

SRAP: Sender-side Receiver-Aware Port selection for High-Speed Multi-Flow TCP

Shingo Hattori¹, Osamu Tatebe² (advisor)

1. Graduate School of Science and Technology, University of Tsukuba
2. Center for Computational Sciences, University of Tsukuba



筑波大学
University of Tsukuba



St. Louis, MO
hpc ignites.

Abstract

Achieving **400 Gbps** requires aggregating **multiple flows** across cores rather than pushing single-flow limits. However, with **16x25 Gbps TCP** flows, random ephemeral ports cause receiver-side packet steering (RSS) to concentrate flows on single CPU cores, degrading throughput from **25 Gbps** to below **5 Gbps**.

Unlike receiver-side solutions that suffer from cache misses and state migration overhead, **SRAP** transparently controls source ports at senders, ensuring **collision-free 1:1 mapping** without runtime remapping costs or application modification.

With **16x25 Gbps** flows, random assignment achieves optimal distribution with probability 1.04%, causing throughput to vary from 44.8 to 395 Gbps, while our approach consistently delivers **23.3-25 Gbps per flow** with guaranteed 1:1 flow-to-core mapping.

1. Introduction

How to achieve 400G throughput with pacing...? Aggregating multiple flows, but there are difficulties.

Single-flow throughput is limited by sequential processing and CPU clock speeds. Aggregating multiple flows across multiple cores provides a scalable path—achieving 400 Gbps requires distributing load across multiple paced flows (e.g., 16x25 Gbps).

Pacing prevents microbursts and packet loss, but when RSS maps multiple paced flows to the same CPU core, maintaining target-rate pacing becomes impossible.

We propose Source Port Assignment Protocol (SRAP), a sender-side approach that prevents flow collisions through intelligent source port selection, achieving 9x better worst-case throughput than random assignment while avoiding dynamic remapping overhead.

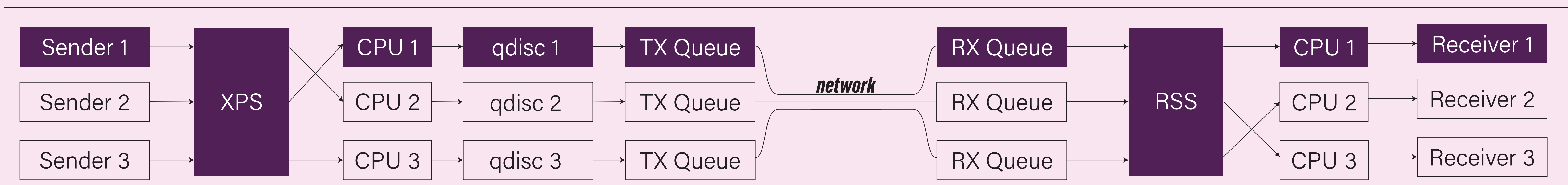


Figure 1: Linux Packet send/recv Flow

2. Problem Analysis

RSS's random flow-to-core mapping destroys packet pacing.

RSS distributes packets across CPU cores by hashing 4-tuples to map flows to receive queues. Random ephemeral ports create unpredictable flow-to-core mappings - with 16 flows and 32 cores, collision-free probability is only 1.04%.

Comparison: RSS++ suffers from cache misses and state migration overhead. ARFS requires kernel modifications. SRAP avoids these issues.

