Understanding O-RAN: Architecture, Interfaces, Algorithms, Security, and Research Challenges

Michele Polese[®], *Member, IEEE*, Leonardo Bonati[®], *Member, IEEE*, Salvatore D'Oro[®], *Member, IEEE*, Stefano Basagni[®], *Senior Member, IEEE*, and Tommaso Melodia, *Fellow, IEEE*

Abstract—The Open Radio Access Network (RAN) and its embodiment through the O-RAN Alliance specifications are poised to revolutionize the telecom ecosystem. O-RAN promotes virtualized RANs where disaggregated components are connected via open interfaces and optimized by intelligent controllers. The result is a new paradigm for the RAN design, deployment, and operations: O-RAN networks can be built with multi-vendor, interoperable components, and can be programmatically optimized through a centralized abstraction layer and data-driven closed-loop control. Therefore, understanding O-RAN, its architecture, its interfaces, and workflows is key for researchers and practitioners in the wireless community. In this article, we present the first detailed tutorial on O-RAN. We also discuss the main research challenges and review early research results. We provide a deep dive of the O-RAN specifications, describing its architecture, design principles, and the O-RAN interfaces. We then describe how the O-RAN RAN Intelligent Controllers (RICs) can be used to effectively control and manage 3GPP-defined RANs. Based on this, we discuss innovations and challenges of O-RAN networks, including the Artificial Intelligence (AI) and Machine Learning (ML) workflows that the architecture and interfaces enable, security, and standardization issues. Finally, we review experimental research platforms that can be used to design and test O-RAN networks, along with recent research results, and we outline future directions for O-RAN development.

Index Terms-Open RAN, O-RAN, cellular, 5G, 6G

I. Introduction

THE COMPLEXITY of cellular networks is increasing [1], with next-generation wireless systems built on a host of heterogeneous technologies and frequency bands. New developments include massive Multiple Input, Multiple Output (MIMO) [2], millimeter wave and sub-terahertz communications [3], network-based sensing [4], network slicing [5], [6], [7], [8], and Machine Learning (ML)-based digital signal processing [9], among others. This will impose increasing capital and operational costs for the networks operators, which will

Manuscript received 1 August 2022; revised 18 November 2022 and 13 January 2023; accepted 17 January 2023. Date of publication 23 January 2023; date of current version 23 May 2023. This work was supported in part by the U.S. National Science Foundation under Grant CNS-1923789 and Grant CNS-2112471, and in part by the U.S. Office of Naval Research under Grant N00014-20-1-2132. (Corresponding author: Michele Polese.)

The authors are with the Institute for the Wireless Internet of Things, Northeastern University, Boston, MA 02115 USA (e-mail: m.polese@northeastern.edu; l.bonati@northeastern.edu; s.doro@northeastern.edu; s.basagni@northeastern.edu; melodia@northeastern.edu).

Digital Object Identifier 10.1109/COMST.2023.3239220

have to continuously upgrade and maintain their infrastructure to keep up with new market trends and technology and customer requirements [10].

Managing and optimizing these new network systems require solutions that open the Radio Access Network (RAN). This makes it possible to expose data and analytics and to enable data-driven optimization, closed-loop control, and automation [11]. Current approaches to cellular networking, however, are far from open. Today, RAN components are monolithic units, all-in-one solutions that implement each and every layer of the cellular protocol stack. They are provided by a limited number of vendors and seen by the operators as black-boxes. Reliance on black-box solutions has resulted in: (i) limited reconfigurability of the RAN, with equipment whose operations cannot be fine-tuned to support diverse deployments and different traffic profiles; (ii) limited coordination among network nodes, preventing joint optimization and control of RAN components; and (iii) vendor lock-in, with limited options for operators to deploy and interface RAN equipment from multiple vendors. Under these circumstances, optimized radio resource management and efficient spectrum utilization through real-time adaptation become extremely challenging [12].

