

Accelerated Packet Processing Based on Nvidia Bluefield DPU

Day Final - Team 6 NTHU_LSALAB

Team Name

Team 6 NTHU_LSALAB

Mentors

- Kevin Chen
- SungTa Tsai

Members

Dr. Jerry (Chi-Yuan) Chou, NTHU
(Large-scale System Lab)



Ivan Ou, Realtek



Aiden Huang, NTHU
(Speaker)

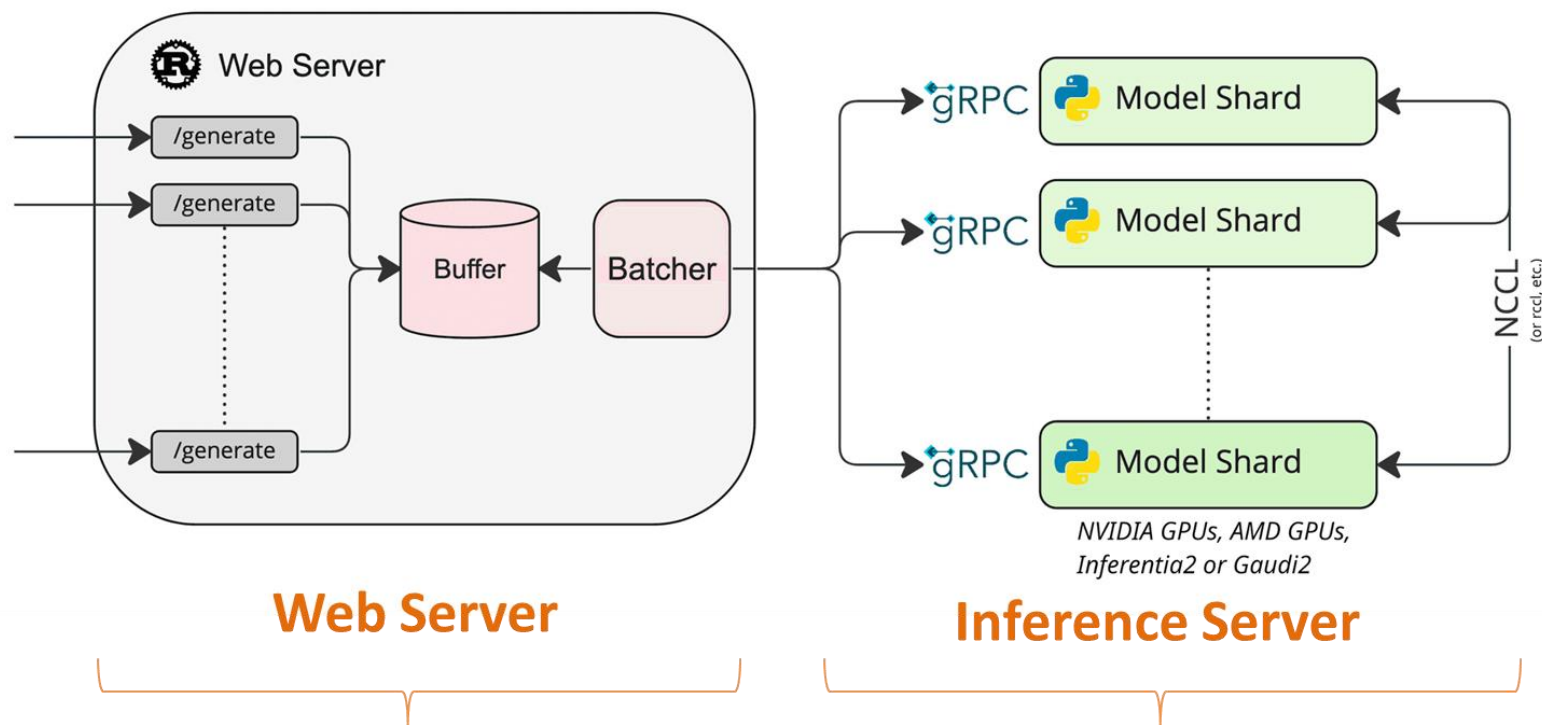


Samuel Ke, CHTTL



Can we accelerate inference server?

