



White paper

## Acceleration technologies: realizing the potential of network virtualization

*June 2019*

Gorkem Yigit and Caroline Chappell

# Contents

<b>1.</b>	<b>Executive summary</b>	<b>2</b>
<b>2.</b>	<b>Drivers for acceleration technologies</b>	<b>3</b>
2.1	Acceleration technologies are needed to move network virtualization forward	3
2.2	Which virtualization use cases require acceleration?	4
<b>3.</b>	<b>Introduction to acceleration technologies</b>	<b>5</b>
3.1	Software-centric acceleration	5
3.2	Hardware-centric acceleration	9
3.3	Acceleration for NFV/SDN security functions: Intel QAT and HPE performance results	12
3.4	Summary comparison of acceleration technologies	13
<b>4.</b>	<b>Acceleration-as-a-service is crucial to abstract the complexity of acceleration resources</b>	<b>14</b>
<b>5.</b>	<b>Conclusion and recommendations</b>	<b>15</b>
<b>6.</b>	<b>About the authors</b>	<b>17</b>
<b>7.</b>	<b>Analysys Mason's consulting and research are uniquely positioned</b>	<b>18</b>
<b>8.</b>	<b>Research from Analysys Mason</b>	<b>19</b>
<b>9.</b>	<b>Consulting from Analysys Mason</b>	<b>20</b>

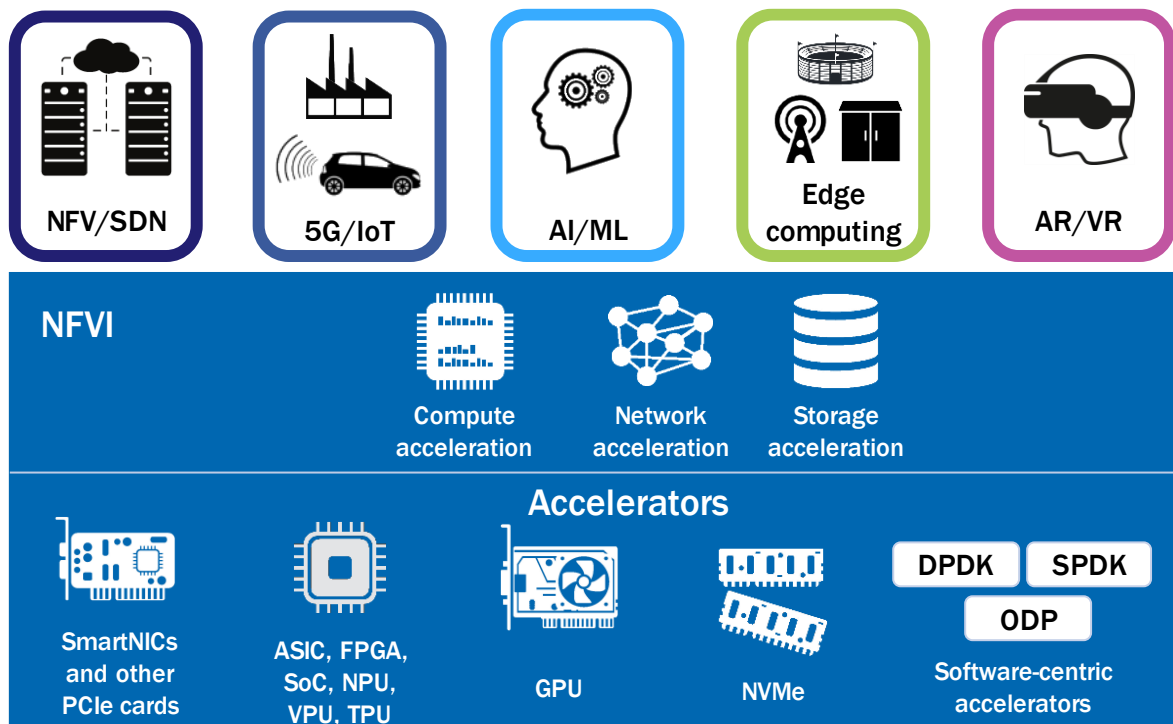
## List of figures

Figure 1.1: Acceleration technologies and use cases [Source: Analysys Mason, 2019].....	2
Figure 2.1: Virtualization use cases and their acceleration requirements [Source: Analysys Mason, 2019]....	4
Figure 3.1: Native OVS and DPDK accelerated OVS [Source: Analysys Mason, Intel 2019] .....	7
Figure 3.2: SR-IOV and PCI- passthrough models [Source: Analysys Mason, 2019] .....	8
Figure 3.3: Comparison of accelerated vSwitch, SR-IOV and PCI passthrough [Source: Analysys Mason, Intel 2019] .....	8
Figure 3.4: VNF and NFVI acceleration using Smart NICs [Source: Analysys Mason, 2019] .....	10
Figure 3.5: Comparison of hardware acceleration options [Source: Analysys Mason, 2019] .....	11
Figure 3.6: Intel QAT acceleration performance results on HPE servers [Source: Intel, HPE, 2019].....	12
Figure 3.7: Summary comparison of acceleration technologies (Source: Analysys Mason, 2019) .....	14

## 1. Executive summary

Operators are using industry-standard, commercial off-the-shelf (COTS) servers with general purpose CPUs to build their network function virtualization infrastructures (NFVIs) as they transition from closed, custom appliances to software-based network functions running on horizontal, open and highly automated platforms. The goal of network function virtualization (NFV) is still to have a homogenous, hyperscale, COTS architecture for all virtual network functions (VNFs). In practice, acceleration technologies (both fixed-function and programmable), in conjunction with tightly coupled general compute resources, can improve the overall networking and application performance of current VNFs to meet the quality of service (QoS), automation and security requirements of existing and future services, and improve capital and operational expenditures.

Figure 1.1: Acceleration technologies and use cases [Source: Analysys Mason, 2019]



This whitepaper defines the various acceleration technologies that are typically explored, in conjunction with general-purpose hardware and software technologies, to maximise the performance and efficiency of NFV/SDN, 5G and artificial intelligence (AI)/machine learning (ML) applications. No single architecture can respond to or satisfy all use cases, but this paper seeks to outline some of the currently deployed acceleration technologies as well as those that are emerging, alongside performance enhancements achieved by general purpose CPUs. It presents performance results for accelerating virtualized security functions as a proof point for emerging acceleration technologies. Finally, it discusses the need for an acceleration technology abstraction layer and a uniform management framework to simplify and automate the operations of heterogeneous acceleration resources with an ‘as-a-service’ model. It highlights industry initiatives such as OpenStack Cyborg and OPNFV DPACC, which aim to help operators to achieve these goals.

## 2. Drivers for acceleration technologies

This section discusses the benefits of having common infrastructure and operations for all VNF workloads based on general purpose hardware and software. It explains the challenge that this poses for VNFs that require high performance and throughput to perform data plane tasks. It suggests that an NFVI based on general-purpose CPUs can be augmented with technologies that accelerate compute-intensive VNF and NFVI workloads, and then identifies the use cases and functions that need this capability.

### 2.1 Acceleration technologies are needed to move network virtualization forward

NFV and software-defined networking (SDN) underpin operators' aspirations to become more efficient, competitive and innovative. Many operators have already passed the early stages of network virtualization, having implemented NFV/SDN within specific domains that are limited in scope and size. Operators are now looking to scale up these deployments to support the delivery of dynamic, on-demand network services (including virtual/universal customer premises equipment (vCPE/uCPE), software-defined wide area networks (SD-WAN) and managed security) and to lay the groundwork for high-capacity and low-latency 5G networks (such as edge, cloud RAN and virtual evolved packet core (vEPC)). Analysys Mason forecasts that operator spending on NFV and SDN will increase from USD3.7 billion in 2017 to USD23.8 billion in 2022, growing at a CAGR of 45%.<sup>1</sup>

Network virtualization is moving from a phase of closed, vertically siloed, custom appliances to one where multiple VNFs from different vendors run on a horizontal, open and highly automated platform. This will require network functions to be fully decoupled from specialised appliances and implemented as cloud-native software components (VNFs) on a common NFVI. The NFVI should consist of standard, off-the-shelf compute and memory resources and a virtualization layer. Eliminating interdependencies between hardware and software is essential if operators are to realise the financial and operational benefits of network virtualization. Operators can lower capex using inexpensive COTS hardware and white boxes, which reduce vendor lock-in and leverage IT market commoditization. They can also shrink opex through the use of low-cost IT automation tools. Together, these capabilities deliver a simpler, more flexible and programmable infrastructure on which to run the network.

