

# Exploring Computing Platforms for Radio Access Networks

Heterogeneous Computing at the Network Edge

[Andy Butcher](#), Technical Staff, Server Advanced Engineering, Dell EMC

[Joseph Boccuzzi](#), Principal Wireless Access Architect, Network & Custom Logic Group, Intel Corporation

## Abstract

A commentary is provided on server workloads pertinent to the wireless communications industry including 4G and 5G cellular. A summary of experiments in the Dell labs details benefits of PowerEdge products with FPGA offloads.

## Executive Summary

Communications service providers are envisioning increased demand for mobile services including media and content delivery, mobile gaming, augmented and virtual reality, and connected vehicles. To satisfy this emerging demand, the buildout of 5G cellular infrastructure has commenced. Computing infrastructure is behind the scenes supporting radio access networks and other core services so this wireless ecosystem can grow. This paper analyzes the current workloads surrounding radio access networks, how they are changing, and where opportunities for offload exist to support enhanced computing performance. This paper briefly discusses “Edge” computing workloads, a closely related topic as infrastructure sites that support the radio network are expected to host multi-access edge computing applications as well.

# Table of Contents

1	Introduction .....	4
2	C-RAN Background .....	4
3	3GPP 4G Network Architecture and Software Stack .....	7
3.1	4G Network Architecture .....	7
3.2	4G Software Stack .....	9
3.3	4G Edge Services Network Architecture .....	10
4	3GPP 5G NR (New Radio) Network Architecture and Software Stack .....	11
4.1	5G Network Architecture .....	11
4.2	Example FH & BH BW Requirements for 4G & 5G .....	13
4.3	5G Software Stack .....	14
4.4	5G Edge Service Network Architecture .....	15
4.5	5G Network Slicing Example .....	16
4.6	Software Offload Functionality .....	17
5	FlexRAN test setup .....	18
5.1	Components and System diagram .....	18
5.2	Results .....	19
6	Platform Architectures .....	21
6.1	Flexible Architectures for the Network Edge .....	21
6.2	Platform Example .....	21
6.3	Products .....	23
6.3.1	Rack Servers .....	23
6.3.2	Modular Data Centers .....	23
7	Conclusions .....	23
8	Acknowledgements .....	23
9	References .....	24

# 1 Introduction

A collective of communications service providers wrote a seminal whitepaper in 2012 describing a desire to transform their infrastructure to utilize standard IT equipment (servers, switches, storage appliances) to perform functions traditionally accomplished by customized equipment like routers, firewalls, and radio access network nodes.<sup>1</sup> The concept of NFV, Network Functions Virtualization, was born. This led to a significant standardization effort headed by the NFV ISG (Industry Standards Group) within ETSI. Naturally, server manufacturers would be excited about this transformation, so how to best optimize products for this market has been an important topic. With the advent of 5G Cellular, this is an ongoing conversation. Indeed, in the original whitepaper, mobile network nodes are specifically mentioned as one of the many applications for NFV.

As the original whitepaper stated, one of the technical concerns associated with executing workloads on general purpose CPUs has always been performance and latency for data plane applications, and a corollary to this concern was avoiding the need for proprietary hardware to alleviate any performance degradation. Despite the mention of avoiding proprietary hardware, a category in which some accelerators might be grouped, potential resources in the toolbox available to solutions providers are FPGA peripherals to augment the CPUs. In fact, it could be argued that these programmable devices are not proprietary at all; they are becoming more widely used in Enterprise IT markets, and for applications where latency and fast operations not directly supported by standard instruction sets are needed, FPGAs can be an effective part of the solution. This is recognized by the ISG as evidenced by the related discussions on hardware abstraction layers. The need is further validated by the development of BBDev in the DPDK open source project.<sup>2</sup>

This paper will explore an application of FPGA acceleration for this market. In particular, C-RAN (Centralized Radio Access Networks) is the primary focus for this investigation. Additional background is provided in anticipation of future workload optimizations in this area. Also, because of the transformative nature of this topic, some highlights of Multi-access Edge Computing are discussed.

## 2 C-RAN Background

A traditional 4G LTE implementation of a radio access network is shown in Figure 1. At the heart of this implementation is the baseband unit (BBU), which provides the back-haul interface to the mobile network core and the front haul interface to the remote radio head (RRH). The software protocol stack in the radio equipment controller is specified by 3GPP.<sup>3</sup> In the spirit of NFV, the BBU is a candidate for implementation by an enterprise server.

Some of the tenets of C-RAN predate the concept of NFV, as evidenced by the China Mobile paper<sup>4</sup> in which the idea of centralizing baseband functionality was discussed. Among the proposed benefits of centralizing this functionality were lowering the total cost of ownership, improving system capacity utilization, and offering a path to performance improvement at the cell edge.

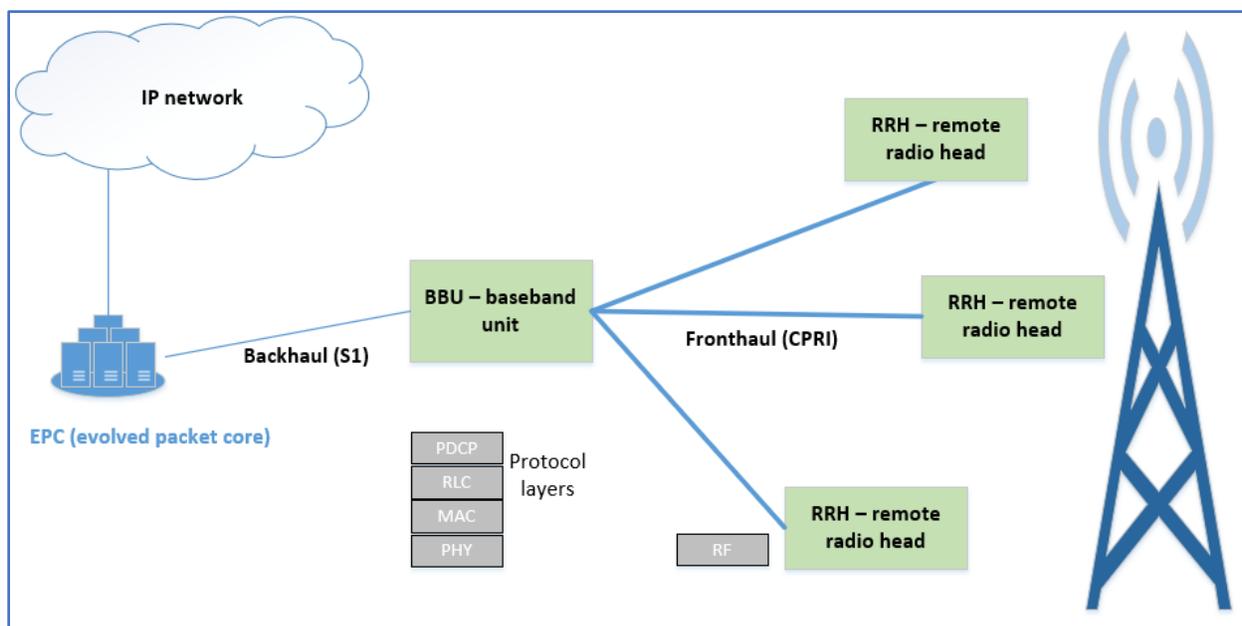


Figure 1: Traditional Radio Access network (4G LTE)

Referring to Figure 2, the following attributes can be associated with a Centralized Radio Access Network:

1. The protocol stack implemented by the BBU can be split in different ways between a Centralized Unit and a Distributed unit, with different implications and tradeoffs for bandwidth and latency. (Note the protocol can also be split in a three-tiered architecture including a CU, DU, and remote radio head RRH.) By moving the PHY layer to a DU, bandwidth on the front haul is greatly relieved compared to the CPRI implementation which entails sending time domain radio IQ samples over the interface.
2. Because of the various possibilities in splitting the protocol between the CU and DU, emerging standards are evolving to define the new front haul interface. The standard defined by the xRAN consortium, absorbed into the ORAN Alliance,<sup>5</sup> is one example.
3. New 5G technology will coexist with 4G LTE radios. In places where CPRI is perpetuated with legacy equipment, controllers for these locations will be co-located with controllers for the new technology. One benefit of servers handling this workload is that the same equipment can handle either scenario with different software.
4. The development of these centralized sites will provide a means for deployment of multi-access edge computing (MEC).<sup>i</sup> This will increase the richness and quality of service of the endpoint use cases.

There are additional benefits envisioned with C-RAN. It enables CoMP, coordinated multipoint transmission, a technique to increase cell edge coverage and throughput through scheduling. By providing connections to several base stations at once, data can be passed through the least loaded base stations for better resource utilization. On the receiver side, performance can be increased by using several sites to selectively use the best signal, which can change rapidly due to multipath fading. Also, centralizing mobility and traffic management decisions can lead to

<sup>i</sup> Formerly a group named “Mobile Edge Computing” within ETSI, European Telecommunications Standards Institute

fewer handoff failures and less network control signaling, which can also be a potential savings with less need for inter-base station networks.

As the concept of C-RAN evolved, the movement to virtualize workloads also pertained to this area, and the term vRAN (Virtualized Radio Access Networks) was created. Also, in parallel, the concept of disaggregated hardware and software led to principles such as SDN (software defined networking), and along with that the term Cloud-RAN was born, and in some cases this created confusion with the overloaded term C-RAN. However, of prime importance remains in the ability to implement the radio access networks using software and standard enterprise servers.