To overcome these limitations, in the last decade several research and standardization efforts have promoted the Open RAN as the new paradigm for the RAN of the future. Open RAN deployments are based on disaggregated, virtualized and software-based components, connected through open and well-defined interfaces, and interoperable across different vendors. Disaggregation and virtualization enable flexible deployments, based on cloud-native principles. This increases the resiliency and reconfigurability of the RAN. Open and interoperable interfaces also allow operators to onboard different equipment vendors, which opens the RAN ecosystem to smaller players. Finally, open interfaces and software-defined protocol stacks enable the integration of intelligent, data-driven closed-loop control for the RAN.

The O-RAN Alliance is an industry organization created in 2018 with the goal of implementing these principles on top of 3GPP LTE and NR RANs. Specifically, the O-RAN Alliance embraces and extends the 3GPP NR 7.2 split for base stations [13]. The latter disaggregates base station functionalities into a Central Unit (CU), a Distributed Unit (DU), and a Radio Unit (RU). Moreover, O-RAN connects them to intelligent controllers through open interfaces that can stream telemetry from the RAN and deploy control actions and policies to it.

The O-RAN architecture includes indeed two RAN Intelligent Controllers (RICs) that perform management and control of the network at near-real-time (10 ms to 1 s) and non-real-time (more than 1 s) time scales [12]. Finally, the O-RAN Alliance is defining a virtualization platform for the RAN, and extending the definition of 3GPP and eCPRI interfaces to connect RAN nodes.

Contributions: The Open RAN paradigm and, specifically, O-RAN networks will drastically change the design, deployment, and operations of the next generations of cellular networks. They will enable, among other things, transformative applications of ML for optimization and control of the RAN [12]. In this paper, we provide a detailed overview of how O-RAN will revolutionize future cellular networks. We do so through a comprehensive analysis of the O-RAN technical specifications, architectural components, of the interfaces connecting them, and of the ML and closed-loop control workflows that O-RAN enables. We also discuss the new security challenges and opportunities introduced by O-RAN, as well as the main publicly available experimental platforms that enable research and development of O-RAN components. Finally, we survey recent results on design and optimization of O-RAN, and discuss the issues that need to be addressed to fully realize the O-RAN vision. The goal is to offer the interested reader a clear picture of the state of the art in O-RAN, and a deep understanding of the opportunities that the Open RAN introduces in the cellular ecosystem.

Other papers [12], [14], [15], [16], [17], [18], [19], [20] introduce the O-RAN building blocks and architecture, with use cases mostly related to the application of machine learning to the RAN. The literature on Open RAN also includes several high-level white papers that summarize different elements of the O-RAN architecture [10], [21], [22], [23], [24], [25], [26], [27], [28], [29], [30], [31], [32]. Differently from these, we introduce here a multi-faceted perspective on O-RAN, which starts from the foundational principles, covers in details the architectural components and the interfaces, and then connects these elements to highlight AI/ML use cases, security issues, deployment options, testbeds, and future research and development challenges. Notably, this is the first paper that describes in detail the full set of O-RAN specifications for the RICs and interfaces, including how O-RAN effectively enables control of 3GPP-defined network elements through custom logic running on the intelligent controllers.

Paper structure: The rest of this paper is organized as shown in Fig. 1. Sections II to V introduce specific components of O-RAN networks; Sections VI to XI discuss topics that are relevant to the overall O-RAN vision and architecture; Section XII concludes this work. In particular, Section II describes the key principles of the O-RAN architecture, and introduces its components and the control loops that O-RAN enables. The near-real-time RIC and RAN control are discussed in Section III, while the non-real-time RIC is presented in Section IV. Section V is a deep dive on the O-RAN interfaces that connect the RAN and the RICs. Section VI describes the Artificial Intelligence (AI)/ML workflow supported in O-RAN networks. Section VII summarizes the main O-RAN use cases and related research

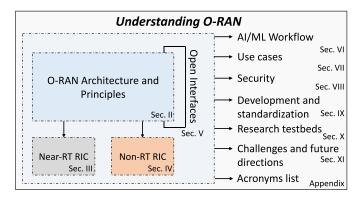


Fig. 1. O-RAN components and paper organization. Sections II—V (left part of the figure) introduce the general architecture of O-RAN, the RICs, and the open interfaces connecting them. Sections VI—XI (right part of the figure) discuss topics that relate to the overall Open RAN architecture, including technical research challenges to be addressed.

results. Section VIII reviews security challenges in O-RAN, and Section IX presents the development efforts and structure of the O-RAN Alliance. Publicly-available research and experimental platforms for O-RAN are discussed in Section X. Finally, Section XI provides an outlook on future directions and challenges for the Open RAN, and Section XII concludes the paper. We also include examples of O-RAN messages and a list of acronyms at the end of the paper.

