



Protocol Buffer DPU Offloading in the RPC Datapath

SC24 Workshop: Communication, I/O, and Storage at Scale on Next-Generation Platforms - Scalable Infrastructures

Raphaël Frantz, Jerónimo Sánchez García, Marcin Copik, Idelfonso Tafur Monroy, Juan José Vegas Olmos, Gil Bloch, Salvatore Di Girolamo

Quantum Terahertz Systems, Electro-Optical Communications, Department of Electrical Engineering

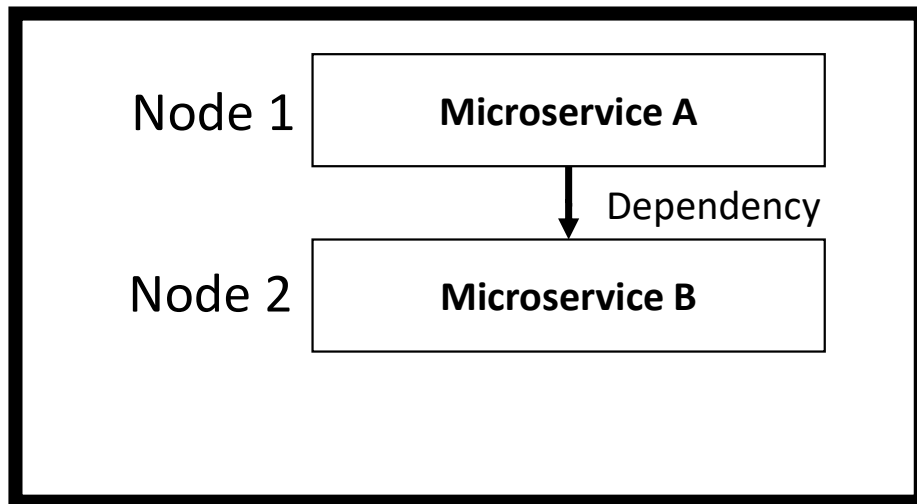
Outline

1. Background
2. Designs
3. Benchmarks
4. Conclusion

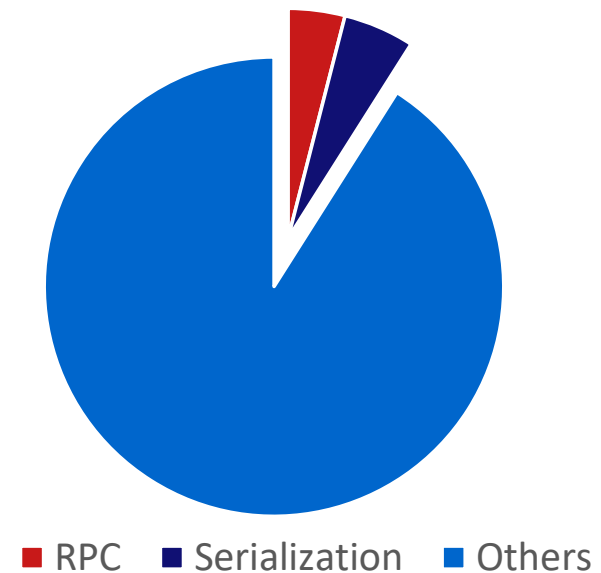
Microservice Architecture

- Today's applications are **split** in independent “microservices”
- Communication is done with RPCs

Data center



% Data Center Tax



Kanev, S. et al., 2016
Profiling a warehouse-scale computer

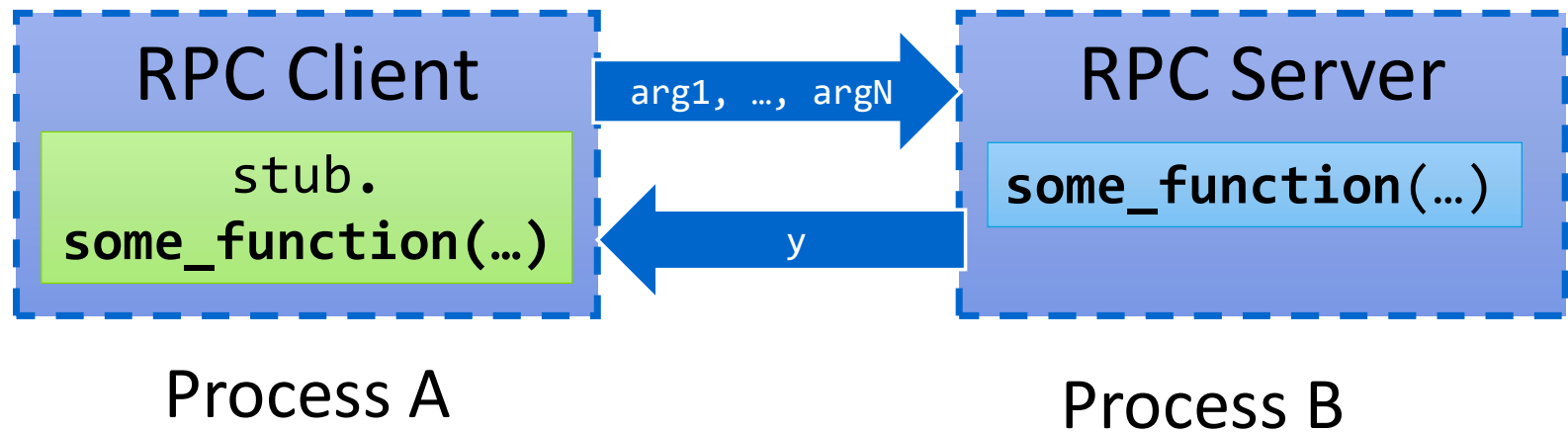
Remote Procedure Calls (RPCs)