Achieving disaggregation of network hardware and software is a key focus for operators during their network virtualization journey, but it is not without challenges. Standard, general-purpose hardware platforms are well-suited to handle many use cases. However, some network functions and workloads (especially those that are related to the data plane) require a high level of QoS (for areas such as throughput, latency and jitter), predictable performance and security. When implementing such VNFs on COTS hardware, operators often expect parity with custom network appliances and perceive a performance gap. This is due to the following reasons.

- The immaturity of many VNFs. These VNFs are based on software architecture that is still optimized for proprietary hardware environments. General purpose CPUs are not as optimized to handle specific VNF workloads. Only a few VNFs have been developed or rearchitected for cloud-native infrastructure using COTS platforms at this point.
- The virtualization layer in the NFVI adds processing overheads which operators want to minimize.

<sup>1</sup> For more information, see Analysys Mason's *Digital infrastructure: worldwide forecast 2018–2022*. Available at: [www.analysysmason.com/Research/Content/Reports/digital-infrastructure-forecast-rma16](http://www.analysysmason.com/Research/Content/Reports/digital-infrastructure-forecast-rma16).



As the scale and scope of virtualized networks grows, adding CPUs and servers may not be the optimal way to tackle this problem due to power, space and cooling constraints. This approach fails to address the technological limitation imposed by virtualization overheads from the multiple layers of packet processing needed as traffic flows from network interface cards (NICs) to VNFs. This limitation makes it difficult to satisfy service level agreement (SLA) and QoS requirements, especially for ultra-low-latency 5G use cases. Various approaches can be studied and implemented.

To match, or even surpass, the performance and latency levels of traditional appliances without negating the capex and opex benefits of virtualization, as well as to avoid virtualization overheads, operators should consider incorporating software and/or hardware acceleration technologies into their NFV infrastructure. This can help operators to address increasing overhead costs (which cannot be billed to end customers) by reducing the cost-per-bit and can also help to achieve the optimal allocation of computing, networking and storage for different types of VNF.

## 2.2 Which virtualization use cases require acceleration?

VNFs responsible for executing data plane functions fulfil a variety of specific networking, security and media-related tasks such as switching, routing, traffic management, cryptography (SSL, IPsec), compression and transcoding. Data plane functions need to process and forward traffic at near-line rates, that is, at the same speed as their network interfaces (for example, 40G or 100G), so they may suffer from NFVI performance bottlenecks and cost-effective scalability if they rely solely on general-purpose hardware and software-based resources to do this. Use cases that involve the service chaining of multiple VNF components can add further load on the infrastructure. In these cases, the throughput requirement is amplified and the tunnel processing associated with overlay networks (GRE, VXLAN) may add further overhead. The virtualization use cases and workloads that may benefit from acceleration are detailed in Figure 2.1 below.

**Figure 2.1: Virtualization use cases and their acceleration requirements [Source: Analysys Mason, 2019]**

Possible use cases	Functions	Potential acceleration points
vBNG/BRAS	VPN, firewall, DPI, multicast and service chains consisting of BNG and other data plane-heavy functions	<ul style="list-style-type: none"> <li>Layer 2/3/4 packet forwarding (high throughput in the long term; for example, 1000+ Gbit/s)</li> <li>Overlay networking, Layers 4–7 traffic management (QoS/traffic shaping), pattern matching</li> </ul>
vCPE, SD-WAN and vSD-WAN (centralised SD-WAN capabilities as VNFs)	Router, VPN, firewall, DPI and service chains	<ul style="list-style-type: none"> <li>IPsec, encryption/decryption, pattern matching, overlay networking (VXLAN/GRE/MPLS), SD-WAN security monitoring</li> <li>Centralised, higher-density vSD-WAN nodes that amplifies requirements for all tasks above</li> </ul>
vRAN	Baseband PHY layer (Layer 1) for signal processing	<ul style="list-style-type: none"> <li>Layer 1 and potentially Layer 2/3 packet processing</li> </ul>
vEPC/5G core	P/S GW, deep packet inspection (DPI); in 5G core UPF	<ul style="list-style-type: none"> <li>Layer 2/3/4 forwarding, 5G core increased throughput requirements (for example, 200+ Gbit/s)</li> <li>Overlay networking (VXLAN/GRE/MPLS)</li> <li>Layers 4–7 traffic management (hierarchical QoS)</li> </ul>
vIPsec	Aggregation points (Wi-Fi hotspot, vCPE, service provider edge) and enodeB backhauling	<ul style="list-style-type: none"> <li>Cryptography (encryption, decryption), IPsec protocol, SSL record layer processing</li> <li>Authentication processing</li> </ul>
vNGFW	Firewall, Intrusion Prevention Systems (IPS),	<ul style="list-style-type: none"> <li>Layer 2/3 forwarding, network address and port translation (NAPT), load balancing, pattern matching</li> <li>Cryptography (encryption, decryption)</li> </ul>

Possible use cases	Functions	Potential acceleration points
	SSL VPN and Deep Packet Inspection (DPI)	<ul style="list-style-type: none"> <li>• Compression/decompression</li> </ul>
vIMS	SBC, MRF	<ul style="list-style-type: none"> <li>• Media/audio transcoding</li> </ul>
AI/ML	Deep learning, neural networks, big data analytics	<ul style="list-style-type: none"> <li>• Parallel processing, image/face recognition, natural language processing</li> <li>• Data compression/decompression for analytics</li> </ul>
Video and graphic processing	IPTV/OTT/cable head-end, edge, cloud gaming and AR/VR	<ul style="list-style-type: none"> <li>• Video encoding/transcoding, compression</li> </ul>
IoT	URLLC use cases: factory automation, robotics, health care, smart transportation	<ul style="list-style-type: none"> <li>• Latency – each application requires a different range of latency, some as low as &lt;1–10ms and can be affected by virtualisation overhead</li> <li>• Most of the above functions/workloads can apply if they are deployed to deliver IoT applications.</li> </ul>

These use cases and functions will be deployed in different parts of operators' networks including 4G/5G core, edge locations (multi-access edge, central offices) and enterprise and video networks. Operators will need to identify the functions/workloads that make most sense to accelerate and choose a suitable acceleration solution for these functions and their deployment scenarios.

## 3. Introduction to acceleration technologies

This section introduces the range of acceleration technologies available to NFVI builders and describes their capabilities. It divides acceleration technologies into two categories: software-centric acceleration and hardware-centric and discusses the advantages and challenges of solutions in each category, as well as the potential of hybrid architectures. It summarizes the suitability of the various acceleration solutions for different VNF and NFVI use cases.

### 3.1 Software-centric acceleration

Software-centric acceleration solutions can be implemented as an additional layer in various parts of virtualized networks, for example, in the CPU, hypervisor and VNFs themselves to augment VNF and NFVI performance. Software-centric acceleration solutions are based on acceleration frameworks, such as the Data Plane Development Kit (DPDK), Open Data Plane (ODP) for SoC (System on Chip) and the Storage Performance Development Kit (SPDK) for storage applications. These frameworks leverage the capabilities of underlying chipsets including CPUs and SoCs and provide a set of libraries, drivers and interfaces that enable developers to build acceleration solutions on top of the underlying hardware for specific network virtualization demands.

#### Data Plane Development Kit (DPDK)

DPDK is a widely used software-centric acceleration framework created by Intel in 2010. Its source code was made available to developers under the Open Source BSD License in 2017. DPDK provides data plane libraries to accelerate Layer 3 packet processing and throughput performance of virtualized resources on various CPU architectures (such as Intel®-based architecture, ARM and POWER) and NICs from multiple vendors. Through

its developer ecosystem, DPDK supports a growing variety of data plane functions and use cases in virtualized networks (including vSwitch, crypto, compression and baseband acceleration). vSwitch acceleration is one of its most common use cases.

