



user density. In Third-Generation Partnership Project (3GPP) long-term evolution (LTE), small cells draw significant attention on both the physical and higher layer [3], [4], where impacts on the RAN protocol and system architecture are discussed.

As networks become denser, intercell interference increases and interference scenarios become more complex due to multiterrier interference. Furthermore, the higher the deployment density, the higher the chance that a RAP will carry no or only low traffic-load due to spatial and temporal traffic fluctuations. Currently, 15–20% of all sites carry about 50% of the total traffic [5]. Centralized processing permits to selectively turn RAPs on and off to address the spatiotemporal traffic fluctuations. In addition, it allows for efficient interference avoidance and cancellation algorithms across multiple cells as well as joint detection algorithms. Centralized RAN (C-RAN) recently attracted attention as one possible way to efficiently centralize RAN processing [6]. In C-RAN, remote radio heads (RRHs) are connected through optical fiber links to a data center where all baseband processing is performed [7], [8]. Thus, by pooling baseband processing in baseband units (BBUs), centralization gains are achieved. However, BBUs are based on specialized hardware platforms utilizing digital signal processors (DSPs) [9]. As a long-term goal, it is beneficial to deploy cloud-computing platforms running on general-purpose hardware, leading to a cloud-RAN system as outlined subsequently in this article.

Only fiber links are capable of supporting the necessary data rates between the RRH and the BBU. This constitutes the main drawback of C-RAN, i.e., it requires very high data rate links to the central BBU. In [8], the authors report a required backhaul (BH) transmission rate of 10 Gbit/s for time-domain LTE (TD-LTE) with eight receive antennas and 20-MHz bandwidth. Due to the use of optical fiber, C-RAN deployments are less flexible as only spots with existing fiber access may be chosen or fiber access must be deployed, which is very cost-intensive. Future mobile networks will deploy heterogeneous BH solutions that are optimized for different scenarios. This mix of BH characteristics will also imply a mix of more C-RAN solutions that require high-capacity BH and more decentralized solutions compatible with BH solutions that introduce high latency and stronger throughput constraints [10].

The RAN as a Service (RANaaS) concept is introduced in [11]. It addresses the deficiencies of C-RAN to allow for a centralization over heterogeneous BH. The main characteristics of RANaaS are the flexible assignment of RAN functionality between the RAPs and the central processor, the deployment of commodity hardware at the central processor, and the tight integration of RAN, BH network, and central processor. In this article, we focus on the challenges and benefits of implementing signal processing algorithms on a cloud-computing platform. Hence, in the following, we refer to the concept of centralization toward commodity cloud-computing platforms as cloud-RAN. More details on the architecture design of the underlying 5G mobile network as well as fundamental concepts from medium access control (MAC) and network layer of the cloud-RAN concept are given in [11]. Further challenges in 5G mobile networks, which are beyond the scope of this article, are introduced in [2] and [12], among others. However,

cloud-RAN will foster approaches currently under discussion for 5G such as massive multiple-input, multiple-output (MIMO) and multiple radio access technologies.

## FLEXIBLE CENTRALIZATION THROUGH CLOUD-RAN

### FLEXIBLE ASSIGNMENT OF RAN FUNCTIONALITY

A flexible assignment of RAN functionality can consider both the cloud-platform resource availability and the small-cell BH characteristics. In addition, cloud-computing platforms allow for the scalability that is required to cope with temporal and spatial traffic fluctuations in mobile networks. This scalability is a fundamental requirement to improve the utilization of mobile networks and to allow for an economically and ecologically sustainable operation of mobile networks.

Cloud-RAN is a disruptive technology in many ways and imposes new challenges on the signal processing in 5G mobile networks. Most importantly, it will exploit standard processor technology [general-purpose processors (GPPs)] to execute RAN functionality. By contrast, currently discussed C-RAN technology considers a baseband pooling approach where a large number of DSPs are provided at a central entity [8], [9]. Although this allows for resource sharing, C-RAN still uses specialized and expensive hardware and software. Hence, it is misleading to consider C-RAN as an example of cloud computing according to the IT definition by the U.S. National Institute of Standards and Technology [13].

Cloud-RAN will further foster scalable algorithms that are designed for cloud-computing environments and leverage massive parallelism. This implies that algorithms should not be simply ported to cloud-computing platforms but rather redesigned to gain from the available computing resources. Cloud-RAN allows for the deployment of algorithms that scale with the need for cooperation and coordination among the individual cells, i.e., depending on the traffic demand and user density, RAPs may be differently grouped or different algorithms may be deployed. In the following sections, this article provides more detailed examples for algorithms that benefit from an application to cloud-computing platforms.

To enable cloud-RAN, it is necessary to have a system architecture that provides the required interfaces without disruptive changes to an existing deployment. This architecture has been introduced in [11]. It does not imply changes to existing interfaces but introduces the concept of a virtual eNodeB (veNB). A veNB is composed of one or more RAPs, a cloud-computing platform, and the necessary BH links between these nodes. It maintains the same interfaces as a 3GPP LTE eNodeB (eNB) to maximize backward-compatibility. This system architecture requires 1) that the functionality at the eNB can be decomposed into reassignable functions and 2) that each function can be assigned either to the central processor or local RAPs. Furthermore, a tight integration of RAN, BH, and central processor is required, e.g., through joint coding as introduced in [14].

### OPPORTUNITIES OF CLOUD-RAN

Cloud computing offers the ability of computational load balancing to RANs. This allows for spending more computational

efforts on critical operations, e.g., in the case of interference scenarios or difficult channel conditions. In these scenarios, more advanced and computationally intense algorithms may be needed and could be executed in a cloud environment. By contrast, traditional implementations are hard real-time systems. Hence, a certain task such as decoding or scheduling is always executed within the same time window.

A flexible assignment of functionality will also allow for shaping the signaling load on the BH connection. For instance, in the case of high-capacity, low-latency BH, the central processor may process directly in-phase/quadrature (I/Q) samples. In the case of higher latency and lower bandwidth on the BH, the central processor may only perform upper-layer functionality. This will require changes to the operation of the BH and the signal processing platform, and it may require changes to the RAN protocol stack.

Cloud-RAN will open the door for many new applications in 5G. It offers the possibility of using signal processing software dedicated to a special purpose based on the actual service. It reflects the diversity of services, use cases, and deployments through flexibility and scalability of the signal processing platform. In addition, it may even take into account the complexity and abilities of terminals during the processing of signals. Finally, cloud-RAN avoids the typical vendor lock-in as in current deployments that follow a similar development observed in the mobile core network, which may be implemented on cloud-platforms [15].

The flexible centralization of RAN functionality will impact the operation of the 3GPP LTE RAN protocol stack and may even be limited by dependencies within the protocol stack. Table 1 provides an overview of promising functions of the 3GPP LTE radio protocol stack, which may be considered for a partial centralization. In general, the lower we place the functional split within the

protocol stack, the higher the overhead and the more stringent BH requirements are. Centralizing functionality on the physical layer (PHY) allows for computational diversity that depends directly on the number of users per RAP. Due to temporal and spatial fluctuations, the computational load can be balanced across cells. Central processing also allows for implementing multicell algorithms to avoid or exploit interference, e.g., intercell interference coordination and cooperative multipoint processing [16].