High-level abstraction of communication with function calls

In code:

```
y = stub.some_function(arg1, ..., argN)
```

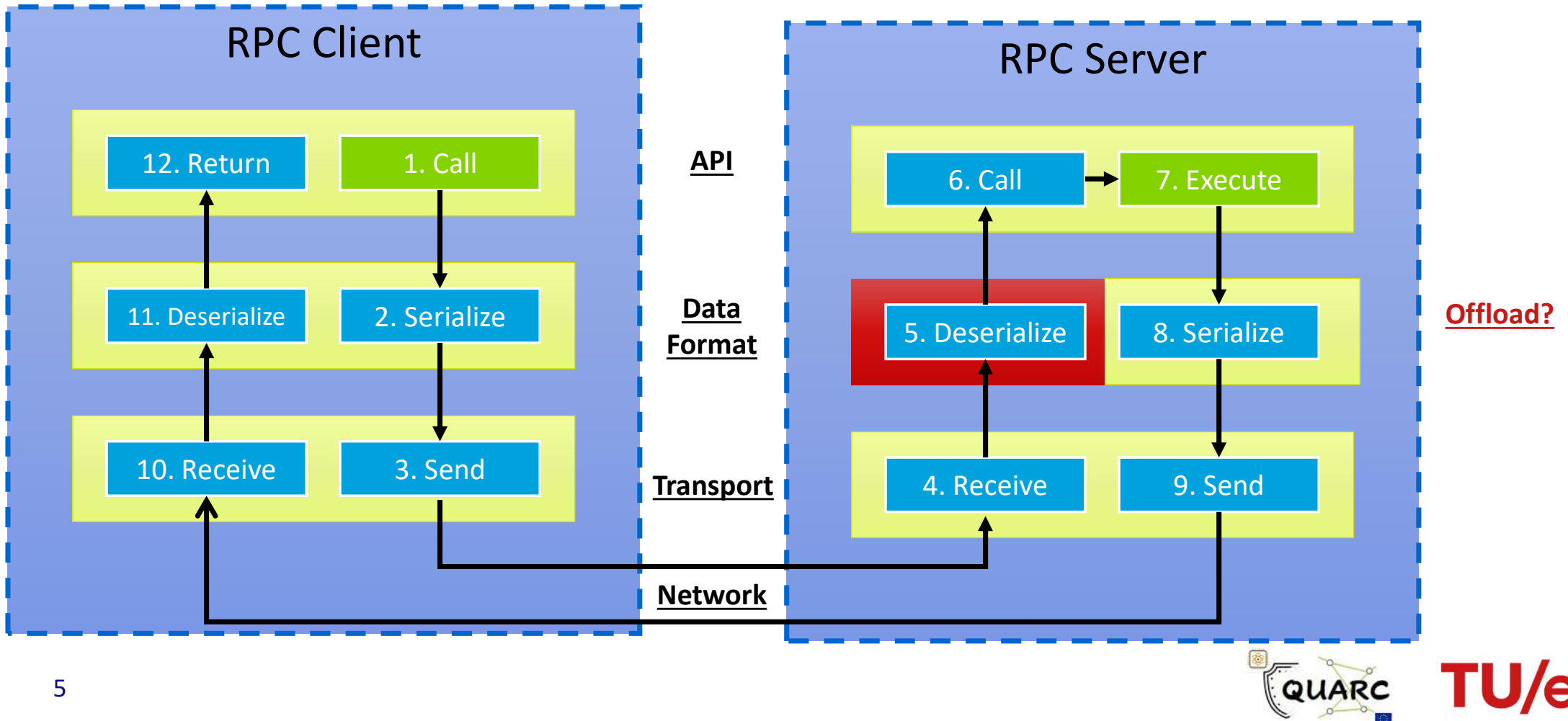
At runtime:



RPC Datapath

Exposed to the User

RPC Framework



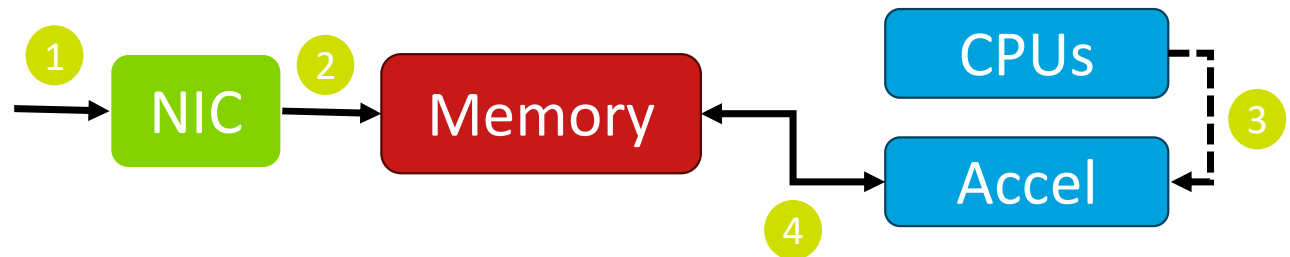
Strategies to Offload RPC Deserialization

● Time

No
Offloading



Hardware
Offloading



Software
Offloading



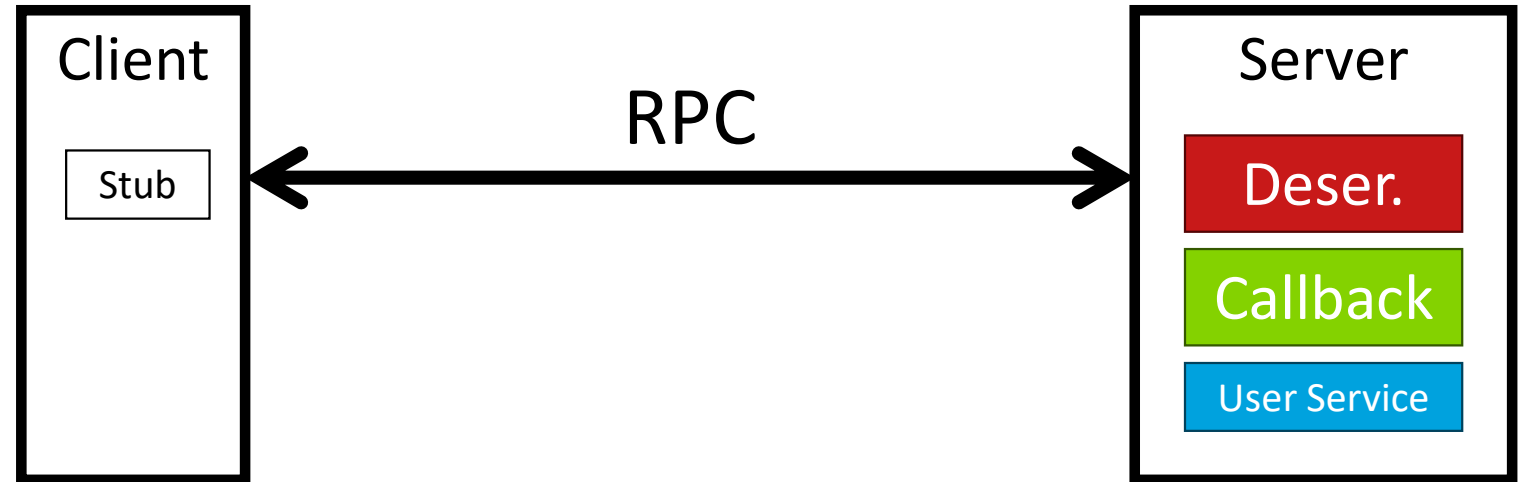
Cao, S. et al., 2022.

Accelerating Data Serialization/Deserialization Protocols with In-Network Compute