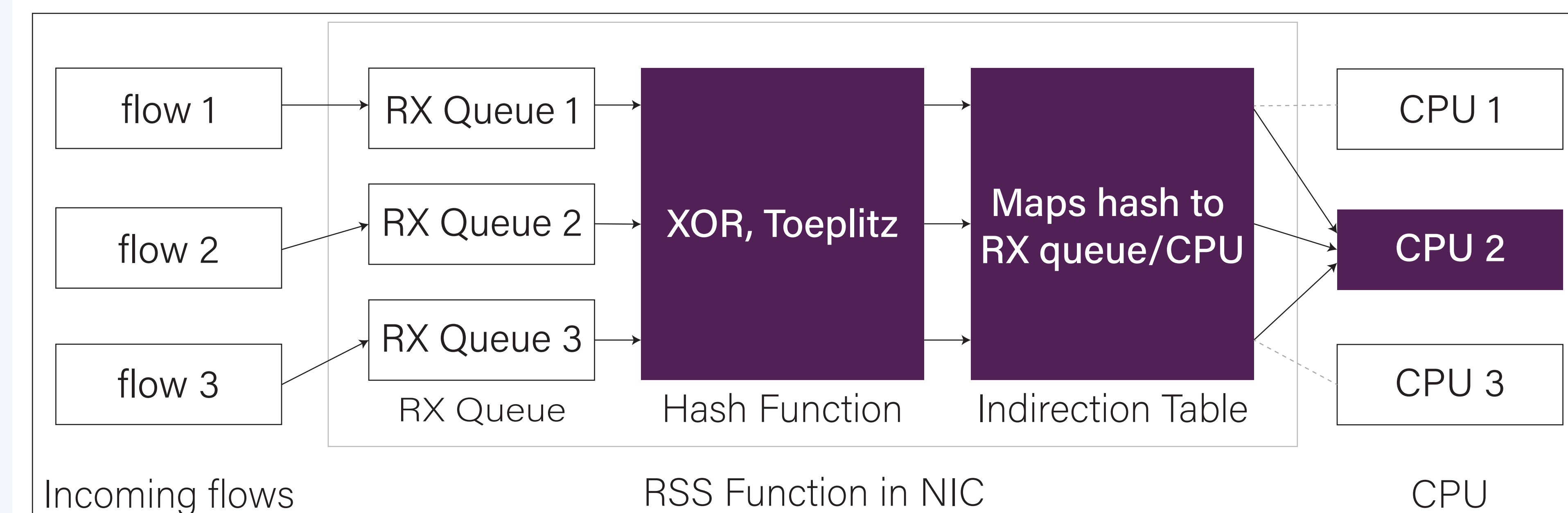


Figure 2: Packet reception flow from NIC to CPU

Flow collisions limit throughput—worst case: 16 flows on one core, <5 Gbps each. Random ports cause unpredictable collisions, making aggregation unreliable. Sender CPUs are controllable but receiver RSS is opaque. We enable predictable mappings via source ports.

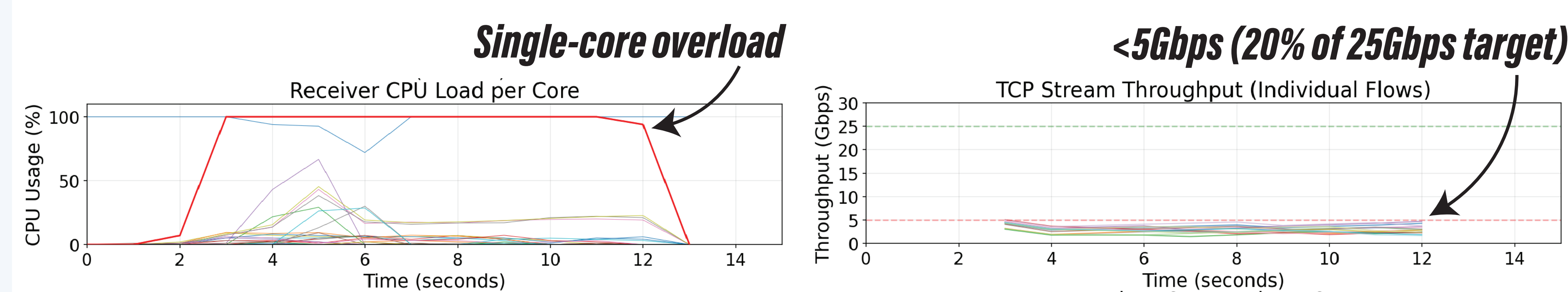


Figure 3: Performance with previous method

3. Proposed Approach

Select source ports at sender to guarantee collision-free mapping.

SRAP controls receiver-side flow distribution by selecting source ports that hash to different CPU cores through three coordinated components.

- RSS Information Exchange:** Exchanges RSS configurations (indirection tables, hash functions) with receivers—no app core binding needed, only RSS table matters.
- Port-to-CPU Predictor:** Computes RSS hashes and returns collision-free ports.
- Transparent Port Binder:** Intercepts connect() calls transparently, requests collision-free ports, and calls bind() to control the RSS 4-tuple without application modification.

By controlling source ports, SRAP prevents receiver-side flow collisions, achieving predictable multi-flow pacing without modifying applications or receivers.