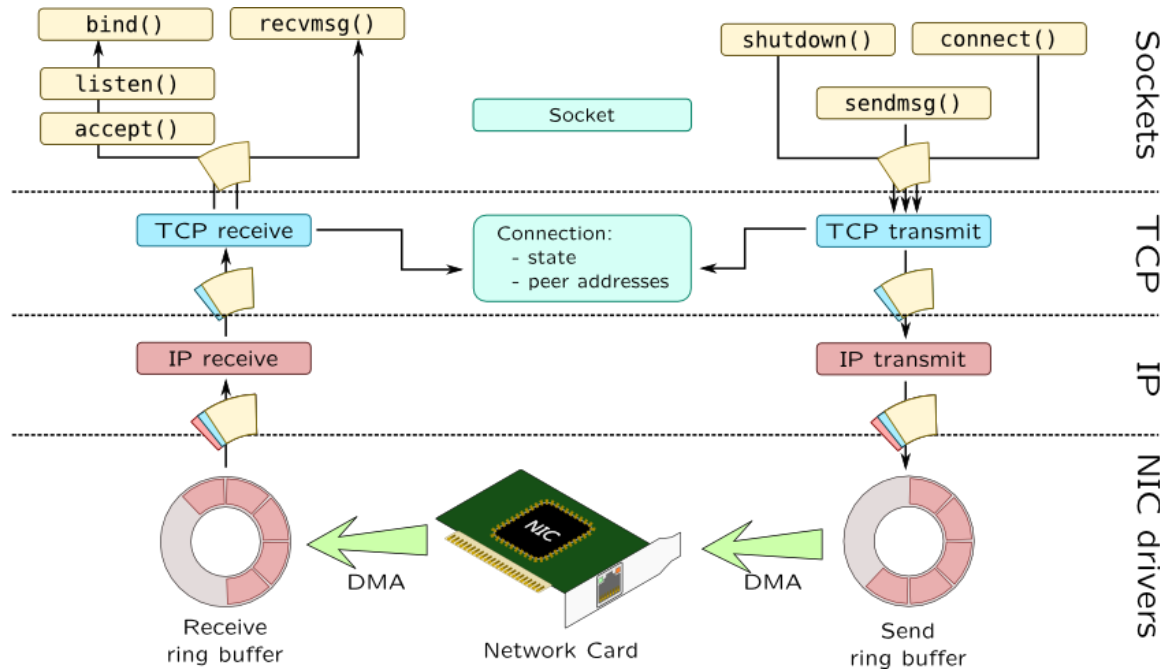
Problem the team is trying to solve.



Can we accelerate the inference process?

CPU Packet Processing

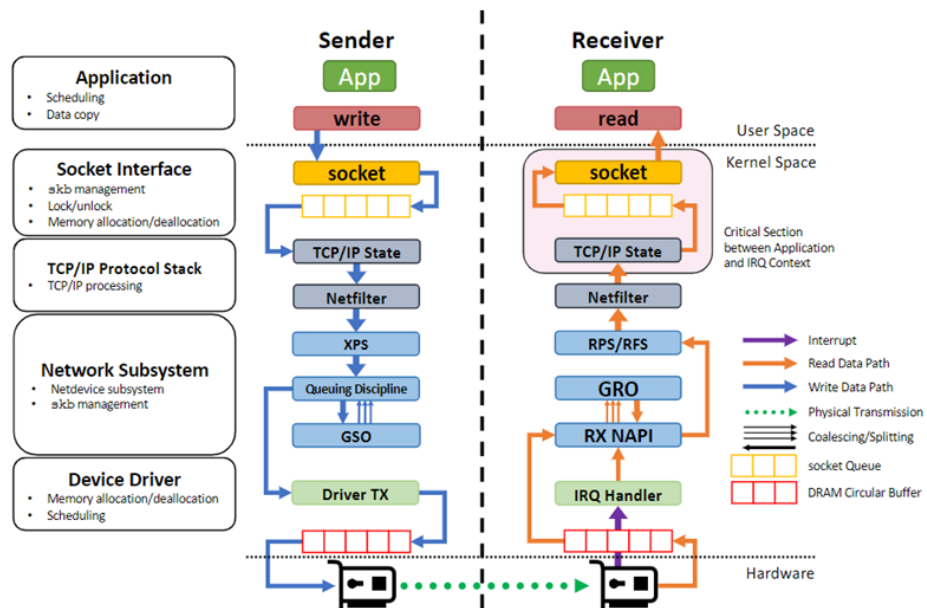
Packet processing refers to the operations performed by a processing unit to handle network packets, including tasks such as packet forwarding, classification, inspection, and transformation based on specific protocols or rules.



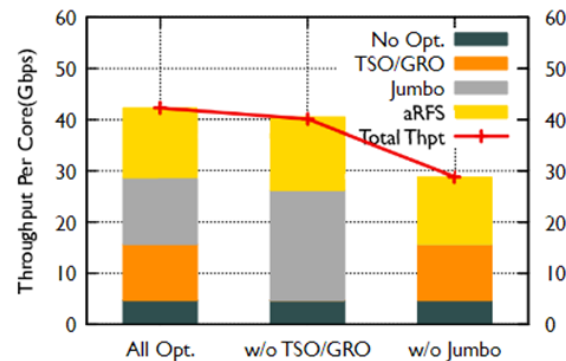
- The most commonly used processing unit for network tasks is the CPU
- Which works in conjunction with the network subsystem.

The limitations of CPU

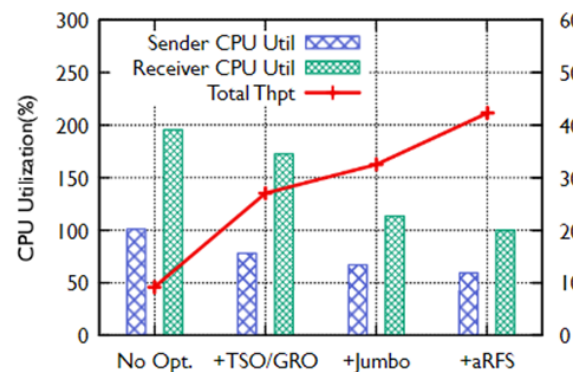
Process a packet is a heavy task.



