

Lauberhorn: a Smart NIC that is part of the OS

Pengcheng Xu
ETH Zurich
Switzerland
pengcheng.xu@inf.ethz.ch

Timothy Roscoe
ETH Zurich
Switzerland
troscoe@inf.ethz.ch

Remote Procedure Calls (RPC) are essential for various data center workloads, including micro-services, serverless, networked filesystems, and many more. Previous work on characterizing RPC workloads [10] demonstrated that short RPC invocations in the ballpark of 1 microsecond make up a significant portion of all RPC workloads.

Despite the high frequency of short RPC workloads, traditional datacenter RPC architecture using PCIe DMA-based NICs incur high latency and CPU overhead. This is due to a deliberate *split of networking states* between the CPU and NIC. This historical and stable split comes from numerous factors. Firstly, there is a fundamental distrust between the OS and NIC, largely coming from the perception that the NIC never does what the OS wants. The introduction of various IOMMUs has further strengthened this view. Secondly, architecture researchers are well-known to ignore the OS, attempting to bypass it by assigning applications to fixed CPU cores to minimize scheduling disruptions. Last but not least, there is a common misconception that fine-grained interaction between the OS and NIC is slow; this might be true in traditional descriptor-ring-based PCIe DMA, but not so much in PCIe MMIO, and even less so in modern cache-coherent interconnects.

Recent works explore alternative packet delivery mechanisms to the host CPU. Many of them embody kernel bypass mechanisms to remove the OS kernel from the critical path, making trade-offs between energy-efficiency, flexibility, and throughput [1, 3, 4, 6, 7, 11]. Others integrate the NIC tightly with the operating system through various on-chip mechanisms and similarly sacrifice flexibility [5, 9]. Tightly integrating the CPU and NIC on the same die also poses challenges in practicality, since compute power and networking speeds do not evolve in lock-step.

We argue that kernel-bypass approaches are at their limit due to various advancements in the data center architecture. Servers can come with over a hundred cores, making obsolete mechanisms for multiplexing CPU time between multiple tasks. Furthermore, cache-coherent interconnects allow for communication between the CPU and NIC lower latency and higher throughput, without the need for huge batch sizes.

Such advancements call for a radically different, OS-centric approach based on new interconnects and sharing of state between the OS and NIC. In our ongoing project, Lauberhorn, we propose an accelerated communication architecture that fuses the NIC and the OS with coherently attached FPGA,

effectively freeing the CPU from all overheads in RPC processing to achieve high efficiency and performance.

We implement a variant of the protocol proposed by Ruzhanskaia, et. al. [8] on top of cache-coherence protocols like CXL to achieve efficient fine-grained communication between the CPU and NIC. We then implement all but handler execution in RPC processing on the NIC, offloading operations such as compression, encryption, and arguments marshaling to the FPGA as hardware accelerators. Finally, a tight integration between the NIC and the OS scheduling system allow for performance as good as kernel-bypass approaches, but with far more flexibility.

A key novelty and contribution of Lauberhorn is how it uses precise scheduling states to dispatch RPC requests. We implement a *multi-level scheduling* approach, putting the NIC completely in charge of deciding where to dispatch an RPC request. The NIC holds authority of which user space RPC application runs on a CPU core. The CPU starts in a kernel thread and polls for a user space process to dispatch on the core. After a process is dispatched on the core, the CPU serves the request and continues polling *in user space* on a different pair of cache lines. The NIC can then dispatch another request to the same process with *zero software overhead* (a simple jump on the CPU core). Alternatively, it can issue a command asking the CPU core to go back to kernel mode to dispatch a different user process. Such an approach allows maximal performance without software overhead when one service is hot, while retaining flexibility switching between processes when necessary.

We set an ambitious target for end-to-end RPC latency of 1 microsecond. Preliminary results on Enzian [2] show that we can send and receive Ethernet frames from the CPU over the ECI cache coherence protocol in around 800 nanoseconds, making a promising case for our latency target.