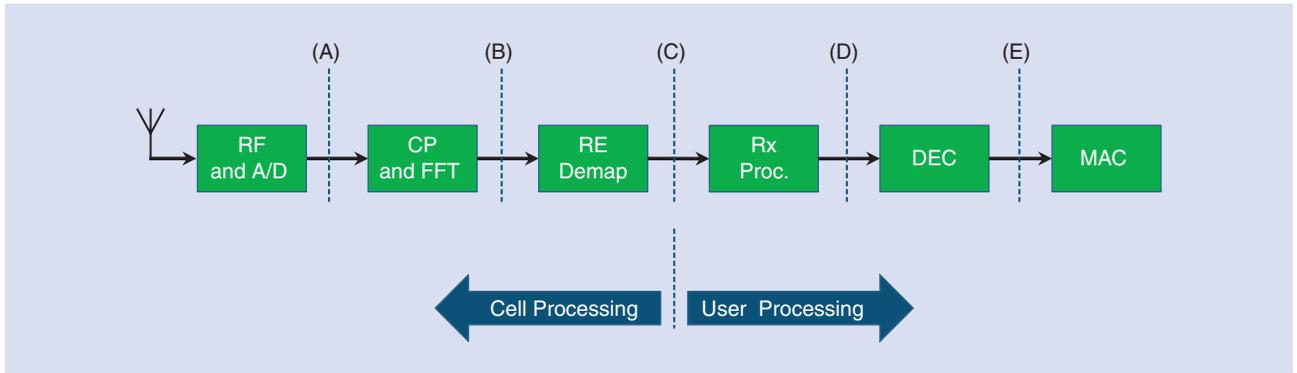
### FUNCTIONAL SPLIT

In this subsection, we introduce several functional split options that determine the execution of processing in the RAP or in the cloud-platform and directly influence the required BH data rate. The discussion is focused on the uplink (UL) since its processing load dominates the downlink (DL) processing. Detailed investigations of such splits have also been conducted in [17], but here we focus more on the opportunities of a flexible split. By relying on GPPs as opposed to dedicated hardware as used in the C-RAN concept, and through extensive use of function virtualization, the envisioned architecture allows us to adapt the functional split flexibly in time (e.g., according to traffic demand) and location (e.g., depending on the density of the deployment). Figure 1 illustrates the principle LTE signal processing chain of an UL receiver and different options of placing a functional split. Notice that similar shifts are also possible for DL processing as considered, e.g., in the context of precoding for massive MIMO systems in [18].

Subsequently, we discuss these split options and give numerical results on the required BH data rates per link between one RAP and the cloud-platform for a simple configuration as specified in Table 2.

**[TABLE 1] THE BENEFITS AND SIGNAL PROCESSING CHALLENGES FOR THE CENTRALIZATION OF SELECTED 3GPP LTE RADIO PROTOCOL FUNCTIONALITY ON THE PHY AND LOWER MAC LAYER.**

CENTRALIZED FUNCTIONALITY	CENTRALIZED REQUIREMENTS	CENTRALIZATION BENEFITS	CHALLENGES FOR SIGNAL PROCESSING
DETECTION AND FEC-DECODING	<ul style="list-style-type: none"> <li>■ DEPENDS ON CONTROL OVERHEAD IN UL</li> <li>■ LATENCY REQ. DEPENDS ON TIMING REQ. IN DL</li> <li>■ STRONG RELIABILITY</li> </ul>	<ul style="list-style-type: none"> <li>■ COOPERATIVE RECEIVER (RX)</li> <li>■ COMPUTATIONAL DIVERSITY</li> </ul>	<ul style="list-style-type: none"> <li>■ PREDETECTION AT RAP TO REDUCE BH OVERHEAD</li> <li>■ OPTIMAL QUANTIZATION OF SIGNALS AND EXCHANGE OVER BH</li> </ul>
FEC-ENCODING AND MODULATION AND PRECODING	<ul style="list-style-type: none"> <li>■ DEPENDS ON CONTROL OVERHEAD IN DL</li> <li>■ STRONG RELIABILITY</li> </ul>	<ul style="list-style-type: none"> <li>■ COOPERATIVE TRANSMITTER (TX)</li> <li>■ ADVANCED PRECODING</li> <li>■ COMPUTATIONAL DIVERSITY</li> </ul>	<ul style="list-style-type: none"> <li>■ SEPARATE PRECODING DECISION AND EXECUTION AT RAP AND CENTRAL PROCESSOR</li> <li>■ OPTIMAL QUANTIZATION OF SIGNALS AND EXCHANGE OVER BH</li> </ul>
LINK RELIABILITY PROTOCOLS (E.G., HARQ)	<ul style="list-style-type: none"> <li>■ DEPENDS ON ENTITY THAT PERFORMS RETRANSMISSION DECISION</li> </ul>	<ul style="list-style-type: none"> <li>■ SIMPLIFIED CENTRALIZATION OF SCHEDULING AND DECODING</li> </ul>	<ul style="list-style-type: none"> <li>■ PREDEFINED TIMING OF (N)ACK MESSAGES</li> <li>■ SEPARATION OF RETRANSMISSION DECISION AND PACKET COMBINING</li> <li>■ STRONG INTERACTION WITH OTHER FUNCTIONS, E.G., SCHEDULER, EN-/DECODER</li> </ul>
SCHEDULING AND INTERCELL RRM	<ul style="list-style-type: none"> <li>■ FLEXIBLE REQUIREMENTS</li> </ul>	<ul style="list-style-type: none"> <li>■ MULTICELL GAINS</li> <li>■ COMPUTATIONALLY EXPENSIVE ALGORITHMS</li> <li>■ GAINS DEPEND ON BH QUALITY</li> </ul>	<ul style="list-style-type: none"> <li>■ SCALABLE LATENCY REQUIREMENTS MUST BE SUPPORTED</li> <li>■ INTERCELL INTERFERENCE COORDINATION (ICIC) BASED ON CHANGING QUALITY OF CHANNEL STATE INFORMATION</li> <li>■ CHANGING COMPUTATIONAL COMPLEXITY</li> </ul>



**[FIG1]** The functional split between RAPs and the cloud-platform for UL transmission.

### I/Q FORWARDING (A)

By immediately forwarding the time-domain receive signals that have been downconverted to the baseband and analog-to-digital (AD) converted (indicated by block RF/AD), the complete receive frame including the cyclic prefix (CP) has to be transmitted over the BH link to the cloud-platform. This approach is usually referred to as *radio-over-fiber (RoF)* and is used in the common public radio interface (CPRI) standard [19]. The main benefit of this split is that almost no digital processing devices are required at the RAPs, potentially making them very small and cheap. If a flexible split varying over time is envisioned, the processing devices would have to be available at the RAPs anyway, nullifying this benefit. Also, the required BH data rate for I/Q forwarding is comparatively high and given as

$$D_{\text{BH}}^{\text{A}} = N_{\text{O}} \cdot f_{\text{S}} \cdot 2 \cdot N_{\text{Q}} \cdot N_{\text{R}} = 2 \cdot 30.72 \text{ MHz} \cdot 2 \cdot 10 \text{ bit} \cdot 2 = 2.46 \text{ Gbit/s}. \quad (1)$$

### SUBFRAME FORWARDING (B)

By removing the CP and transforming the Rx signal to frequency-domain using fast Fourier transformation (FFT), guard subcarriers can be removed (block CP/FFT). Since the number of guard subcarriers in LTE is  $\approx 40\%$ , this decreases the required BH data rate significantly.

$$D_{\text{BH}}^{\text{B}} = N_{\text{Sc}} \cdot T_{\text{S}}^{-1} \cdot 2 \cdot N_{\text{Q}} \cdot N_{\text{R}} = 1,200 \cdot (66 \mu\text{s})^{-1} \cdot 2 \cdot 10 \text{ bit} \cdot 2 = 720 \text{ Mbit/s}. \quad (2)$$

**[TABLE 2]** EXEMPLARY TRANSMISSION PARAMETERS FOR CALCULATING THE IMPACT OF FUNCTIONAL SPLIT CHOICES ON THE BH DATA RATE.