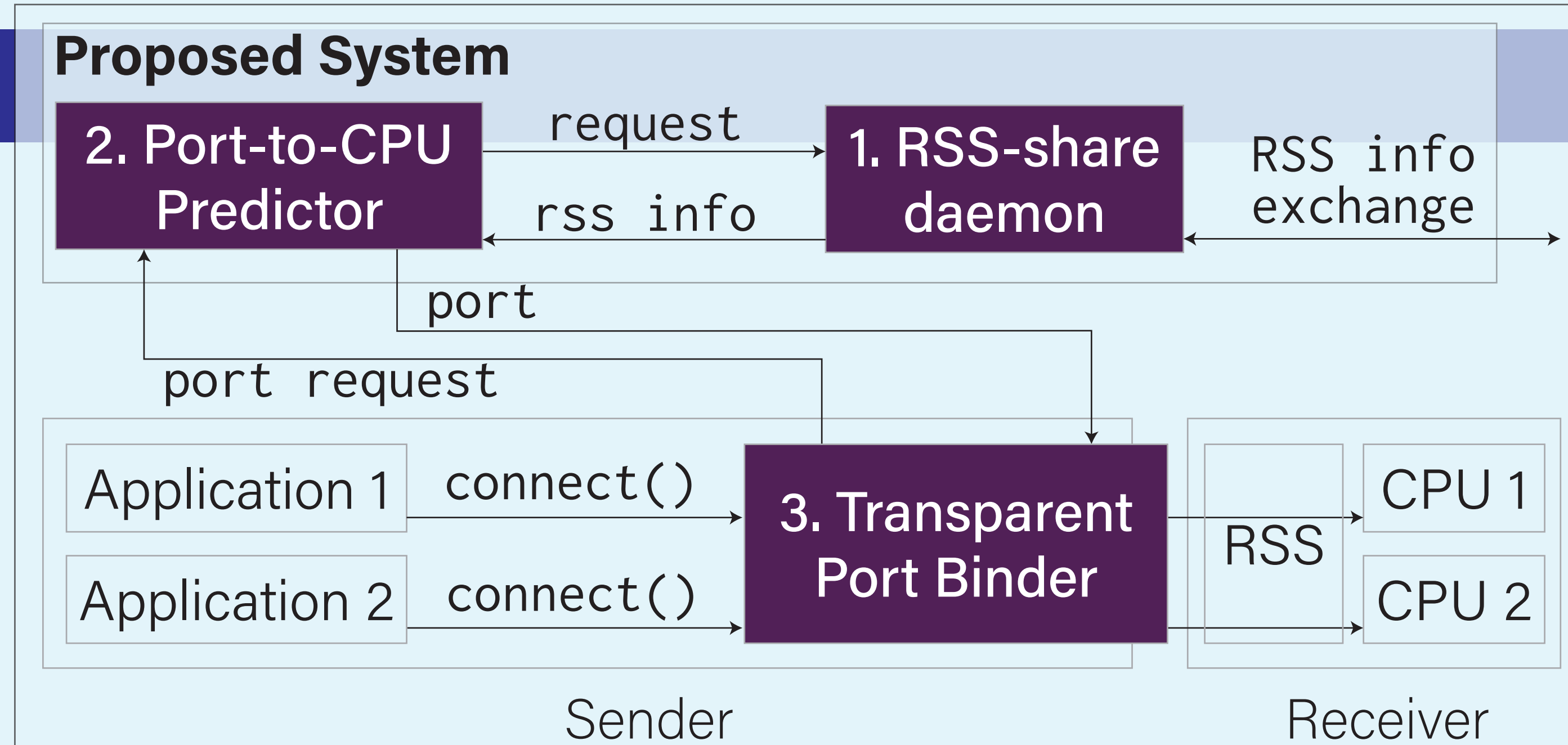


Figure 4: Proposed system architecture

4. Evaluation

Successfully achieve 400 Gbps aggregate with pacing.

Setup : AMD EPYC 9354P, 32 cores with DDR5 256GB and ConnectX-7 400GbE NICs running Linux 6.8.0. 16x25 Gbps TCP CUBIC flows with Linux FQ₃ pacing, iperf3 zero-copy optimizations, and 100MB windows.