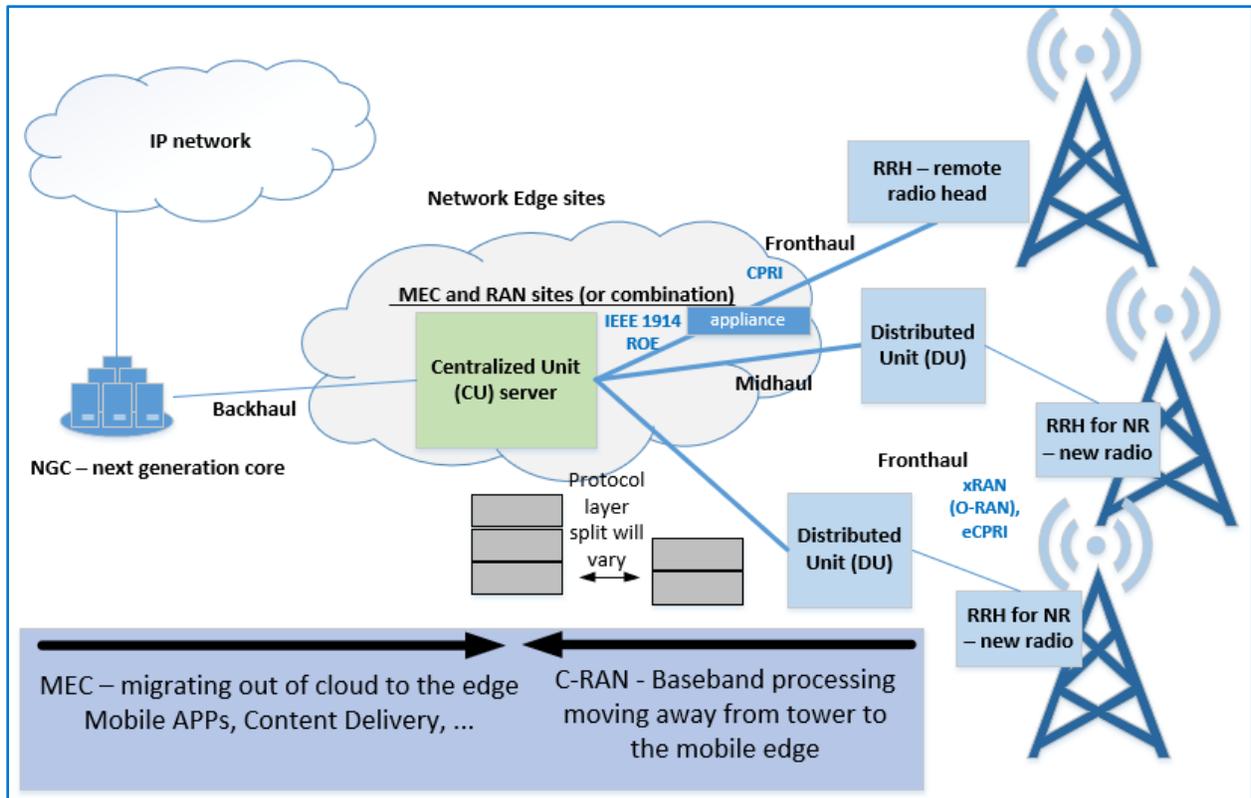


Figure 2: Radio access network employing centralized and distributed units

User traffic is growing exponentially as operators are forced to deliver these demands at lower costs and faster deployments that scale and enable new applications. 4G has been successful in delivering a data network, and 5G is expected to not only extend these capabilities but also deliver a services network, promising to create an ecosystem that delivers a very wide range of vertical markets and their respective applications. Which will have diverse requirements in terms of data rates, latency, mobility, connection density, reliability, spectrum efficiency and energy efficiency. An incremental improvement over 4G does not satisfy the expected degree of flexibility and scalability needed. vRAN helps create a network that is flexible, scalable, achieves lower overall cost and provides high service availability. It also allows for heterogeneous networks, such as the convergences of fixed and mobile networks, to coexist more naturally. It should also be noted that 4G and 5G cellular will coexist, so each architecture will be examined.

## 3 3GPP 4G Network Architecture and Software Stack

### 3.1 4G Network Architecture

To do this workload study, it is necessary to describe the architectural elements and software stack in further detail before diving into performance optimization. In this section the LTE network (RAN + ePC) architecture elements are identified. Referring to Figure 3, the access network is represented by a collection of eNBs (evolved NodeB), and the core network is represented by the ePC (Evolved Packet Core), which connects the PDN (Public Data Network) over the SGi interface. Also shown are the eNBs connected to each other via the X2 interface and to the ePC via the S1 interface.

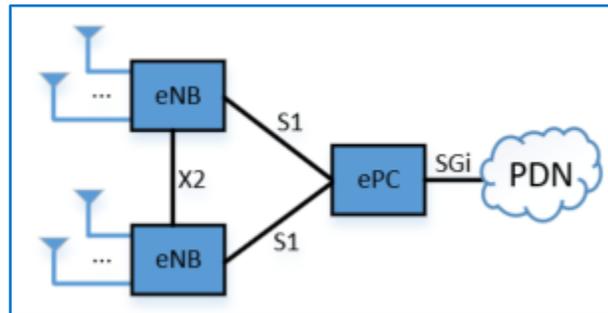


Figure 3: LTE Network Architecture high level block diagram

A more detailed network architecture is provided in Figure 4, where the eNB is split into the Remote Radio Head (RRH) and Baseband Unit (BBU). The functionality performed within the RRH is the RF spectral up and down conversion, ADC/DAC, Digital Front End (DFE), Power Amplifier Linearization techniques such as (ET, DPD), and Beamforming. The connection between the RRH and BBU is called the Front Haul (FH). The FH interface has been traditionally based on the Common Public Radio Interface (CPRI) standard<sup>6</sup> and is evolving to a more open and flexible approach through the efforts of O-RAN organization and the IEEE1914.3 standardization.<sup>7</sup>

The BBU will perform the L1 (PHY), L2 (MAC, RLC, PDCP) and L3 (RRC, etc.) functionality. The connection to the ePC is called the Back Haul (BH). The ePC contains the Mobility Management Entity (MME) which will perform Mobility Management, S-GW/P-GW selection and Tracking Area list management functions. The Serving Gateway (S-GW) will be the mobility anchor point for inter-eNB handover and Packet routing and forwarding. The Packet Data Network Gateway (P-GW) will perform UE IP address allocation, Service Level charging, Rate Enforcement, and PDN (public data network) connectivity.<sup>8</sup>

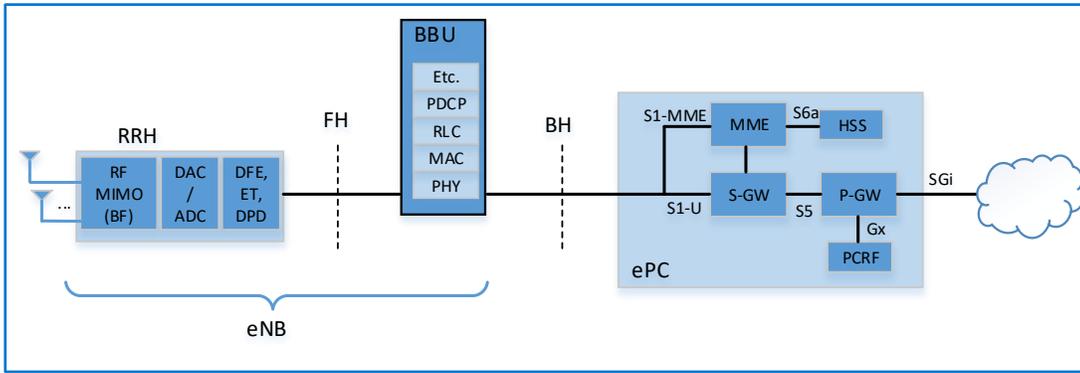


Figure 4: LTE Network Architecture block diagram

The FH bandwidth utilization has been increasing with the cell capacity. The required bandwidth, now and in the future, brings forth complexity and costly solutions. Figure 5 illustrates network element partitioning options for the RAN to help mitigate this increase in FH capacity requirements. In the variation labelled (a), the classical functionality is shown in which the bulk of the software stack is implemented by the BBU. Progressing to option (b), the PHY layer is split into upper-PHY and lower-PHY partitions, where lower-PHY is moved to the RRH, allowing the data stream to traverse the FH before it is transformed for radio consumption, alleviating some of the bandwidth required. Finally, in option (c), where the BBU is split into a Distributed Unit (DU) and a Centralized Unit (CU), the architecture has evolved to move even more of the functionality toward the edge closer to the antennas.

The benefits of a RAN split architecture are to support flexible SW & HW implementations to allow for scalable and cost-effective solutions. The choice of which functional split to select depends on deployment scenario, target services, and availability of transport network. The FH split provides a reduction in the required FH bandwidth, as well as a relaxation in the latency requirements. The functional split shown is also referred to as Option 7.2.<sup>ii</sup> A motivation for this approach is to provide the ability to separate functions that scale with user data rates from those functions that scale with RF bandwidth and number of antennas.