PARAMETER	SYMBOL	VALUE
BANDWIDTH	$B$	20 MHz
SAMPLING FREQUENCY	$f_{\text{S}}$	30.72 MHz
OVERSAMPLING FACTOR	$N_{\text{O}}$	2
NUMBER OF USED SUBCARRIERS	$N_{\text{Sc}}$	1,200
SYMBOL DURATION	$T_{\text{S}}$	66.6 $\mu\text{s}$
QUANTIZATION/SOFT BITS PER I/Q	$N_{\text{Q}}$	10
RX ANTENNAS	$N_{\text{R}}$	2
SPECTRAL EFFICIENCY	$S$	3 bit/cu
ASSUMED RB UTILIZATION	$\eta$	50%

As an FFT can be implemented on dedicated hardware very efficiently, the implementation in the RAP is worthwhile compared to the split option I/Q forwarding (A). As the per-cell based processing does not depend on the actual load of the RAP, load balancing gains can be only achieved if RAPs are completely turned off.

### RX DATA FORWARDING (C)

If only a part of the resource elements (REs) are actually utilized by the user equipment (UE) in a cell, only these REs remain after RE demapping (block RE Demap) and have to be forwarded to the cloud-platform. The required BH data rate is directly given by the fraction of utilized RE and thus, the subsequent splits can profit from load balancing gains.

$$D_{\text{BH}}^{\text{C}} = D_{\text{BH}}^{\text{B}} \cdot \eta = 720 \text{ Mbit/s} \cdot 0.5 = 360 \text{ Mbit/s}. \quad (3)$$

To allow for a joint processing of received signals from multiple RAPs, it has to be ensured that only REs of UE not considered for joint processing are removed, even if they are not (primarily) associated with the current RAP.

### SOFT-BIT FORWARDING (D)

The receive processing (block Rx Proc) per user consists of equalization in frequency domain, inverse discrete Fourier transformation (IDFT), MIMO receive processing, and demapping. In a MIMO scheme utilizing receiver diversity, the signals of multiple antennas are combined during channel equalization, thus removing the dependency on the number of receive antennas. This results in a reduced BH load of  $D_{\text{BH}}^{\text{D}} = D_{\text{BH}}^{\text{C}}/N_{\text{R}} = 180 \text{ Mbit/s}$ . In contrast, for spatial multiplexing with  $N_{\text{S}}$  layers per UE, the BH would correspond to  $D_{\text{BH}}^{\text{D}} = D_{\text{BH}}^{\text{C}} \cdot N_{\text{S}}/N_{\text{R}}$ . By this split, only joint decoding of soft bits forwarded by several RAPs is possible in the cloud-platform. Also note that usually the number of soft bits per symbol would depend on the modulation scheme (e.g., three soft-bits per information bit), and thus  $N_{\text{Q}}$  and the BH data rate would depend directly on the modulation order, which in turn depends on the access channel quality due to radio resource management (RRM).

### MAC (E)

During forward error correction (FEC) decoding (block DEC), data bits are recovered from the received symbols and redundant

bits are removed, resulting in the pure MAC payload at the decoder output. The resulting BH data rate depends largely on the used modulation and coding scheme (MCS), which is reflected here by the exemplary spectral efficiency  $S = 3$  bit/cu.

$$D_{\text{BH}}^E = N_{\text{sc}} \cdot T_s^{-1} \cdot \eta \cdot S$$

$$= 1,200 \cdot (66 \mu\text{s})^{-1} \cdot 0.5 \cdot 3 \text{ bit/cu} = 27 \text{ Mbit/s.} \quad (4)$$

FEC decoding is a complex task that is commonly performed on dedicated hardware and hence a centralized decoding on GPPs has not been considered in C-RAN. However, as outlined later in this article, recent results show that it can be performed on GPPs. On the other hand, performing decoding in the RAPs according to the split option MAC (E) terminates the possibility for joint PHY-layer processing in the cloud-platform and only cooperation on higher layers, e.g., joint scheduling, remains possible. As PHY-layer cooperation mainly revolves around interference mitigation, this option is beneficial in scenarios where RAPs are well separated, e.g., for indoor deployments or in narrow street canyons.

Obviously, the required BH data rate and the required processing power in the cloud decreases significantly when the functional split is shifted to the higher PHY processing layers or even to the MAC. However, this is traded off with lower centralization gains in terms of spectral efficiency and computational load balancing. The advantage of a flexible split is that we can reap the benefits of both extremes: load balancing for low traffic situations and high spectral efficiency by cooperative processing for high traffic. Since current BH standards like CPRI only support a very specific functional split, new and more flexible standards will have to be defined to enable cloud-RAN architectures.

The huge BH bandwidth requirements of functional shifts on the lower PHY layers also shows that improved and optimized BH technologies are required. While technologies offering sufficient bandwidth are already available [10], a joint design of radio access and BH links should be also considered to use the deployed capacity as efficiently as possible. Additionally, to further limit the BH rate between the RAPs and the cloud-platform, cooperative processing strategies could be used to directly exploit lower-layer interaction between RAPs. This would allow the use of heterogeneous BH technologies to interconnect the

RAPs and implement joint distributed detection techniques as depicted in Figure 2 and discussed in the next section.

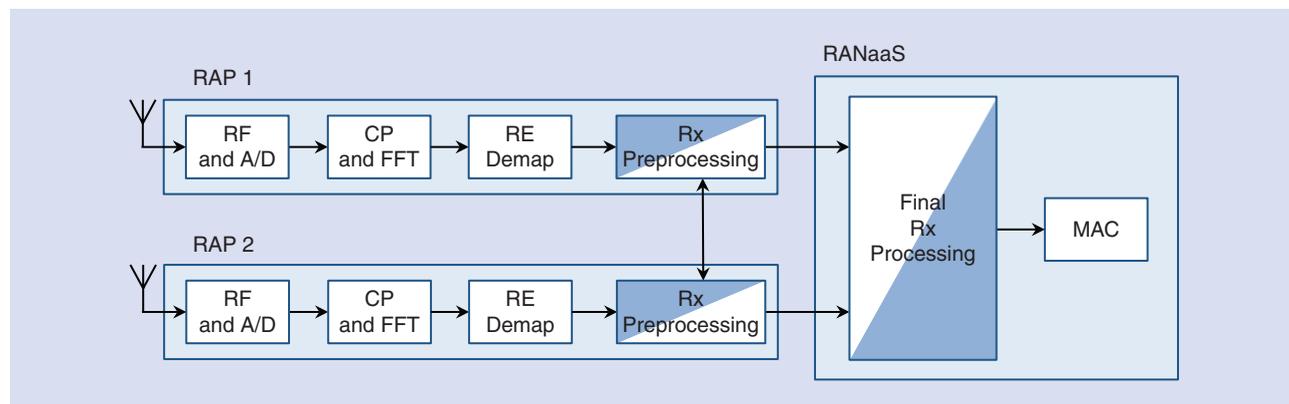
### SIGNAL PROCESSING IN THE CLOUD

The difficulty of implementing RAN functionality in a cloud-platform lies in the tight constraints caused by the 3GPP LTE protocol stack. This implies that individual tasks need to finish within a predefined time window. Figure 3 shows relevant parts of the 3GPP LTE protocol stack and two exemplary functional splits that correspond to options (C) and (D) in Figure 1. In the following, we discuss the benefits and challenges of a cloud implementation of three representative parts of the signal processing chain.