II. O-RAN KEY ARCHITECTURAL PRINCIPLES

The Open RAN vision is based on years of research on open and programmable networks. These principles have been at the center of the Software-defined Networking (SDN) transformation in wired networks [33] in the past 15 years, and have started moving into the wireless domain more recently. For example, the xRAN Forum—an initiative led by operators—has proposed a standardized fronthaul interface, and introduced the idea of open, standardized interfaces for the integration of external controllers in the RAN [30]. In parallel, the Cloud RAN (C-RAN) architecture (promoted, among others, by the operator-led C-RAN Alliance [34]) has emerged as a solution to centralize most of the baseband processing for the RAN in virtualized cloud data centers [35], [36], connected to remote radio units through high speed fronthaul interfaces. C-RAN enabled more refined signal processing and load balancing techniques by leveraging centralized data and control paths, while reducing costs by multiplexing computational resources. In 2018, these two initiatives joined forces to launch the O-RAN Alliance with the overall goal of specifying and eventually standardizing an architecture and a set of interfaces to realize an Open RAN [34]. In just four years, the O-RAN Alliance has scaled up to more than 300 members and contributors. Its specifications are expected to drive 50% of RAN-based revenues by 2028 [10].

Overall, it is possible to identify four foundational principles for the Open RAN in the literature and in the O-RAN specifications, as discussed next. These include disaggregation; intelligent, data-driven control with the RICs; virtualization; and open interfaces [37].

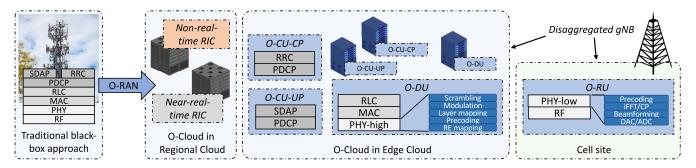


Fig. 2. Evolution of the traditional black-box base station architecture (left) toward a virtualized gNB with a functional split (right, including the CU and DU at the edge, and the RU at the cell site). The functional split distributes the higher layers of the stack in the CU, which features RRC, SDAP, and PDCP. The DU features the RLC, MAC, and the higher part of the physical layer. This is distributed according to the 3GPP 7.2x split, which features frequency-domain functionalities in the DU (including scrambling, modulation, layer mapping, part of precoding, and mapping into physical resource blocks), and the time-domain functionalities in the RU (with precoding, FFT and CP addition/removal, beamforming, and the RF components).

A. Disaggregation

As shown in Fig. 2, RAN disaggregation splits base stations into different functional units, thus effectively embracing and extending the functional disaggregation paradigm proposed by 3GPP for the NR Next Generation Node Bases (gNBs) [38]. The gNB is split into a Central Unit (CU), a Distributed Unit (DU), and a Radio Unit (RU) (called O-CU, O-DU, and O-RU in O-RAN specifications). The CU is further split into two logical components, one for the Control Plane (CP), and one for the User Plane (UP). This logical split allows different functionalities to be deployed at different locations of the network, as well as on different hardware platforms. For example, CUs and DUs can be virtualized on white box servers at the edge (with hardware acceleration for some of the physical layer functionalities) [14], [39], while the RUs are generally implemented on Field Programmable Gate Arrays (FPGAs) and Application-specific Integrated Circuits (ASICs) boards and deployed close to RF antennas.