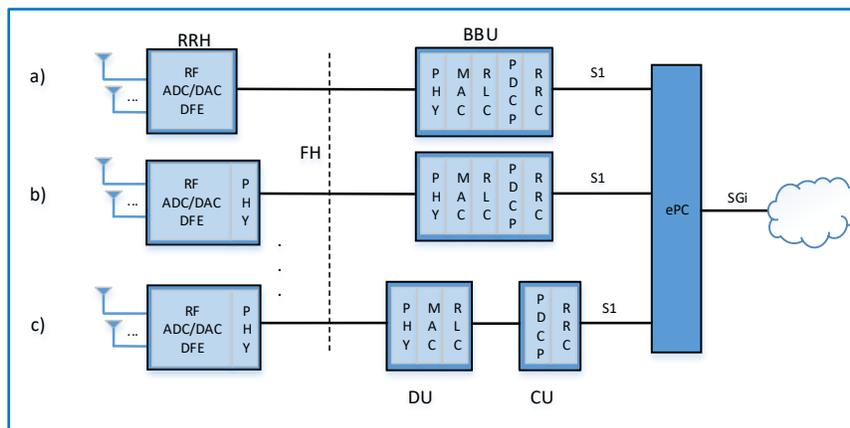


Figure 5: Front Haul split options for the LTE Network.

<sup>ii</sup> 3GPP defines many functional splits with various advantages and disadvantages. The 7.2 split entails a split between the “high PHY” and “low PHY” between the BBU and RRH.

### 3.2 4G Software Stack

To continue with developing the context for the workload optimization study, this section discusses how the software components are distributed across the 4G network elements. First, Figure 6 depicts the software components for the control plane. An example of a Mobility Management (MM) transaction flow is provided between the UE (user endpoint) and the ePC. The interface between the eNB and MME is defined by the S1-MME protocol (control plane).

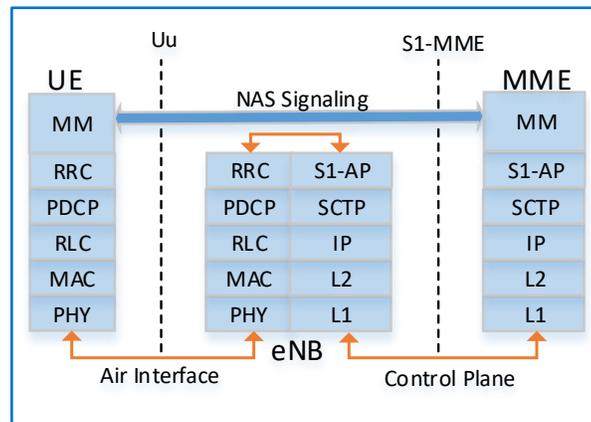


Figure 6: LTE control plane functionality.

The user plane software components are distributed across the network elements, as shown in Figure 7. An example of a data transaction flow is provided between the UE and a server located in a data center outside the Core Network (ePC). The interface between the eNB and MME is defined by the S1-U protocol (user plane).

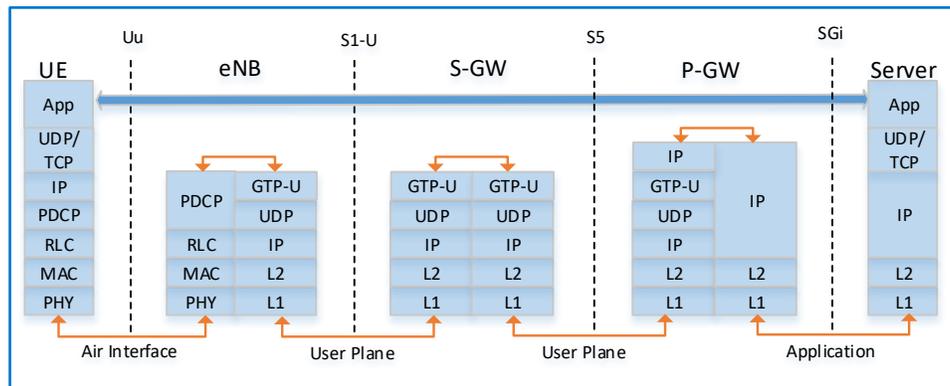


Figure 7: LTE user plane functionality.

This diagram shows a data packet exiting the eNB must traverse two entities (S-GW and P-GW) before leaving the core network. Each of the functions that operate on a packet introduce more delay (latency) that packet will encounter. This is an area that was addressed in defining the 5G core network, discussed in the sections that follow, also suggesting the benefit of computing offloads that reduce latency.

### 3.3 4G Edge Services Network Architecture

The trend with 4G has been to adopt the cloud concept by centralizing compute and storage resources deep in the network. This approach works well when applications are not demanding higher data rates and lower latency. To address the overall end-to-end latency, 3GPP Release 14 introduced the concept of Control and User Plane Separation (CUPS) to allow a more flexible and scalable network. Standardization of all the interfaces is required to guarantee success and vendor interoperability. CUPS will allow for separate evolution, in terms of features supported and achievable performance, between the Control and User Planes. A block diagram showing the network entities split between control and user planes is provided in Figure 8.

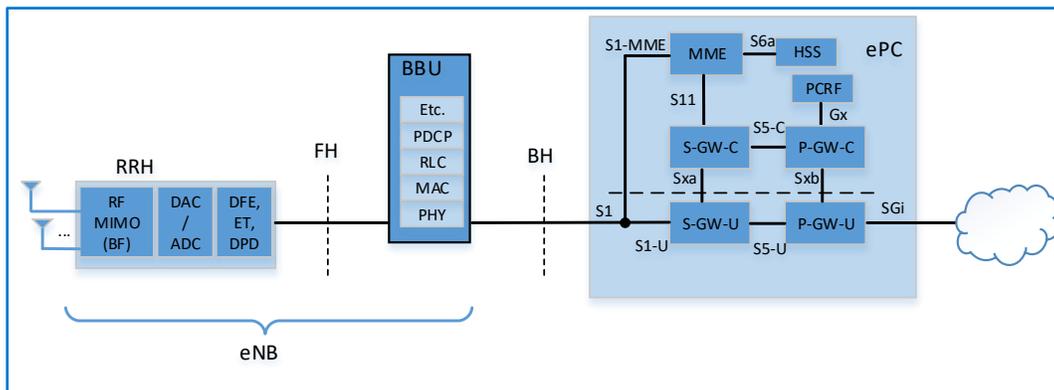


Figure 8: LTE Control and User Plane Separation.

With the control and user plane separated and the interfaces clearly defined, the network can be re-architected to address latency sensitive applications. Figure 9 separates the user plane from the ePC to allow placement of these functions closer to the edge of the network, sometimes referred to as Edge Cloud, mobile Edge Computing will occur to support services such as immersive video applications, analytics, autonomous driving, and remote control. MEC moves applications and network functions much closer to the edge of the network or closer to the users. Edge services can be collocated with the RAN as well as the Core network to provide greater scalability. Moreover, access to the PDN via the Edge Cloud is shown, offering lower latency path to critical applications. In addition to latency reduction, Edge Cloud deployments reduce the back and forth traffic through the network since the operations needed to be performed in the packet are moved closer to the source (i.e. handset, or UE). To support various network requirements, the P-GW functionality is also kept within the ePC domain, along with the PDN access and included MEC functionality. Lastly, Edge Computing can deliver an enhanced QoE to the end user simply by the improvement observed in the network reliability.

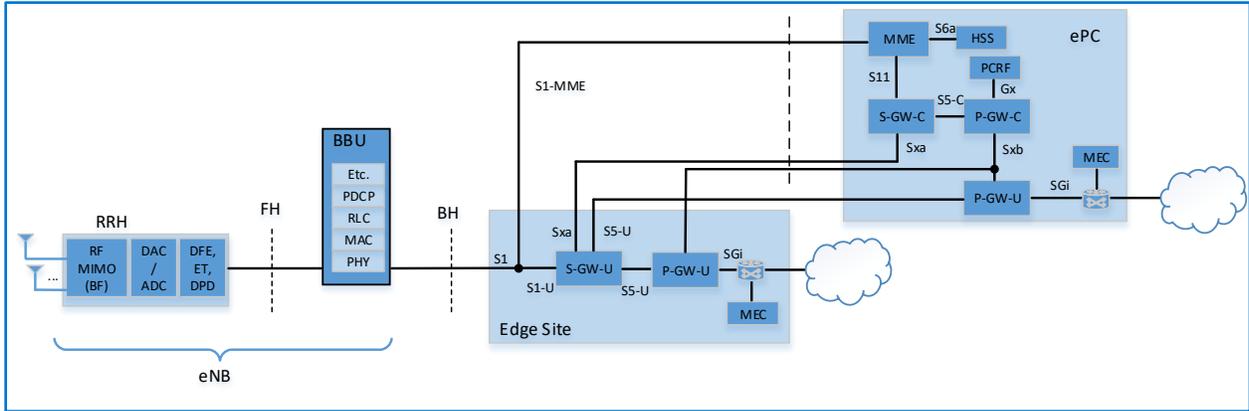


Figure 9: 4G Network Architecture Utilizing Edge Services.

The well-defined interfaces between the eNB and ePC should not be barriers to collocating Edge functionality needed to satisfy service requirements. The Edge Cloud enables a new world of applications to cellular communications such as augmented reality, virtual reality, and more. Careful selection of edge functions is required to best support targeted services; a hierarchical architecture is needed.

## 4 3GPP 5G NR (New Radio) Network Architecture and Software Stack

### 4.1 5G Network Architecture

Figure 10 illustrates some of the tenets previously shown that will become an inherent part of the 5G cellular architecture. The gNB is capable of being separated into the Distributed Unit (DU) and Centralized Unit (CU), whose interconnection is defined by F1 (by the 3GPP standards organization). When compared to Figure 3, the following changes in nomenclature are apparent: eNB --> gNB, X2 --> Xn, S1 --> N2/N3, ePC --> 5G CN and SGI --> N6. Lastly, the F1 interface was introduced.

