# Towards Faster Columnar Data Transport Using RDMA
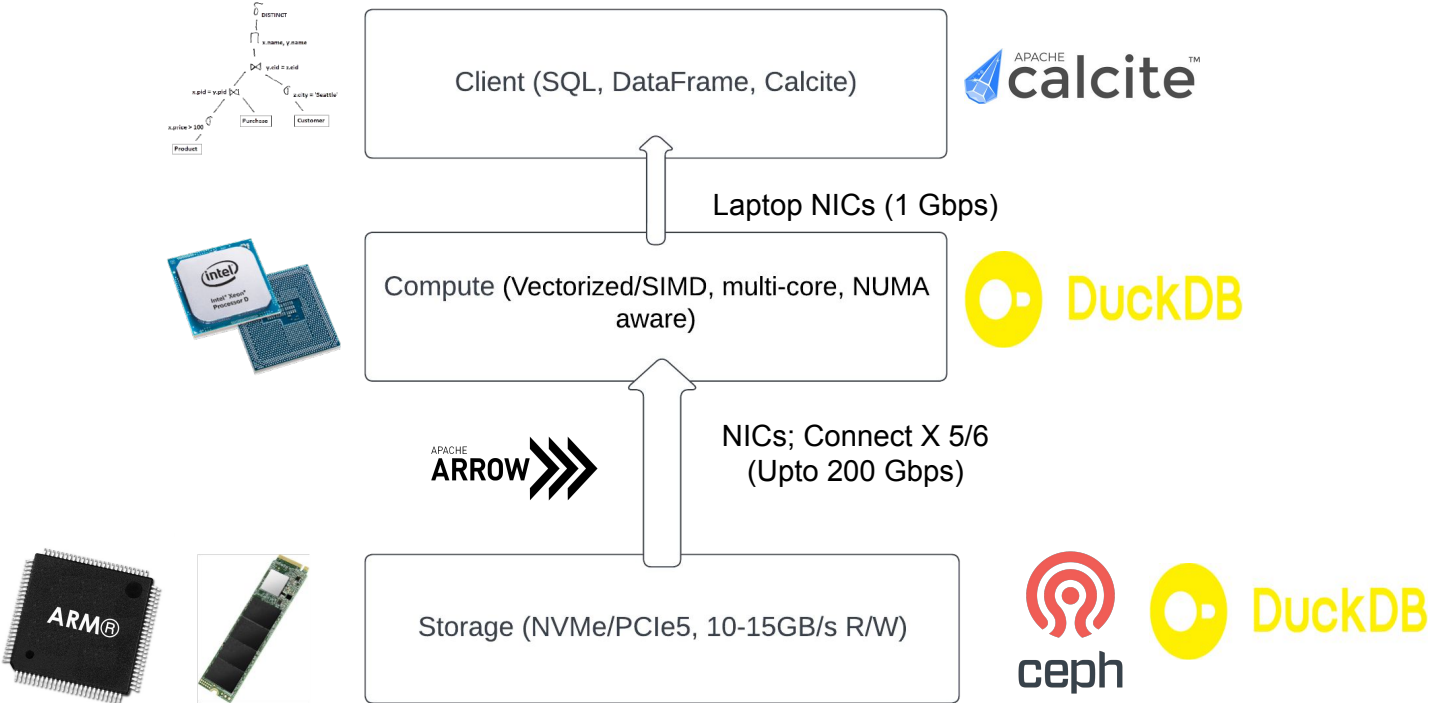