The O-RAN Alliance has evaluated the different RU/DU split options proposed by the 3GPP, with specific interest in alternatives for physical layer split across the RU and the DU [13]. The selected 7.2x split strikes a balance between simplicity of the RU and the data rates and latency required on the interface between the RU and DU. In split 7.2x, the RU performs time-domain functionalities, with precoding, Fast Fourier Transform (FFT), cyclic prefix addition/removal, and Radio Frequency (RF) operations, which makes the RU inexpensive and easy to deploy. The DU then takes care of the remaining functionalities of the physical layer, and of the Medium Access Control (MAC) and Radio Link Control (RLC) layers [40], [41], [42], including scrambling, modulation, layer mapping, part of precoding, and mapping into physical resource blocks. The operations of these three layers are generally tightly synchronized, as the MAC layer generates Transport Blocks (TBs) for the physical layer using data buffered at the RLC layer. Finally, the CU units (CP and UP) implement the higher layers of the 3GPP stack, i.e., the Radio Resource Control (RRC) layer, which manages the life cycle of the connection [43]; the Service Data Adaptation Protocol (SDAP) layer, which manages the Quality of Service (QoS) of the traffic flows (also known as bearers) [44]; and the Packet Data Convergence Protocol (PDCP) layer, which takes care

of reordering, packet duplication, and encryption for the air interface, among others [45].

B. RAN Intelligent Controllers and Closed-Loop Control

The second innovation is represented by the RICs, which introduce programmable components that can run optimization routines with closed-loop control and orchestrate the RAN. Specifically, the O-RAN vision includes two logical controllers that have an abstract and centralized point of view on the network, thanks to data pipelines that stream and aggregate hundreds of Key Performance Measurements (KPMs) on the status of the network infrastructure (e.g., number of users, load, throughput, resource utilization), as well as additional context information from sources outside of the RAN. The two RICs process this data and leverage AI and ML algorithms to determine and apply control policies and actions on the RAN. Effectively, this introduces data-driven, closed-loop control that can automatically optimize, for example, network and RAN slicing, load balancing, handovers, scheduling policies, among others [12]. The O-RAN Alliance has drafted specifications for a non-real-time RIC, which integrates with the network orchestrator and operates on a time scale longer than 1 s, and a near-real-time RIC, which drives control loops with RAN nodes with a time scale between 10 ms and 1 s. Figure 3 provides an overview of the closed-loop control that the RICs enable throughout the disaggregated O-RAN infrastructure. The figure also shows its real-time extensions which, among others, are considered by the next Generation Research Group (nGRG), a task force founded by the O-RAN Alliance in June 2022 whose goal is to identify and steer the evolution of O-RAN specifications to support 6G and beyond. In the next paragraphs, we will discuss the role of each RIC and related control loops.

Non-real-time RIC and Control Loop: The non-real-time (or non-RT) RIC is a component of the Service Management and Orchestration (SMO) framework, as illustrated in Fig. 4, and complements the near-RT RIC for intelligent RAN operation and optimization on a time scale larger than 1 second [46], [47], [48]. Using the non-real-time control loop, the non-RT RIC provides guidance, enrichment information, and management of ML models for the near-RT RIC [49]. Additionally, the non-RT RIC can influence SMO operations, which gives

Control and learning objective	Scale (devices)	Input data	Timescale	Architecture	Challenges and limitations	
Policies, models, slicing	> 1000	Infrastructure KPMs	Non-real-time > 1 s	Non-real-time RIC	Orchestration of large-scale deployments	Support
User Session Management e.g., load balancing, handover	> 100	CU KPMs e.g., number of sessions, PDCP traffic	Near-real-time 10-1000 ms	A1 gNB Near-real-time E2 CU	Process streams from multiple CUs and sessions	ed by O-
Medium Access Management e.g., scheduling policy, RAN slicing	> 100	MAC KPMs e.g., PRB utilization, buffering	Near-real-time 10-1000 ms	RIC CZ F1	Small time scales, control many DUs/UEs	RAN
Radio Management e.g., scheduling, beamforming	~10	MAC/PHY KPMs e.g., PRB utilization, channel estimation	Real-time < 10 ms	Mobile devices Open FH	Custom real-time loops not supported	For furth
Device DL/UL Management e.g., modulation	1	I/Q samples	Real-time < 1 ms	RU	Device- and RU-level standardization	າer study

Fig. 3. Closed-loop control enabled by the O-RAN architecture, and possible extensions, adapted from [12]. The control loops are represented by the dashed arrows over the architectural diagram.