### NFVI virtual networking acceleration: accelerated vSwitch, SR-IOV and PCI passthrough

Operators are presented with various NFVI networking options for forwarding traffic from north to south - from network interfaces (for example, NIC) to virtual machines (VMs), as well as east to west and between VNFs and their component VMs or containers. The main approaches to traffic forwarding use one or more of the following: a software-based virtual switch (vSwitch), single root I/O virtualization (SR-IOV) and/or PCI passthrough. There are distinct advantages and disadvantages associated with each approach to NFVI networking and these approaches can be combined in hybrid architectures depending on use case demands.

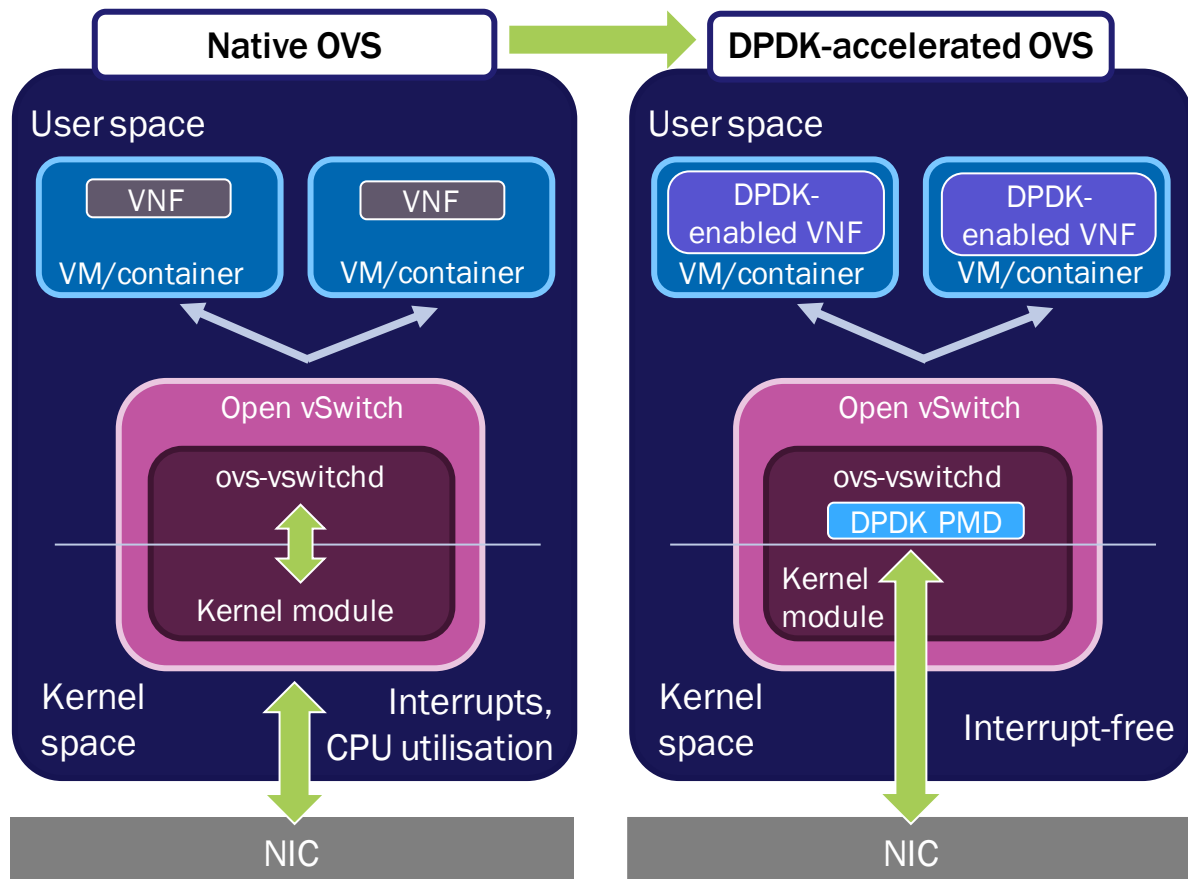
vSwitch is a software layer that sits in the hypervisor, where it aggregates and controls NFVI traffic. It provides isolation from the underlying hardware for network I/O and manages network traffic to and from VNFs. The main advantage of vSwitch is that it is abstracted from the hardware layer and, when coupled with an SDN controller and orchestration, provides a highly flexible and programmable infrastructure. However, its native software implementation without acceleration may not meet the requirements of high-packet processing and throughput use cases. DPDK plays a crucial role in addressing this problem by optimising the way vSwitches process and forward packets. An example of this is DPDK Accelerated Open vSwitch (OVS).<sup>2</sup>

OVS is a vSwitch for Linux-based hypervisors (such as KVM) and consists of two main components: kernel space and user space. Native deployment of OVS introduces performance bottlenecks due to ‘interrupts’ that occur when a packet received from the NIC is first processed in kernel space and then moved to the application in the user space. This process increases CPU utilization and creates overhead. An OVS with DPDK removes this overhead by bypassing the kernel space with an interrupt-free approach. Packets are polled (using poll-mode driver (PMD)) and moved directly to the application in user space. DPDK also accelerates the forwarding between VNFs and vSwitches. Figure 3.1 below illustrates OVS acceleration with DPDK. Intel tests show that DPDK can improve OVS performance by around 12 times and continues to improve.<sup>3</sup>

<sup>2</sup> DPDK supports various vSwitch offerings in the market including open-source (OVS, VPP) and vendor-specific solutions (VMware, Cisco). OVS is a widely adopted open-source vSwitch solution and plays an important role in OpenStack and OpenDayLight implementations.

<sup>3</sup> Testing conducted by Intel as of March 2016. Hardware configurations: Server platform: Supermicro X10DRH-I server with dual-integrated 1 GbE ports via Intel® Ethernet Controller I350-AM2 gigabit Ethernet; Chipset: Intel® C612 chipset; Processor: 1x Intel® Xeon® processor E5-2695 v4 @ 2.10 GHz, 120 W, 45 MB cache processor, 18 cores, 36 hyper-threaded cores per processor; Memory: 64 GB total: Samsung 8GB 2Rx8 PC4-2400MHz, 8GB per channel, 8 Channels; Local storage: 500 GB HDD Seagate SATA Barracuda 7200.12 (SN:Z6EM258D; PCIe: 2 x PCI-E 3.0 x8 slot; NICs: 2 x Intel® Ethernet Converged Network Adapter X710-DA4, total: 8 Ports; 2 ports from each NIC used in tests; BIOS: AMIBIOS Version: 2.0 Release Date: 12/17/2015. Software: Host OS: Fedora 23 x86\_64 (Server version), Kernel version: 4.2.3-300.fc23.x86\_64; VM OS: Fedora 23 x86\_64 (Server version), Kernel version: 4.2.3-300.fc23.x86\_64; QEMU-KVM: QEMU-KVM version 2.5.0, libvirt version: 1.2.18.2-2.fc23.x86\_64. Open vSwitch: Open vSwitch 2.4.9 Commit ID: 53902038abe62c45ff46d7de9dcec30c3d1d861e, Intel® Ethernet Drivers: i403-1.4.25, Intel® Ethernet Converged Network Adapters X710-DA4; DPDK: DPDK version 2.2.0. For more information, see Intel Software Developer Zone (19 December 2016), *Open vSwitch\* with DPDK Overview*. Available at: <https://software.intel.com/en-us/articles/open-vswitch-with-dpdk-overview>.

