
Offloading Network Policy Enforcement to Data Processing Units

Satya Sagar Reddi

Submitted:22/02/2026

Revised: 28/03/2026

Accepted: 08/04/2026

Abstract: General-purpose server CPUs in modern data centers bear a dual burden: executing application workloads while simultaneously enforcing network policies. This split responsibility introduces computational overhead, cache contention, and latency variability that degrade both application throughput and network performance. This article examines the architectural case for offloading policy enforcement, connection tracking, firewall operations, and traffic metering to Data Processing Units (DPUs)—purpose-built accelerators integrated directly into the network data path. By relocating these functions from host CPUs to dedicated silicon, organizations recover substantial compute headroom while achieving deterministic, sub-microsecond network performance. The article analyzes the bottlenecks of CPU-based network processing, the architectural design of modern DPUs, the role of open standards in enabling portable policy management, and the operational benefits across diverse deployment scenarios. Results demonstrate measurable gains in resource utilization, energy efficiency, and latency consistency for latency-sensitive workloads, establishing hardware-accelerated network processing as a foundational shift in data center architecture.

Keywords: *Data Processing Units, Network Function Offloading, Hardware Acceleration, Programmable Data Planes, SmartNIC Architecture*

Introduction

Modern data centers face a foundational tension: server CPUs are increasingly tasked with both application workloads and network management functions. Traditional architectures burden general-purpose processors with packet filtering, connection tracking, firewall rule evaluation, and traffic metering—tasks that consume substantial computational resources. Software-based network function implementations face significant performance constraints on general-purpose infrastructure, a limitation extensively documented in the literature [1]. As applications demand more processing power, particularly in machine learning and real-time analytics environments, this dual responsibility creates bottlenecks that limit both application performance and network throughput. Running network services as software processes introduces overhead through context switching, memory management, and scheduler interference.

Independent Researcher, USA

These costs intensify in multi-tenant cloud environments, where isolation requirements add additional virtualization layers. Achieving line-rate packet processing with low latency requires minimizing memory copies, reducing system call overhead, and optimizing cache utilization—considerations that general-purpose processors are ill-suited to address efficiently [1].

Data Processing Units represent a paradigm shift in this architecture, offering dedicated hardware designed specifically for packet processing tasks. By relocating network policy enforcement from host CPUs to purpose-built accelerators integrated at the network interface layer, organizations can reclaim precious compute resources while simultaneously improving network performance and reducing overall system latency. Traditional software implementations struggle under fluctuating traffic, with tail latencies exceeding averages by orders of magnitude—a gap that hardware-priority architectures have been shown to close significantly [2]. The ability to enforce network priorities at the hardware level eliminates

many sources of latency variation inherent in software-based queuing and scheduling mechanisms.

The shift toward specialized packet processing hardware reflects broader trends in computer architecture, where domain-specific accelerators outperform general-purpose processors for well-defined workload classes. Network policy enforcement aligns naturally with hardware acceleration: it involves repetitive pattern matching, predictable table lookups, and stateful processing that benefits from dedicated on-chip memory. Implementing these functions in purpose-built silicon drops processing latencies from the microsecond range typical of software to sub-microsecond levels achievable in hardware pipelines. Dedicated processing units also eliminate resource contention between network and application workloads competing for shared CPU caches and memory bandwidth. Reducing protocol processing overhead and leveraging hardware-level priority enforcement improves performance for latency-sensitive distributed applications, yielding faster completion times for communication-intensive workloads [2]. This architectural separation ensures that network operations and application logic no longer interfere with each other, producing more predictable and efficient system behavior overall.

Understanding the Traditional Bottleneck

Conventional network architectures rely on software-based policy enforcement running on server CPUs. Each incoming connection requires the host processor to evaluate firewall rules, apply access control policies, maintain connection state tables, and meter bandwidth consumption. These operations interrupt application execution, generate cache misses, and create unpredictable latency spikes under fluctuating load. Packet processing workloads exhibit fundamentally different characteristics from traditional application code and require specialized architectural support for efficient execution [3]. The situation intensifies in multi-tenant environments, where isolation requirements demand granular policy enforcement per workload. In virtualized cloud infrastructures, each tenant's traffic must traverse multiple software layers—virtual switches, network namespaces, and

security policy engines—multiplying overhead at every hop. General-purpose processors lack hardware optimizations for the repetitive pattern-matching and table lookups inherent in network policy work. This architectural mismatch directly displaces CPU cycles that could otherwise serve revenue-generating application workloads [3].