Results : Figure 4 demonstrates SRAP's effectiveness—all 16 flows maintain stable 23-25 Gbps throughput while CPU load is evenly distributed at 50-65% across cores. This contrasts sharply with random assignment where flows dropped below 5 Gbps and CPUs reached 100% utilization due to collisions.

SRAP achieves perfect 1:1 flow-to-core mapping, eliminating the 99% collision probability inherent in random ports. Aggregate throughput reaches 380-395 Gbps, compared to 44.8-395 Gbps variance without control, enabling reliable multi-flow aggregation.

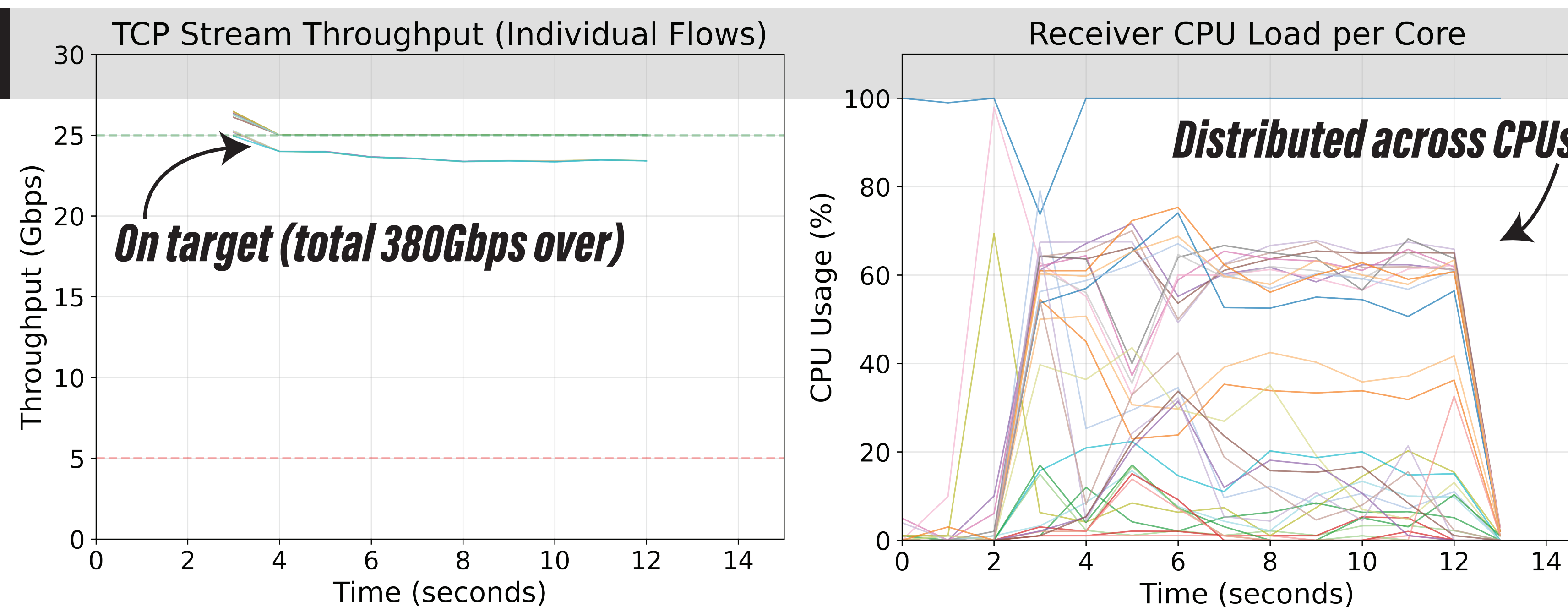


Figure 5: Performance with SRAP

5. Conclusion & Discussion

Scalable pacing for 400G aggregate and beyond.

We demonstrated that random port assignment causes severe RSS flow collisions, with only 1.04% probability of achieving optimal distribution for 16 flows on 32 cores, resulting in throughput varying from 44.8 to 395 Gbps.

SRAP (Sender-side Receiver-Aware Port selection) transparently controls source ports, achieving consistent 1:1 mapping and 380-395 Gbps throughput without modifying applications or receivers.

This demonstrates that source port engineering transforms unpredictable multi-flow performance into reliable high-throughput communication essential for 400G+ environments, achieving 9x faster worst-case transfers while preventing network-wide congestion from pacing breakdown.