The era of 100Gbps is coming.



(a) Throughput-per-core (Gbps)

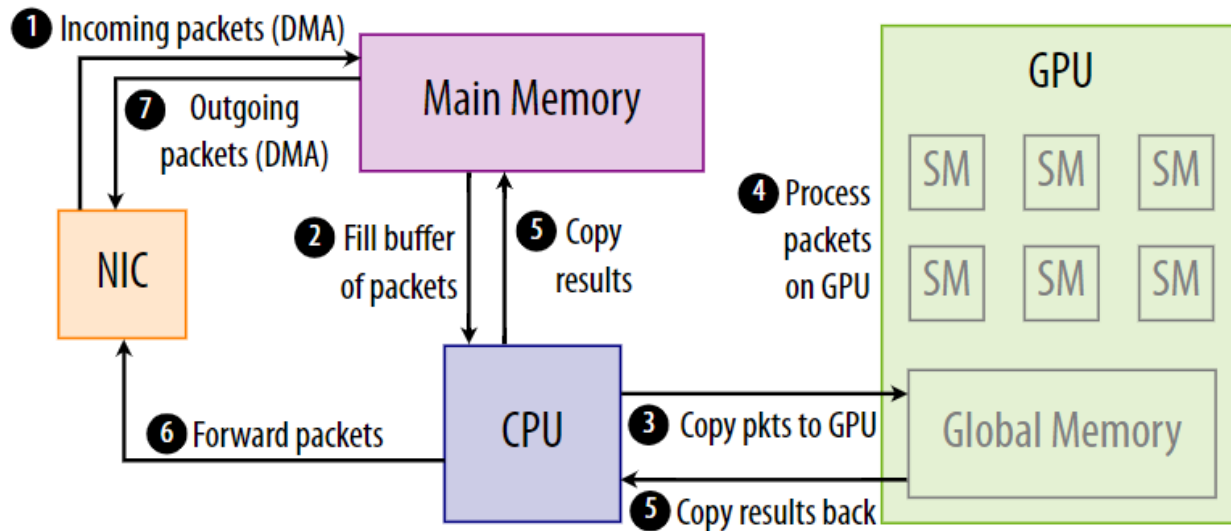


(b) CPU utilization (%)

⇒ A single core is no longer sufficient.

⇒ Receiver-side CPU is the bottleneck.

Packet Processing on the GPU

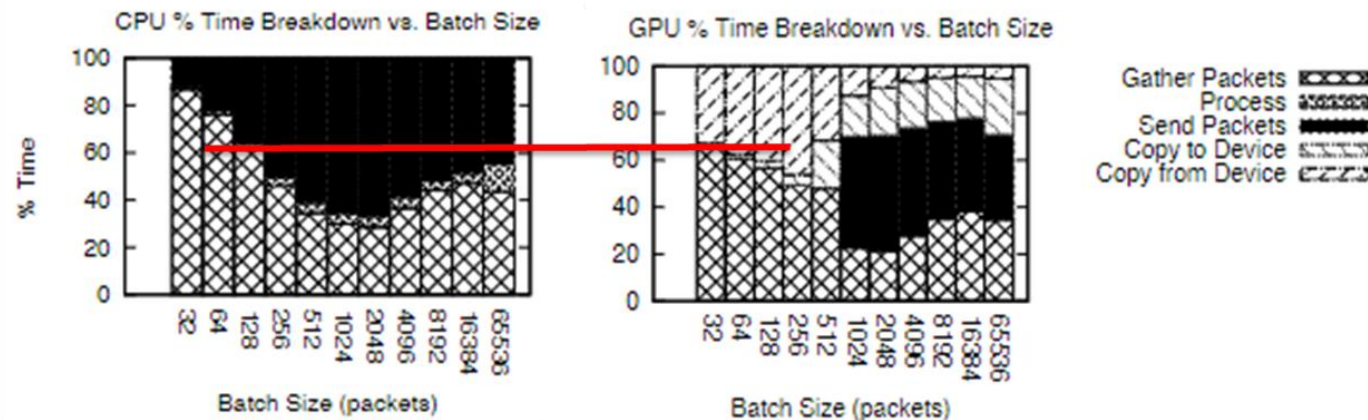
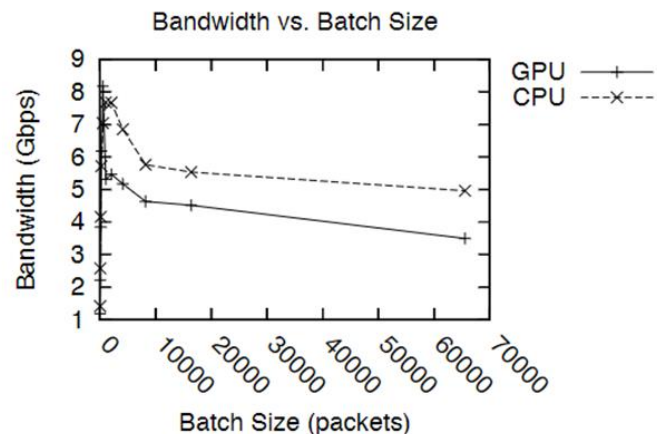


Since GPU are strong in processing, the throughput are willing to be raised.