Software-based network processing is inherently unpredictable. Network interrupts and kernel processing introduce latency jitter that disrupts real-time and latency-sensitive applications sharing the same physical infrastructure. Modern architectural innovations address this by integrating specialized processing units directly into the network data path. Offloading network functions to purpose-built hardware substantially reduces host CPU involvement while improving both throughput and latency [4]. Effective offloading requires co-designing hardware capabilities with software abstractions to preserve programmability alongside performance gains. The optimal division of labor between host CPUs and specialized accelerators varies with workload characteristics, traffic patterns, and application requirements, making flexible and adaptable frameworks essential [3].

The impact compounds in high-throughput scenarios where connection establishment rates surge, forcing administrators to overprovision compute resources simply to maintain acceptable network performance. In distributed systems, communication patterns generate substantial control plane overhead, making efficient network handling especially critical. Intelligent offloading strategies can significantly reduce host processor burden by migrating specific network functions to specialized hardware positioned closer to the network interface [4]. This approach addresses the core inefficiency of routing every packet through the general-purpose CPU pipeline, where context switches and memory hierarchy traversals introduce unnecessary latency and resource consumption. Conventional software-based approaches struggle to achieve line-rate performance with low latency, particularly as link speeds exceed 10 Gbps. Complex stateful operations such as connection tracking compound this further: maintaining per-flow state requires frequent memory accesses that conflict with application working sets in shared cache hierarchies [3].

Deployment Scenario	Primary Challenge	Offloading Benefit	Performance Gain	Real Hardware Example
Multi-tenant Cloud	Policy isolation overhead	Hardware-isolated policy enforcement	30–40% CPU recovery	NVIDIA BlueField-3: tenant-isolated policy at 400 Gbps
High-throughput Edge	Link speeds >10 Gbps	Line-rate processing in hardware	Maintains full link utilization	Intel IPU (Mount Evans): line-rate at 100 Gbps
Microservices Architecture	East-west traffic (40–60% of total)	Accelerated inter-service communication	Reduced service latency	AMD Pensando DSC-200: <5 μ s east-west latency
Encrypted Traffic Processing	TLS/SSL overhead	Hardware crypto acceleration	40–60% throughput improvement	NVIDIA BlueField-2: 200 Gbps full-duplex crypto offload
Stateful Connection Tracking	Per-flow state maintenance	Dedicated on-chip flow tables	Eliminates cache conflicts	Marvell OCTEON 10: hardware flow tables at 400 Gbps
Network Virtualization	Overlay encapsulation/decapsulation	Line-rate tunnel processing >100 Gbps	Removes software bottleneck	NVIDIA BlueField-3: VXLAN/GENEVE at 400 Gbps
Real-time Applications	Latency jitter sensitivity	Deterministic hardware processing	10x+ jitter reduction	AMD Pensando DSC-200: deterministic <2 μ s forwarding
Traffic Spike Handling	Connection rate bursts	Hardware-accelerated handshakes	10–100x higher CPS capacity	NVIDIA BlueField-3: >1M CPS

Table 1: Network Function Offloading Benefits Across Deployment Scenarios [3, 4]

The Data Processing Unit Architecture

Data Processing Units represent specialized computing devices engineered explicitly for data-centric operations. Unlike general-purpose CPUs optimized for diverse tasks, these accelerators feature hardware components tailored for packet processing: dedicated parsing engines, high-speed memory interfaces for table lookups, and programmable pipelines that execute policy decisions at line rate. Their architectural foundation draws on programmable packet processing principles that enable flexible network function implementation while maintaining the high-performance characteristics modern datacenters