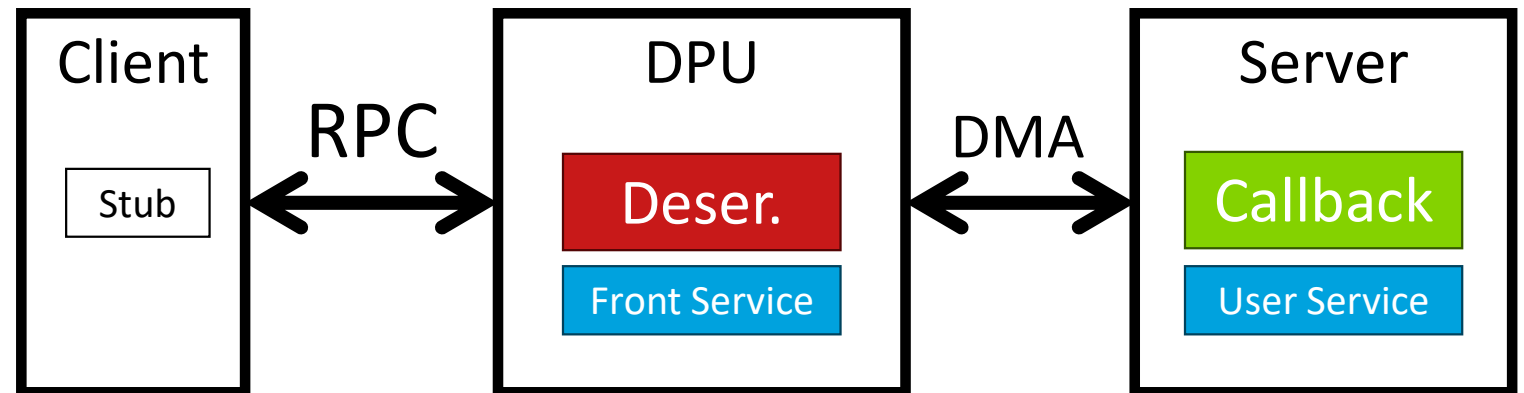


Our Solution to Offload RPC Deserialization

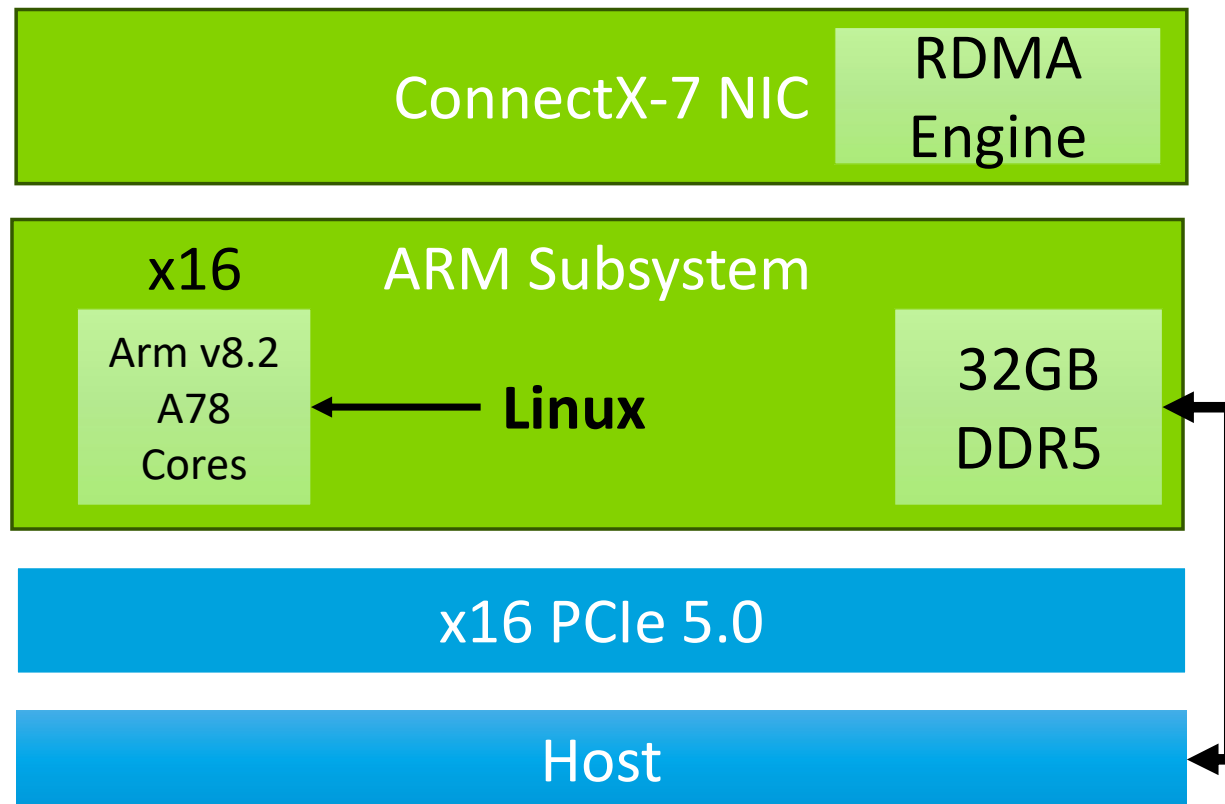
No offloading



Offloading

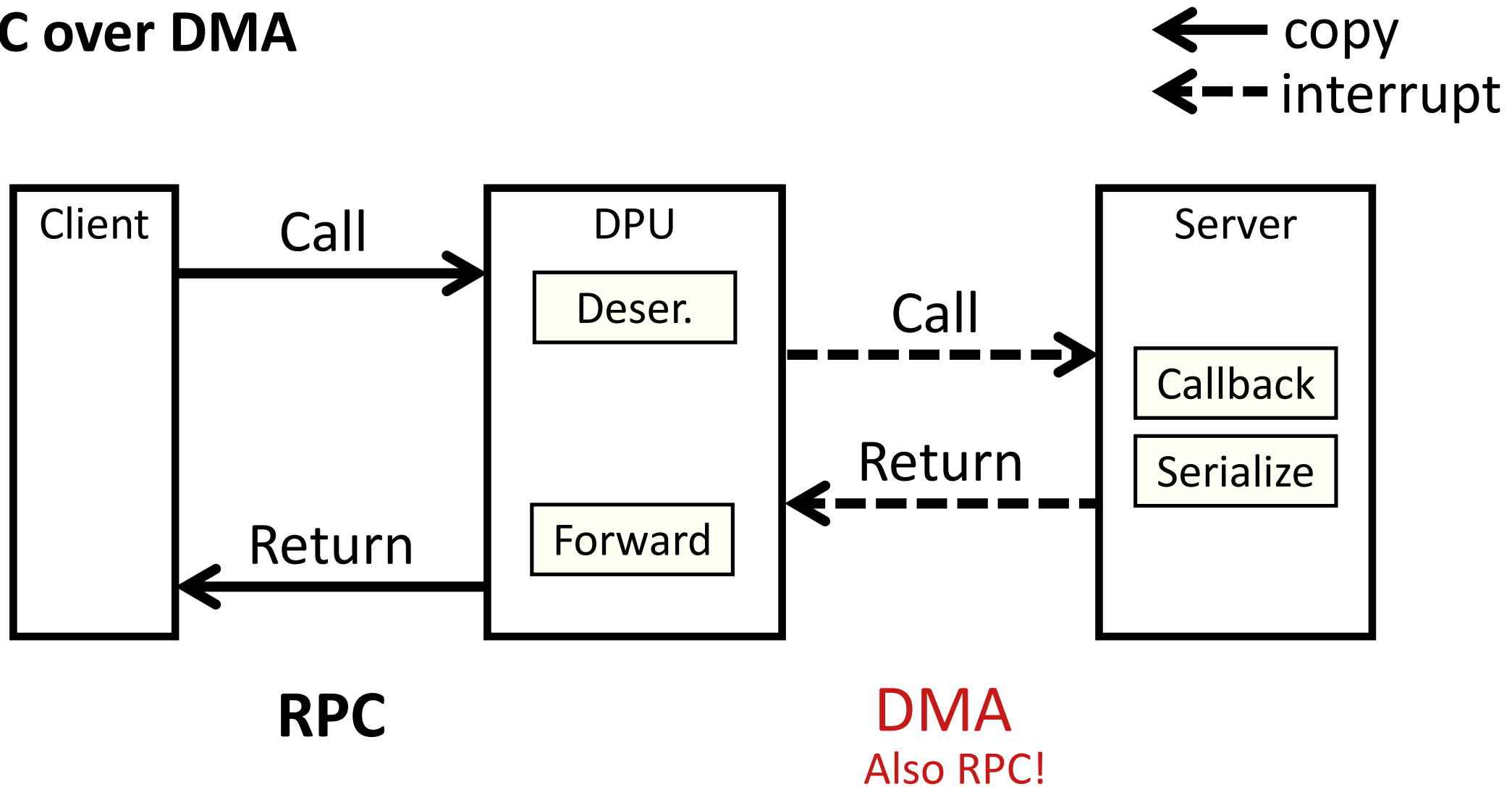


BlueField-3 Data Processing Units (DPUs)