CPU-GPU Based Processing

The devil is in the experiments.



Although it spends less time processing, **it has to copy packets to the GPU for processing and copy the results back.**



What if, we can transfer packets without CPU?

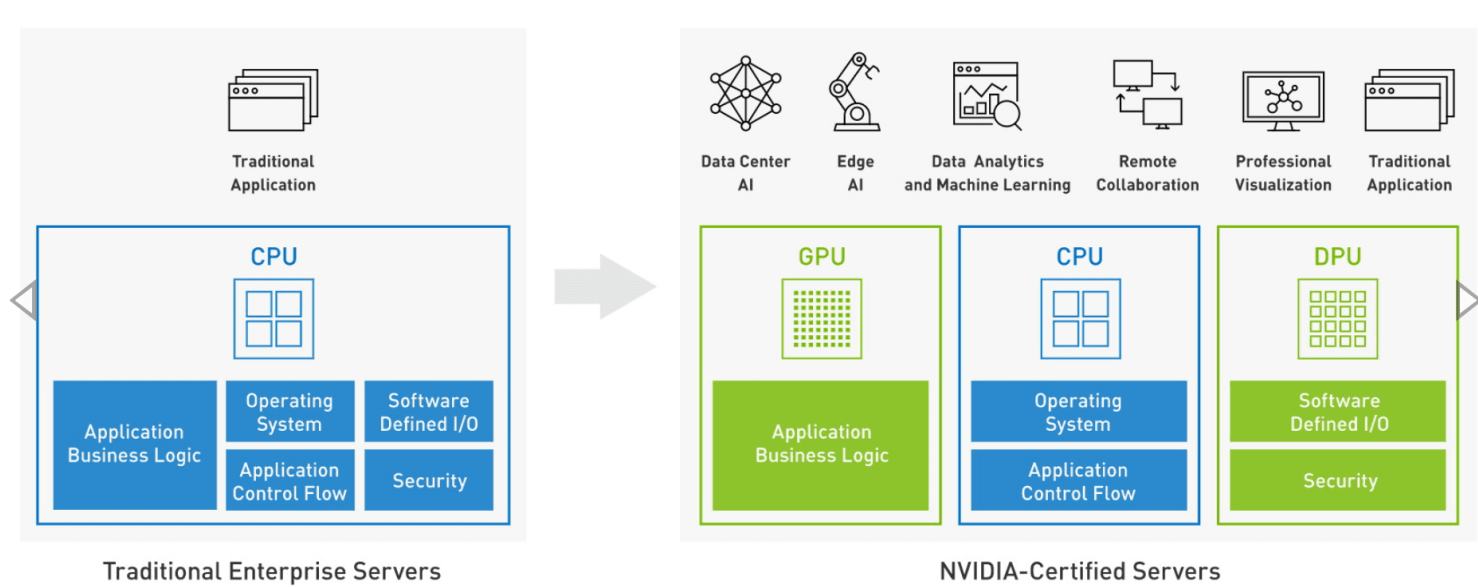


Source: https://www.cs.cmu.edu/~bvavala/misc/project740/15-740_Project_files/Report.pdf

Data Process Unit

What was your initial strategy?

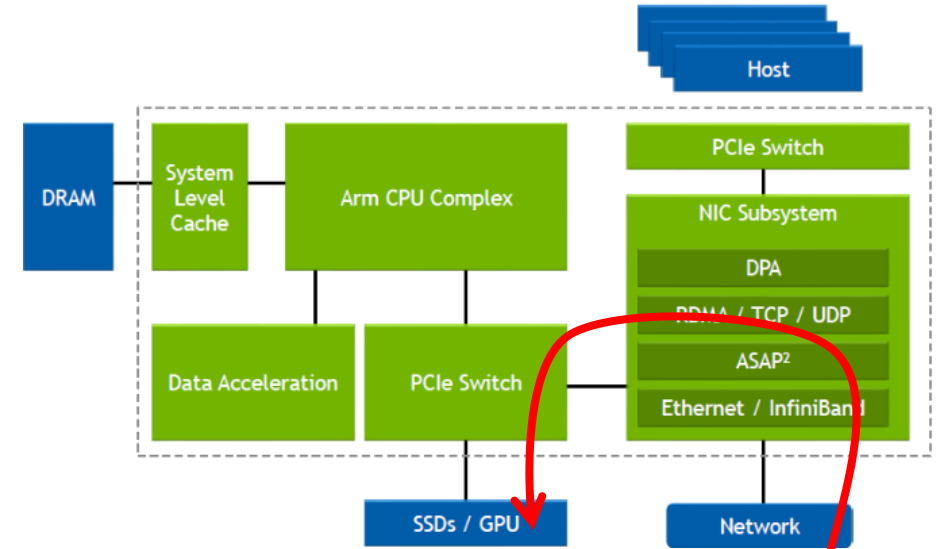
- ❑ Data Processing Unit (DPU) is a pivotal component in the realm of advanced computing architectures.
- ❑ A new class of programmable processor and will join CPUs and GPUs as one of the three pillars of computing.