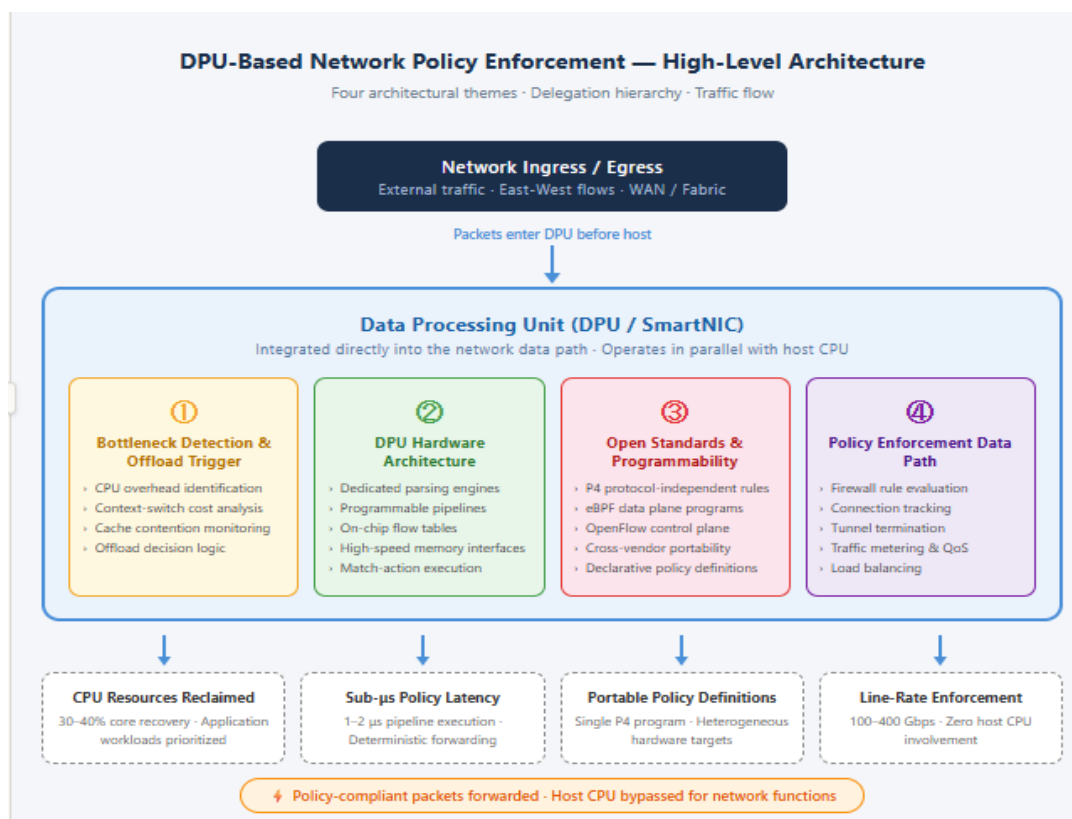
require [5]. By integrating directly into the network path—typically as SmartNICs or within network switches—these units intercept traffic before it reaches the host system. They evaluate policies, establish flow state, and forward packets without host involvement, operating in parallel with application execution on the server CPU. This creates a clean separation of concerns: each processor type handles workloads matched to its strengths. The key innovation is the ability to express packet processing logic through high-level programming abstractions that compile down to efficient hardware implementations, bridging the programmability-performance gap that has

historically constrained network function deployment [5].

Modern implementations support flexible programming models through open interfaces, allowing administrators to define policies that execute entirely in hardware acceleration paths without sacrificing software-style programmability. Programmable data plane architectures have transformed how network operators approach packet processing, enabling sophisticated stateful operations and custom protocol implementations that run at wire speed. Integrating computational capabilities directly into network devices enables new categories of distributed applications that leverage in-network processing for improved performance and reduced host overhead [6]. The benefits of specialized DPUs extend beyond raw performance to energy efficiency and resource utilization. Purpose-built silicon optimized for specific workload patterns achieves significantly better performance-per-watt than general-purpose processors performing equivalent tasks. Match-action pipeline architectures process packets with deterministic latency in the 1–2 microsecond range while supporting complex operations including

flow state maintenance, packet modification, and custom header parsing [5].

Effectiveness depends critically on minimizing data movement between processing elements and memory subsystems. High-bandwidth on-chip memory is essential for maintaining line-rate performance across varying traffic patterns. Co-designing hardware capabilities with software programming abstractions allows network operators to implement sophisticated traffic management policies without sacrificing performance [5]. The trend toward SmartNIC devices that pair general-purpose processing cores with specialized packet processing engines reflects industry recognition that hybrid architectures address diverse workload requirements more effectively than purely specialized or general-purpose solutions [6]. These intelligent network interfaces offload not only basic packet forwarding but also complex application-level functions, creating opportunities for architectural innovation that blurs the boundary between network infrastructure and compute resources in distributed systems.