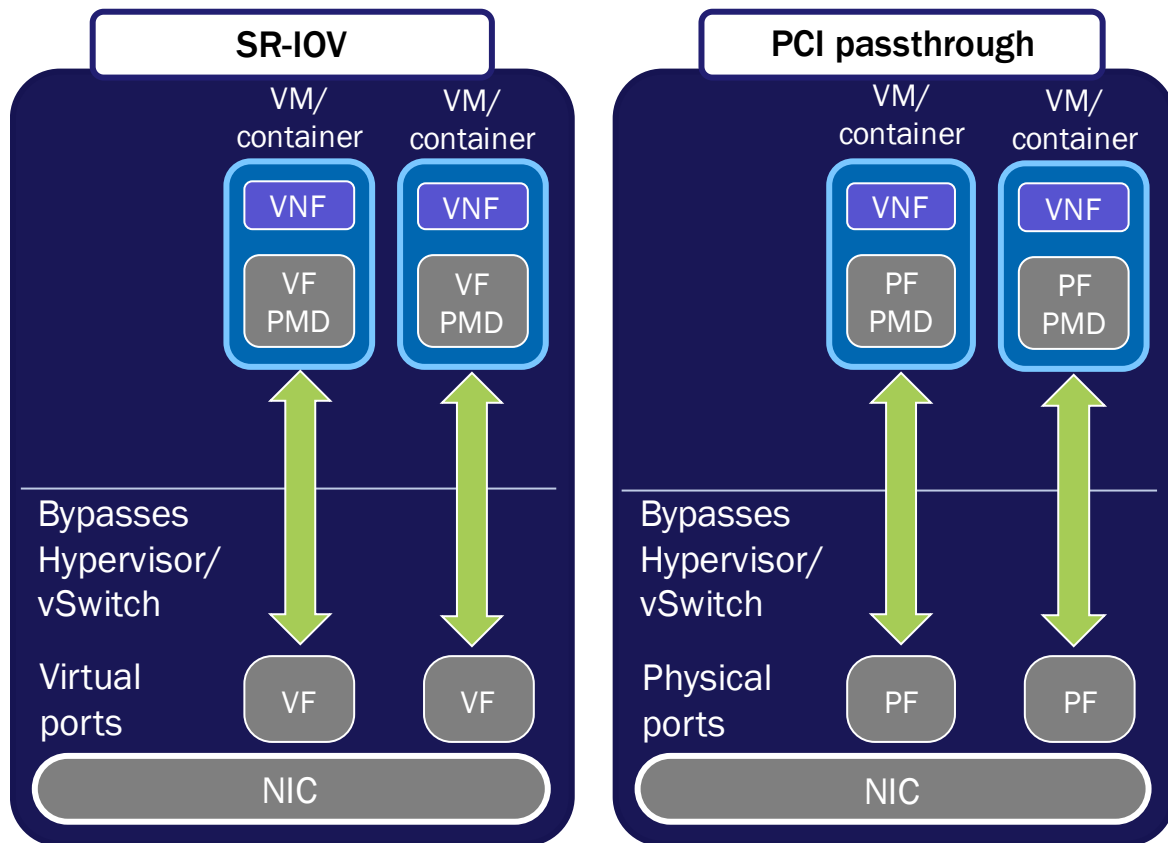
Figure 3.1: Native OVS and DPDK accelerated OVS [Source: Analysys Mason, Intel 2019]



PCI passthrough and SR-IOV are alternative options to vSwitch. They are similar technologies that enable operators to assign NIC resources directly to specific VMs and containers, bypassing the hypervisor and vSwitch. **PCI passthrough** allows an entire single physical NIC port (PF) to be dedicated to a specific VM. **SR-IOV** is a more-advanced standard than PCI passthrough. It exposes the ports of a single NIC device as shared resources, in other words, as virtual ports (VFs) that can be shared by multiple VMs and containers. Both technologies require hardware drivers that sit inside VMs. Figure 3.2 below provides a description of these approaches.



Figure 3.2: SR-IOV and PCI-passthrough models [Source: Analysys Mason, 2019]



SR-IOV and PCI passthrough have been widely used in virtualized networks because they enable high throughput (near line-rate) and low-latency packet processing at low or no CPU utilization. However, these benefits may come at the expense of hardware independence, flexibility and portability, which are the main goals of NFV/SDN deployments. This is because these approaches tie VNFs rigidly to the infrastructure equipped with these capabilities and bypassing the virtualization layer introduces additional complexity in the management layer. Figure 3.3 below summarises the main advantages and disadvantages of accelerated vSwitch and SR-IOV and PCI passthrough.

Figure 3.3: Comparison of accelerated vSwitch, SR-IOV and PCI passthrough [Source: Analysys Mason, Intel 2019]

Criteria	Accelerated OVS	SR-IOV/PCI passthrough
North/south traffic (maximum throughput to VM/containers)	Lower performance than SR-IOV and PCI passthrough	Higher (near-line rate performance)
East/west traffic in the same node (VM to VM, for example, service chaining)	High performance	Lower performance than accelerated vSwitch because packets need to traverse the network interface in each hop, adding overhead
Latency	Higher than SR-IOV and PCI passthrough	Low
CPU utilization	Higher than SR-IOV and PCI passthrough but significantly lower than native OVS	Low or no impact

Criteria	Accelerated OVS	SR-IOV/PCI passthrough
Hardware dependency/lock-in	No, decoupled from underlying hardware with common interfaces (such as VirtIO)	High, tightly coupled; may require vendor-specific network hardware interfaces and drivers to reside in VMs
SDN programmability	High – thanks to centralised management	Limited – despite being supported by several SDN controllers and orchestrators, its management and control is complex unlike vSwitch
VNF portability	High, supports live migration	Limited – VMs are rigidly tied to specific hardware ports

OVS-DPDK has other limitations beyond providing lower performance than SR-IOV. For instance, because it bypasses the kernel space with its in-built maturity mechanisms, it presents certain security challenges. OVS-DPDK has been found to reduce performance in cases where overlay networking (VXLAN, GRE) is needed because these packets still need to traverse kernel space, which results in inefficiencies.<sup>4</sup> Other solutions have been introduced to tackle these issues in accelerating OVS, including TC flower-based OVS, which uses the TC flower classifier in the Linux subsystem to accelerate packet processing, as well as vendor-specific offerings.

Both types of software-centric acceleration co-exist in the market. Hybrid acceleration architectures that combine vSwitch and SR-IOV/PCI passthrough in the same NFVI create additional complexity, however. Operators may prefer to use the former for use cases that involve service chaining/high inter-VM traffic (such as vBNG, vCPE and Gi-LAN) and use the latter for use cases that require best possible performance and lowest latency (such as vRouters). Some operators are adopting accelerated vSwitch in their NFV/SDN deployments in conjunction with hardware-centric accelerators (for example, SmartNIC) to further enhance NFVI virtual networking performance and achieve a more complete acceleration of the OVS data path.

### 3.2 Hardware-centric acceleration

Hardware-centric acceleration solutions consist of specialised hardware components that enable higher performance and efficiency than general-purpose computing, networking and storage resources in an NFVI. Such components provide pure hardware acceleration by augmenting general-purpose resources (working in a complementary way or independently) and can be used in combination with software acceleration technologies (including DPDK, SR-IOV and TC flower) as discussed in section 3.1.

Examples of domains that may benefit from accelerating tasks to more-efficient hardware-based acceleration components include:

- NFVI networking acceleration, for example, vSwitch packet processing
- VNFs with complex, compute-intensive network and non-network functionalities (including IPsec, encryption, transcoding and tunnel processing)
- storage (storage virtualization, networking, software-defined storage and hyperconverged infrastructure).