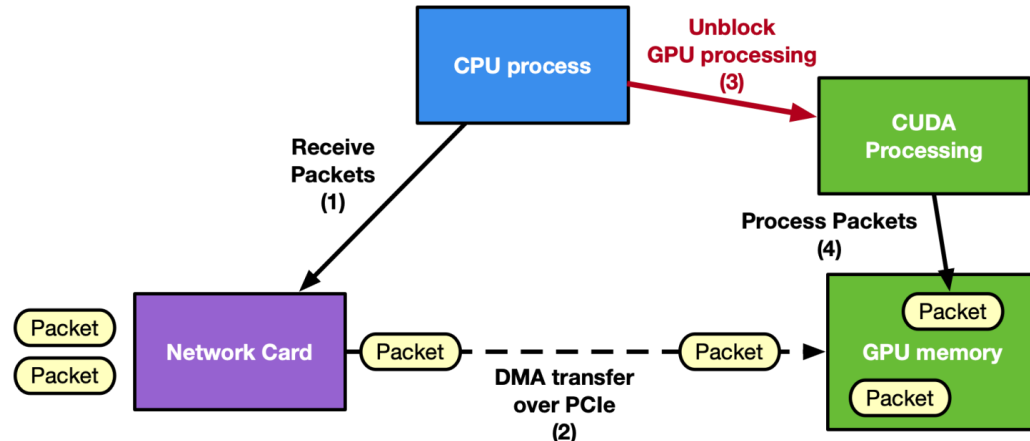
- An enhanced SmartNIC for Network Processing.
- A bridge between NIC to CPU/GPU.

DPU System Architecture

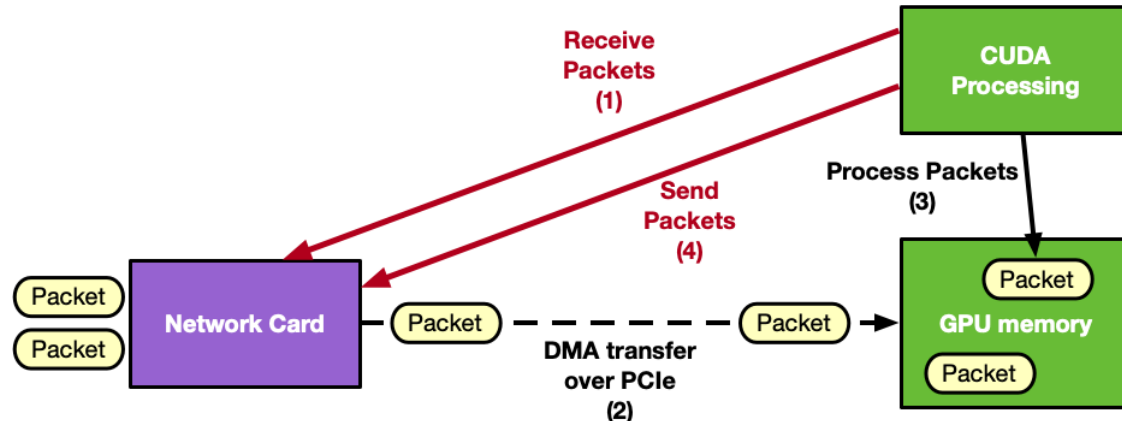
- **DPU utilized the same subsystem from existing kernel.**
- **The best part of DPU is programmable.**



Before: CPU-Centric



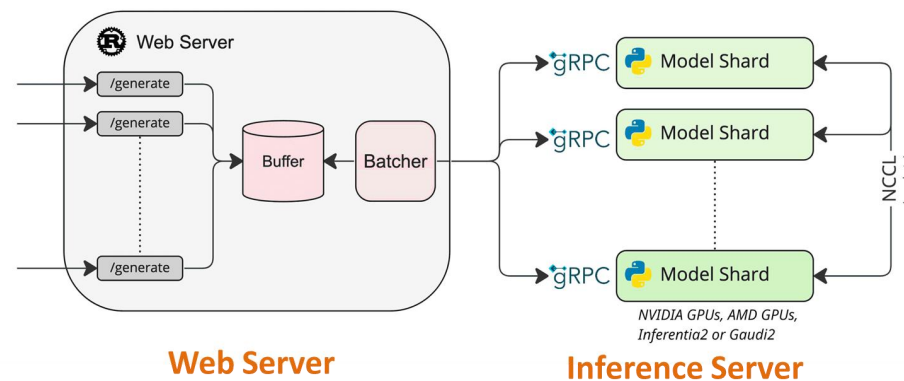
After: GPU-Centric



The devil is in the experiments again.

What was your initial strategy?

Problem: Web Server consumes CPU resource on inferencing.