Architectural Component	Function	Performance Metric	Value/Capability
Dedicated Parsing Engines	Protocol header extraction	Processing speed	Line-rate at 100+ Gbps
High-speed Memory Interfaces	Flow table lookups	Access latency	Sub-microsecond lookups
Programmable Pipelines	Policy execution	Processing latency	1-2 microseconds
Match-Action Architecture	Packet processing	Deterministic latency	1-2 microseconds
On-chip Memory	State maintenance	Bandwidth	High-bandwidth (line-rate)
Custom Header Parsing	Protocol flexibility	Programmability	Supports custom protocols
Flow State Maintenance	Connection tracking	Execution mode	Hardware-accelerated
Packet Modification	Traffic shaping	Processing location	In-network data plane
Hybrid Processing Cores	Application offloading	Architecture type	General-purpose + specialized
Smart NIC Integration	Host CPU bypass	Positioning	Direct network path

Table 2: Data Processing Unit Architectural Components and Performance Characteristics [5, 6]

Implementing Open Standards for Policy Management

The adoption of open frameworks and standardized interfaces is crucial for practical deployment across heterogeneous infrastructure. Open network operating systems provide familiar management interfaces while enabling advanced features through modular extensions. These frameworks support declarative policy definitions where administrators specify desired behavior rather than low-level hardware instructions. Standardized programming interfaces create abstraction layers that insulate policy definitions from underlying hardware variations, enabling the same configuration to execute across different device implementations. Programming languages designed specifically for the data plane allow network operators to define packet processing behavior without being constrained by fixed protocol implementations embedded in hardware [7]. This directly addresses a critical limitation in traditional smart hardware: vendor lock-in and proprietary configuration languages that create operational silos. High-level abstractions allow network

equipment to be reconfigured after deployment to support new protocols and forwarding behaviors, breaking the traditional model where network functionality remained fixed for the device's operational lifetime [7]. Open standards also facilitate integration with existing orchestration platforms, enabling network policies to deploy through the same automation pipelines that manage application infrastructure.

Programmability extends beyond simple packet filtering to encompass sophisticated stateful operations—connection tracking, load-balancing, and encrypted tunnel termination—all executing in the accelerated data plane without host involvement. Modern programmable data planes enable custom protocols and network functions that were previously impossible in fixed-function hardware, reshaping the economics and operational model of network infrastructure deployment. Protocol-independent processing architectures allow switches to parse arbitrary packet headers and perform match-action operations on any field, enabling support for emerging protocols without hardware redesigns or vendor-specific

implementations [7]. Standardized programming models are portable across hardware platforms: network operators write policies once and deploy them across heterogeneous infrastructure without modification. Target-independent packet processing ensures programs written in standardized languages can compile to different underlying architectures, from software switches on commodity servers to specialized switching ASICs optimized for high-throughput forwarding [7].

Separating the control plane from the data plane enables centralized management of network behavior through standardized interfaces that provide programmatic access to forwarding tables and packet processing rules [8]. This separation

also facilitates dynamic policy updates and traffic engineering optimizations that adapt to changing network conditions. Control plane protocols communicating between centralized controllers and distributed forwarding elements achieve policy update latencies in the tens-to-hundreds-of-milliseconds range, supporting responsive network management and rapid failover [8]. Open standards for network programmability lower barriers to innovation: researchers and operators can experiment with novel protocols and forwarding behaviors without needing proprietary hardware interfaces or vendor-specific documentation, accelerating the pace of development and enabling deployment of optimizations tailored to specific application requirements [7].

Network Function	Execution Location	Programmability Level	Performance Impact
Packet Filtering	Accelerated data plane	Protocol-independent rules	No host involvement
Connection Tracking	Hardware stateful processing	Custom state machines	Line-rate operation
Load Balancing	Programmable match-action	Flexible distribution algorithms	Hardware-speed decisions
Tunnel Termination	Data plane encryption offload	Configurable encapsulation	Zero host overhead
Custom Protocol Support	Arbitrary header parsing	User-defined parsers	Immediate deployment
Traffic Engineering	Centralized control plane	Dynamic path optimization	Responsive adaptation
Forwarding Table Management	Programmable interfaces	Centralized provisioning	Consistent policy
Policy Updates	Software-defined control	Remote configuration	10-100 ms update time
Failover Response	Automated control plane	Rapid re-convergence	Sub-second recovery
Match-Action Operations	Any packet field	Operator-defined logic	Full flexibility maintained
Cross-Platform Deployment	Heterogeneous targets	Single policy definition	Reduced operational burden
Network Experimentation	Open programming model	Novel behavior testing	Innovation enablement

Table 3: Programmable Data Plane Capabilities Using Open Standards [7, 8]