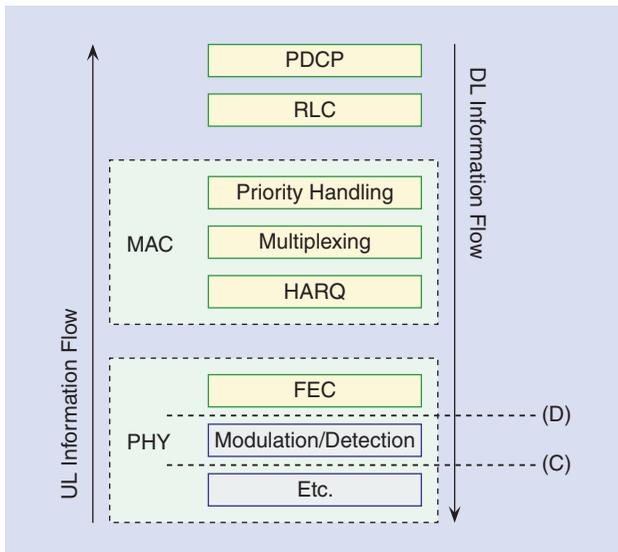
### HYBRID AUTOMATIC REPEAT REQUEST

Among all the timers defined in LTE, the one associated to the acknowledgment (ACK) of a UL physical frame at the MAC layer is the most critical one. The reception status of any frame sent through the air interface needs to be fed back to the transmitter, to proceed to the transmission of a new frame ACK or to attempt a retransmission negative ACK (NACK). This hybrid automatic repeat-request (HARQ) operation is performed at the MAC level, after all the physical processing of a codeword is done (detection, demodulation, and FEC decoding). In LTE, each frame sent at subframe  $n$  needs to be acknowledged (ACK or NACK) at subframe  $n + 4$  in both UL and DL directions, a subframe lasting 1 ms [20]. Hence, the overall receive process has to finish in 3 ms to stay compliant with the 3GPP LTE HARQ timing. This timing includes the processing at the RAPs of the physical blocks located before the split (see Figure 3 and both functional split options therein), the processing at the cloud-platform of the physical blocks located after the split and the round-time trip through the BH. However, some algorithms such as turbo-decoders underly a computational jitter which implies that the decoding time may vary. Hence, it may happen that packets are retransmitted even though they would have been decoded with more computational resources, i.e., either more time or more parallel processors. This computational jitter also adds up to the overall delay that needs to be considered.

To relax the timing constraint for the receive processing, we may adapt the HARQ process. The authors in [17] suggest for



**[FIG2]** The cooperative Rx preprocessing among RAPs with final Rx processing in the central processor.



**[FIG3]** The LTE protocol stack and exemplary functional splits.

example to suspend the HARQ process until the end of the receive processing. In the case that the receive processing is not finished in time, an ACK is sent after 3 ms to meet the timing requirements while receive processing is continued. If, at the end, successful decoding is not possible, a NACK is sent. As the UE does not immediately drop out a package when receiving an ACK to cope with transmission errors on the feedback channel, a retransmission of the particular packet can be scheduled later. However, this approach halves the achievable UE peak rate [17]. This drawback can be avoided by a preliminary HARQ process, where the initial feedback message is determined by estimating the decoding success based on the quality of the received signals (e.g., using models from link level simulations [21]). If correct decoding is likely, a preliminary ACK is sent to the UE, otherwise a preliminary NACK. Again, the standard techniques capturing feedback errors automatically handle erroneous preliminary feedback messages. This approach relaxes the timing constraints for the receive processing chain. It separates the most complex processing parts and the most latency-critical parts but still allows for high data rates depending on the reliability of preliminary ACK/NACK.

### FORWARD ERROR CORRECTION

The tight requirement of finishing the overall detection within 3 ms poses a significant challenge for executing FEC decoding within the cloud-platform due to its high complexity. Usually, FEC decoders are implemented in specialized hardware, such as application-specific integrated circuit (ASIC) designs or field-programmable gate array (FPGA) implementations [22]. However, the introduction of many-core architectures opens new perspectives for massively parallel implementations. To meet stringent requirements on data rates, cloud-based FEC decoders will need to fully exploit the available parallelism of a cloud-computing platform. In this context, low-density parity check (LDPC) [23] and turbo codes [24] are two promising candidates because both allow for accommodating various degrees of parallelization.

From a high-level perspective, two main approaches can be used to exploit parallelism in multicore platforms. The first approach parallelizes the decoder itself through decomposition of the decoding algorithms into multiple threads that run in parallel. Second, multiple codewords may be decoded in parallel. The first approach decreases the latency per codeword but introduces more synchronization overhead across different threads. By contrast, the second approach uses less synchronization objects and therefore increases the parallelization gain. However, it may introduce a higher latency per codeword compared to the first approach.

For very high throughput applications, LDPC codes are known to compare favorably against turbo codes because LDPC decoding allows for a higher degree of parallelism [25], [26]. Hence, LDPC codes are suitable for the first approach of decoder parallelization. However, software-based parallel LDPC decoders barely achieve throughputs of a few tens of Mbit/s, as reported in [27] for graphical processing units (GPUs), or in [28] for the signal processing on-demand architecture (SODA). In both cases, the main reason is the need for synchronization across different threads to access shared objects that results in scalability issues [27].

By contrast, parallelizing multiple codewords eliminates the need for synchronizing objects. This results in better scalability properties and the throughput of the multicodeword decoder is known to increase almost linearly with the number of cores [29]. Furthermore, it allows for different codes, algorithms, and configurations running in parallel. Multicodeword LDPC decoders have been reported to achieve throughputs up to 80 Mbit/s on the IBM CELL Broadband Engine [27], [30], with 24–96 codewords decoded in parallel. Recently, central processing unit (CPU) and GPU implementations of multicodeword turbo decoders have also been reported in [31] with a peak throughput from 55 Mbit/s to 122 Mbit/s, as the number of decoding iterations decreases from eight to four.

Figure 4 shows experimental results for spectral efficiency and required computational complexity of an 3GPP LTE UL decoder. To obtain these results, the turbo-decoder has been implemented on a default VMWare ESXi server with Ubuntu Linux host operating system, GNU C++ compiler, and codeword multithreading to account for the virtualization overhead. We measured the required CPU time to decode one codeword and determined the average CPU time within the 90% confidence interval.

Figure 4(a) shows the achievable spectral efficiency for a given signal-to-noise ratio (SNR) (additive white Gaussian noise, no fading). We illustrate the results for two cases: maximum throughput (high number of iterations possible) and low complexity (number of iterations limited to two). Reducing the complexity of the decoding process results in a performance penalty of 1–2 dB. In Figure 4(b), we show the required computational resources for a 10-MHz 3GPP LTE system. The required complexity strongly depends upon the SNR. First, it increases linearly with the number of information bits, which implies a logarithmic increase of complexity in SNR. Second, the complexity increases with the number of iterations that are necessary to decode a codeword. As shown in [32], the complexity increases superlinearly with decreasing SNR (in decibels) for a fixed MCS. In Figure 4(b),

markers show the SNR where the next higher MCS has been chosen. We notice at each of these markers an increase of the computational demand, which is then quickly decreasing in SNR.