References

- [1] BELAY, A., PREKAS, G., KLIMOVIC, A., GROSSMAN, S., KOZYRAKIS, C., AND BUGNION, E. IX: a protected dataplane operating system for high throughput and low latency. In *Proceedings of the 11th USENIX Conference on Operating Systems Design and Implementation* (USA, 2014), OSDI'14, USENIX Association, p. 49–65.
- [2] COCK, D., RAMDAS, A., SCHWYN, D., GIARDINO, M., TUROWSKI, A., HE, Z., HOSSLE, N., KOROLJA, D., LICCIARDELLO, M., MARTSENKO, K., ACHERMANN, R., ALONSO, G., AND ROSCOE, T. Enzian: an open, general, CPU/FPGA platform for systems software. In *ASPLOS '22: Proceedings of the Twenty-Seventh International Conference on Architectural Support for Programming Languages and Operating Systems* (February 2022).

- [3] FRIED, J., RUAN, Z., OUSTERHOUT, A., AND BELAY, A. Caladan: Mitigating Interference at Microsecond Timescales. pp. 281–297.
- [4] HUMPHRIES, J. T., NATU, N., CHAUGULE, A., WEISSE, O., RHODEN, B., DON, J., RIZZO, L., ROMBAKH, O., TURNER, P., AND KOZYRAKIS, C. ghOS: Fast & Flexible User-Space Delegation of Linux Scheduling. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles* (Virtual Event Germany, Oct. 2021), ACM, pp. 588–604.
- [5] IBANEZ, S., MALLERY, A., ARSLAN, S., JEPSEN, T., SHAHBAB, M., KIM, C., AND McKEOWN, N. The nanoPU: A Nanosecond Network Stack for Datacenters. pp. 239–256.
- [6] KAFFES, K., CHONG, T., HUMPHRIES, J. T., BELAY, A., MAZIÈRES, D., AND KOZYRAKIS, C. Shinjuku: Preemptive Scheduling for {usecond-scale} Tail Latency. pp. 345–360.
- [7] PETER, S., LI, J., ZHANG, I., PORTS, D. R. K., WOOS, D., KRISHNAMURTHY, A., ANDERSON, T., AND ROSCOE, T. Arrakis: The Operating System is the Control Plane. In *11th Symposium on Operating Systems Design and Implementation (OSDI’14)* (Broomfield, Colorado, USA, October 2014).
- [8] RUZHANSKAIA, A., XU, P., COCK, D., AND ROSCOE, T. Rethinking Programmed I/O for Fast Devices, Cheap Cores, and Coherent Interconnects, Sept. 2024. arXiv:2409.08141 [cs].
- [9] SCHUH, H. N., KRISHNAMURTHY, A., CULLER, D., LEVY, H. M., RIZZO, L., KHAN, S., AND STEPHENS, B. E. CC-NIC: a Cache-Coherent Interface to the NIC. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 1* (La Jolla CA USA, Apr. 2024), ACM, pp. 52–68.
- [10] SEEMAKHUP, K., STEPHENS, B. E., KHAN, S., LIU, S., WASSEL, H., YEGANEH, S. H., SNOEREN, A. C., KRISHNAMURTHY, A., CULLER, D. E., AND LEVY, H. M. A Cloud-Scale Characterization of Remote Procedure Calls. In *Proceedings of the 29th Symposium on Operating Systems Principles* (Koblenz Germany, Oct. 2023), ACM, pp. 498–514.
- [11] ZHANG, I., RAYBUCK, A., PATEL, P., OLYNYK, K., NELSON, J., LEIJA, O. S. N., MARTINEZ, A., LIU, J., SIMPSON, A. K., JAYAKAR, S., PENNA, P. H., DEMOULIN, M., CHOUDHURY, P., AND BADAM, A. The Demikernel Datapath OS Architecture for Microsecond-scale Datacenter Systems. In *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles* (Virtual Event Germany, Oct. 2021), ACM, pp. 195–211.