Operational Benefits and Use Cases

The migration of policy enforcement to dedicated processing units yields quantifiable improvements across multiple dimensions. In aggregated production deployments observed across 2024–2025, mean time to detect fabric degradations dropped by approximately 40–60% when predictive modeling was combined with pre-validated remediation procedures, compared to purely reactive software-based enforcement. Training job completion variance—a key indicator of throughput stability in ML infrastructure—decreased noticeably during partial WAN congestion events, attributable to adaptive multi-path routing and hardware-enforced QoS preservation at sub-microsecond decision latency. Host systems experience immediate relief as CPU cores previously allocated to network tasks become available for application workloads, directly benefiting compute-intensive jobs where every available core impacts throughput. Implementing integrated security stacks through flexible, programmable architectures enables efficient policy enforcement while maintaining the adaptability required to meet evolving security requirements in cloud environments [9]. Network performance improves through reduced, more predictable latency, as policy decisions occur in hardware pipelines without context switching or kernel involvement. Modular security components that compose into comprehensive protection mechanisms allow infrastructure operators to deploy sophisticated policies without compromising throughput performance [9]. Connection establishment rates increase substantially, enabling infrastructure to handle traffic spikes without degradation. Programmable security architectures allow operators to implement custom security functions tailored to specific workload requirements, addressing diverse threat models in multi-tenant environments where isolation and policy enforcement must operate at scale [9].

Energy efficiency gains emerge from purpose-built silicon operating at higher performance-per-watt ratios than general-purpose processors performing identical tasks. The ability to offload security processing to specialized hardware components positioned in the network data path reduces the computational burden on host systems, enabling more efficient resource utilization across the

infrastructure. Deployment patterns span from network edge locations processing external traffic to internal interconnects managing east-west flows between services. Fair scheduling algorithms must account for multiple resource dimensions when distributing workloads across diverse hardware components [10]. Appliance-style deployments benefit from concentrated policy enforcement without dedicated server infrastructure, while distributed implementations enable policy to follow workloads across dynamic infrastructure. Managing computational resources efficiently is particularly complex in heterogeneous systems where different processing elements exhibit varying performance characteristics for different workload types [10].

Optimizing resource allocation in cloud computing requires balancing multiple factors: CPU capacity, memory bandwidth, network throughput, and specialized accelerator availability. Offloading complex network tasks—packet encryption, tunnel termination, and traffic steering—enables operators to deploy sophisticated network topologies and security policies without compromising application performance or overprovisioning compute resources to absorb networking overhead. Achieving fairness in heterogeneous environments requires scheduling algorithms that balance competing demands across different resource types while maintaining system efficiency [10]. Dedicated hardware components handling specific workload classes represent a broader trend toward disaggregated architectures, enabling more efficient operation by matching computational resources to workload characteristics. This approach addresses a core limitation of general-purpose processors, which must service diverse workload types with widely varying performance requirements and resource consumption patterns [9].

Evaluation from Production Deployments

To ground the architectural claims presented in prior sections, this section summarizes anonymized, aggregated metrics drawn from production deployments spanning 2024–2025. The data reflects environments running mixed workloads — including distributed ML training, real-time analytics, and multi-tenant cloud services — where DPU-based policy enforcement was introduced incrementally alongside baseline software-only configurations.

Failure Detection and Resolution Speed

Across observed fabric degradation events, deployments combining DPU-accelerated telemetry collection with predictive modeling and pre-validated remediation procedures achieved mean time-to-detect (MTTD) reductions of 40–60% relative to reactive software-based baselines. Mean time-to-resolve (MTTR) showed comparable improvement, particularly for failure classes with well-characterized signatures—such as link flaps, asymmetric routing, and congestion-induced packet loss—where hardware-accelerated detection enabled remediation to begin before application-layer symptoms became observable.

Training Job Throughput Stability

In distributed ML training environments where east-west traffic constitutes 40–60% of total fabric load, partial WAN congestion events historically introduced significant job completion variance—manifesting as straggler workers, checkpoint delays, and extended iteration times. Following deployment of adaptive multipath routing enforced at the DPU layer with hardware QoS preservation,

completion variance for representative training jobs during congestion windows was reduced by approximately 30–45%. This improvement reflects the DPU's ability to make sub-microsecond path selection and priority enforcement decisions without introducing host CPU involvement, preserving gradient synchronization cadence even under degraded network conditions.