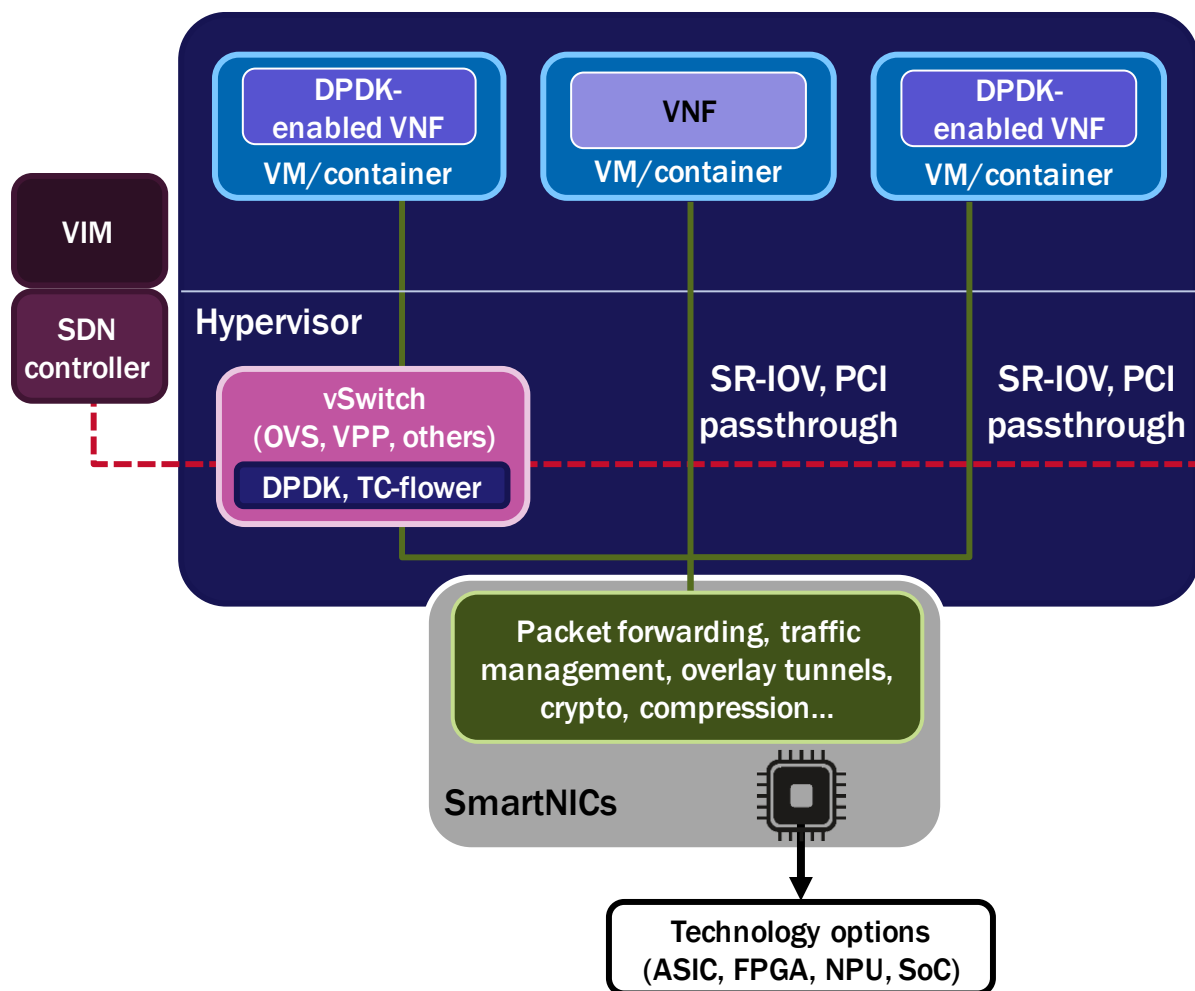
Many hardware-centric acceleration solutions are available for virtualized networks based on different types of chipset, including application-specific integrated circuit (ASIC), field programmable gate array (FPGA),

<sup>4</sup> European Telecommunications Standards Institute (December 2015), *ETSI GS NFV-IFA 001: Network Functions Virtualisation (NFV); Acceleration Technologies; Report on Acceleration Technologies & Use Cases*. Available at: [https://docbox.etsi.org/ISG/NFV/Open/Publications\\_pdf/Specs-Reports/NFV-IFA%20001v1.1.1%20-%20GS%20-%20Acceleration%20-%20UCs%20report.pdf](https://docbox.etsi.org/ISG/NFV/Open/Publications_pdf/Specs-Reports/NFV-IFA%20001v1.1.1%20-%20GS%20-%20Acceleration%20-%20UCs%20report.pdf).

network processing unit (NPU) and multicore processors, which may be implemented – for performance and low-power reasons – as systems-on-a-chip (SoCs). Hardware-centric acceleration solutions are deployed in a wide variety of use cases and scenarios but also often for scenario-specific accelerating purposes. GPUs, vision processing units (VPU) and tensor processing units (TPUs) can all give excellent results for specific acceleration cases including video processing, vision systems and AI/ML respectively. These hardware solutions come in different form factors and can reside in various parts of the NFVI.

SmartNICs are evaluated by operators because they provide a natural evolution path from conventional NICs by adding programmable acceleration. SmartNIC solutions can be based on different chipsets and architectures (including ASIC, FPGA, SoC and NPU) and there is a large ecosystem of vendors that provide these solutions including Broadcom, Ethernet Networks, Marvell, Mellanox, Napatech, Netronome and Xilinx. Other options include PCI-linked devices (for example, Intel® QuickAssist Technology (Intel® QAT) accelerator, discussed in section 3.3.), network attached devices or integrated architectures (such as Intel® Xeon® D processors with integrated Intel® QAT).

Figure 3.4: VNF and NFVI acceleration using Smart NICs [Source: Analysys Mason, 2019]



Hardware-centric accelerators boost VNF and NFVI performance, but they can also add to NFVI costs and operational complexity without standardization and management abstraction. Operators need to carefully evaluate the cost/benefit analysis of additional hardware acceleration in NFVI and should be prepared to handle the additional management complexity of heterogenous, multi-vendor acceleration devices.

No single technology option for hardware-centric acceleration prevails; there are trade-offs in terms of performance, cost and programmability in each of the various options, as detailed in Figure 3.5.

**Figure 3.5: Comparison of hardware acceleration options [Source: Analysys Mason, 2019]**

Hardware acceleration	Performance	Price–performance	Programmability and Flexibility
ASIC	●●●● Specifically designed for an application/workload	●●●● Excellent cost-performance	●○○○ Limited by the initial design capabilities
FPGA	●●●○ High performance but small trade-off with flexibility	●●●○ High performance at high cost; prices are decreasing	●●●○ Can be programmed on the fly, but requires specialist programming skills or vendor support
NPU	●●●○ High performance but small trade-off with flexibility	●●●○	●●○○ Can be programmed but has limited focus on network I/O, solutions are typically vendor-proprietary
SoC	●●○○ High flexibility - performance trade-off	●●●○ Good performance at moderate cost	●●●● General purpose, C language, Linux

Operators' choice of technology will be dependent on use case demands. For example, ASIC-based acceleration is more suitable to scenarios where acceleration requirements do not change often, and maximum cost–performance is desired. FPGA would suit specific edge deployments where application requirements are dynamic and require on-the-fly programmability (for example, perform DPI during the day and switch to transcoding to support a sport events broadcast in the evening). SoCs can be deployed for security applications such as encryption acceleration or transcoding use cases.

### AI/ML, next-generation video and AR/VR services will need acceleration

Graphics Processing Units (GPUs), such as those provided by NVIDIA and AMD, are commonly used for multimedia processing workloads such as video encoding/transcoding and compression. New generations of codecs such as H.265, VP9 and AV1 are increasingly being adopted to deliver high-quality video services (4K/8K) and low-latency, high-bandwidth AR/VR services. Virtualized video networks and edge computing services would benefit from GPU acceleration and accelerate to meet the requirements of these services.

GPUs have become a popular solution for AI/ML demands in the past 2 years. The parallel processing capabilities of GPUs bring significant improvement over CPUs; GPUs typically comprise thousands of cores, which can perform millions of calculations in parallel as demanded by tasks such as deep learning, neural networks, image/face recognition and natural language processing. Other hardware acceleration options such as FPGAs, ASICs and SoCs, as well as specialised solutions such as Google's TPUs are also increasingly used to accelerate computationally intensive AI/ML workloads.

Vision processing units (VPUs), such as Intel® Movidius™ Myriad™ X VPUs, are emerging as purpose-built processors for computer vision applications. They process real-world elements and images including AR/VR, robotics and drones. VPUs are specifically designed for low power and high performance for mobility requirements, which may not be met with existing GPU architectures.

### 3.3 Acceleration for NFV/SDN security functions: Intel QAT and HPE performance results

NFV/SDN applications such as vCPE, SD-WAN, vADC, vIPsec, vGi-Lan and vNGFW need to perform compute-intensive security tasks (such as cryptography including encryption/decryption (SSL/TLS), authentication, public key functions, compression and decompression) with high throughput and low latency. The performance, QoS and CPU utilization of these applications can be significantly improved by using hardware-centric accelerators, such as Intel QuickAssist Technology (Intel QAT),<sup>5</sup> used in conjunction with CPU to assist with encryption and compression workloads. Figure 3.6 provides an overview of performance improvements achieved using Intel QAT acceleration in several NFV/SDN use case tests.