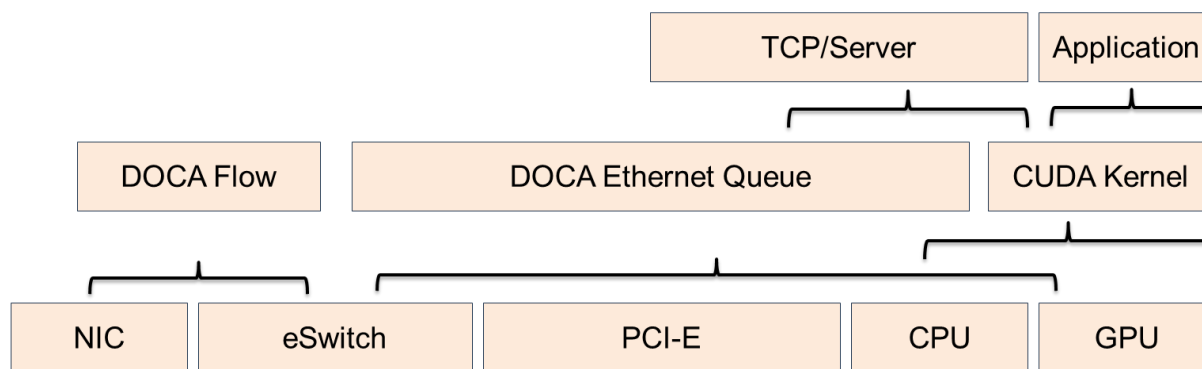


Easy, DPU handles everything ?

Application

Middleware

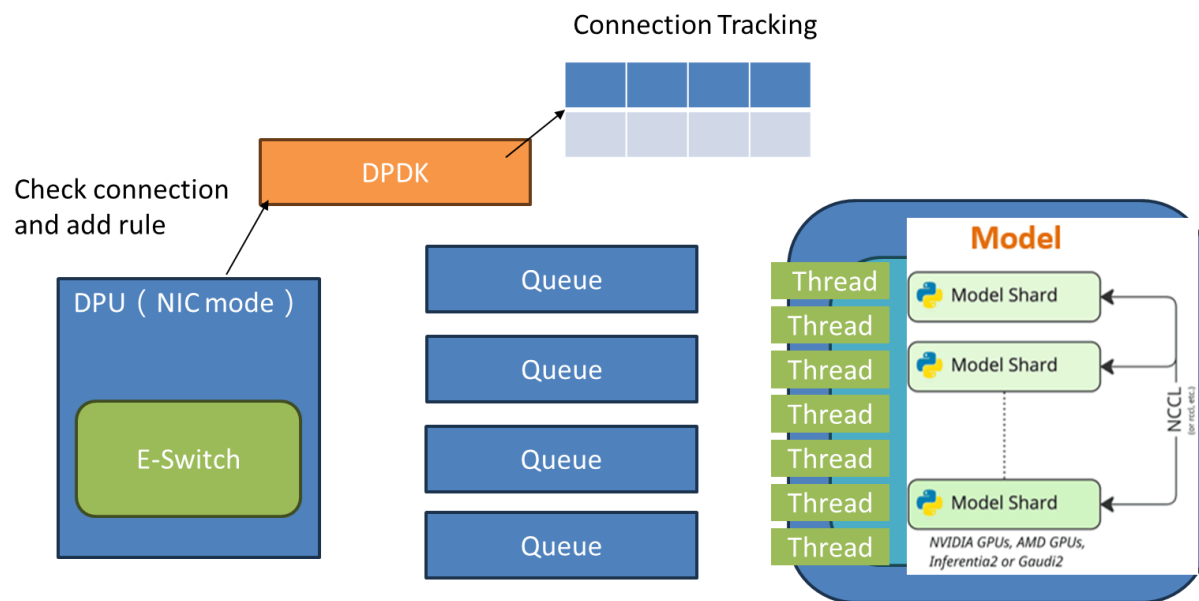
Infrastructure



DPU gives you only the possibility.

The devil is in the experiments again.