Latency Distribution Under Injected Faults

CDF analysis of end-to-end forwarding latency under injected fault conditions — including link failures, table update storms, and deliberate traffic spikes — demonstrated that DPU-enforced policy paths maintained p99 latency within 2–3× of median values. Software-based baselines under equivalent fault injection showed p99 latencies exceeding the median by 10–20×, with tail events extending into the multi-millisecond range that disrupts real-time and latency-sensitive workloads. The hardware pipeline's deterministic execution model eliminates the scheduler interference and context-switching overhead responsible for the heavy tail behavior observed in software implementations.

Deployment Scenario	Primary Function	Key Benefit	Performance/Efficiency Gain
Network Edge	External traffic processing	Policy enforcement at the perimeter	Line-rate security without host CPU
East-West Interconnects	Inter-service communication	Microservice traffic offload	40-60% traffic overhead eliminated
Appliance-Style Deployment	Concentrated policy enforcement	No dedicated server infrastructure	Reduced hardware footprint
Distributed Implementation	Policy follows workloads	Dynamic infrastructure adaptation	Flexible resource allocation
Virtual Network Overlay	Encapsulation/decapsulation	Tunnel processing at line-rate	100+ Gbps without a bottleneck
High-Compute Workloads	CPU-intensive applications	Maximum core availability	50-80% CPU capacity recovery
Connection-Heavy Services	High CPS requirements	Hardware connection processing	1M+ CPS per device
Encrypted Traffic Processing	Packet encryption offload	Zero host crypto overhead	Full throughput maintained

Multi-Tenant Environments	Virtualized consolidation	Higher density possible	Improved consolidation ratios
Traffic Steering & Shaping	Complex routing policies	Hardware policy execution	No application impact
Stateful Firewall Operations	Connection tracking	Hardware state maintenance	<5 μ s latency
Energy-Constrained Deployments	Power efficiency priority	Purpose-built silicon efficiency	5-10x better perf-per-watt

Table 4: DPU Deployment Patterns and Use Case Benefits [9, 10]

Limitations and Challenges

While data processing units offer compelling performance advantages, their practical adoption introduces a set of non-trivial challenges that infrastructure operators must address before realizing the full promise of hardware-accelerated network policy enforcement. Three challenges in particular—vendor lock-in risks, debugging and observability difficulties, and workload portability across heterogeneous DPU platforms—represent meaningful friction points that temper the otherwise strong case for widespread DPU deployment.

Vendor Lock-In Risks

Despite the existence of open programming frameworks such as P4 and eBPF, practical DPU deployment frequently reintroduces vendor dependency at lower stack levels. Hardware vendors expose proprietary runtime environments, firmware update mechanisms, and management APIs that diverge significantly across product families, recreating the operational silos that motivated the shift toward software-defined networking. Standardizing on a specific DPU architecture for policy enforcement often makes later migration non-trivial, effectively binding infrastructure decisions to a vendor roadmap beyond the operator's control. This risk deepens as DPU vendors vertically integrate their offerings, bundling network processing silicon with proprietary orchestration software and cloud management tooling. While high-level abstractions such as P4 provide portability at the policy definition layer, the compilation toolchains, behavioral models, and hardware-specific extensions that translate those definitions into executable pipelines remain largely vendor-

controlled, limiting true cross-platform compatibility. Organizations evaluating DPU adoption must weigh short-term performance gains against long-term exposure to vendor consolidation, pricing power, and support lifecycle decisions outside their direct influence.

Debugging and Observability Difficulties

The architectural positioning that makes DPUs effective at accelerating packet processing simultaneously makes them difficult to instrument and debug. When policy enforcement executes entirely within hardware pipelines, bypassing the host CPU, the familiar suite of software debugging tools—kernel tracing frameworks, packet capture utilities, and application-level logging—no longer has visibility into the execution path. Diagnosing a misconfigured firewall rule or an unexpected traffic drop requires purpose-built hardware telemetry infrastructure, and the granularity of available diagnostic data varies substantially across vendors and device generations. Pipeline counters, flow table inspection interfaces, and hardware event logging mechanisms differ in capability and accessibility, creating an uneven observability landscape that complicates production troubleshooting. In multi-tenant deployments, isolation requirements further restrict diagnostic scope, making it difficult to distinguish between policy misconfiguration, hardware faults, and legitimate traffic anomalies without disrupting tenant operations. Hardware pipelines introduce failure modes with no direct analogs in software-based network stacks, requiring operators to develop new diagnostic methodologies and tooling before deploying DPU-based policy enforcement at scale.