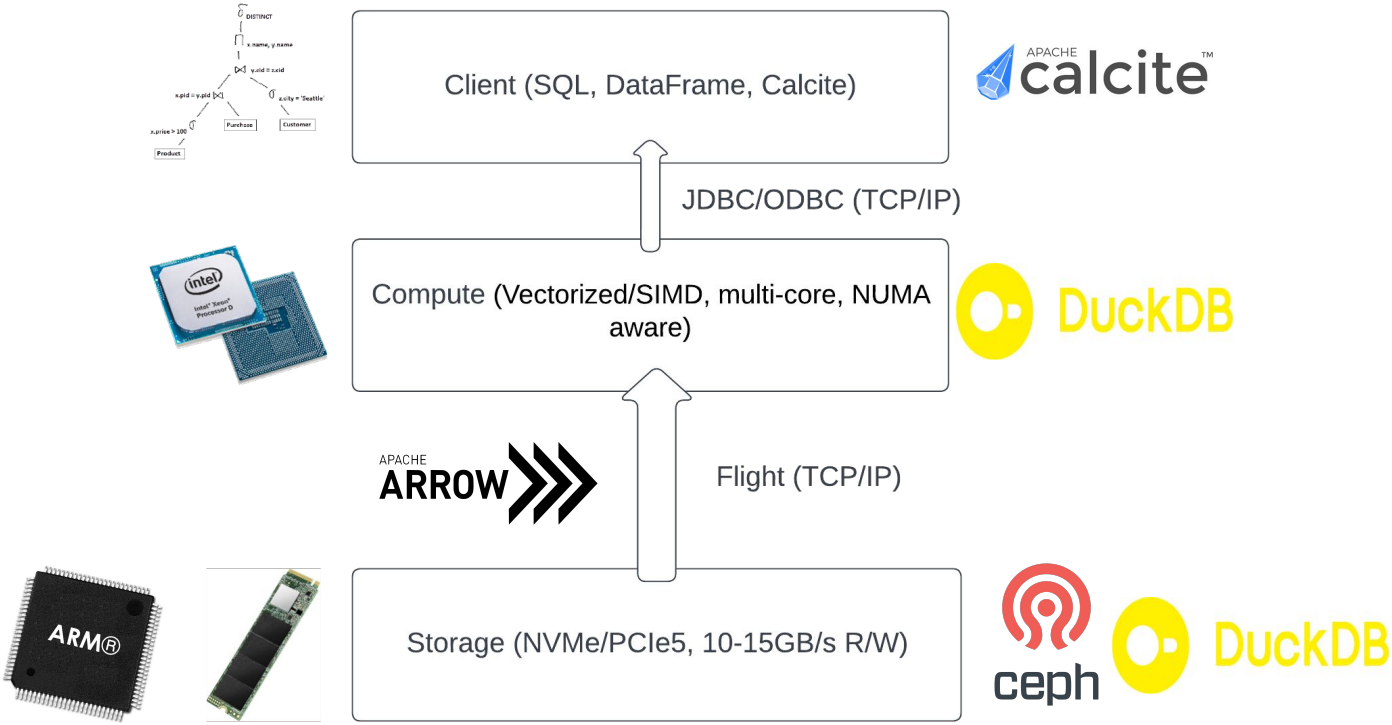Jayjeet Chakraborty