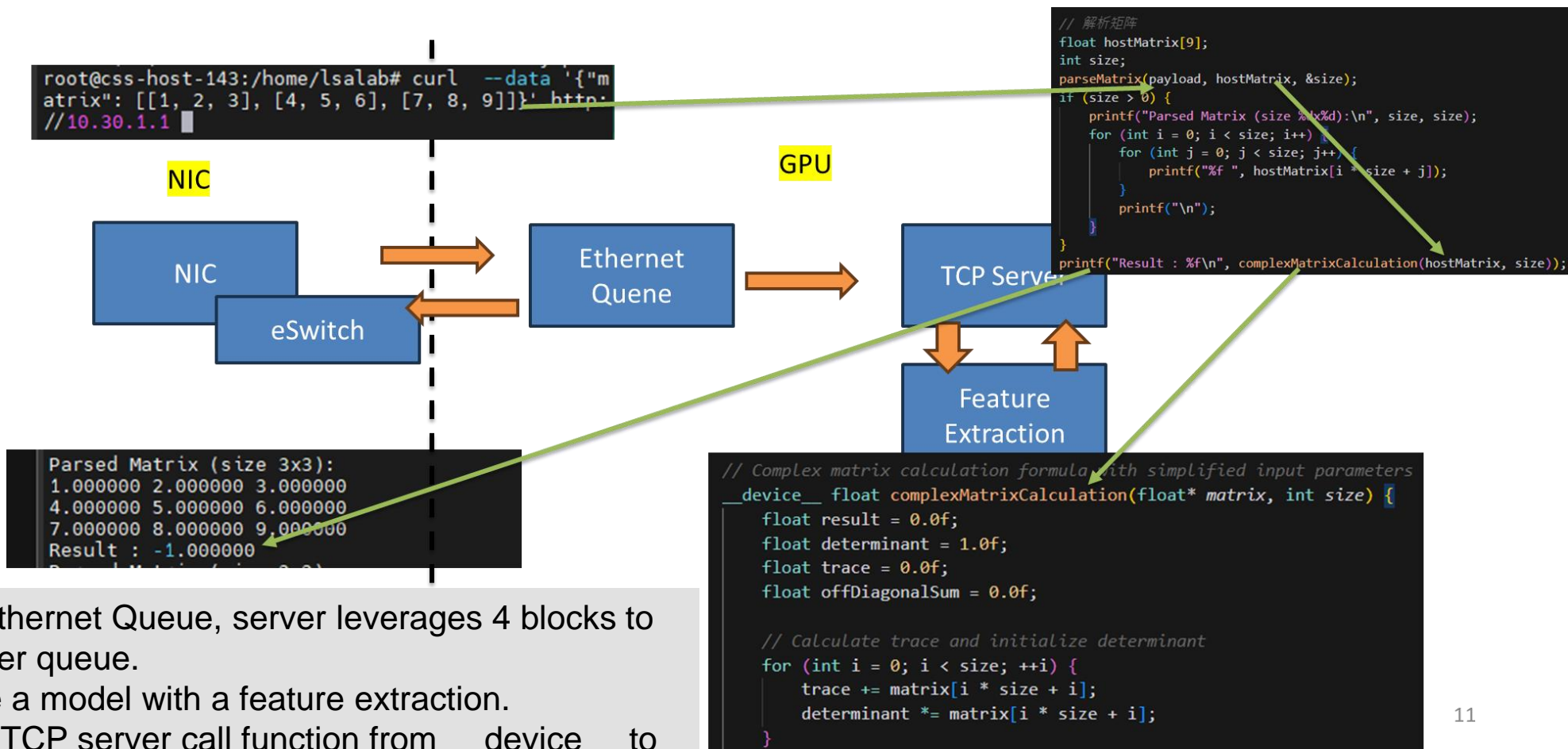
How did this strategy change?



1. Make GPU receive TCP/UDP based packets.
2. Conduct a model function on GPU
3. Transfer the TCP payload to the GPU and measure its processing efficiency

Build from Scratch

What were you able to accomplish?

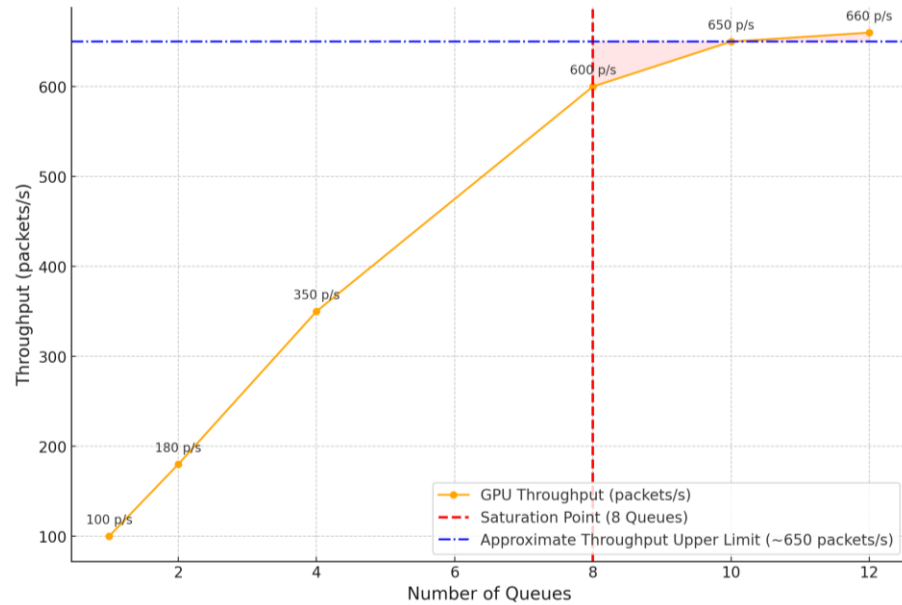


- For each Ethernet Queue, server leverages 4 blocks to consume per queue.
- We replace a model with a feature extraction.
- On CUDA, TCP server call function from __device__ to __device__ without Dynamic Parallelism