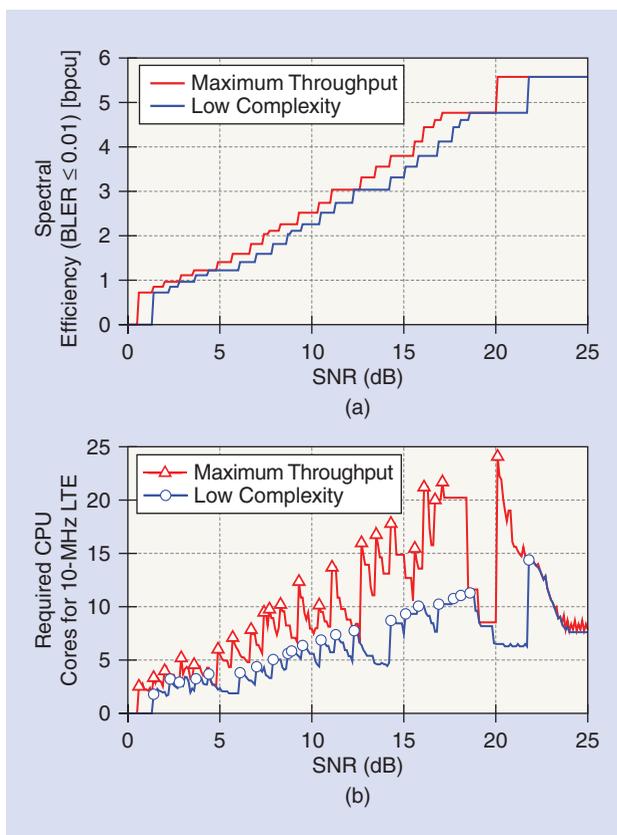
Apparently, this strongly varying computational demand allows for the exploitation of multiuser computational diversity at the centralized processor. For instance, the central processor can perform computational load balancing across multiple users to reduce the ratio of peak-to-average computational efforts. Furthermore, the central processor can actively shape the computational demand by selecting MCS to satisfy a computational constraint, e.g., in the case of a traffic burst the computational requirements may significantly increase and may exceed the available resources if MCSs are chosen based on maximum throughput. Finally, the computational load can be actively shaped by adjusting the number of quantization bits  $N_Q$  used for forwarding the Rx signals from the RAP to the cloud-platform over the BH. Figure 5 shows the tradeoff between number of turbo iterations and quantization bits  $N_Q$  for different modulation schemes at a target bit error rate (BER) of  $10^{-4}$ . Obviously, the decoding latency can significantly be reduced by increasing the number of quantization bits  $N_Q$  on the cost of a higher BH transmission rate.

### MULTIUSER DETECTION

Consider again the functional split option, Rx Data Forwarding (C) in Figure 3. In this case, I/Q samples are forwarded over high-capacity BH links to the central processor that performs joint multiuser detection (MUD) using the Rx signals of several RAPs. The joint processing of many RAPs implements a virtual MIMO architecture and the huge computational power offered by the cloud-platform allows for aggressive RRM across the RAPs. However, due to the heterogeneous nature of BH networks, it is also beneficial to use a mix of local processing at RAPs, cooperative processing among RAPs, and central processing in the cloud-platform. Promising techniques that are adaptable to changing BH and radio access parameters are, among others, multipoint turbo detection (MPTD) and in-network processing (INP).

The underlying idea of MPTD [33] is to schedule (edge) users attached to different RAPs on the same resource. Then, a joint detection of these users through a turbo processing approach is performed [21], [34]. Such processing could be done either centrally on the cloud-platform or locally in each RAP. If it is fully centralized, MPTD benefits from high degree of spatial diversity due to the different locations of the involved RAPs. Due to this spatial diversity increase, the centralization gain can be quite significant compared to a classical distributed detection.

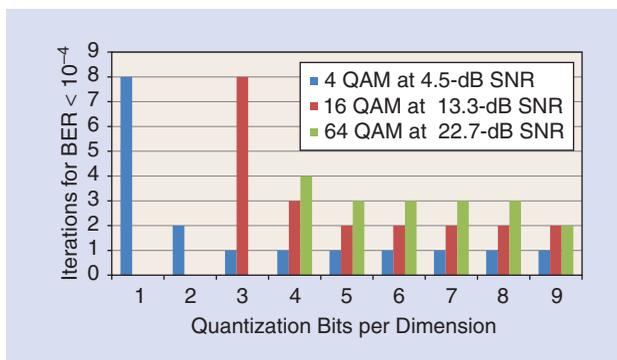
This split of functionality may offer significant centralization gains compared to distributed detection methods. This is illustrated in Figure 6 for an UL scenario with  $N_{UE} = 2$  users each equipped with  $N_T = 1$  transmit antenna. Both users interfere with each other at  $N_{RAP} = 2$  RAPs each equipped with  $N_R = 2$  receive antennas. We assume the worst case of identical path-losses. In addition, these results consider Rayleigh channel fading and LTE-compliant MCSs [20]. Figure 6 shows that at a frame error rate (FER) of 0.01 a centralization



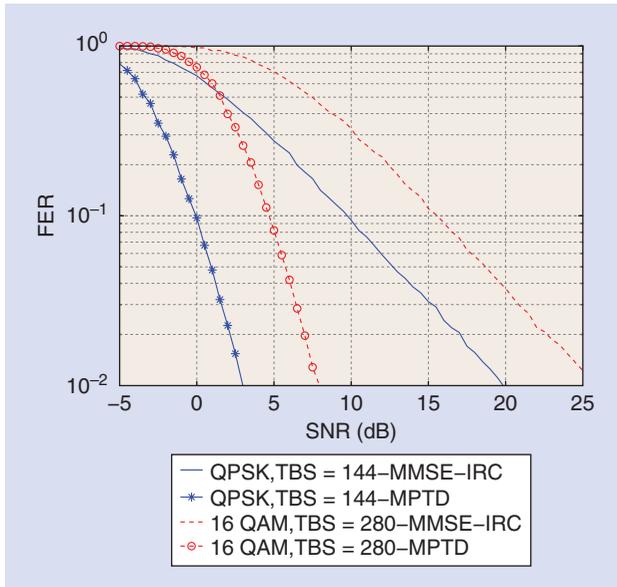
**[FIG4] Throughput and computational complexity results for turbo-decoding using an out-of-the-box cloud-computing platform and 3GPP LTE MCSs. (a) Spectral efficiency. (b) Required CPU cores.**

gain of about 17 dB can be achieved by MPTD compared to a linear minimum mean square error (MMSE) filter with interference rejection combining (IRC) [35].