**Figure 3.6: Intel QAT acceleration performance results on HPE servers [Source: Intel, HPE, 2019]**

Use case	VNF	Configuration	Performance Improvements (w/Intel QAT)
Virtual Gi-LAN Firewall (vGi-Lan)	F5 BIG-IP Virtual Edition v14.1	VNF using 8 vCPUs (Second Generation Intel® Xeon® Gold 6230N processor-based server) with Intel® QAT providing public key encryption for TLS/SSL.	<p>Higher TLS/SSL performance<sup>6</sup>:</p> <ul style="list-style-type: none"> <li>• <b>5x</b> greater TLS/SSL transactions per second (TPS) (TLS1.2 AES128-SHA 2K Key): from 7034 TPS to 34173 TPS</li> <li>• <b>2.4x</b> greater TLS/SSL bulk encryption throughput: from 7744 Mbit/s to 15360 Mbit/s</li> </ul> <p>The performance boost with Intel® QAT is achieved at 56% CPU utilization while the baseline is at 99.9%.</p>
Virtual SD-WAN (vSD-WAN)	Nuage Networks SD-WAN 5.3.3 NSGv (Virtual Network Service Gateway)	VNF using 4 vCPUs (Second Generation Intel® Xeon® Gold 5218N processor-based server) with Intel® QAT providing encryption/decryption acceleration for IPsec.	<p>Higher IPsec performance<sup>7</sup>:</p> <ul style="list-style-type: none"> <li>• <b>2x</b> greater IPsec throughput</li> <li>• <b>2x</b> greater VNF capacity for more applications on the same server node.</li> </ul>
Virtual Next Generation Firewall (vNGFW)	Fortinet Fortigate	VNF using 4 vCPUs (Second Generation Intel® Xeon® Gold	<p>Higher VPN Performance<sup>8</sup></p> <ul style="list-style-type: none"> <li>• <b>3x</b> greater VPN throughput</li> </ul>

<sup>5</sup> Intel QAT can take various forms, for example, integrated into CPU (Intel Xeon processor), in a System on Chip package (Intel Atom® SoC), or as PCI linked accelerator cards (Intel QuickAssist Adapter).

<sup>6</sup> Testing conducted by F5 as of March 2018. Configurations: Intel® QuickAssist Technology and Second Generation Intel® Xeon® Scalable Processors (QTY 2) 6230N with 192 GB total memory (12 slots / 16GB / DDR4 2667MHz), µcode 0x4000019, Bios: PLYXCRB 1.86B.0568.D10.1901032132, µcode: 0x4000019 on CentOS 7.5 with Kernel 3.10.0-862, KVM Hypervisor; 1 x Dual Port 40GbE Intel® Ethernet Network Adapter XL710; ; 1 x Intel® QuickAssist Adapter 8970, TLS1.2: AES128-GCM-SHA256 2K key with 3x QAT Physical Functions; Application: BIG-IP Virtual Edition (VE) v14.1 with Intel QAT Enabled.

<sup>7</sup> Testing conducted by Intel as of February 2018. Configurations: Intel QuickAssist Technology and Second Generation Intel Xeon Scalable processors (QTY 2) 5218N with 192 GB total memory (12 slots / 16GB / DDR4 2667MHz), µcode 0x4000019, Bios: PLYXCRB 1.86B.0568.D10.1901032132, µcode: 0x4000019 on CentOS 7.5 with Kernel 3.10.0-862, KVM Hypervisor; 1x Intel QuickAssist Adapter 8970, Cipher: AES-128 SHA-256; Intel Ethernet Converged Network Adapter X520-SR2; Application: Nokia Nuage SDWAN NSGv 5.3.3U3.

<sup>8</sup> Testing conducted by Intel as of March 2018. Configurations: Intel QuickAssist and Second Generation Intel Xeon Scalable processors (QTY 2) 6230N with 192 GB total memory (12 slots / 16GB / DDR4 2667MHz), µcode 0x4000019, Bios: PLYXCRB 1.86B.0568.D10.1901032132, µcode: 0x4000019 on CentOS 7.5 with Kernel 3.10.0-862, KVM Hypervisor; 1 x Intel QuickAssist Adapter 8970, IPsec: AES128-SHA256; 1 x Dual Port 40GbE Intel Ethernet Network Adapter XL710; Application: FortiGate VM64-KVM.



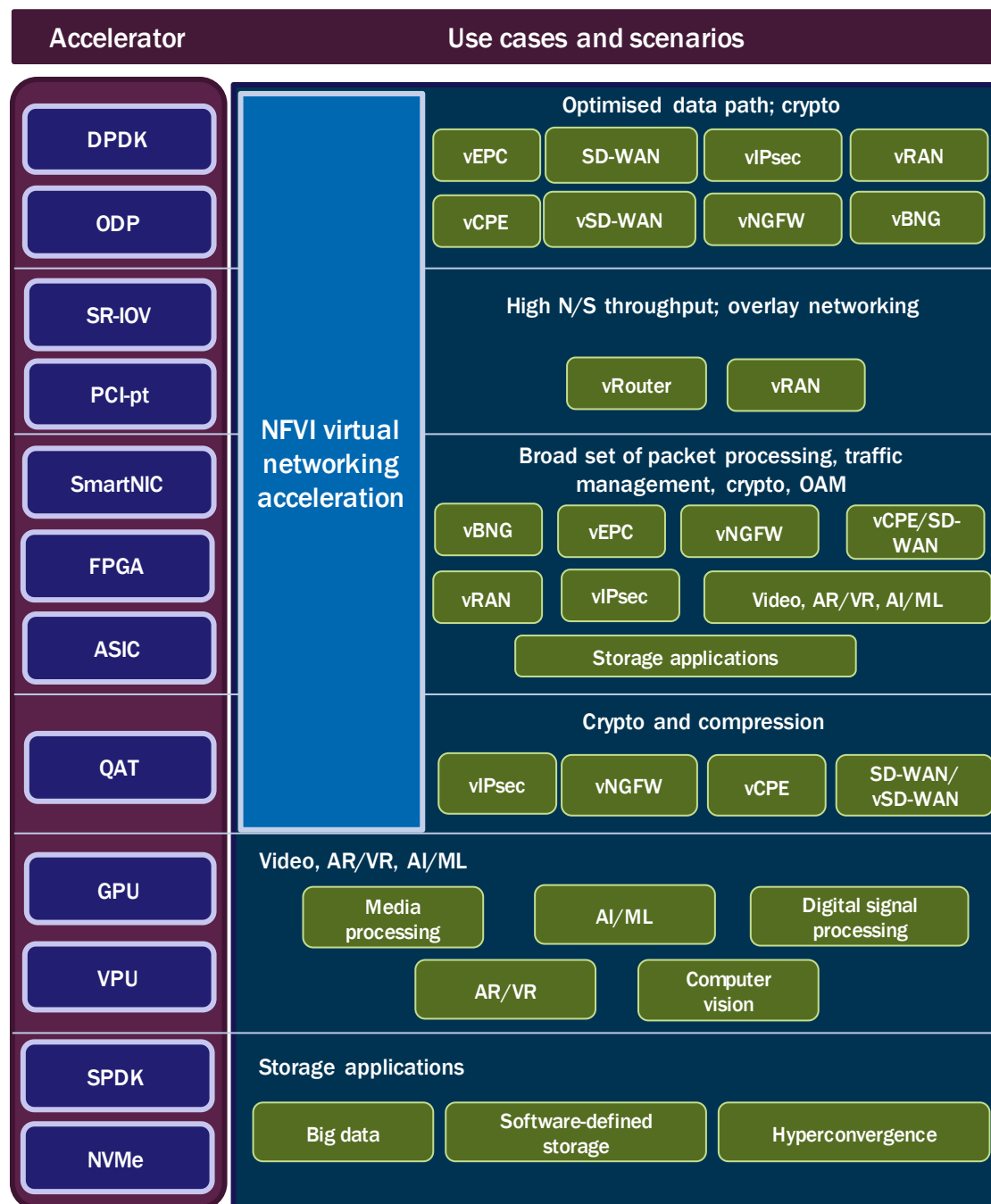
Use case	VNF	Configuration	Performance Improvements (w/Intel QAT)
		6230N processor-based server) with Intel® QAT providing encryption/decryption acceleration for VPN.	