Direct Memory Access

RPC over DMA



Auto-generated Code

Domain-Specific Language (DSL)

```
message Foo {  
  int32 i = 1;  
  string s = 3;  
}
```

protoc
Compiler

C++ Class

Foo.pb.h
Foo.pb.cc

serialize()
deserialize()
getters
setters

Our
plugin

Application Info

Class Info

Class Info

Default instance (bytes)

Field Info

Field Info

Offset (int)

Type (enum)

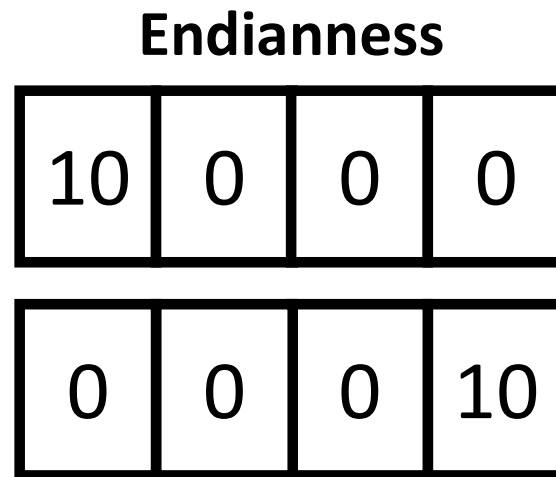
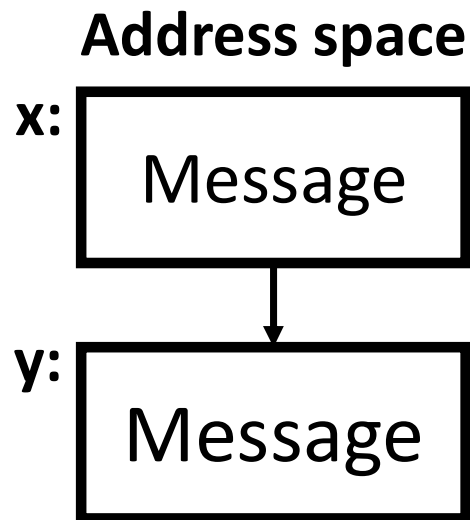
Class (ptr)

C++ Reflection

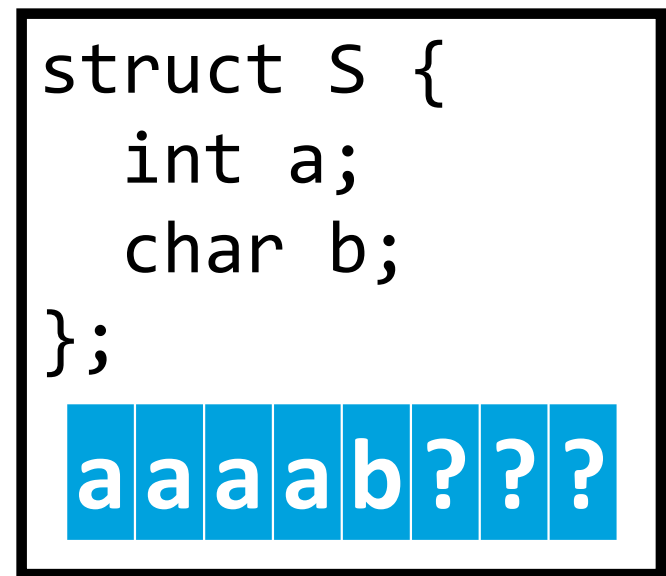
Foo.adt.pb.h
Foo.adt.pb.cc

Challenges with offloading Deserialization

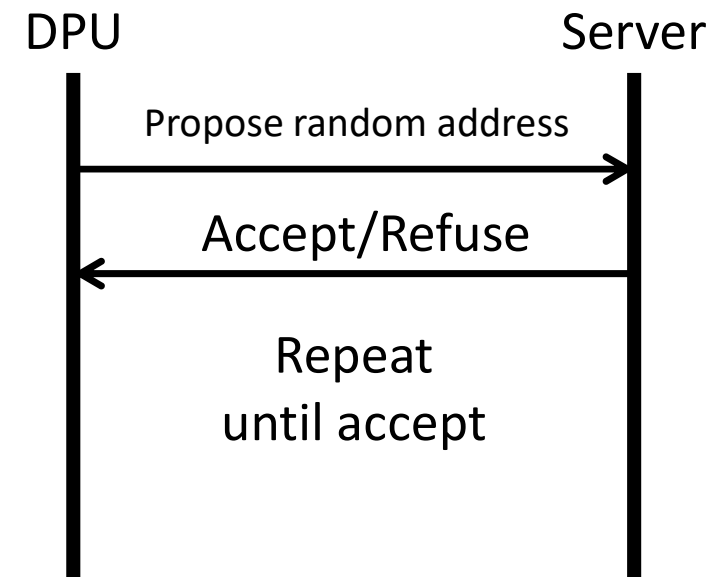
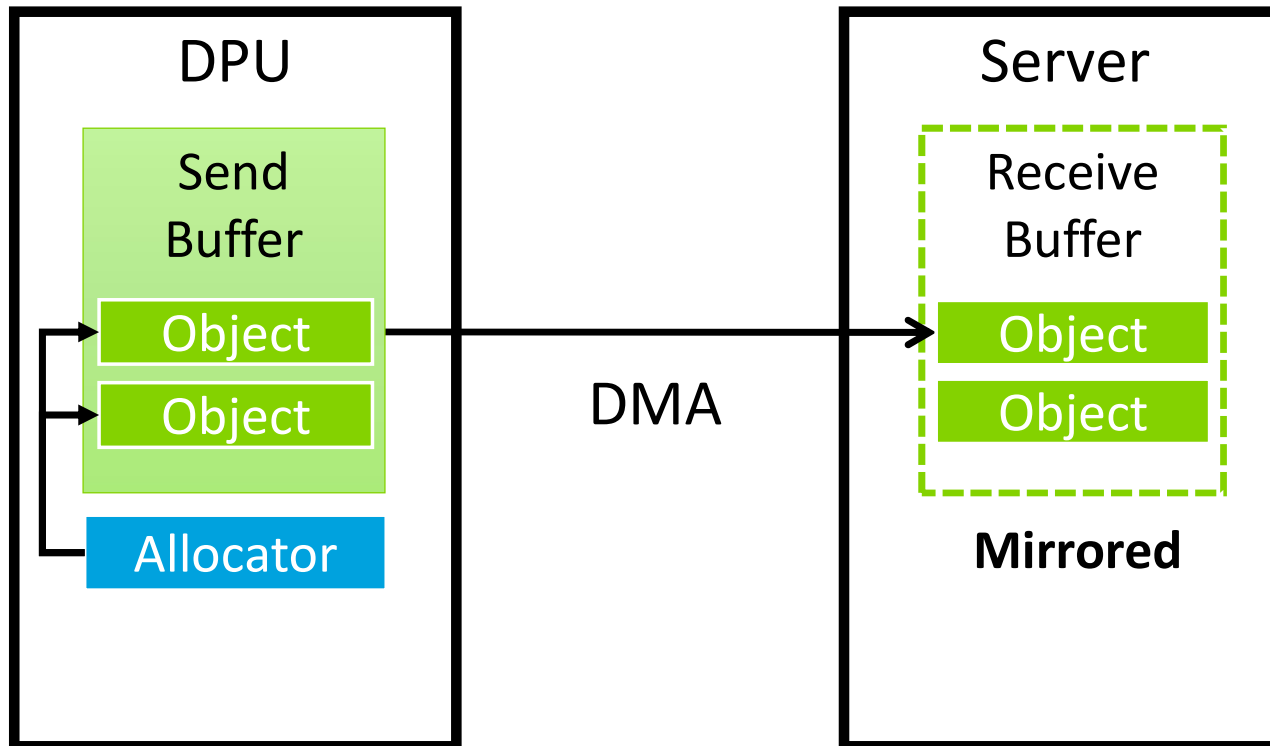
1. Address space ☹️
2. Application Binary Interface (ABI)
 - DPU ABI = Server ABI = Itanium 😊