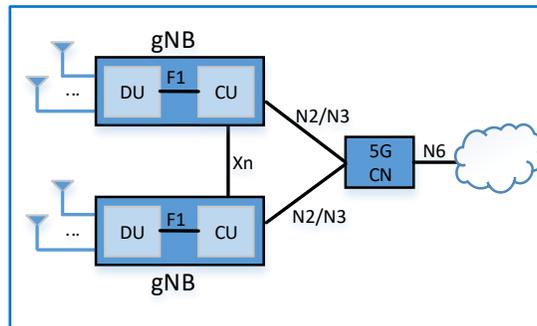


Figure 10: 5G NR Network Architecture high level block diagram

A more detailed network architecture block diagram is provided in Figure 11, where the gNB is split into the RRH and BBU. The functionality performed within the RRH is the RF spectral up and down conversion, ADC/DAC, Digital Front End (DFE), Power Amplifier Linearization

techniques such as (ET, DPD), and Beamforming. The connection between the RRH and BBU is again called the FH. The FH interface is evolving from CPRI to eCPRI, IEEE1914.3, as well as O-RAN.

The connection to the 5G CN is called the BH. The 5G CN will contain the Access and Mobility Function (AMF) which will perform Mobility Management, Access Authentication/Authorization, Terminate N2 interface, and Connection Management. The Session Management Function (SMF) will perform UE IP address allocation & Management and Traffic steering to UPF. The User Plane Function (UPF) will perform the PDN connection, Packet routing and forwarding, QoS handling of User Plane (UL/DL rate enforcement), and Anchor point for Intra-Inter-RAT mobility. The United Data Management (UDM) will perform User ID handling, and Subscription Management.<sup>9</sup> The 5G CN was designed based on the Service Based Architecture (SBA) principles. This direction was chosen to support cloud native based implementations and to allow for easy integration of MEC applications.

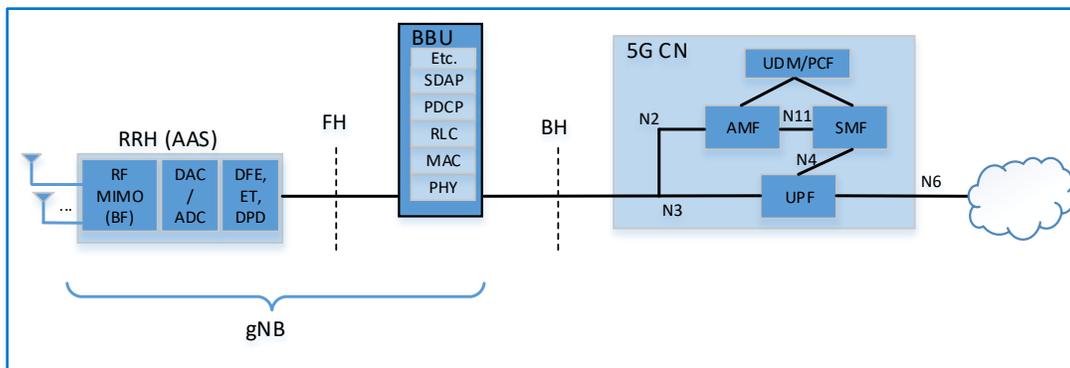


Figure 11: 5G NR Network Architecture block diagram

An example FH functional split is provided in Figure 12. Note again the DU/CU interface is now defined by 3GPP and called the F1 interface. Depending on the deployment scenario, a FH Gateway can be inserted between the RRH and DU functions. The purpose of this FH Gateway would be to aggregate traffic, perform any protocol conversions, distribute timing, support FH split options, and more. A motivation of the particular DU/CU functional split shown is security. The PDCP layer provides Integrity and Confidentiality protection to both the user and control planes.

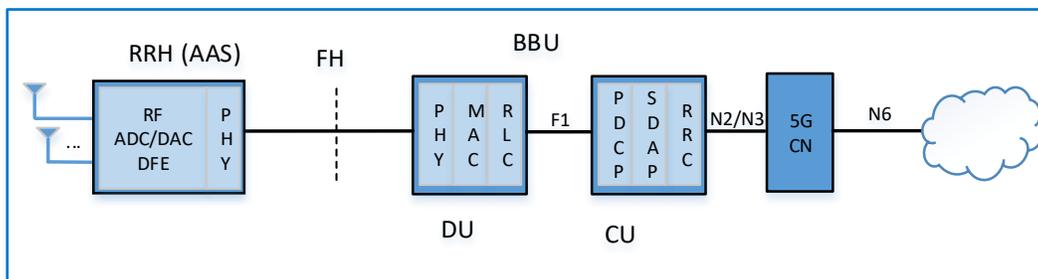


Figure 12: Front Haul split options for the 5G NR Network.

Figure 13 provides a more detailed diagram of the physical layer RAN FH functional splits. This example corresponds to the Option 7.2 provided by O-RAN front haul specification. The left side of the dashed line (denoted as FH) represents the DU/CU functionality performed by the

physical layer (PHY). The right side is the functionality represents the RRH functionality. The O-RAN Alliance front haul specification defines User Plane, Control Plane, Synchronization Plane, and Management Plane multiplexing to allow for flexible and scalable access over the FH.

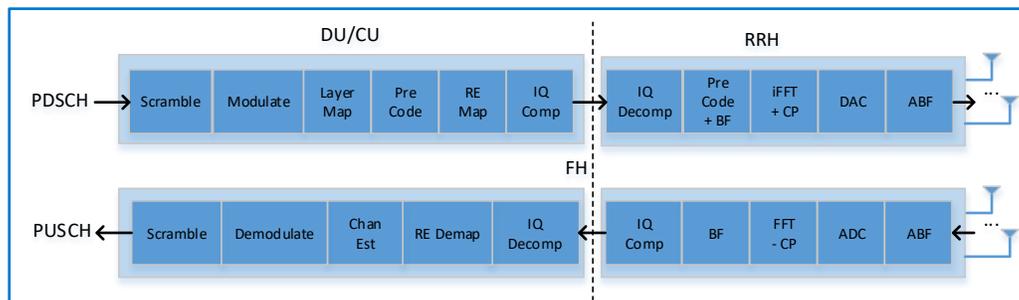


Figure 13: RAN FH Functional Split (Option 7.2).

## 4.2 Example FH & BH BW Requirements for 4G & 5G

In this section, two cell configuration examples (4G and 5G) are provided to demonstrate bandwidth requirements, which will impact Ethernet device and port selection. The O-RAN functional split (Option 7-2) is assumed in this estimation. Assumed in the calculations are a 30% packet overhead for control messages, a 60% compression ratio (when enabled), a 80%/20% TDD Ratio for DL/UL allocation, and a total of 3 cells to be supported. Downlink uses a greater bandwidth than uplink, hence it determines the front haul and back haul overall requirement.

### An LTE deployment scenario:

- ⇒ RF Transmitted BW = 20MHz (FFT size = 2048, using 1200 Subcarrier)
- ⇒ DL Configuration = 4x4 (4 Layers) @ 256QAM and UL Configuration = 2x2 (2 Layers) @ 64QAM
- ⇒ 3x Carrier Aggregation (CA): Total BW = 20MHz \* 3 = 60MHz
- ⇒ FDD duplex
- ⇒ fs (sampling frequency) = 30.72MHz per carrier (SCS = 15KHz)
- ⇒ DL Bit Rate per Cell = 400Mbps \* 3 = 1.2Gbps and UL Bit Rate per Cell = 150Mbps \* 3 = 450Mbps
- ❖ The required per Cell FH BW = 8.4Gbps<sup>iii</sup> (w/o IQ Compression) or 5Gbps (w/ IQ Compression enabled). The required BH BW = 4.7Gbps (aggregated over 3 Cells).

### A 5G deployment scenario:

- ⇒ RF Transmitted BW = 100MHz (FFT size = 4096, using 3276 subcarrier)
- ⇒ DL Configuration = 4x4 (4 Layers) @ 256QAM and UL Configuration = 2x2 (2 Layers) @ 64QAM
- ⇒ TDD duplex
- ⇒ fs (sampling frequency) = 122.88MHz per carrier (SCS = 30KHz)
- ⇒ DL Bit Rate per Cell = 2.14Gbps and UL Bit Rate per Cell = 800Mbps

<sup>iii</sup> Cell FH BW = (# of CA) x (16bit I or Q) x (Compression Ratio) x (2 IQ Complex) x (4 Antennas) x (# of Subcarriers) x (14 symbols/TTI) / (1x10<sup>6</sup>)

- ❖ The required per Cell FH BW = 15.3Gbps (w/o IQ compression or 9.2Gbps w/ IQ compression enabled). The required BH BW = 6.7Gbps (aggregated over 3 Cells).

These configurations and results are illustrated in Figure 14: a) for LTE and b) for the 5G example. The respective FH rates have increased proportional to the user data rates. In order to design a network to support 4G and 5G as well as have room for growth, 25GbE ports can be used for the interfaces in the scenarios shown below, and devices of greater bandwidth can be selected when more capacity is required.