**UC Santa Cruz**

# Modern Datacenter Hardware

- **Fast memory devices**
  - NVMe
  - PCIe5, DDR5, CXL
- **Fast networking infrastructure**
  - ConnectX-5/6 NICs
  - Upto 200 Gbps bandwidth
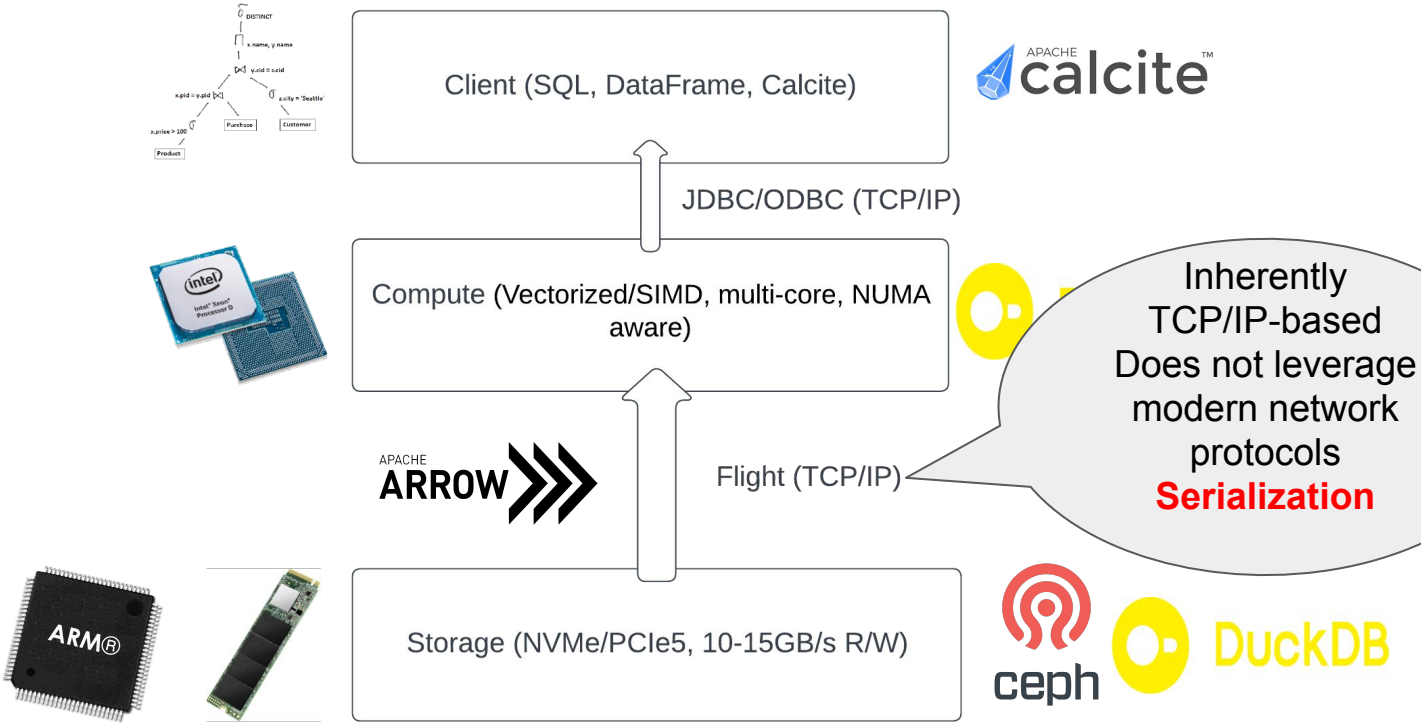- **Fat CPUs**
  - Intel Xeon
  - Intel Sapphire Rapids