ABI

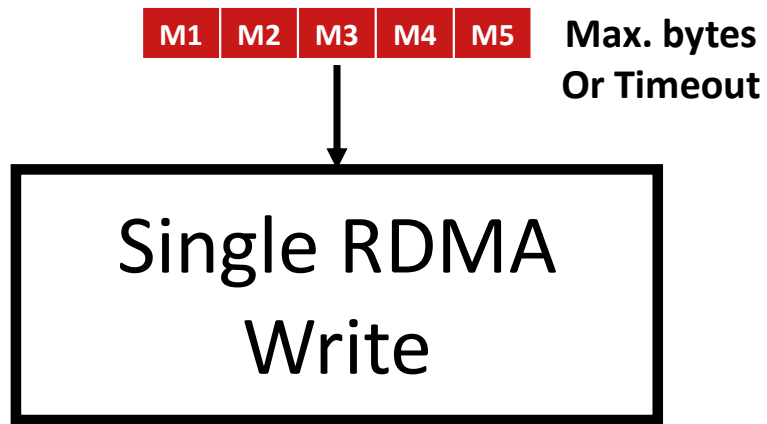


Shared Address Space

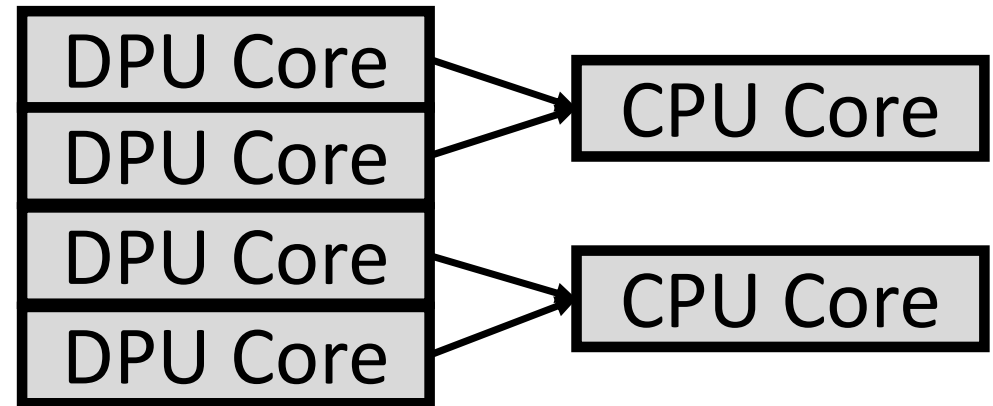


DMA-based RPC Protocol

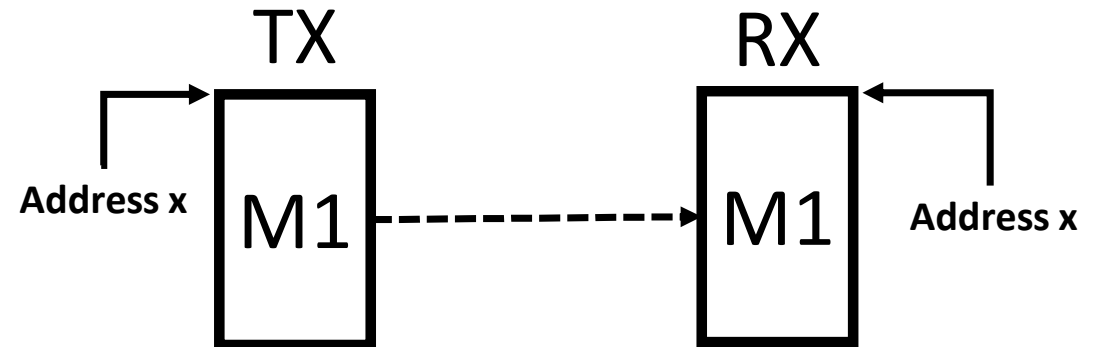
Pack Small Messages



Parallelized Queues



Shared Address Space



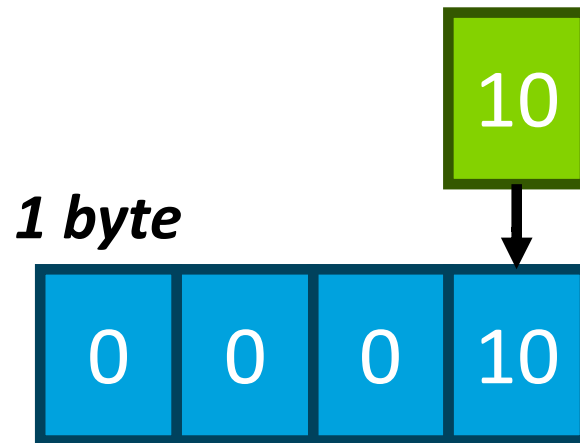
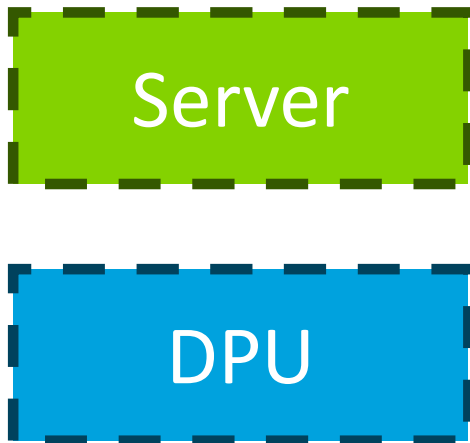
Benchmarking: Deserialization Workload Types

Discriminate two types of costs:

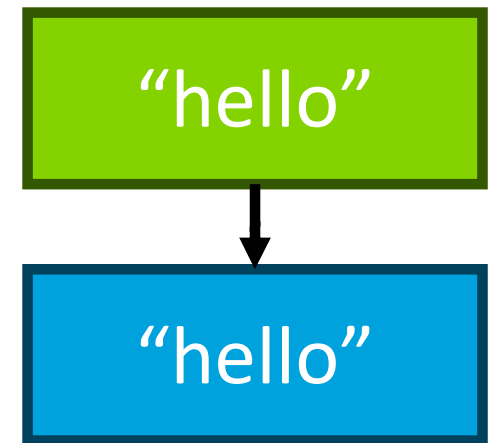
- High **Compute** Cost
- High **Copy** Cost