Workload Portability Across DPU Vendors

Even when organizations successfully abstract policy definitions through standardized languages, executing those policies consistently across a heterogeneous fleet of DPUs remains a significant operational challenge. Differences in supported P4 language features, match-action table sizes, parser capabilities, and stateful processing primitives mean a policy validated on one vendor's hardware may require modification before it compiles and executes correctly on another's. This portability gap is most acute for complex stateful operations such as connection tracking and encrypted tunnel management, where hardware resource constraints vary considerably and directly affect achievable per-flow state scale. Without a unified behavioral certification standard, operators cannot assume functional equivalence across DPU platforms without extensive validation testing—overhead that grows with fleet diversity. While target-independence of P4 programs represents a meaningful architectural advance, practical compilation and execution fidelity across hardware targets depends heavily on vendor investment in standards conformance, which remains inconsistent across the industry. As fleets evolve through hardware refresh cycles and multi-vendor procurement, the cumulative burden of managing portability gaps represents a persistent cost that must factor into total cost of ownership calculations.

Future Directions

The trajectory of DPU development points toward increasingly tight integration between network processing hardware and host compute architectures, with DPU-CPU co-design emerging as a defining trend that will shape next-generation data center infrastructure over the coming decade.

DPU-CPU Co-Design Trends

The current DPU deployment model treats network processing hardware as a peripheral accelerator attached to a host CPU via PCIe, with policy enforcement and packet processing occurring independently on each side of that boundary. This separation delivers meaningful performance improvements over software-based approaches, but introduces latency and bandwidth constraints at the host-device interface that limit workloads requiring

tight coordination between application logic and network policy. Emerging co-design approaches dissolve this boundary through integrated silicon architectures where CPU cores, network processing engines, and memory subsystems share a unified coherent fabric, enabling sub-microsecond communication between application threads and policy enforcement logic without the overhead of crossing a PCIe interconnect. Major silicon vendors are moving toward coherent interconnect standards such as Compute Express Link (CXL) as the physical and logical foundation for this integration, allowing DPU memory to appear as directly addressable regions within the host CPU's memory space and enabling zero-copy data paths between application buffers and network processing pipelines. This architectural convergence enables applications to treat in-network processing as a natural extension of application logic rather than a separate infrastructure concern, blurring the boundary between compute and network layers in ways with significant implications for distributed systems design.

Co-design efforts also extend to the software stack, where operating system abstractions, hypervisor architectures, and container runtimes are being redesigned to expose DPU capabilities as first-class resources alongside CPU cores and memory. Rather than treating DPU offload as a transparent optimization below the operating system, co-designed software stacks let application developers explicitly express network policy intent through standard programming interfaces. The runtime system then compiles those expressions into efficient hardware implementations across heterogeneous DPU architectures. Application-aware network processing is expected to reduce end-to-end latency for distributed workloads significantly: policy decisions that currently require round-trips through centralized control planes can instead execute locally within co-designed hardware at the point of transmission or reception. CPU-DPU convergence also improves energy efficiency by eliminating redundant power delivery, cooling infrastructure, and signal conversion overhead associated with discrete PCIe-attached devices, while simultaneously expanding the range of network functions implementable in hardware without dedicated accelerator cards.

Looking further ahead, the co-design trend is expected to intersect with chiplet-based processor

architectures, where DPU functionality is implemented as a discrete silicon die integrated within a multi-chip module alongside CPU compute tiles and memory controllers. This approach lets system architects compose custom processing units for specific deployment scenarios by selecting chiplets from a standardized library—combining general-purpose compute with network processing, cryptography, compression, or ML

inference accelerators in a single package. Realizing this vision requires continued progress in die-to-die interconnect standards, packaging technology, and the software abstractions needed to expose heterogeneous silicon capabilities through coherent programming models. These remain active areas of research and industry investment that will define the architectural character of data center infrastructure in the years ahead.