An alternative approach that faces the joint MUD problem from an optimization perspective is INP. It allows for the solution of general estimation problems in a distributed, decentralized way within a network. The special class of consensus-based algorithms achieves this by iteratively reaching consensus of the estimates among the processing nodes [36], [37]. The adaptation of INP for an iterative distributed MUD has recently been presented in [38]. Due to its generic structure, INP can also be implemented with



**[FIG5] The number of turbo iterations required for BER  $< 10^{-4}$  versus number of quantization bits  $N_Q$  per I/Q dimension.**

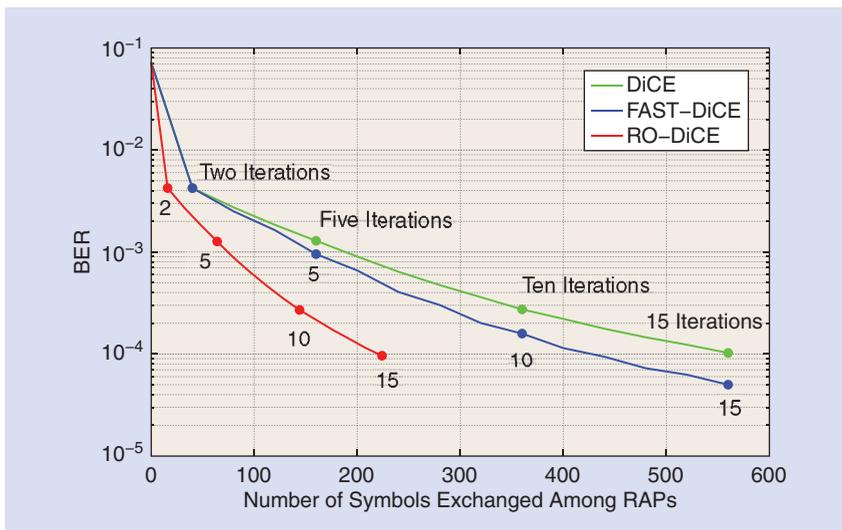


[FIG6] The FER over SNR for MPTD (centralized) and MMSE-IRC (distributed).

the desired mix of local, cooperative, and central processing within the distributed architecture in Figure 2 allowing for shifting the BH traffic flexibly within the network.

By combining for each RE the  $N_T$  transmit signals of all  $N_{UE}$  users into the signal vector  $x$ , the receive signal vector at RAP  $j$  is given by  $y_j = H_j x + n_j$  with  $H_j$  denoting the effective channel matrix and  $n_j$  representing the additive noise vector. In case of a fully centralized solution, the receive signals  $y_j$  of all  $N_{RAP}$  RAPs have to be forwarded to the central processing node and can be collected into the receive signal vector  $y = [y_1^H \dots y_{N_{RAP}}^H]^H = Hx + n$ , where  $H$  and  $n$  denote the stacked channel matrix and the stacked noise vector. The solution of a centralized least squares (LS) problem

$$\hat{x} = \underset{x}{\operatorname{argmin}} \|y - Hx\|^2 \quad (5)$$



[FIG7] The BER over BH overhead at SNR of 10 dB.

is given by  $\hat{x} = H^+ y$  with the Moore–Penrose pseudo-inverse  $H^+ = (H^H H)^{-1} H^H$ . For a distributed calculation of this central solution, local estimates  $\hat{x}_j$  per node  $j$  are introduced to reformulate the LS problem by a set of local optimization problems

$$\hat{x}_j = \underset{\{x_i\}_{i \in \mathcal{J}}}{\operatorname{argmin}} \sum_{j=1}^{N_{RAP}} \|y_j - H_j x_j\|^2 \quad (6a)$$

$$\text{s.t. } x_j = x_i \quad \forall j \in \mathcal{J}, i \in \mathcal{N}_j, \quad (6b)$$

where  $\mathcal{J}$  denotes the set of all RAPs and  $\mathcal{N}_j$  the set of all RAPs connected with RAP  $j$ . The consensus constraint (6b) directly couples estimates of neighboring nodes guaranteeing that the estimates of all nodes converge to the central LS solution (5). In [37], the distributed consensus-based estimation (DiCE) algorithm has been introduced, which allows for parallel processing across the involved RAPs. Furthermore, the required information exchange is reduced by the reduced-overhead-DiCE (RO-DiCE) [39] approach and the fast-DiCE implementation improves the convergence speed [40].

Figure 7 shows the BER for uncoded binary phase shift keying transmission with  $N_{UE} = 2$  users with  $N_T = 2$  transmit antennas to  $N_{RAP} = 4$  RAPs with  $N_R = 4$  over Rayleigh-fading channels and fully connected mesh network of RAPs. It further compares the three different DiCE implementations for a fixed SNR of 10 dB versus the number of signals exchanged among the RAPs. Obviously, with an increasing number of iterations the BER performance improves at the cost of an increased communication overhead. In particular, the Fast-DiCE approach allows for a faster convergence and the RO-DiCE reduces the overhead by 60% at the same BER.

MUD imposes new challenges on signal processing within a cloud-computing environment. Among other challenges, synchronization needs to be maintained and taken into account. Furthermore, the data exchange between virtual machines needs to be orchestrated to allow for low delays during the MUD process. Scalability and resource pooling are two major advantages of cloud computing. This requires a hypervisor that takes into account requirements and constraints from the RAN functionality and distributes the work load accordingly, e.g., resources per virtual machine, assignment of users to virtual machines, mapping of communication clusters to virtual machines, and massive parallelization across multiple virtual machines and possible different hardware racks.

## CONCLUSIONS

This article discussed benefits and challenges that may be implied by cloud-computing platforms on signal processing algorithms. The novel RANaaS concept was introduced, which realizes cloud technologies in 5G mobile networks and allows for a flexible functional split between RAPs and the centralized cloud-platform. This allows

for centralization benefits, but also introduces challenges due to the strict timing constraints imposed by the 3GPP LTE protocol stack. These challenges were identified and enabling technologies were discussed.

## ACKNOWLEDGMENTS

The research leading to these results has received funding from the European Union's Seventh Framework Programme (FP7/2007–2013) under grant agreement number 317941. We would like to acknowledge our iJOIN colleagues' contributions, although the views expressed in this article are our own and do not necessarily represent the project.

## AUTHORS

**Dirk Wübben** (wuebben@ant.uni-bremen.de) received the Dipl.-Ing. (FH) degree in electrical engineering from the University of Applied Science Münster, Germany, in 1998, and the Dipl.-Ing. (Uni) degree and the Dr.-Ing. degree in electrical engineering from the University of Bremen, Germany, in 2000 and 2005, respectively. In 2001, he joined the Department of Communications Engineering, University of Bremen, Germany, where he is currently a senior researcher and lecturer. His research interests include wireless communications, signal processing for multiple antenna systems, cooperative communication systems, and channel coding. He has published more than 80 papers in international journals and conference proceedings. He is a board member of the Germany Chapter of the IEEE Information Theory Society and an editorial board member of Elsevier's journal, *Physical Communication*.

**Peter Rost** (peter.rost@ieee.org) received his Ph.D. degree from Technische Universität Dresden, Germany, in 2009, and his M.Sc. degree from the University of Stuttgart, Germany, in 2005. From 1999 to 2002 he was with the Fraunhofer Institute for Beam and Material Technologies, Dresden, Germany, and from 2002 to 2005 he was with IBM Deutschland Entwicklung GmbH, Böblingen, Germany. In June 2005 he joined the Vodafone chair of Prof. Gerhard Fettweis at Technische Universität Dresden and focused on different aspects of relaying in the context of mobile communications systems. Since April 2010, he has been a member of the Mobile and Wireless Networks group at NEC Laboratories Europe, where he is working as a senior researcher in business unit projects, 3GPP RAN2 as active delegate, and the European Union's Seventh Framework Programme projects FLAVIA and iJOIN, which he currently manages as technical manager. He was the Technical Program Committee chair at the Spring 2013 IEEE Vehicular Technology Conference and a member of the IEEE Communications Society GLOBECOM/ICC Technical Committee (GITC) and *IEEE Transactions on Wireless Communications* Executive Editorial Committee. He has published more than 30 scientific publications and authored 18 patents and patent applications.