# Data Processing Architecture using CS

# Data Processing Architecture using CS

# Data Processing Architecture using CS

# What is Serialization ?

**The process of converting 2D tables/record batches into network transferable format**
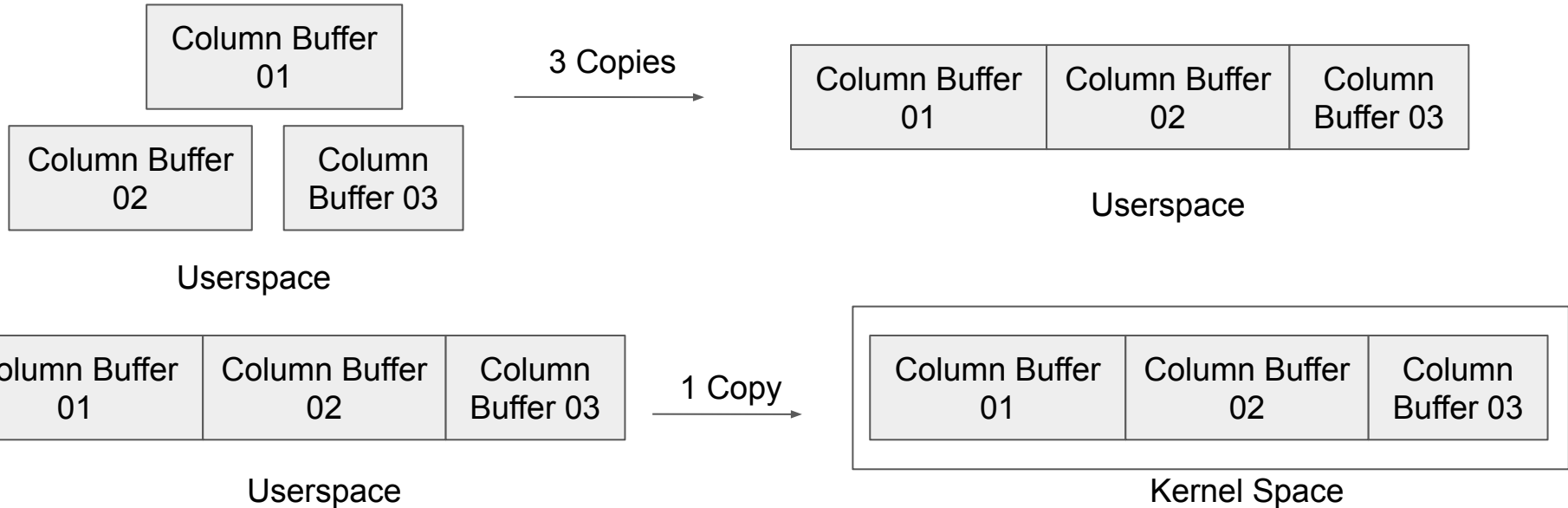
# What is Serialization ?

**Copy the individual buffers holding tabular data into a single-contiguous buffer as required by TCP/IP**

# What is Serialization ?

**Copy the individual buffers holding tabular data into a single-contiguous buffer as required by TCP/IP**

# Why is Serialization bad ?

- Unwanted memory copies
- Wastage of CPU cycles
- Added overhead for Computational storage
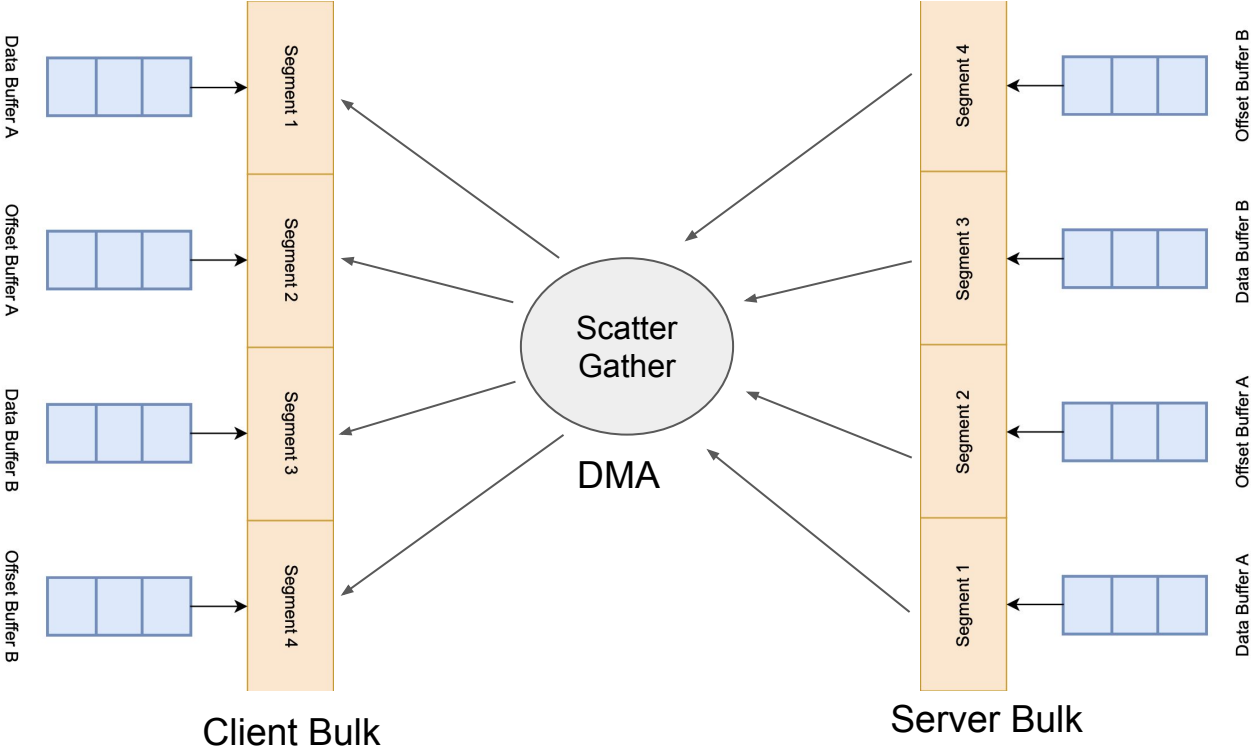
# Why is Serialization bad ?

