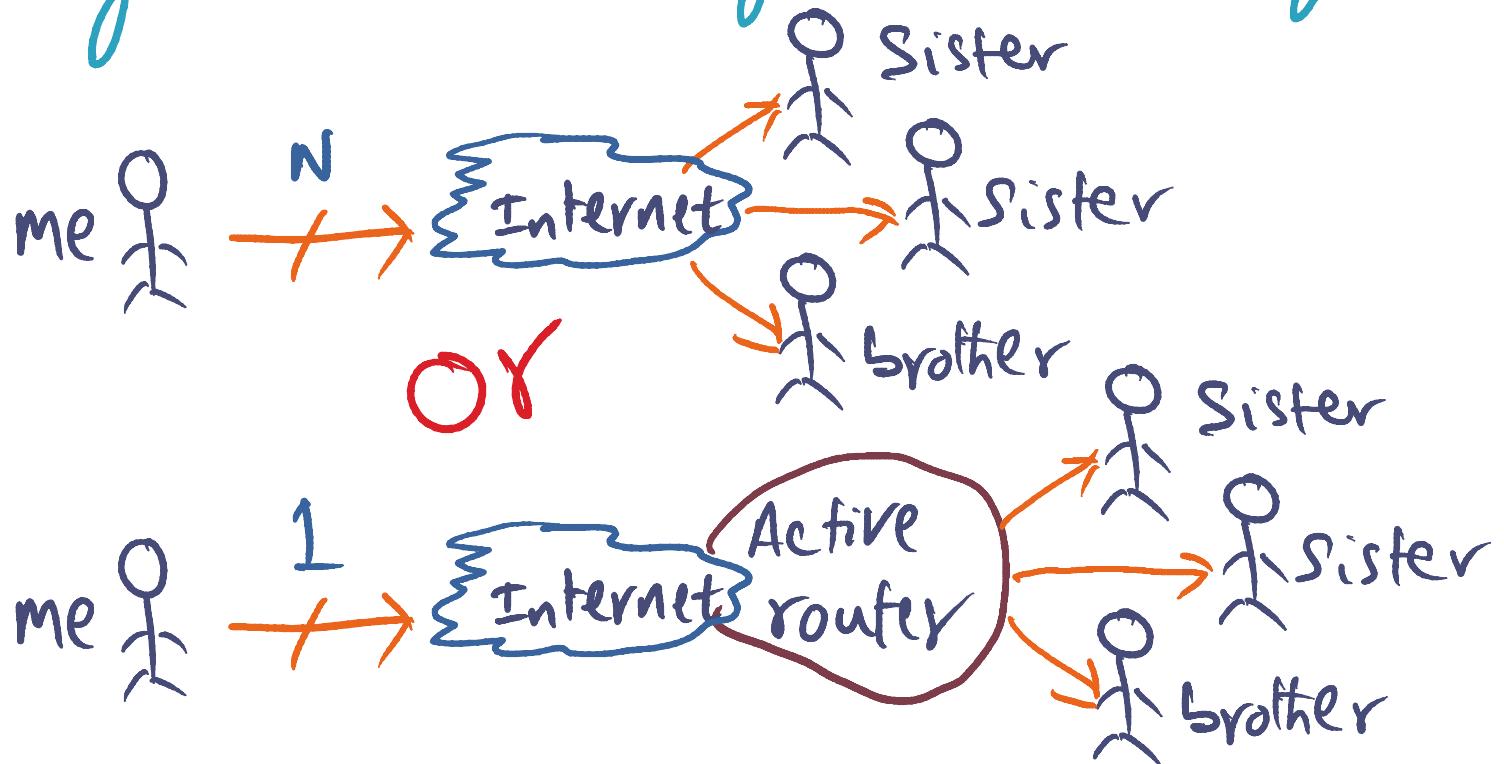
An Example

Me sending Diwali Greetings to my Siblings



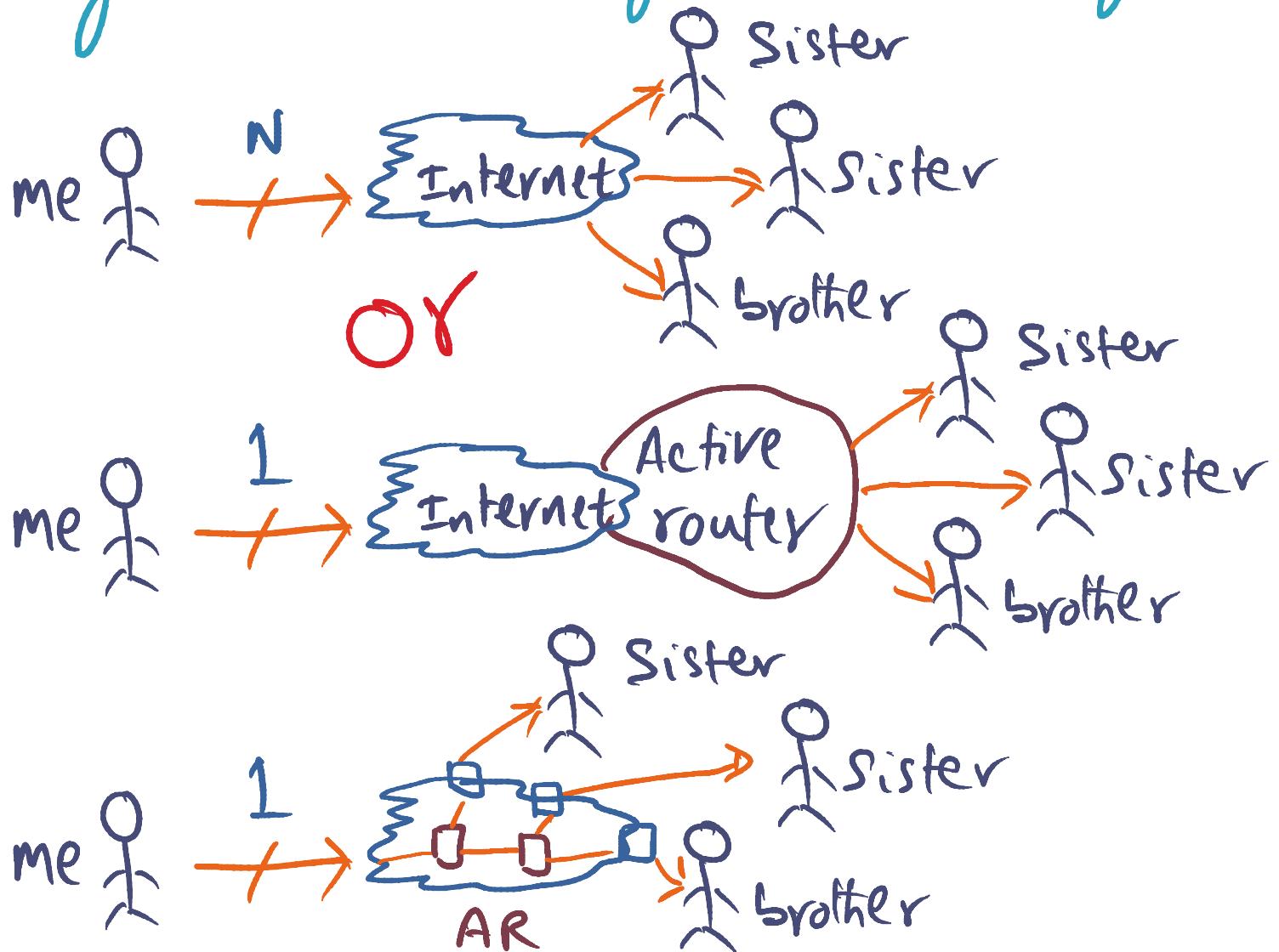
An Example

Me sending Diwali Greetings to my Siblings



An Example

Me sending Diwali Greetings to my Siblings



How to implement the vision?