Performance continues to improve for commercial VNF providers that enable Intel QuickAssist Technology to accelerate security and compression workloads, especially with the recently launched Second Generation Intel® Xeon® Scalable processors, including the Intel Xeon Scalable Gold processors used in this testing.

### 3.4 Summary comparison of acceleration technologies

Operators implementing NFVIs that will support multiple VNFs and digital services, both today and in future, will need to understand the range of available acceleration technologies and their applicability to individual NFV and SDN use cases. Figure 3.6 summarizes the different types of acceleration technology and the use cases for which they are suitable. It should be noted that this is not an exhaustive list of use cases/scenarios for each of the acceleration options. In addition, these are not mutually exclusive (multiple options can be combined for a specific use case).

Figure 3.7: Summary comparison of acceleration technologies (Source: Analysys Mason, 2019)



## 4. Acceleration-as-a-service is crucial to abstract the complexity of acceleration resources

As the number and type of acceleration technologies needed in an NFVI grows, operators face a substantial challenge in managing both the technologies themselves and VNF access to them. As explained previously, acceleration technologies come in different forms (software, hardware, hybrid), are implemented in different

devices (smart NICs, crypto cards) in different components of the NFVI (PCI linked, integrated with the CPU, or even remote) and are supplied by a broad set of vendors. The requirement for a heterogeneous NFVI that is able to support multiple use cases and scenarios with high-performance demands is spurring a strong, parallel need for an acceleration technology abstraction layer. The NFVI would use such a layer to manage a complex set of acceleration resources and expose their capabilities through a common API in an ‘as-a-service’ model to virtualized networks.

Such an acceleration abstraction layer should:

- ensure the complete independence of VNFs from the NFVI. It should provide a common interface across all acceleration technologies, regardless of type or vendor, to prevent vendor lock-in and eliminate the need to re-write/modify VNFs to work with a specific acceleration technology.
- streamline and automate the lifecycle management of acceleration resources (identify/discover, provision, configuration, maintenance and monitoring), and acceleration scenarios, including assigning the right acceleration resources dynamically to adapt to workloads/use cases in edge infrastructure.

The industry is making significant efforts to achieve these goals. For example, ETSI NFV ISG defines an acceleration abstraction layer (AAL), which enables the decoupling of VNFs from the underlying hardware accelerators and provides abstraction interfaces for VNFs and supports for the management of these acceleration resources in the VIM.<sup>9</sup> Open-source initiatives are also following this course. OPNFV has the Data Plane Acceleration Project (DPACC), which aims to provide a common suite of abstraction APIs to enable VNF portability and resource management for integrated SoC accelerators. OpenStack has responded to industry demands with its Cyborg project. OpenStack Cyborg<sup>10</sup> is an emerging general-purpose management framework for the automated lifecycle management of software and hardware acceleration resources (including FPGA, GPU, crypto cards, DPDK, ODP and storage). Currently, management of these resources can be complex and time-consuming as each may require specific tuning and configuration. Network engineers need to develop their own Python scripts or Ansible playbooks to carry out these processes. OpenStack Cyborg aims to simplify these processes with a set of software components, RESTful APIs and generic drivers, which enable dynamic resource discovery, attachment/detachment of acceleration devices and driver management, and work in conjunction with OpenStack Nova or standalone in bare metal (for example, OpenStack Ironic).

Acceleration-as-a-service is a key step for the movement towards complete hardware and software independence in virtualized networks. Building open, unified management platforms with common, standardised interfaces will be critical and it will require close collaboration between operators, vendors and industry groups to achieve it.

## 5. Conclusion and recommendations

Operator networks are transforming into software-defined, automated, cloud-native NFV infrastructures enhanced by advanced analytics, AI/ML techniques to manage the rapid growth in traffic and bandwidth demands and to achieve dynamic, on-demand delivery of network services. 5G and edge computing will usher in the new era of cloud-based services such as AR/VR, robotics and many other URLLC applications that will

<sup>9</sup> ETSI GS NFV-IFA 001, ETSI GS NFV-IFA 002, ETSI GS NFV-IFA 004.

<sup>10</sup> For more information, see <https://wiki.openstack.org/wiki/Cyborg>.

be developed and delivered through these new infrastructures. Acceleration technologies will be pivotal to this transformation by enabling operators to meet the performance, latency, QoS, subscriber density and security requirements of existing and future applications with optimum TCO.

As presented in this paper, many network virtualization, AI/ML, 5G and IoT use cases will require operators to take advantage of acceleration technologies. There are many software- and hardware-centric acceleration solutions available for these use cases, which can be deployed in different parts of the NFVI in different form factors and can be supplied by different vendors or through open-source initiatives. Operators will need to identify and deploy the most-suitable acceleration and accelerating technologies for their application/workload needs, networks (including fixed, wireless, core and access), deployment scenarios and operational capabilities. Moreover, operators need to carefully evaluate the cost, performance, programmability and flexibility of the acceleration solutions because there is no one-size-fits-all approach as they often represent a trade-off. For example, fixed-function (ASIC) accelerators are preferable in scenarios where acceleration requirements are constant and maximum price-performance is desired (for example, 24/7 active security workloads), while FPGA or SoC-based accelerators provide lower price-performance than ASIC but they are more programmable and can adapt to dynamic workloads (for example, in edge infrastructure).

Potential performance and cost benefits of acceleration technologies should not come at the expense of hardware and software independence in virtualized networks. The sheer multitude of use cases and technologies, as well as the current lack of standardization across these technologies and their suppliers, will inevitably increase the heterogeneity in NFVI, which can lead to additional operational complexity and negate the automation, cost and agility benefits of virtualization on COTS hardware. To mitigate these risks, operators need to:

- consider carefully when and where to use acceleration
- use an acceleration technology abstraction layer in the NFVI to decouple VNFs from underlying acceleration resources in order to independently accommodate all their potential use case needs and to automate the lifecycle management of acceleration resources and scenarios with an ‘as-a-service’ model.

Strong industry collaboration to set a common framework for acceleration standardization and support for open-source initiatives such as OpenStack Cyborg and OPNFV DPACC will be crucial in order to achieve operators’ goal to create horizontal and highly automated NFVIs.

## 6. About the authors



**Gorkem Yigit** (Senior Analyst) is the lead analyst for the Video and Identity Platforms programme and a contributor to the Digital Infrastructure Strategies and Network Automation and Orchestration programmes, focusing on producing market share, forecast and research collateral. He has published research on NFV/SDN services business cases, identity management in the digital economy, and has been a key part of major consulting projects including Telco Cloud Index and IPTV/OTT procurement. He holds a cum laude MSc degree in Economics and Management of Innovation and Technology from Bocconi University (Milan, Italy).



**Caroline Chappell** (Research Director) is the lead analyst for Analysys Mason's Digital Infrastructure Strategies research programme. Her research focuses on service provider adoption of cloud, and the application of cloud technologies to fixed and mobile networks. She is a leading exponent of SDN and NFV and the potential that these technologies have to enhance business agility and enable new revenue opportunities for service providers. Caroline investigates key cloud and network virtualization challenges and helps telecoms customers to devise strategies that mitigate the disruptive effects of cloud and support a smooth transition to the era of software-controlled networks.

This white paper was commissioned by HPE. Analysys Mason does not endorse any of the vendor's products or services.



**Hewlett Packard Enterprise** Hewlett Packard Enterprise advances the way people live and work. Learn more at [hpe.com/dsp/infrastructure](http://hpe.com/dsp/infrastructure).

Published by Analysys Mason Limited • Bush House • North West Wing • Aldwych • London • WC2B 4PJ • UK  
Tel: +44 (0)20 7395 9000 • Email: [research@analysysmason.com](mailto:research@analysysmason.com) • [www.analysysmason.com/research](http://www.analysysmason.com/research)

Registered in England No. 5177472

© Analysys Mason Limited 2019

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means – electronic, mechanical, photocopying, recording or otherwise – without the prior written permission of the publisher.

Intel, the Intel logo, Intel Atom, Xeon, and Movidius and Myriad are trademarks of Intel Corporation or its subsidiaries in the U.S. and/or other countries.