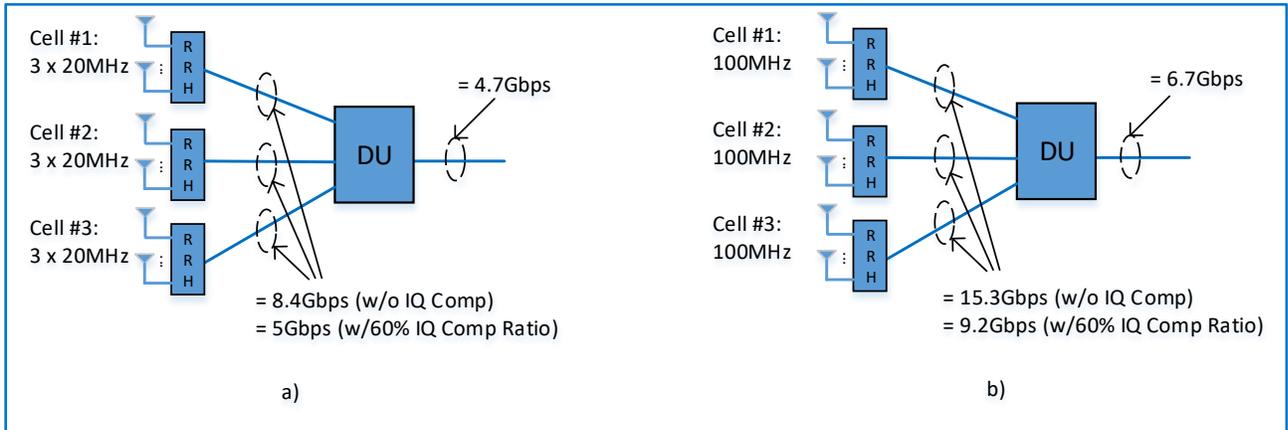


Figure 14: Example 4G (a) & 5G (b) FH and BH BW Requirements

### 4.3 5G Software Stack

Figure 15 shows the SW functionality partitioning between the DU and CU. The DU & CU interface is defined by the F1 interface, which consists of both control (F1-C) and user (F1-U) plane signaling. The CU & 5G CN interfaces are defined by the N2 and N3 interfaces.

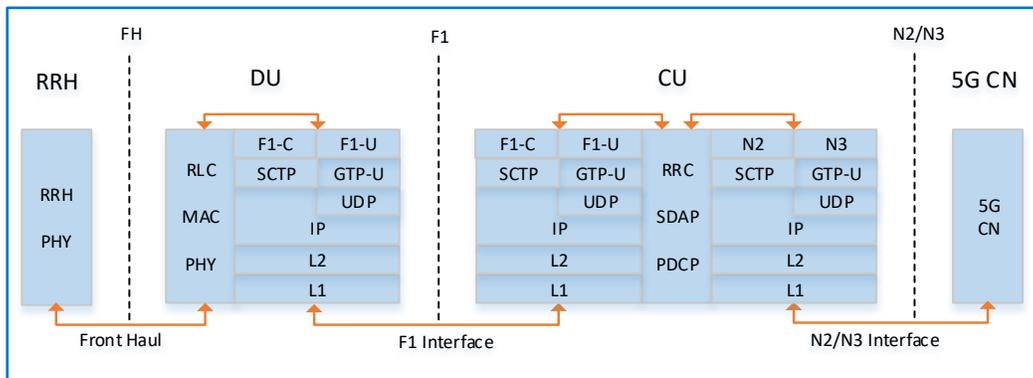


Figure 15: Example DU and CU functionality

The distribution of control plane software components across network elements is illustrated in Figure 16, which depicts a mobility Management (MM) transaction flow. Compared to Figure 6, the 5G stack introduces the Service Data Application Protocol (SDAP) layer. The SDAP is responsible for handling QoS (quality of service) flows, essentially mapping a QoS flow to a radio bearer. It also marks the packets with QoS Flow IDs (QFI) to allow for appropriate treatment in the 5G network.

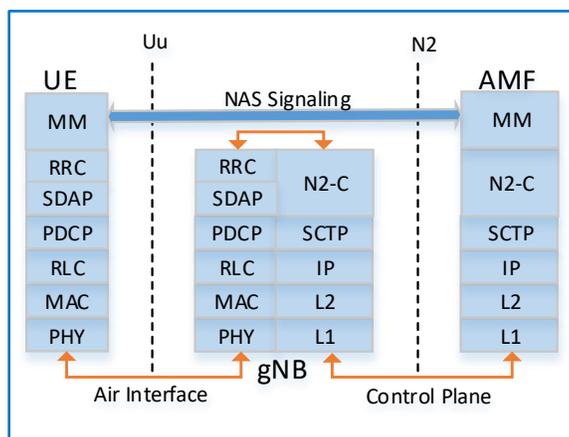


Figure 16: 5G NR control plane functionality

The user plane software components are distributed across the network elements as shown in Figure 17. An example of a data transaction flow is provided. To reduce latency contribution from the network, the two elements in the ePC (namely S-GW and P-GW) have been reduced to one element in the 5G CN (called UPF). This reduces the amount of time the packet needs to traverse the stack and thus the overall end-to-end system latency.

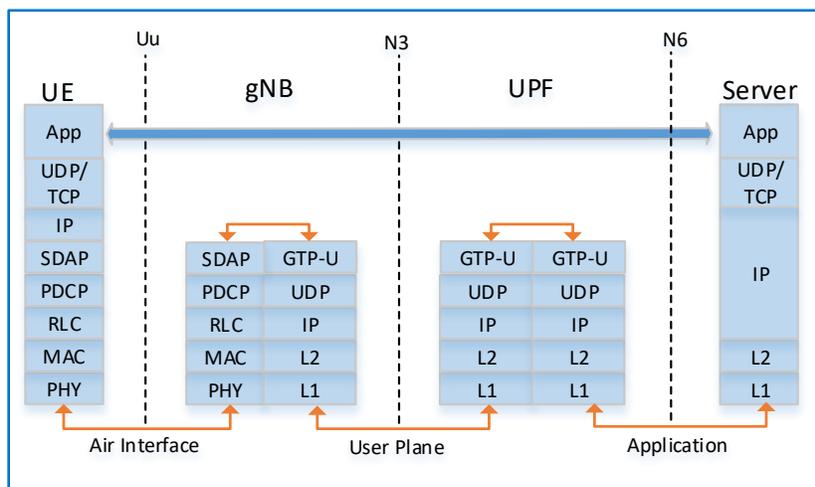


Figure 17: 5G NR user plane functionality.

#### 4.4 5G Edge Service Network Architecture

Since the 5G CN has partitioned the control plane and user plane functionality, the network can be re-architected to address latency sensitive applications. Much like Figure 9, Figure 18 separates the user plane from the 5G CN to allow placement of these functions closer to the edge of the network, also referred to as Edge Cloud. As previously mentioned, and as described in an ETSI whitepaper,<sup>10</sup> this architecture supports emerging mobile edge computing applications.

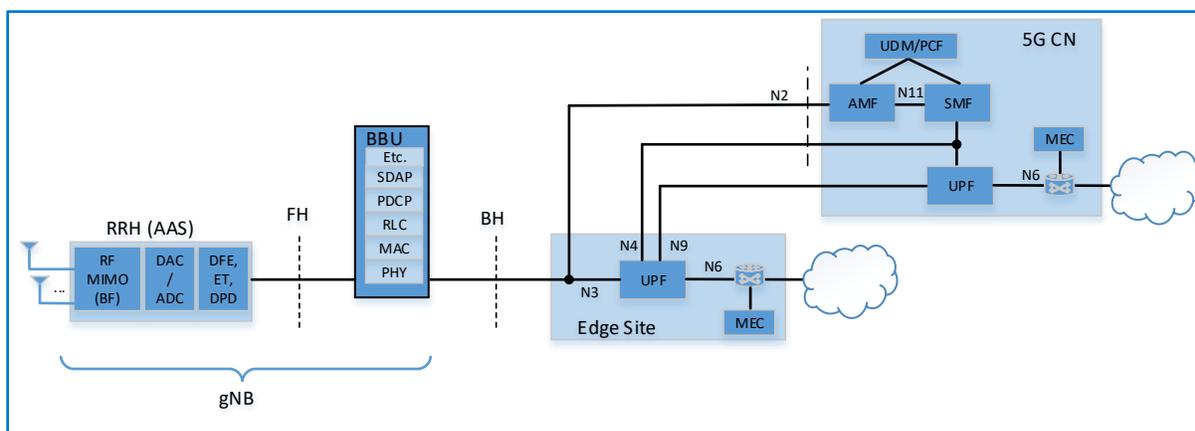


Figure 18: 5G Network Architecture Using Edge Services.

Edge computing benefits users when the time delay of transporting packets is significant. When compute functionality is placed closer to the edge, there is significant reduction in latency. A point worth noting is as the transport latency is reduced, there will be more attention on the remaining compute latency since it may provide the lower limit on the achievable gains. In addition to improving latency the Edge Cloud user encounters, a reduction in network traffic will be achieved, which will also reduce the opportunities for dropped packets and will improve the overall customer QoE.

#### 4.5 5G Network Slicing Example

With the flexibility and scalability of this network architecture, along with the application of SDN (software defined networking) principles, the concept of Network Slicing can be introduced. A slice is defined as a logical arrangement of network functions needed to address a specific service. Figure 19 shows how three distinct services (mMTC, uR-LLC and eMBB) can be accommodated by utilizing network functions in different sections of the network. A software-based network allows for the reallocation and collocation of functions. The mMTC devices (Smart Cites, Health Monitoring, etc.) are not expected to require low latency services, so their transactions can be serviced by either the CN or application located in the PDN. The eMBB devices (AR/VR, Enhanced Video, etc.) are expected to require lower latency than the present 4G networks will allow, so it is reasonable to have their transactions addressed within the Edge Cloud. Lastly, the uR-LLC devices (Autonomous vehicles, industrial automation, etc.) are expected to require very low latency services, so their transactions should be addressed by the CU and Edge Cloud. These three services are highlighted in Figure 19 by their respective dashed lines. Also shown are the relevant network functions at the specific locations. The previous sections have provided the background information needed to support the overall 5G network flexibility.