Compressed Bytes	Max. Value	Generated Proportion
1	127	50%
2	16383	25%
3	2 millions	12.5%
4	2 billions	0.075%
5	max	0.05%

Variable-length Integers

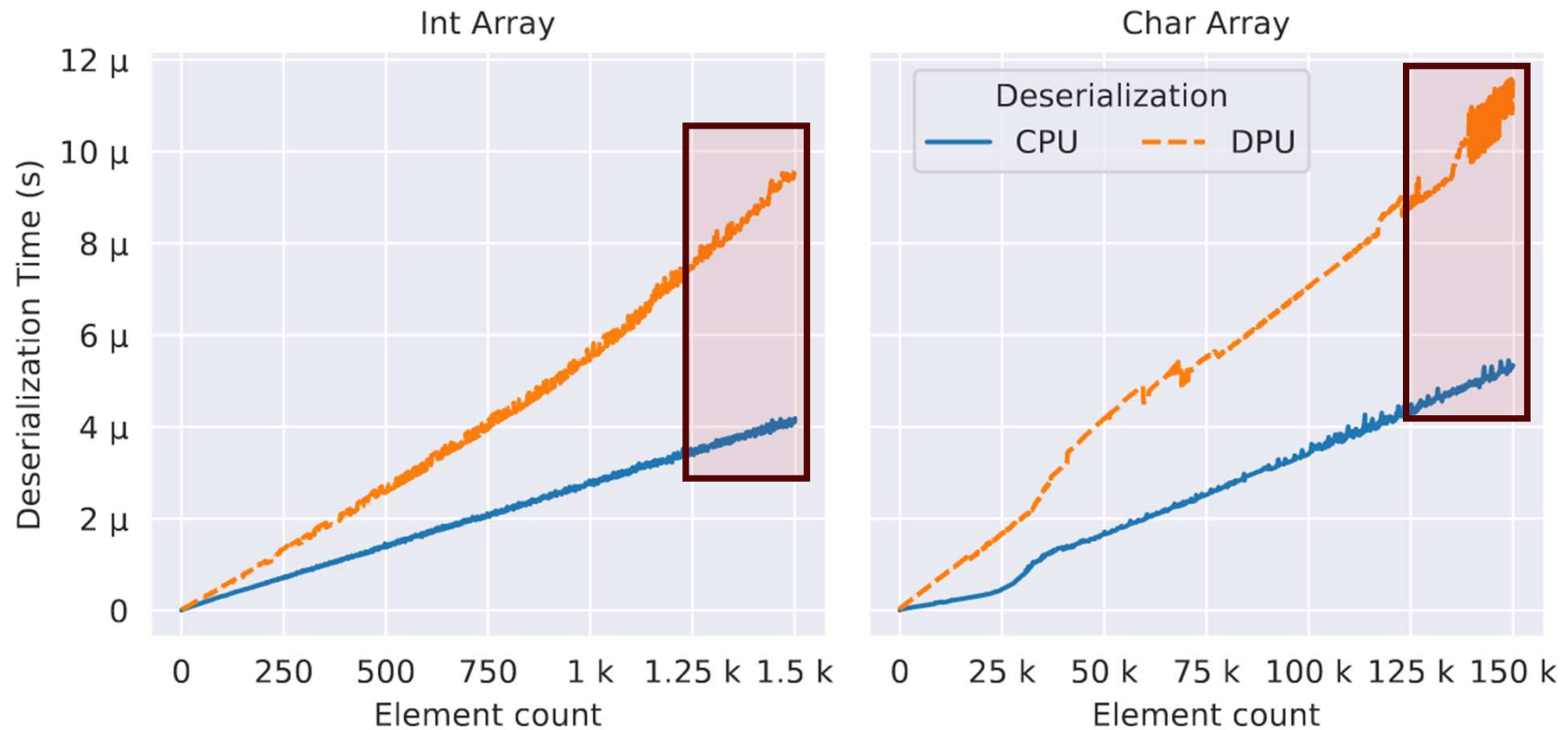


Strings

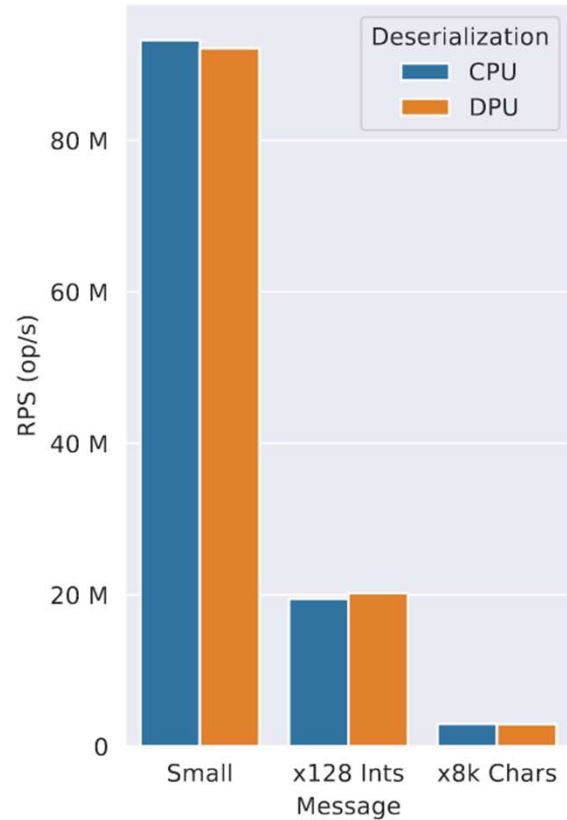


Benchmarking Deserialization

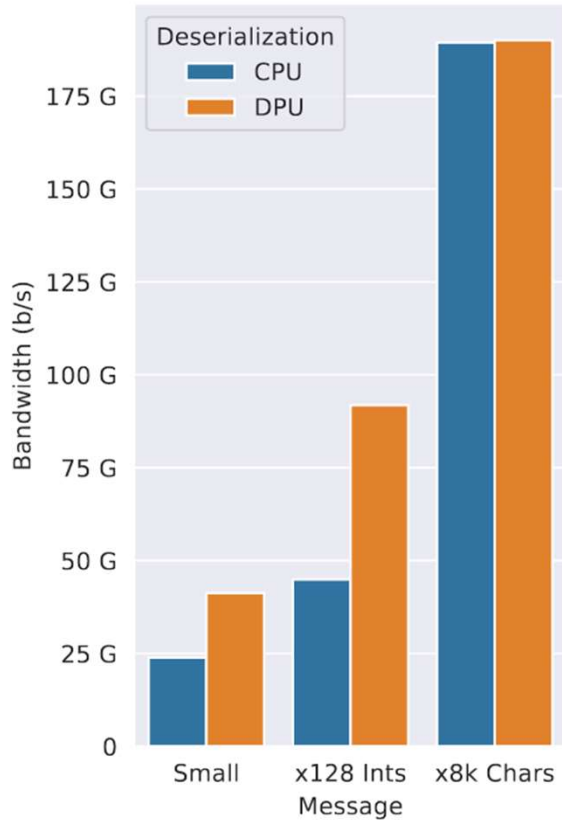
2 DPU cores match the performance of 1 CPU Cores