**Jens Bartelt** (Jens.Bartelt@tu-dresden.de) received his Dipl.-Ing. (M.S.E.E.) degree from Technische Universität Dresden, Germany, in 2012. In 2011–2012 he worked as an intern for Rohde & Schwarz in Munich. Since 2013, he has been a research associate at

the Vodafone Chair Mobile Communications Systems at Technische Universität Dresden, Germany, working toward his Ph.D. degree. His research interests include cloud-based mobile networks, millimeter-wave communication, and channel coding.

**Massinissa Lalam** (massinissa.lalam@sagemcom.com) received his M.S. degree from the Institut National Polytechnique de Grenoble (ENSIMAG, INPG, France) in 2002. He received his Ph.D. degree from the Ecole Nationale Supérieure de Télécommunications de Bretagne (Telecom Bretagne, France) in 2006. In 2007, he took a postdoctoral position with Telecom Bretagne, and from 2008 to mid-2009 he was with Orange Labs (France). He is now with Sagemcom (France) where he works on wireless (Wi-Fi) and cellular (third generation/long-term evolution) technologies. His expertise includes link- and system-level performance evaluation, heterogeneous network, network modeling, and radio resource management. He is the author or coauthor of over 15 international publications.

**Valentin Savin** (valentin.savin@cea.fr) received his M.S. degree in mathematics from École Normale Supérieure de Lyon, France, in 1997 and his Ph.D. degree in mathematics from J. Fourier Institute, Grenoble, France, in 2001. He also holds an M.S. degree in cryptography, security, and coding theory from the University of Grenoble 1. Since 2005, he has been with the Digital Communications Laboratory of CEA-LETI, working on the analysis and design of binary and nonbinary low-density parity check codes for physical- and upper-layer applications. He has published more than 40 papers in international journals and conference proceedings, holds six patents, and is currently participating in or coordinating several French and European Union Seventh Framework Programme research projects in information and communications technology.

**Matteo Gorgolione** (Matteo.Gorgolione@cea.fr) received his M.Sc. degree in telecommunications engineering in 2009 from the Polytechnic University of Turin (Italy) and his Ph.D. degree in information and communication sciences in 2012 from the University of Cergy-Pontoise (France). His Ph.D. dissertation was conducted in collaboration with the Digital Communications Laboratory of CEA-LETI, Grenoble (France). He is currently with CEA-LETI as a postdoctoral research fellow. His research interests are mainly related to the design of binary and nonbinary low-density parity check codes for cooperative communications.

**Armin Dekorsy** (dekorsy@ant.uni-bremen.de) received his Dipl.-Ing. (FH) (B.Sc.) degree from Fachhochschule Konstanz, Germany, Dipl.-Ing. (M.Sc.) degree from the University of Paderborn, Germany, and Ph.D. degree from the University of Bremen, Germany, all in communications engineering. From 2000 to 2007, he worked as research engineer at Deutsche Telekom AG and as a distinguished member of technical staff at Bell Labs Europe, Lucent Technologies. In 2007 he joined Qualcomm GmbH as a European research coordinator conducting Qualcomms' internal and external European research projects like ARTIST4G, BeFemto, and WINNER+. He has held the chair position of the Department of Communications Engineering, University of Bremen, since April 2010. His current research interests include resource management, transceiver design and digital

signal processing for wireless communications systems in health care, automation, and mobile communications. He is a member of the Information Technology Society (ITG) expert committee "Information and System Theory" of the Association for Electrical, Electronic, and Information Technologies (VDE) and the IEEE Communications Society and IEEE Signal Processing Society.

**Gerhard Fettweis** (fettweis@ifn.et.tu-dresden.de) earned his Ph.D. degree from RWTH Aachen in 1990. After one year at IBM Research in San Jose, California, he moved to TCSI Inc., Berkeley, California. Since 1994, he has been the Vodafone chair professor at TU Dresden, Germany, where currently 20 companies from Asia/Europe/United States sponsor his research on wireless transmission and chip design. He coordinates two DFG centers, the Center for Advancing Electronics Dresden (cfaED) and Highly Adaptive Energy-Efficient Computing (HAEC), at Technische Universität Dresden, Germany. He is an IEEE Fellow, member of the German Academy of Science and Engineering (acatech), has an honorary doctorate from Tampere University of Technology, and has received multiple awards. He has helped organize IEEE conferences, was the Technical Program Committee chair of the 2009 IEEE International Conference on Communications as well as the 2012 IEEE Technology Time Machine, and general chair of the Spring 2013 IEEE Vehicular Technology Conference.