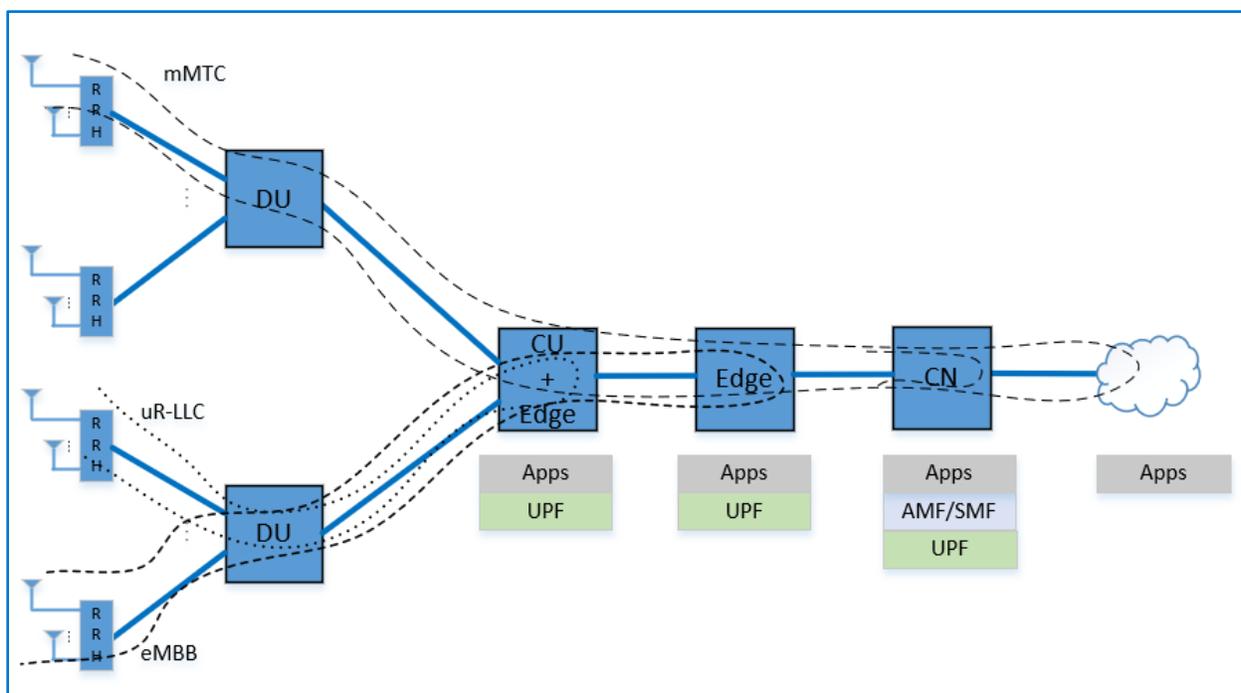


Figure 19: 5G Network Highlighting Network Slicing Example.

## 4.6 Software Offload Functionality

In viewing the functionality required in the RAN (CU & DU), various functions lend themselves to hardware acceleration. The best candidates are in the PHY layer and include forward error correction (FEC) and I/Q sample compression.

The FEC functionality can be accelerated to free up compute resources for other functions. These FEC functions are Turbo Codes (both encode & decode) for LTE<sup>11</sup> and LDPC Codes (both encode & decode) for 5G.<sup>12</sup> These functions are defined by the 3GPP standard, and when accelerated in hardware, free up compute resources (as discussed in the next section) to provide product differentiation opportunities and services. Accelerating the FEC functionality will remove a significant hurdle to deploy C-RAN.

The proposed 4G & 5G FEC look-aside acceleration is depicted in Figure 20. The associated HW/SW functionality partitioning is shown for the encoder direction. As shown, functionality located above the line will be performed in SW, while functionality within the dashed box will be accelerated in HW. Intel's FPGA devices provide an effective way to support acceleration on PCIe cards (as well as solder down applications). The FEC acceleration device will utilize the SR-IOV feature of PCIe to allow support of up to 16 VFs. The proposed 4G & 5G FEC acceleration solution will perform the following operations:

- ⇒ FEC Encoding
  - Code Block CRC Generation, Turbo Encoding (4G), LDPC Encoding (5G) and Rate Matching (Sub-block interleaving, Bit Selection/Collection).
- ⇒ FEC Decoding
  - HARQ Combining, Rate De-Matching (Sub-block de-interleaving, Turbo Decoder (4G), LDPC decoder (5G) and Code Block CRC check.

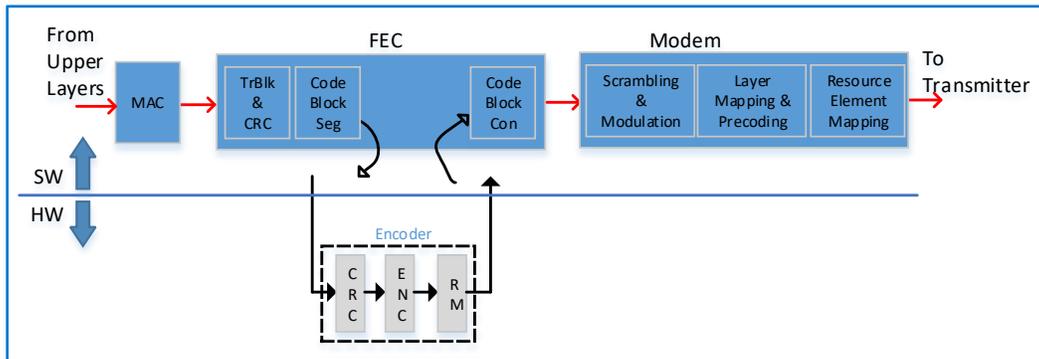


Figure 20: 4G & 5G FEC Acceleration - Encoding Example.

## 5 FlexRAN test setup

Intel's FlexRAN is a software reference solution that implements Layer 1 of the BBU function of an LTE network (and includes a stack for 5G gNB as well). At Dell, an end-to-end test platform was created using FlexRAN as the basis for the Layer 1 software of a radio equipment controller baseband unit. The system running FlexRAN was a PowerEdge R740 equipped with a predecessor version of the N3000 network card, an Intel X520 Ethernet network adapter, two Intel® Xeon® Gold 6148 CPUs at 2.4GHz (microcode version 0200004D), and 256MB of RAM. The server was running CentOS (7.5.1804) with a real time patch to the kernel (3.10.0-rt56), and FlexRAN version 18.12, which required DPDK version 18.08. For the FPGA offload, patches released in versions 19.03 were introduced along with DPDK patches for BBDev.

### 5.1 Components and System diagram

Figure 21 depicts the setup staged to test and characterize C-RAN software. Included in the figure is a software pipeline for layer 1, showing the modules that were identified as beneficial for an FPGA offload. Tests were run before and after the offload was introduced. The Turbo Encoder and Decoder are computationally intensive tasks because of the bit-level interleaving.

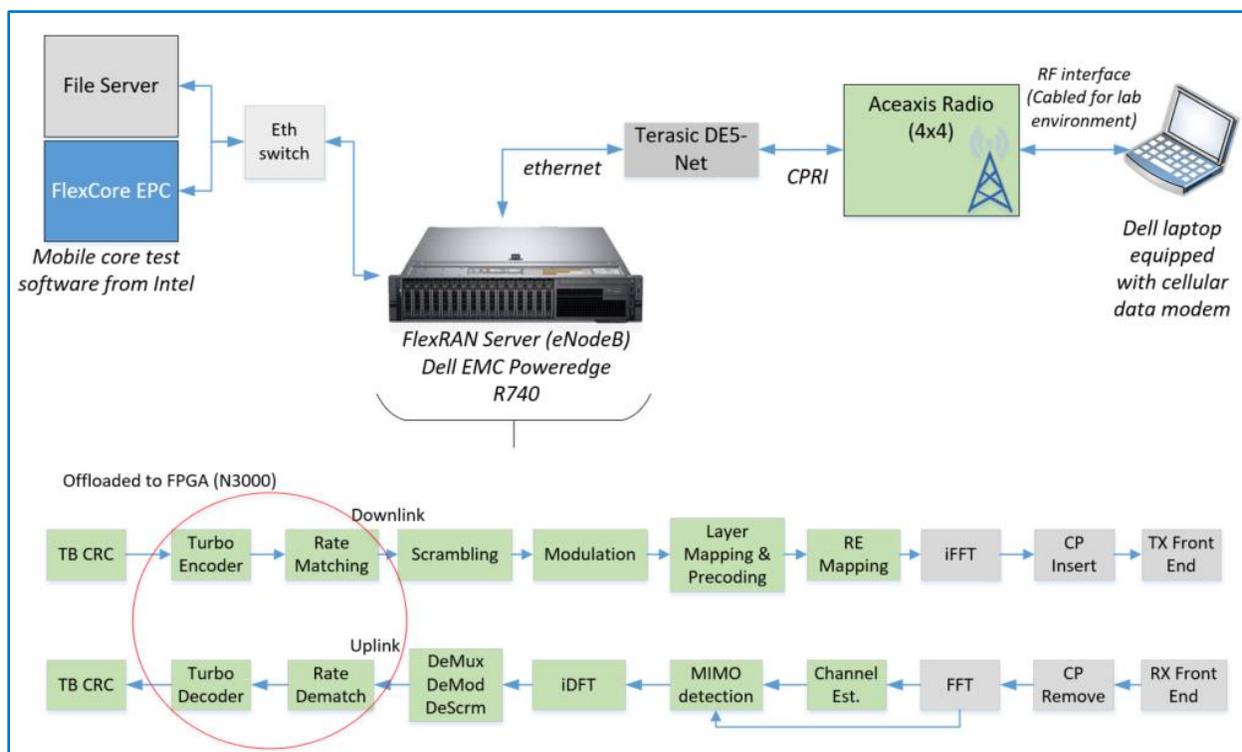


Figure 21: Test setup in Dell lab, with Layer 1 software modules

## 5.2 Results

The benefit of offloading the Turbo operations to an FPGA was explored. The device chosen for the offload was a prototype version of the N3000 network card<sup>13</sup> because it allowed a unique combination of “inline” and “lookaside” operations. Turbo FEC, because of its position in the stack, is not suited for an inline offload. The N3000 was not being used as a network card in this case, but it became part of the test setup because of the vision for inline offloads in future 5G radio test setups. It should be noted that this offload would be equally achievable by a card such as the programmable acceleration card that does not include a NIC ASIC,<sup>14</sup> currently offered by Dell as a standard orderable option.