- Unwanted memory copies
- Wastage of CPU cycles
- Added overhead for Computational Storage

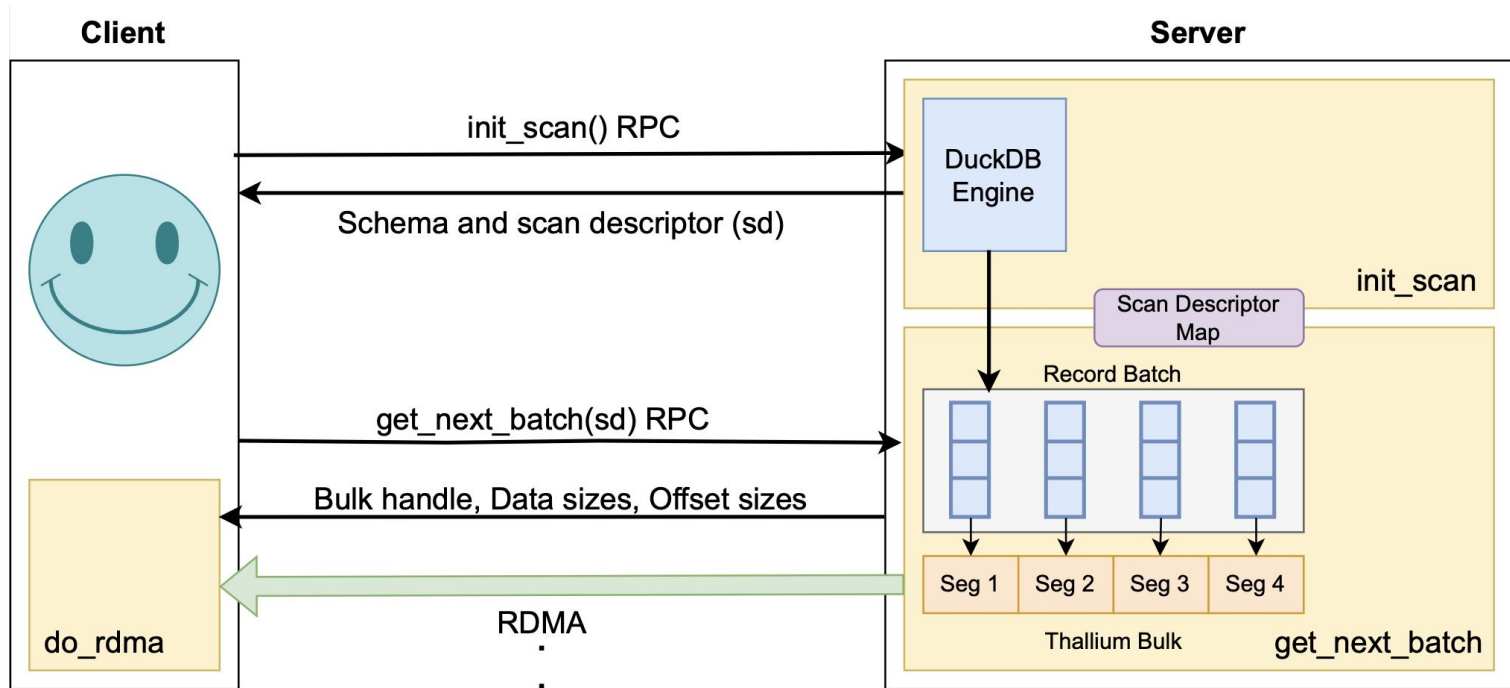## How much can we eliminate the serialization overhead ?