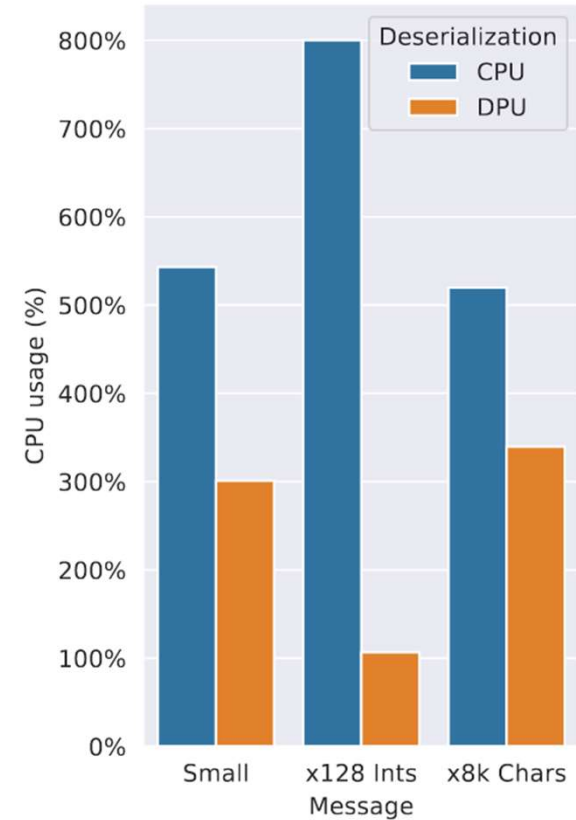
Benchmarking Deserialization in the RPC Datapath



Requests/Sec

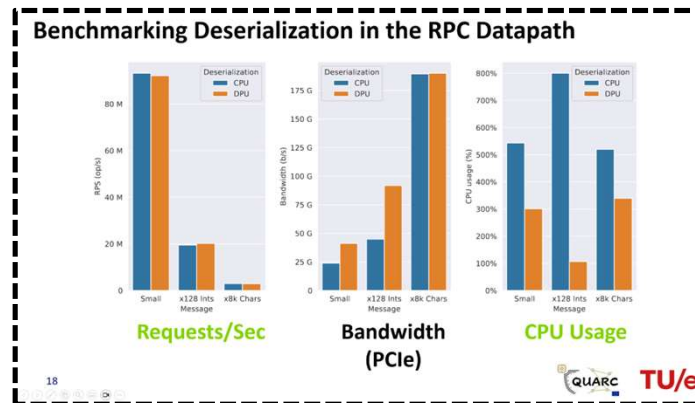
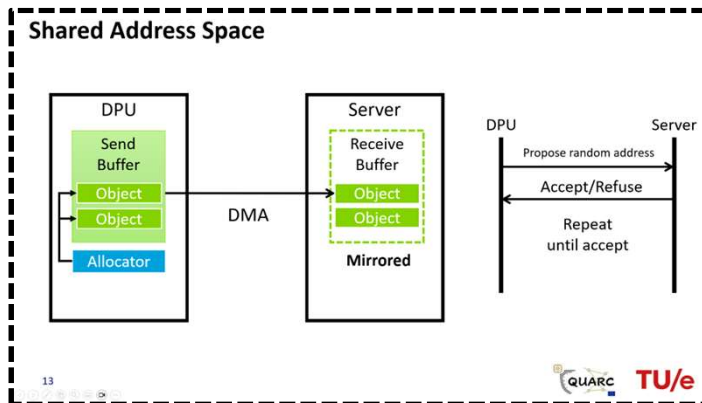
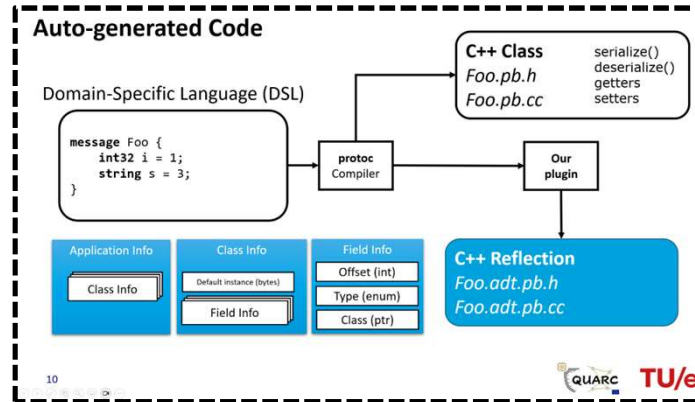
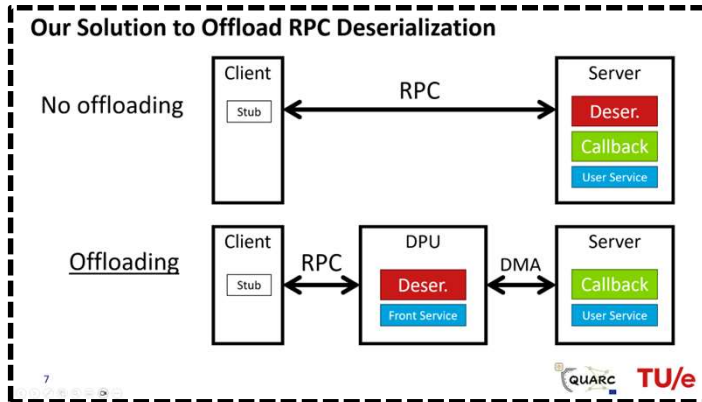


Bandwidth
(PCIe)



CPU Usage

Conclusion



Repository



r.r.a.frantz@tue.nl

This work was partly funded by the European Union Horizon Europe grant number 101073355, and by the grant PID2021-123041OB-I00 funded by MCIN/AEI/10.13039/501100011033 and by "ERDF A way of making Europe".