An instrumented software test mode was used to measure CPU cycles consumed when Turbo FEC was performed by the CPU versus by the FPGA peripheral. These results are illustrated in Figure 22. The actual benefit depends on the radio and cell site configuration; it will vary with the number of cells, transmission buffer sizes, and other parameters.<sup>iv</sup> The graph represents the total processing time for the Layer 1 software as a normalized “1.0”, one component of which is the Turbo processing for the uplink and downlink simulated streams. These cases were performed in duplex mode with concurrent uplink and downlink. A full view of this contrast can

<sup>iv</sup> For example, the number of iterations on the uplink processing of the Turbo FEC, affected by bit errors in the stream, will have an impact on the required processing. These offload benefits increase further with errors in the uplink stream. Taking the last cell of table 1 as an example, there is an 86.5% improvement in cycles required for the operation. With multiple iterations, this offload benefit increases.

be seen in Table 1, which shows normalized cycles attributable to the uplink decode and the downlink encode.

As with any real-time computing task, another important timing consideration is the deviation in processing time. In datacenter applications, this has been called the “long tail” in latency. This test setup was not equipped to measure the standard deviation, however the maximum duration was captured for the thread that performs the interleaving portion of the Turbo algorithm. (See reference<sup>15</sup> for a description brief description of the algorithm and the IP.) This was measured over a million iterations (time transmission intervals, which is 1ms for LTE) in the test setup of Figure 22. The variation in processing time between the average and maximum thread duration in the offloaded case was only 15% that of the software-only case.

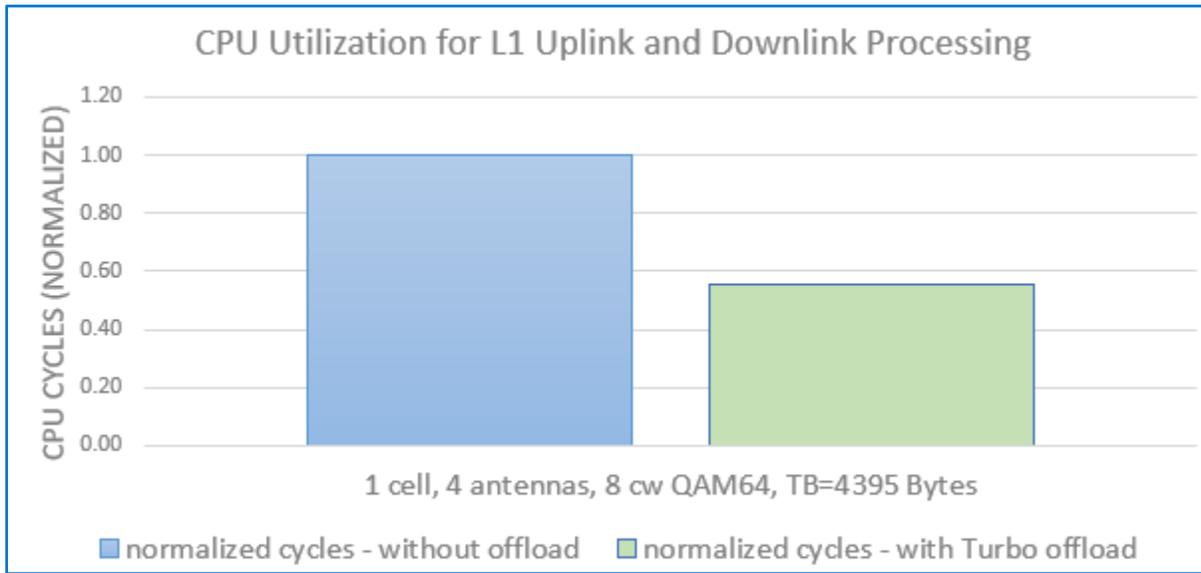


Figure 22: Turbo FEC execution times (results with and without FPGA offload)

test case description	normalized cycles - without offload			normalized cycles - with Turbo offload			% hardware assisted cycles vs. software-only	
	total L1 UL+DL	DL Turbo FEC	UL Turbo FEC	total L1 UL+DL	DL Turbo FEC	UL Turbo FEC	DL Turbo FEC	UL Turbo FEC
1 cell, 4 antennas, 8 CW QAM64, TB=4395 Bytes	1.00	0.27	0.25	0.55	0.03	0.03	10.2%	13.5%

Table 1: CPU cycle comparison on Turbo operations with and without FPGA offload

## 6 Platform Architectures

### 6.1 Flexible Architectures for the Network Edge

Using the 5G NR network architecture split, Figure 23 can be used to conceptualize server platforms that utilize FPGA peripherals. In the DU area, the FPGA will be used, as demonstrated, for look aside acceleration such as 4G/5G FEC. In the CU & 5G CN areas, although not the primary subject of this paper, the FPGA can be used to accelerate machine learning operations, video transcoding, and other tasks in support of Multi-access Edge Computing. The FPGA based solution can be paired with any IA (Intel® architecture) based CPUs, such as the Intel® Xeon® Scalable processor product line available in PowerEdge Servers. Intel offers FPGA and eASIC<sup>16</sup> based solutions to accommodate the HW acceleration needs. FPGA provides great flexibility through its ability to be reprogrammed, while the eASIC based solution is a fixed function device aiming to bridge the performance gap between an ASIC and FPGA implementations.

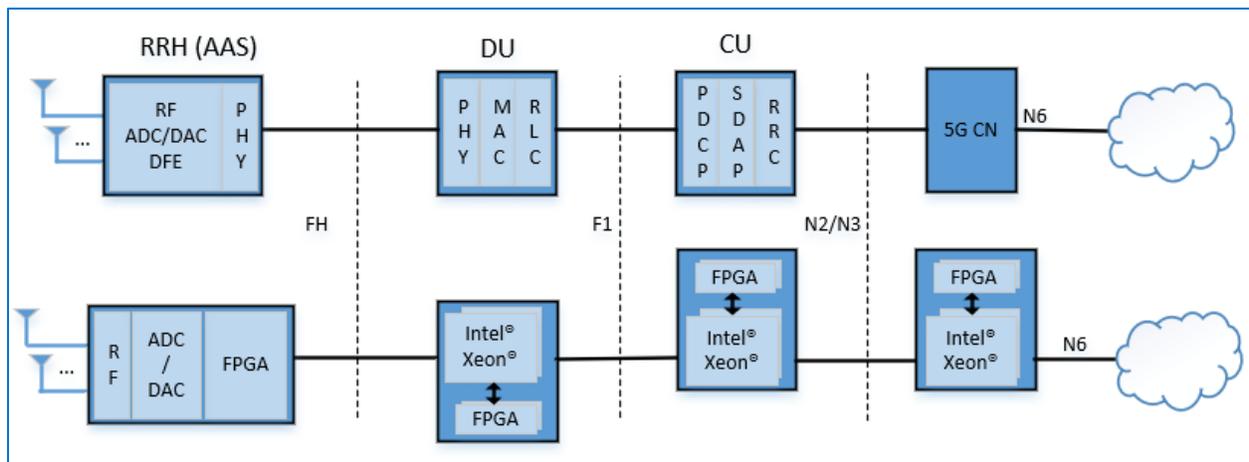


Figure 23: Hardware Accelerator Offload Possibilities

Another emerging application in this arena is indoor wireless. As the trend for more rich and numerous mobile applications expands, the buildout of more (and smaller) cells is inevitable. To implement the controllers for the increasing amount of radio equipment, standard datacenter racks can be deployed indoors. A small computing cluster can implement all the functions described, and where applicable, FPGA peripherals can be used for direct radio interfaces, including legacy interfaces that use CPRI.

### 6.2 Platform Example

More detailed block diagrams of the CU/DU/Edge solutions are provided in this section. An Intel® Xeon® Scalable processor-based solution utilizing up to 8 DDR channels of host memory is shown in Figure 24. A FPGA is placed on a PCIe card to provide various acceleration functionality depending on its location in the network. The FPGA based PCIe card is connected via 16 PCIe lanes for acceleration. Ethernet capability is provided by the NIC PCIe cards.