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors.

Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information visit [www.intel.com/benchmarks](http://www.intel.com/benchmarks).

Performance results are based on testing as of the dates in the configurations and may not reflect all publicly available security updates. See configuration disclosure for details. No product or component can be absolutely secure.

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. No product or component can be absolutely secure. Check with your system manufacturer or retailer or learn more at [intel.com](http://intel.com).

Intel does not control or audit third-party data. You should review this content, consult other sources, and confirm whether referenced data are accurate.

Cost reduction scenarios described are intended as examples of how a given Intel-based product, in the specified circumstances and configurations, may affect future costs and provide cost savings. Circumstances will vary. Intel does not guarantee any costs or cost reduction.

Figures and projections contained in this report are based on publicly available information only and are produced by the Research Division of Analysys Mason Limited independently of any client-specific work within Analysys Mason Limited. The opinions expressed are those of the stated authors only.

Analysys Mason Limited recognises that many terms appearing in this report are proprietary; all such trademarks are acknowledged and every effort has been made to indicate them by the normal UK publishing practice of capitalisation. However, the presence of a term, in whatever form, does not affect its legal status as a trademark.

Analysys Mason Limited maintains that all reasonable care and skill have been used in the compilation of this publication. However, Analysys Mason Limited shall not be under any liability for loss or damage (including consequential loss) whatsoever or howsoever arising as a result of the use of this publication by the customer, his servants, agents or any third party.



## 7. Analysys Mason's consulting and research are uniquely positioned

Analysys Mason is a trusted adviser on telecoms, technology and media. We work with our clients, including communications service providers (CSPs), regulators and end users to:

- design winning strategies that deliver measurable results
- make informed decisions based on market intelligence and analytical rigour
- develop innovative propositions to gain competitive advantage.

We have around 220 staff in 14 offices and are respected worldwide for the exceptional quality of our work, as well as our independence and flexibility in responding to client needs. For over 30 years, we have been helping clients in more than 110 countries to maximise their opportunities.

### Consulting

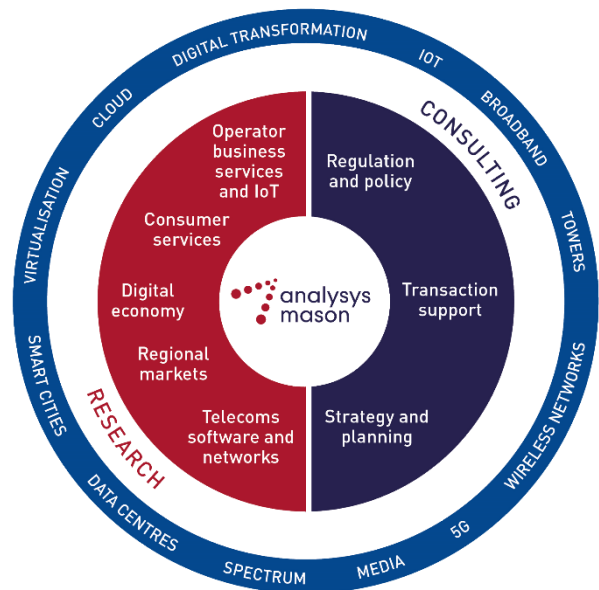
- We deliver tangible benefits to clients across the telecoms industry:
  - communications and digital service providers, vendors, financial and strategic investors, private equity and infrastructure funds, governments, regulators, broadcasters, and service and content providers.
- Our sector specialists understand the distinct local challenges facing clients, in addition to the wider effects of global forces.
- We are future-focused and help clients understand the challenges and opportunities that new technology brings.

### Research

Our dedicated team of analysts track and forecast the different services accessed by consumers and enterprises.

We offer detailed insight into the software, infrastructure and technology delivering those services.

Clients benefit from regular and timely intelligence, and direct access to analysts.



## 8. Research from Analysys Mason

We provide dedicated coverage of developments in the telecoms, media and technology (TMT) sectors, through a range of research programmes that focus on different services and regions of the world.

The division consists of a specialised team of analysts, who provide dedicated coverage of TMT issues and trends. Our experts understand not only the complexities of the TMT sectors, but the unique challenges of companies, regulators and other stakeholders operating in such a dynamic industry.

Our subscription research programmes cover the following key areas.



Each subscription programme provides a combination of quantitative deliverables, including access to more than 3 million consumer and industry data points, as well as research articles and reports on emerging trends drawn from our library of research and consulting work.

**Our custom research service offers in-depth, tailored analysis that addresses specific issues to meet your exact requirements.**

Alongside our standardised suite of research programmes, Analysys Mason's Custom Research team undertakes specialised, bespoke research projects for clients. The dedicated team offers tailored investigations and answers complex questions on markets, competitors and services with customised industry intelligence and insights.

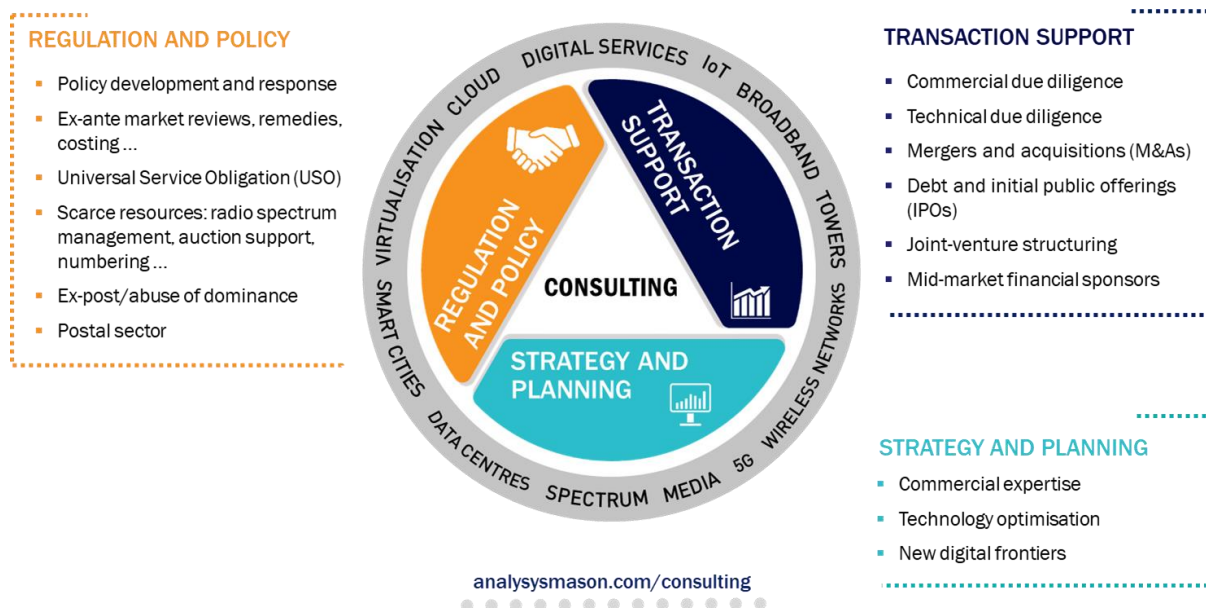
For more information about our research services, please visit [www.analysysmason.com/research](http://www.analysysmason.com/research).

## 9. Consulting from Analysys Mason

For more than 30 years, our consultants have been bringing the benefits of applied intelligence to enable clients around the world to make the most of their opportunities.

Our clients in the telecoms, media and technology (TMT) sectors operate in dynamic markets where change is constant. We help shape their understanding of the future so they can thrive in these demanding conditions. To do that, we have developed rigorous methodologies that deliver real results for clients around the world.

Our focus is exclusively on TMT. We advise clients on regulatory matters, help shape spectrum policy and develop spectrum strategy, support multi-billion dollar investments, advise on operational performance and develop new business strategies. Such projects result in a depth of knowledge and a range of expertise that sets us apart.



We look beyond the obvious to understand a situation from a client's perspective. Most importantly, we never forget that the point of consultancy is to provide appropriate and practical solutions. We help clients solve their most pressing problems, enabling them to go farther, faster and achieve their commercial objectives.

For more information about our consulting services, please visit [www.analysismason.com/consulting](https://www.analysismason.com/consulting).