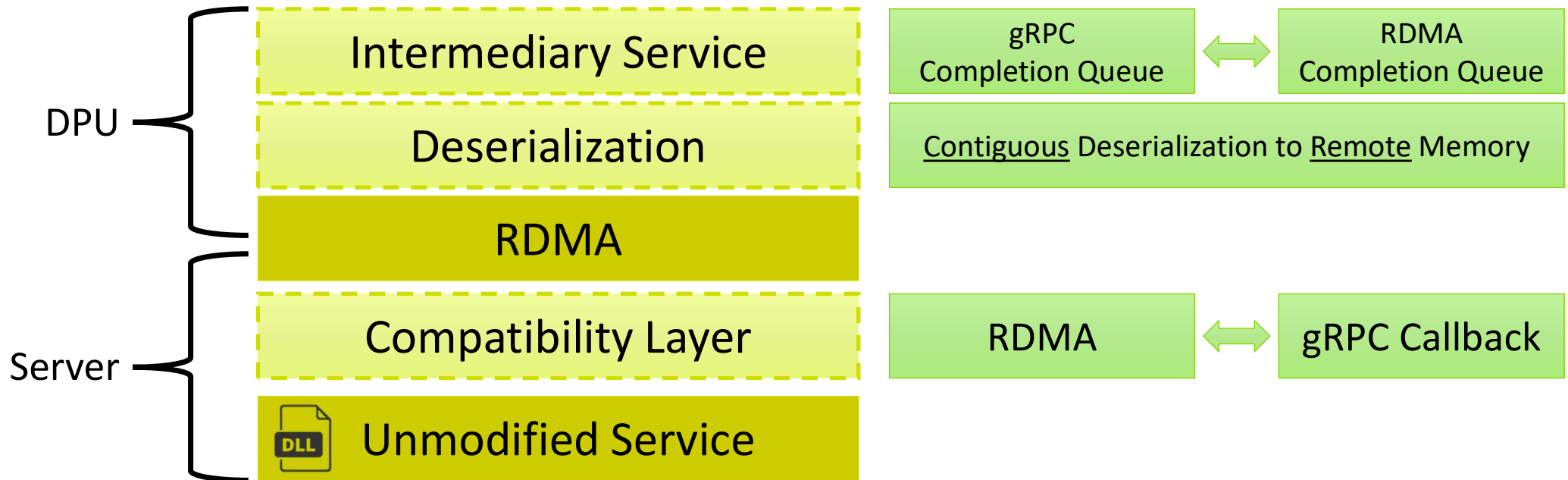
Questions

Protocol Buffer DPU Offloading in the RPC Datapath

Raphaël Frantz – r.r.a.frantz@tue.nl

Integrate the proposed Solution with gRPC

Implement once

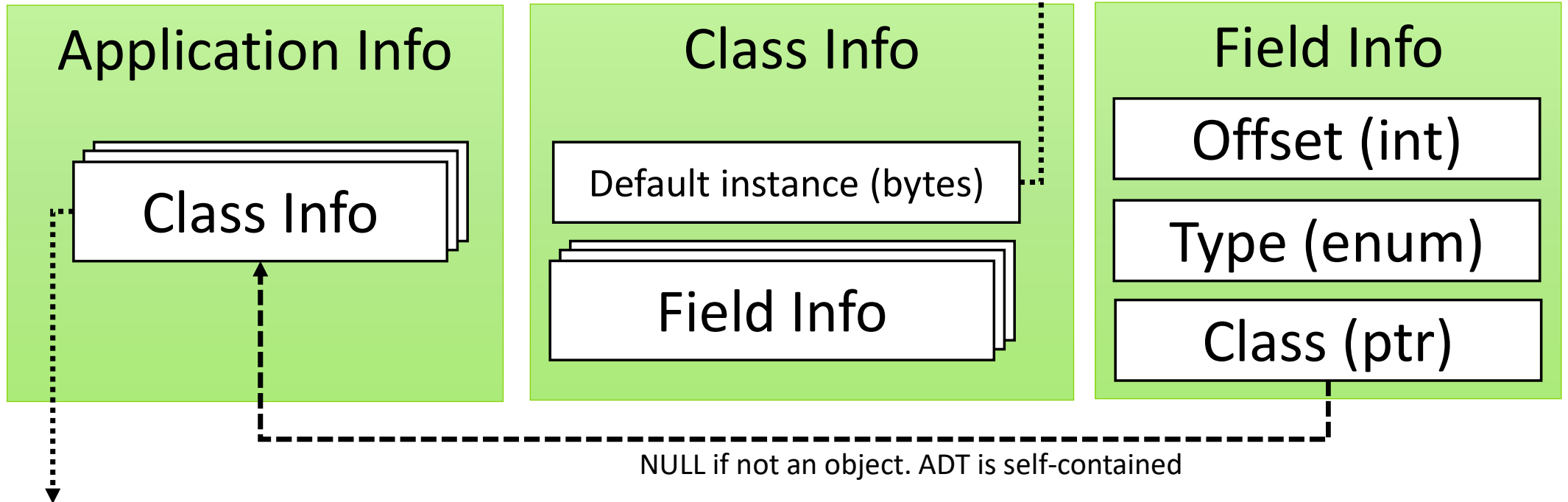


Accelerator Description Table (ADT)

1. Necessary information to deserialize from the DPU
2. Automatically generated by a **protoc** plugin

Defined in server memory space

vptr	fields
------	--------

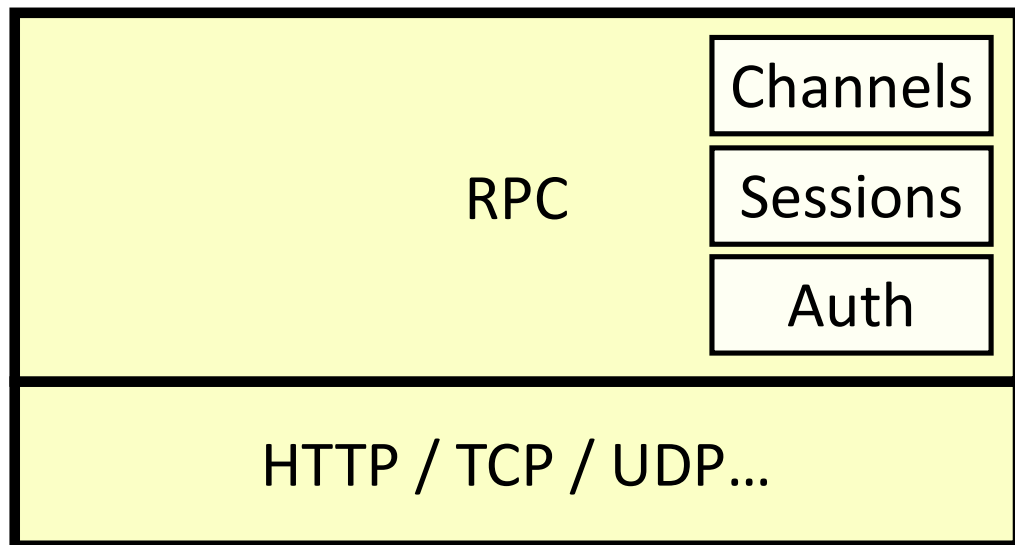


NULL if not an object. ADT is self-contained

All message classes used
by the application

How to communicate between services?

1. **Slow** to develop: TCP, HTTP...
2. **Fast** to develop: Remote Procedure Calls (RPCs)
 - Client/Server → Function call/Callback



client code

```
client.login();
```

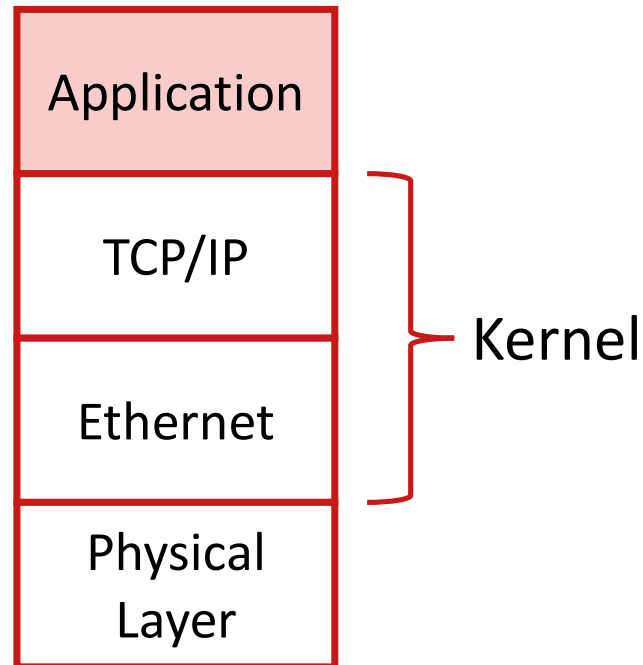
server code

```
login() { ... }
```

Offloading RPC Deser. – RDMA

1. *Kernel Bypass*
 - High bandwidth
 - Low latency
2. Competing Standards
 - Infiniband
 - RoCE
 - iWarp
3. Common API
 - libibverbs

Traditional
network layers



Network layers
with *RDMA*