## REFERENCES

- [1] M. Dohler, R. Heath, A. Lozano, C. Papadias, and R. A. Valenzuela, "Is the PHY layer dead?" *IEEE Commun. Mag.*, vol. 49, no. 4, pp. 159–165, Apr. 2011.
- [2] N. Bhushan, J. Li, D. Malladi, R. Gilmore, D. Brenner, A. Damnjanovic, R. T. Sukhvasi, C. Patel, and S. Geirhofer, "Network densification: The dominant theme for wireless evolution into 5G," *IEEE Commun. Mag.*, vol. 52, no. 2, pp. 82–89, Feb. 2014.
- [3] 3GPP, "Small cell enhancements for E-UTRA and E-UTRAN—Physical layer aspects," Tech. Rep. TR.36.872, Sept. 2013.
- [4] 3GPP, "Small cell enhancements for E-UTRA and E-UTRAN—Higher layer aspects," Tech. Rep. TR.36.842, May 2013.
- [5] H. Guan, T. Kolding, and P. Merz, "Discovery of Cloud-RAN," in *Proc. Cloud-RAN Workshop*, Beijing, China, Apr. 2010.
- [6] NGMN, "Suggestions on Potential Solutions to C-RAN by NGMN Alliance," Tech. Rep., Jan. 2013.
- [7] D. Wake, A. Nkansah, and N. J. Gomes, "Radio over fiber link design for next generation wireless systems," *IEEE/OSA J. Lightwave Technol.*, vol. 28, no. 16, pp. 2456–2464, Aug. 2010.
- [8] K. Chen, C. Cui, Y. Huang, and B. Huang, "C-RAN: A green RAN framework," in *Green Communications: Theoretical Fundamentals, Algorithms and Applications*, J. Wu, S. Rangan, and H. Zhang, Eds. Boca Raton, FL: CRC Press, pp. 279–304, 2013.
- [9] G. Li, S. Zhang, X. Yang, F. Liao, T. Ngai, S. Zhang, and K. Chen, "Architecture of GPP based, scalable, large-scale C-RAN BBU pool," in *Proc. Int. Workshop Cloud Base-Station Large-Scale Cooperative Communications, IEEE GLOBECOM 2012 Workshops*, Anaheim, CA, Dec., pp. 267–272.
- [10] J. Bartelt, G. Fettweis, D. Wübben, M. Boldi, and B. Melis, "Heterogeneous backhaul for cloud-based mobile networks," in *Proc. 1st Int. Workshop Cloud Technologies Energy Efficiency Mobile Communication Networks (CLEEN 2013), IEEE VTC2013-Fall Workshops*, Las Vegas, NV, Sept.
- [11] P. Rost, C. J. Bernardos, A. De Domenico, M. Di Girolamo, M. Lalam, A. Maeder, D. Sabella, and D. Wübben, "Cloud technologies for flexible 5G radio access networks," *IEEE Commun. Mag.*, vol. 52, no. 5, pp. 68–76, May 2014.
- [12] F. Boccardi, R. W. Heath, A. Lozano, T. L. Marzetta, and P. Popovski, "Five disruptive technology directions for 5G," *IEEE Commun. Mag.*, vol. 52, no. 2, pp. 74–80, Feb. 2014.
- [13] P. Mell and T. Grance, "The NIST definition of cloud computing," Tech. Rep., National Inst. Standards Technology, Special Publication 800-145, Sept. 2011.
- [14] J. Bartelt and G. Fettweis, "Radio-over-radio: I/Q-stream backhauling for cloud-based networks via millimeter wave links," in *Proc. 1st Int. Workshop Cloud Technologies Energy Efficiency Mobile Communication Networks, IEEE GLOBECOM 2013 Workshops (IWCPM 2013)*, Atlanta, GA, Dec., pp. 778–783.
- [15] J. Kempf, B. Johansson, S. Pettersson, H. Luning, and T. Nilsson, "Moving the mobile evolved packet core to the cloud," in *Proc. IEEE Int. Conf. Wireless Mobile Computing, Networking, Communications*, Barcelona, Spain, Oct. 2012, pp. 784–791.
- [16] R. Irmer, H. Droste, P. Marsch, M. Grieger, G. Fettweis, S. Brueck, H.-P. Mayer, L. Thiele, and V. Jungnickel, "Coordinated multipoint: Concepts, performance, and field trial results," *IEEE Commun. Mag.*, vol. 49, no. 2, pp. 102–111, Feb. 2011.
- [17] U. Dötsch, M. Doll, H. P. Mayer, F. Schaich, J. Segel, and P. Schier, "Quantitative analysis of split base station processing and determination of advantageous architectures for LTE," *Bell Labs Tech. J.*, vol. 18, no. 1, pp. 105–128, May 2013.
- [18] S. Park, C.-B. Chae, and S. Bahk, "Before/after precoded massive MIMO in cloud radio access networks," *J. Commun. Netw.*, vol. 15, no. 4, pp. 398–406, Aug. 2013.
- [19] CPRI. (2013, Aug.). Common public radio interface (CPRI); Interface specification (V6.0). Tech. Rep. [Online]. Available: <http://www.cpri.info/>
- [20] 3GPP, "Radio access network; evolved universal terrestrial radio access; (E-UTRA); physical channels and modulation (Release 10)," Tech. Rep. TS.36.211, Dec. 2012.
- [21] R. Visoz, A.O. Berthet, and M. Lalam, "Semi-analytical performance prediction methods for iterative MMSE-IC multiuser MIMO joint decoding," *IEEE Trans. Commun.*, vol. 58, no. 9, pp. 2576–2589, Sept. 2011.
- [22] F. Kienle, N. Wehn, and H. Meyr, "On complexity, energy- and implementation-efficiency of channel decoders," *IEEE Trans. Commun.*, vol. 59, no. 12, pp. 3301–3310, Dec. 2011.
- [23] R. G. Gallager, *Low-Density Parity Check Codes* (Research Monograph Series). Cambridge, U.K.: MIT Press, 1963.
- [24] C. Berrou and A. Glavieux, "Near optimum error correcting coding and decoding: Turbo-codes," *IEEE Trans. Commun.*, vol. 44, no. 10, pp. 1261–1271, Oct. 1996.
- [25] M. M. Mansour and N. R. Shanbhag, "High-throughput LDPC decoders," *IEEE Trans. Very Large Scale Integration (VLSI) Syst.*, vol. 11, no. 6, pp. 976–996, Dec. 2003.
- [26] C. Zhang, Z. Wang, J. Sha, L. Li, and J. Lin, "Flexible LDPC decoder design for multigigabit-per-second applications," *IEEE Trans. Circuits Syst. I*, vol. 57, no. 1, pp. 116–124, Jan. 2010.
- [27] G. Falcao, L. Sousa, and V. Silva, "Massively LDPC decoding on multicore architectures," *IEEE Trans. Parallel Distrib. Syst.*, vol. 22, no. 2, pp. 309–322, Feb. 2011.
- [28] S. Seo, T. Mudge, Y. Zhu, and C. Chakrabarti, "Design and analysis of LDPC decoders for software defined radio," in *Proc. IEEE Workshop Signal Processing Systems*, Shanghai, China, Oct. 2007, pp. 201–215.
- [29] A. Diavastos, P. Petrides, G. Falcao, and P. Trancoso, "LDPC decoding on the Intel SCC," in *Proc. IEEE 20th Euromicro Int. Conf. Parallel, Distributed Network-Based Processing (PDP)*, Garching, Germany, Feb. 2012, pp. 57–65.
- [30] G. Falcao, V. Silva, L. Sousa, and J. Marinho, "High coded data rate and multicodeword WiMAX LDPC decoding on Cell/BE," *Electron. Lett.*, vol. 44, no. 24, pp. 1415–1417, Feb. 2008.
- [31] M. Wu, G. Wang, B. Yin, C. Studer, and J. R. Cavallaro, "HSPA+/LTE-A turbo decoder on GPU and multicore CPU," in *Proc. 47th IEEE Asilomar Conf. Signals, Systems, Computers (ASILOMAR)*, Pacific Grove, CA, Nov. 2013, pp. 824–828.
- [32] P. Grover, K. A. Woyach, and A. Sahai, "Towards a communication-theoretic understanding of system-level power consumption," *IEEE J. Sel. Areas Commun.*, vol. 29, no. 8, pp. 1744–1755, Sept. 2011.
- [33] iJOIN. (2013, Nov.). State-of-the-art of and promising candidates for PHY layer approaches on access and backhaul network. Tech. Rep. [Online]. Available: [www.ict-ijoin.eu/wp-content/uploads/2014/01/D2.1.pdf](http://www.ict-ijoin.eu/wp-content/uploads/2014/01/D2.1.pdf)
- [34] G. Caire and R. Müller, "The optimal received power distribution of IC-based iterative multiuser joint decoders," in *Proc. 39th Annu. Allerton Conf. Communications, Control Computing, Monticello, IL, Oct. 2001*.
- [35] R4-132017, "WF on NAICS receiver terminology," 3GPP TSG-RAN WG4#66bis, Apr. 2013.
- [36] H. Zhu, A. Cano, and G. B. Giannakis, "Distributed consensus-based demodulation: Algorithms and error analysis," *IEEE Trans. Wireless Commun.*, vol. 9, no. 6, pp. 2044–2054, June 2010.
- [37] H. Paul, J. Fliege, and A. Dekorsy, "In-network-processing: Distributed consensus-based linear estimation," *IEEE Commun. Lett.*, vol. 17, no. 1, pp. 59–62, Jan. 2013.
- [38] H. Paul, B.-S. Shin, D. Wübben, and A. Dekorsy, "In-network-processing for small cell cooperation in dense networks," in *Proc. 1st Int. Workshop Cloud Technologies Energy Efficiency Mobile Communication Networks, IEEE VTC2013-Fall Workshops (CLEEN 2013)*, Las Vegas, NV, Sept.
- [39] B.-S. Shin, H. Paul, D. Wübben, and A. Dekorsy, "Reduced overhead distributed consensus-based estimation algorithm," in *Proc. 1st Int. Workshop Cloud Technologies Energy Efficiency Mobile Communication Networks, IEEE GLOBECOM Workshops 2013 (IWCPM 2013)*, Atlanta, GA, Dec., pp. 784–789.
- [40] G. Xu, H. Paul, D. Wübben, and A. Dekorsy, "Fast distributed consensus-based estimation for cooperative wireless sensor networks," in *Proc. 18th Int. ITG Workshop Smart Antennas (WSA 2014)*, Erlangen, Germany, Mar.