# Possible Solution

- Eliminate the multiple rounds of `memcpy`
- Use user-space networking libraries
- **Leverage HPC communication frameworks that leverage faster networking protocols**
  - [Mochi Thallium](#) (Argonne National Labs)
    - Supports Infiniband; VPI-enabled ConnectX cards has both Ethernet and Infiniband modes
    - Uses `user-space` RDMA libraries; libfabric and libibverbs
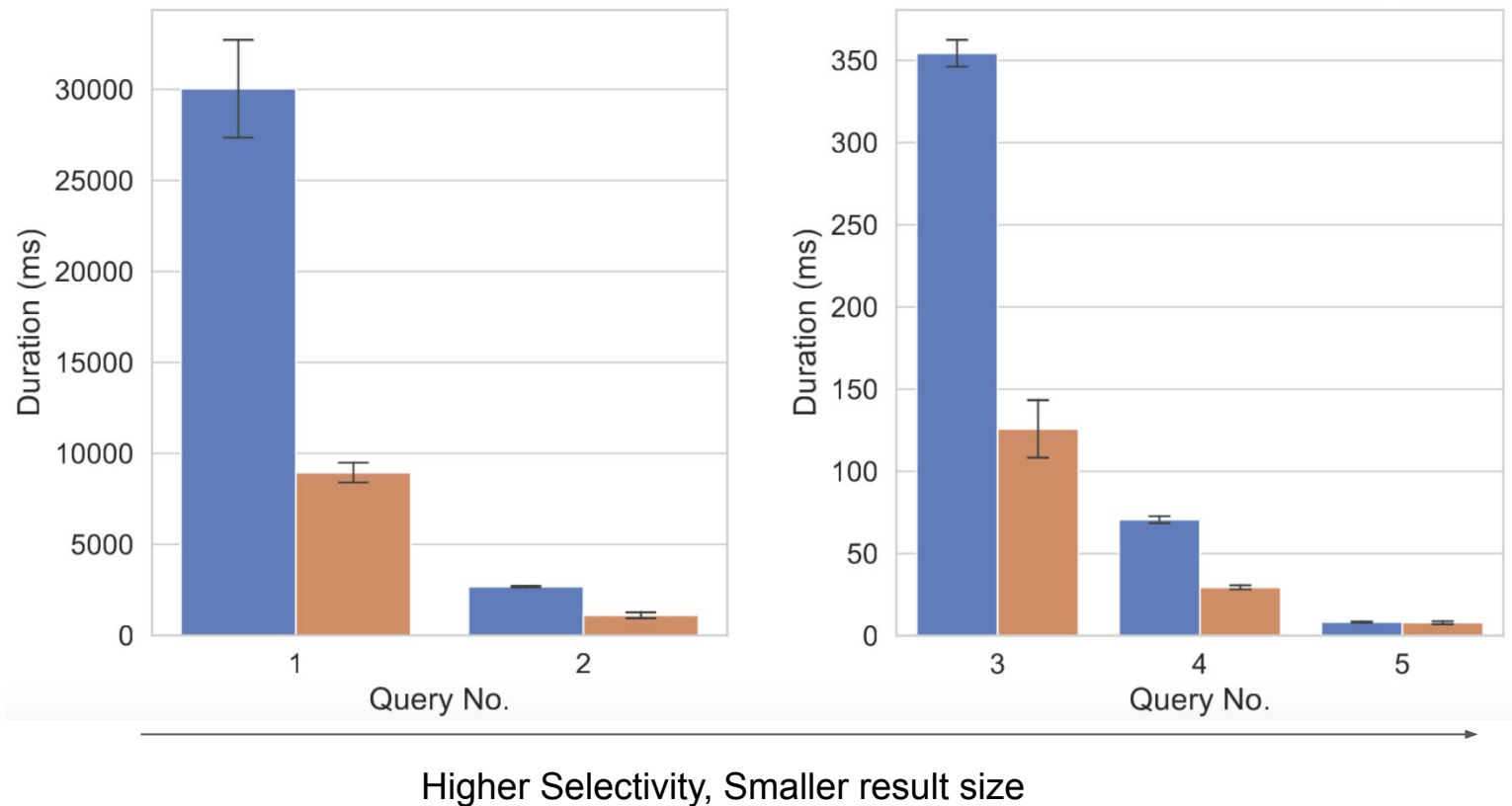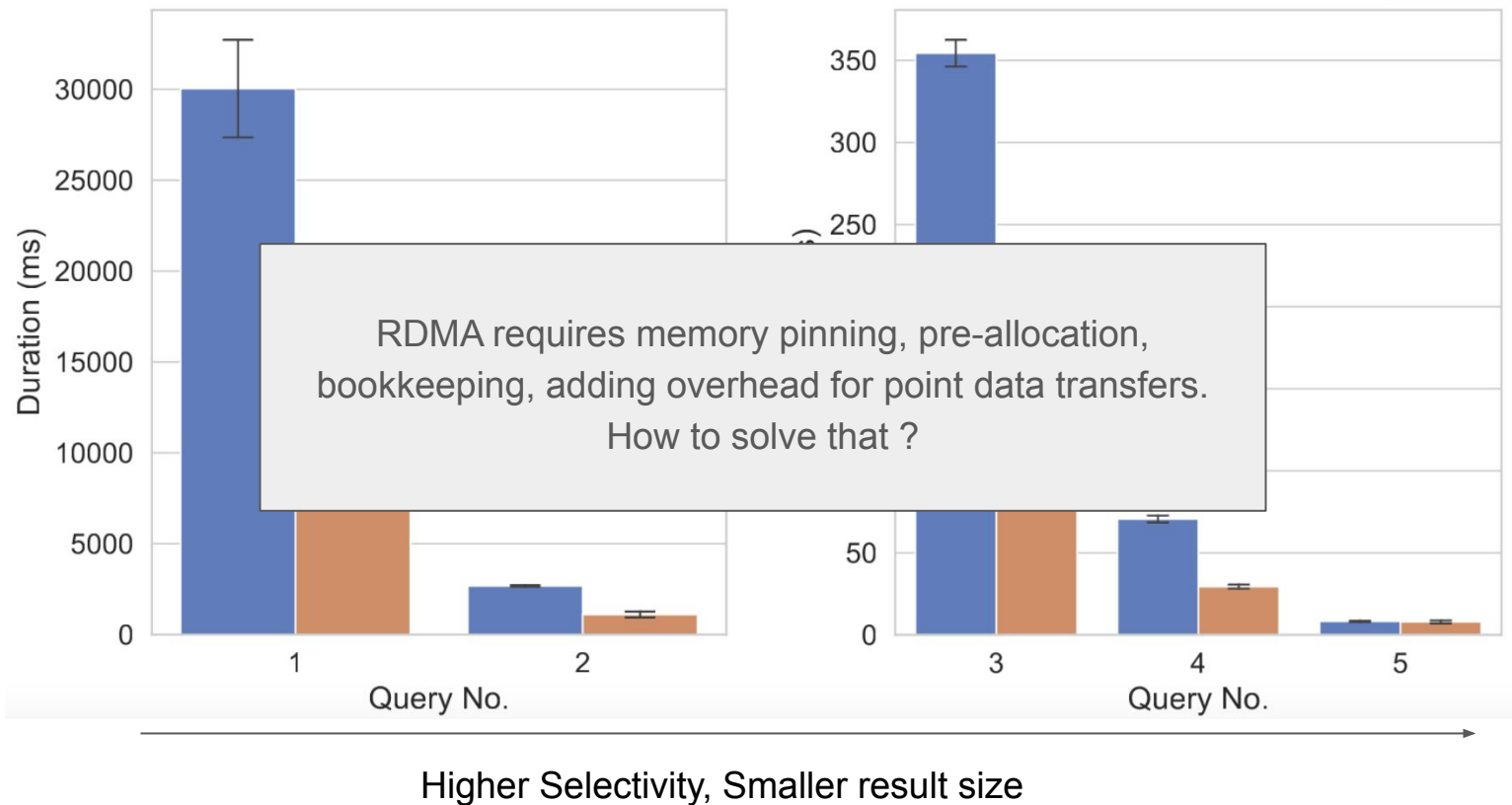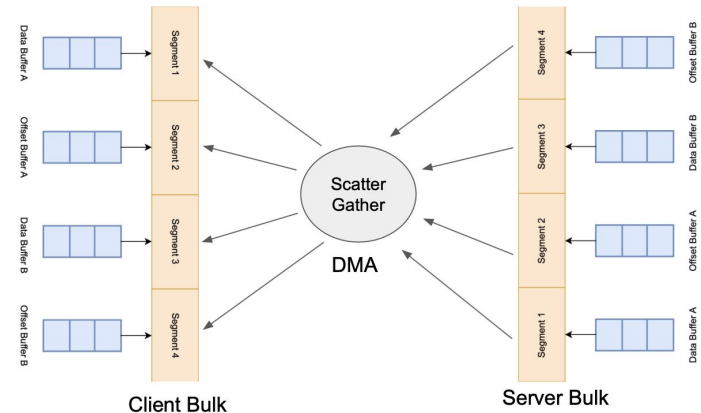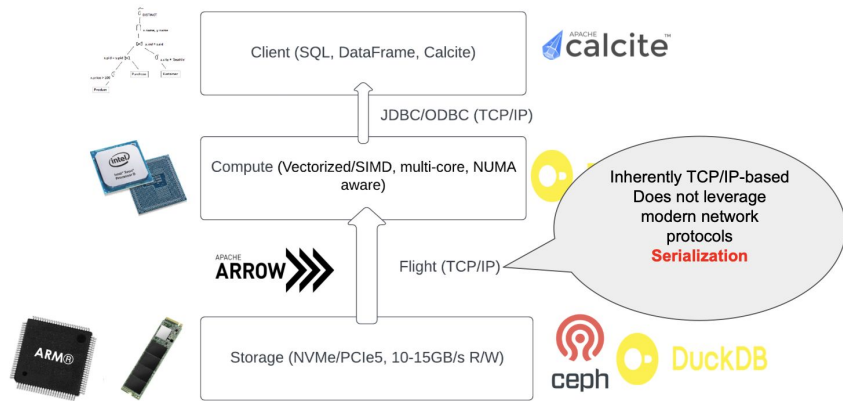
# Mochi Thallium

# Protocol Design

# Initial Evaluations (with DuckDB engine)



Higher Selectivity, Smaller result size

# Initial Evaluations (with DuckDB engine)



RDMA requires memory pinning, pre-allocation, bookkeeping, adding overhead for point data transfers. How to solve that ?

Higher Selectivity, Smaller result size

# Thank You ! (jayjeetc@ucsc.edu)