Accelerated Results

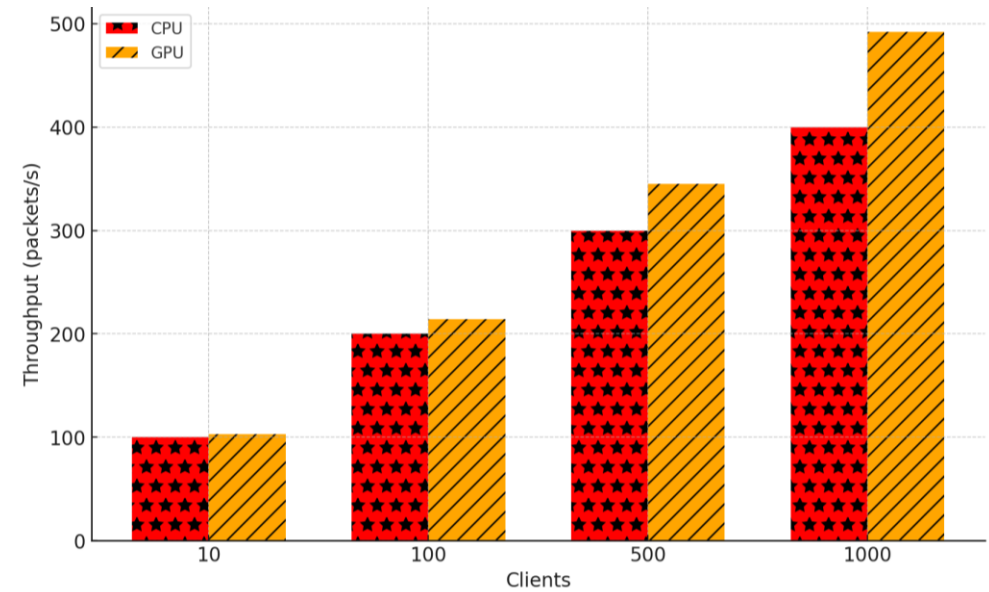
Did you achieve a speed up?

GPU Packets processing ability



GPUNetIO reaches the highest throughput with 8 packets. (no optimization)

Packet Processing Results

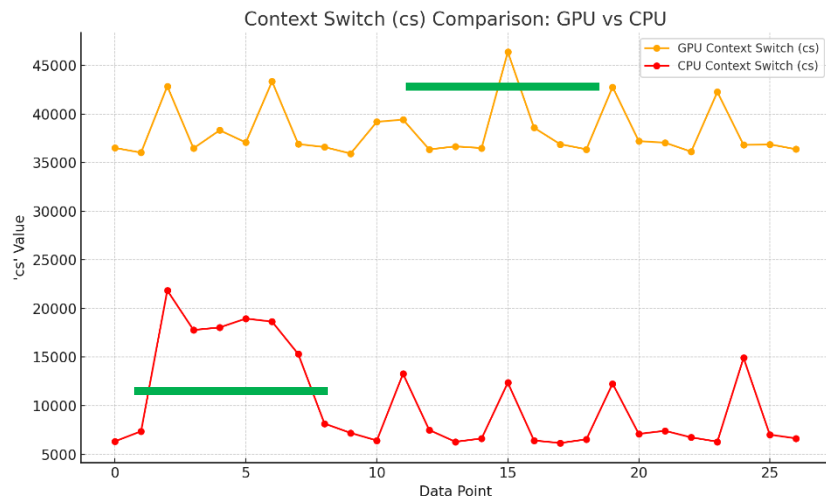
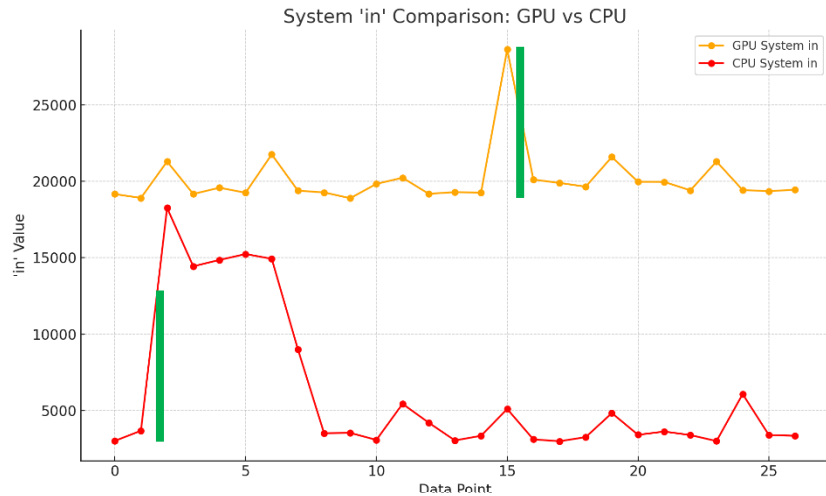


With 4 queues, speed up to ~23%

Accelerated Results

Did you achieve new scientific goals?

CPU Usage (from interrupt rate)



- GPUNetIO triggers interrupts during the middle of the processing flow, whereas the CPU triggers them at the beginning.
- GPUNetIO-based packet processing occupies a CPU core for DPDK, which seems to consume more resources.
- However, we can clearly observe that the throughput is lower compared to CPU-based packet processing.
- The context switch rate also observed higher CPU, in our estimation, we save 110w generally.

Conclusion

- Due to environmental factors, this experiment primarily used DPU BlueField-2 and DOCA 2.5. However, BlueField-3 is already available on the market, and GPU Direct transfer technology has become more mature and diverse.
- Although there are still some overheads in the Hackathon results, it is undoubtedly an architecture worth focusing on in the foreseeable future.
- Shout out to Kevin and Sungta for their amazing support and guidance!
- Thank you to Nvidia for providing this opportunity to communicate and discuss with frontline professionals.



Thank You

OpenACC
More Science. Less Programming