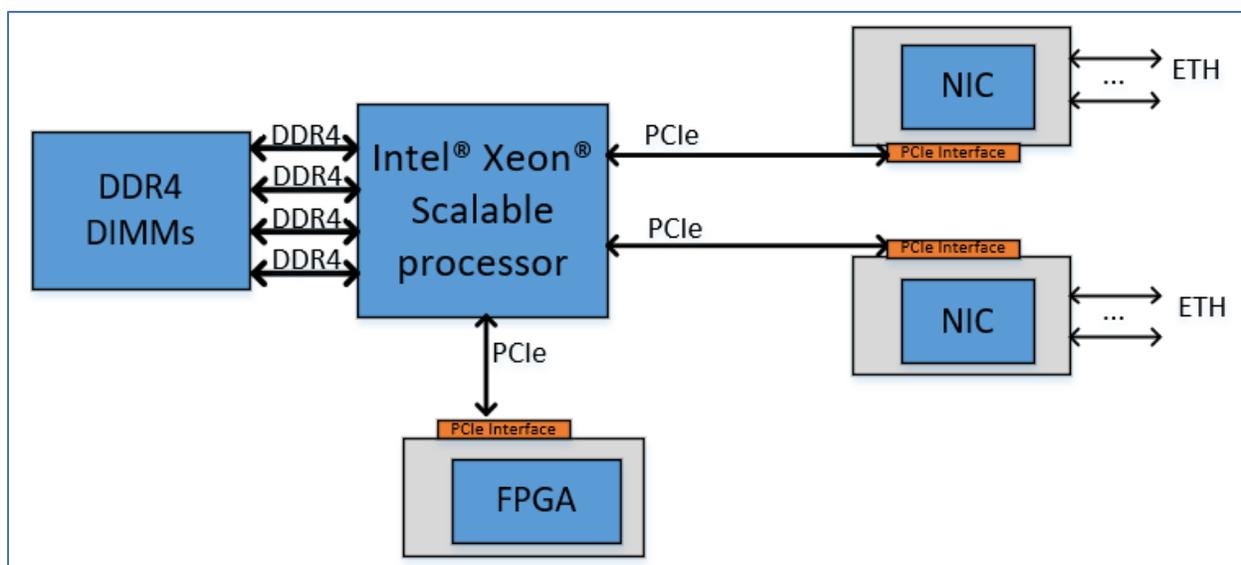


Figure 24: Intel® Xeon® Scalable processor-based solution

This solution is available with a single mother board design where the PCIe card functions can be reconfigured depending on this platform location and function in the network. For use in DU, CU and Edge Cloud deployments, the following examples are given. Depending on space availability, various platform configuration optimization can be used.

- ⇒ **For the DU case, a platform is comprised of:**
  - A: Ethernet connections to interface to the FH direction
  - B: Ethernet connections to interface to the F1 direction
  - C: PCIe card to perform the FEC acceleration functions described above
- ⇒ **For the CU case, a platform is comprised of:**
  - A: Ethernet connections to interface to the F1 direction
  - B: Ethernet connections to interface to the N2/N3 direction
  - C: PCIe card to perform application layer acceleration, if applicable
- ⇒ **For the Edge Cloud case, a platform is comprised of:**
  - A: Ethernet connections to interface to the N2/N3 direction
  - B: Ethernet connections to interface to the N9/N6/N4 direction
  - C: PCIe card to perform application acceleration, if applicable
- ⇒ **For a combined DU & CU case, a platform can be comprised of:**
  - A: Top Ethernet connections to interface to the FH direction
  - B: Bottom Ethernet connections to interface to the N2/N3 direction
  - C: PCIe card to perform FEC acceleration and application layer offloads
- ⇒ **Lastly for a combined DU, CU, and Edge Cloud case, a platform can be comprised of:**
  - A: Top Ethernet connections to interface to the FH direction
  - B: Bottom Ethernet connections to interface to the N2/N4/N9/N6 direction
  - C: PCIe card to perform FEC acceleration and application layer offloads such as AI acceleration, Video Coding, Analytics, etc.

## 6.3 Products

### 6.3.1 Rack Servers

Mainstream servers in the PowerEdge portfolio from Dell EMC can be used to implement radio access networks, and as mentioned as one of the key tenets of NFV workloads, the economy of scale that is achievable by using off-the-shelf equipment is attractive.

The R640 has been validated by the Dell OEM team for NEBS Level-3 compliance.<sup>17</sup> It is a 1U two-socket rack server, offering a high degree of computing density. The carrier grade server also offers extended temperature operation with special thermal tables and fan control algorithms. The R740 has also been certified, and in a 2U space, offers a high degree of expandability for networking bandwidth and acceleration. The XR2 platform is a smaller 1U ruggedized server and can be used in locations where additional environmental constraints are present.<sup>18</sup>

### 6.3.2 Modular Data Centers

Server products alone are not enough to leverage in creating network edge installations. The ESI (Extreme Scale Infrastructure)<sup>19</sup> team at Dell EMC has experience in installing datacenters in a variety of remote locations; this includes micro modular data centers well-suited for deployment in locations where preexisting server facilities may not exist.<sup>20</sup> One of the compelling attributes of these MDCs, featuring built-in cooling and security, is that they can house standard servers in these undeveloped locations.

## 7 Conclusions

Dell EMC has capabilities and expertise across multiple teams to deliver the next generation of computing infrastructure to support the world's existing and emerging wireless networks. From performance optimization expertise across today's Enterprise workloads to research and development in the emerging field of accelerators, Dell EMC has a wide range of deployable solutions that will make Edge Computing a reality. Together with Intel's partnership and their deep wireless expertise, Dell EMC is equipped to be a leading supplier to Communications Service Providers. In this paper, the benefit of wireless protocol offload onto FPGA peripherals has been demonstrated, and the path and relevance to 5G cellular has been shown. All of this can be accomplished with Dell EMC equipment available today and in the future.

## 8 Acknowledgements

This paper would not have been possible without the capable efforts of Imran Masud (Dell EMC) and Sam Diep (Intel).

## 9 References

- <sup>1</sup> [https://portal.etsi.org/nfv/nfv\\_white\\_paper.pdf](https://portal.etsi.org/nfv/nfv_white_paper.pdf)
- <sup>2</sup> [https://doc.dpdk.org/guides/prog\\_guide/bbdev.html](https://doc.dpdk.org/guides/prog_guide/bbdev.html)
- <sup>3</sup> [www.3gpp.org](http://www.3gpp.org)
- <sup>4</sup> China Mobile paper  
<https://pdfs.semanticscholar.org/ea3/ca62c9d5653e4f2318aed9ddb8992a505d3c.pdf>
- <sup>5</sup> <https://www.o-ran.org/>
- <sup>6</sup> <http://www.cpri.info/index.html>
- <sup>7</sup> [https://standards.ieee.org/standard/1914\\_3-2018.html](https://standards.ieee.org/standard/1914_3-2018.html)
- <sup>8</sup> 3GPP TS29.244 (v15.5.0), Interface between the Control Plane and User Plane Nodes, [www.3gpp.org](http://www.3gpp.org)
- <sup>9</sup> 3GPP TS23.501 (v16.0.2), System Architecture for the 5G System, [www.3gpp.org](http://www.3gpp.org)
- <sup>10</sup> ETSI MEC paper [https://portal.etsi.org/Portals/0/TBpages/MEC/Docs/Mobile-edge\\_Computing\\_-\\_Introductory\\_Technical\\_White\\_Paper\\_V1%2018-09-14.pdf](https://portal.etsi.org/Portals/0/TBpages/MEC/Docs/Mobile-edge_Computing_-_Introductory_Technical_White_Paper_V1%2018-09-14.pdf)
- <sup>11</sup> 3GPP TS36.212 (15.5.0), Evolved Universal Terrestrial Radio Access (E-UTRA); Multiplexing and Channel Coding, [www.3gpp.org](http://www.3gpp.org).
- <sup>12</sup> 3gpp TS38.212 (15.5.0), New Radio (NR); Multiplexing and Channel Coding, [www.3gpp.org](http://www.3gpp.org).
- <sup>13</sup> N3000 Network Card  
<https://www.intel.com/content/www/us/en/wireline/products/programmable/applications/nfv.html>
- <sup>14</sup> Programmable Acceleration Card – Arria10  
[https://www.intel.com/content/www/us/en/programmable/products/boards\\_and\\_kits/dev-kits/altera/acceleration-card-arria-10-gx/overview.html](https://www.intel.com/content/www/us/en/programmable/products/boards_and_kits/dev-kits/altera/acceleration-card-arria-10-gx/overview.html)
- <sup>15</sup> Turbo FEC Encode/Decode  
<https://www.intel.com/content/dam/www/programmable/us/en/pdfs/literature/an/an505.pdf>
- <sup>16</sup> <https://www.easic.com>
- <sup>17</sup> Dell OEM – Powerededge solutions <https://www.dellemc.com/en-us/oem/oem-powerededge.htm#scroll=off>
- <sup>18</sup> XR2 tech specs [https://topics-cdn.dell.com/pdf/oth-r440-xr\\_reference-guide2\\_en-us.pdf](https://topics-cdn.dell.com/pdf/oth-r440-xr_reference-guide2_en-us.pdf)
- <sup>19</sup> ESI <https://www.dellemc.com/en-us/solutions/extreme-scale-infrastructure.htm#scroll=off>
- <sup>20</sup> Edge blog (ESI)- micro MDCs <https://blog.dellemc.com/en-us/micro-modular-data-centers-taking-computing-to-edge/>

### Disclaimer

Testing was conducted by Dell as of 6/11/19. See the Test Setup section for configurations. 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 6/11/19 and may not reflect all publicly available security updates. No product or component can be absolutely secure. Check with your system manufacturer or retailer to learn more at intel.com. Intel, the Intel logo, and Xeon are trademarks of Intel Corporation or its subsidiaries in the U.S. and/or other countries.