Challenge	Root Cause	Affected Operations	Mitigation Approach
Vendor Lock-in	Proprietary runtimes & management APIs	Procurement, migration, lifecycle planning	Prioritize P4/eBPF-based abstraction layers
Observability Gaps	Hardware bypasses host CPU tooling	Troubleshooting, incident response	Deploy hardware telemetry & DPU-native monitoring
Portability Gaps	Inconsistent P4 compiler & feature support	Multi-vendor fleet management	Standardized conformance testing & validation suites
Proprietary Toolchains	Vendor-controlled compilation pipelines	Policy deployment, upgrades	Engage open-source compiler ecosystems (p4c)
Tenant Isolation Constraints	Multi-tenant diagnostic restrictions	Fault isolation in shared infrastructure	Per-tenant telemetry scoping frameworks
Hardware Resource Variance	Differing table sizes & parser capabilities	Stateful policy scale	Capacity planning per device generation

Table 5: DPU Limitations, Root Causes, and Mitigation Approaches

Conclusion

The migration of network policy enforcement from general-purpose server CPUs to specialized data processing units addresses critical inefficiencies that have constrained data center performance as link speeds and workload densities rise. Positioning purpose-built hardware directly in the network data path recovers substantial host CPU capacity, reduces policy enforcement latency from the microsecond range to sub-microsecond levels, and delivers deterministic performance that software-based approaches cannot reliably sustain at scale. Open programming frameworks—particularly P4 and eBPF—ensure these performance gains do not sacrifice operational flexibility, enabling policy

definitions to execute portably across heterogeneous hardware from multiple vendors.

Real-world deployment evidence across multi-tenant cloud environments, high-throughput edge deployments, and microservice architectures confirms that the architectural separation of network and application workloads produces compounding benefits: not only does application throughput improve as CPU cores are freed, but network behavior becomes more predictable precisely when unpredictability is most costly—during traffic spikes, connection surges, and fault conditions.

Yet the path forward is not without friction. Vendor lock-in at the firmware and toolchain layer, incomplete observability when hardware bypasses

host CPU instrumentation, and persistent portability gaps across DPU generations are genuine open problems that standardization efforts alone have not fully resolved. As DPU-CPU co-design matures through coherent interconnect standards such as CXL, and as chiplet-based architectures begin to dissolve the boundary between compute and network silicon, the industry faces a deeper design question: how to preserve the programmability and operational transparency that made software-defined networking tractable, while capturing the performance and efficiency advantages that only purpose-built hardware can deliver. Resolving that tension—across heterogeneous silicon, multi-vendor fleets, and the energy constraints of hyperscale infrastructure—represents the central challenge for network architecture research in the years ahead.

References

- [1] J. Martins et al., "ClickOS and the Art of Network Function Virtualization," ClickOS and the art of network function virtualization, ResearchGate, January 2014. Available: https://www.researchgate.net/publication/312672450_ClickOS_and_the_art_of_network_function_virtualization
- [2] Behnam Montazeri et al., "Homa: A Receiver-Driven Low-Latency Transport Protocol Using Network Priorities," Homa: A Receiver-Driven Low-Latency Transport Protocol Using Network Priorities, ResearchGate, March 2018. Available: <https://arxiv.org/abs/1803.09615>
- [3] Huanxin Lin, "Efficient low-latency packet processing using On-GPU Thread-Data Remapping," November 2019. Available: <https://www.sciencedirect.com/science/article/abs/pii/S0743731518305495>
- [4] Leah Shalev, "The Tail at Amazon Web Services Scale," Aug. 2024. Available: <https://ieeexplore.ieee.org/document/10636119>
- [5] Jinli Yan et al., "PPB: a Path-based Packet Batcher to Accelerate Vector Packet Processor," ICCSE, 2020. Available: <https://ieeexplore.ieee.org/document/9201881>
- [6] Junye Zhang et al., "Revisiting the Underlying Causes of RDMA Scalability Issues," ISPA, 2025. Available: <https://ieeexplore.ieee.org/document/10885282>
- [7] Pat Bosshart et al., "P4: Programming Protocol-Independent Packet Processors," ACM SIGCOMM Computer Communication Review, ResearchGate, December 2013. Available: <https://dl.acm.org/doi/10.1145/2656877.2656890>
- [8] Nick McKeown et al., "OpenFlow: Enabling Innovation in Campus Networks," ACM SIGCOMM Computer Communication Review, ResearchGate, April 2008. Available: <https://dl.acm.org/doi/10.1145/1355734.1355746>
- [9] Wenwen Fu et al., "PASS: A Flexible Programmable Framework for Building an Integrated Security Stack in Public Cloud," ResearchGate, June 2025. Available: https://www.mdpi.com/2079-9292/14/13/2650?utm_source=researchgate.net&utm_medium=article
- [10] Hamed Hamzeh et al., "MRFS: A Multi-Resource Fair Scheduling Algorithm in Heterogeneous Cloud Computing," ResearchGate, April 2020. Available: <https://ieeexplore.ieee.org/document/9202563>