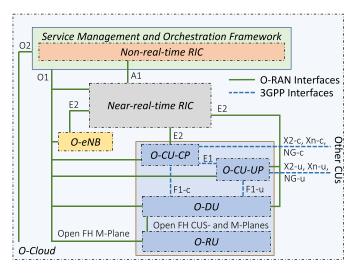


Fig. 4. O-RAN architecture, with components and interfaces from O-RAN and 3GPP. O-RAN interfaces are drawn as solid lines, 3GPP ones as dashed lines.

the non-RT RIC the ability to indirectly govern all the components of the O-RAN architecture connected to the SMO, thus making decisions and applying policies that influence thousands of devices. This presents scalability challenges, as shown in Fig. 3, which need to be addressed through efficient process and software design. Further details on the non-RT RIC and SMO will be given in Section IV.

Near-real-time RIC and Control Loop: The near-real-time (or near-RT) RIC is deployed at the edge of the network and operates control loops with a periodicity between 10 ms and 1 s [50]. As shown in Fig. 3 and Fig. 4, the near-RT RIC interacts with DUs and CUs in the RAN, as well as with legacy O-RAN-compliant LTE evolved Node Bases (eNBs) [51]. The near-RT RIC is usually associated to multiple RAN nodes, thus the near-RT closed-loop control can affect the QoS of hundreds or thousands of User Equipments (UEs).

The near-RT RIC consists of multiple applications supporting custom logic, called xApps, and of the services that are required to support the execution of the xApps. An xApp is a microservice that can be used to perform radio resource management through specific interfaces and service models. It receives data from the RAN (e.g., user, cell, or slice KPMs, as

shown in Fig. 3) and (if necessary) computes and sends back control actions. To support xApps, the near-RT RIC includes (i) a database containing information on the RAN (e.g., list of connected RAN nodes, users, etc.) and serving as a shared data layer among xApps; (ii) messaging infrastructure across the different components of the platform, also supporting the subscription of RAN elements to xApps; (iii) terminations for open interfaces and Application Programming Interfaces (APIs); and (iv) conflict resolution mechanisms to orchestrate control of the same RAN function by multiple xApps. We will further discuss characteristics and functionalities of the xApps in Section III.

Future Extensions to Real-Time Control Loops: Figure 3 also includes loops that operate in the real-time domain, i.e., below 10 ms, for radio resource management at the RAN node level, or even below 1 ms, for device management and optimization. Typical examples of real-time control include scheduling, beam management, and feedback-less detection of physical layer parameters (e.g., modulation and coding scheme, interference recognition) [9]. These loops, which have a limited scale in terms of devices being optimized, are not part of the current O-RAN architecture, but are mentioned in some specifications [49] as for further study.

C. Virtualization

The third principle of the O-RAN architecture is the introduction of additional components for the management and optimization of the network infrastructure and operations, spanning from edge systems to virtualization platforms. According to [37], all the components of the O-RAN architecture shown in Fig. 4 can be deployed on a hybrid cloud computing platform called O-Cloud. Specifically, the O-Cloud is a set of computing resources and virtualization infrastructure that are pooled together in one or multiple physical datacenters. This platform combines physical nodes, software components (e.g., the operating system, virtual machine hypervisors, etc.), and management and orchestration functionalities [52], and specializes the virtualization paradigm for O-RAN [53]. It enables (i) decoupling between hardware and software components; (ii) harmonization and definition of the hardware capabilities for the O-RAN infrastructure; (iii) sharing of the hardware among different tenants; and (iv) automated deployment and instantiation of RAN functionalities.

The O-RAN Alliance Working Group (WG) 6 is also developing hardware acceleration abstractions (called Acceleration Abstraction Layers (AALs)) that define common APIs between dedicated hardware-based logical processors and the O-RAN softwarized infrastructure, e.g., for channel coding/decoding and Forward Error Correction (FEC) [54], [55]. These efforts also reflect into commercial hardware-accelerated, virtualized RAN implementations that can support the requirements of 3GPP NR use cases (e.g., Ultra Reliable and Low Latency Communications (URLLC) flows [56]) also on commercial hardware (e.g., the NVIDIA Aerial platform [57], NEC Nuberu [58], and [59] from Intel). The authors of [60] discuss FPGA-based acceleration of the physical layer decoding with a prototype based on OpenAirInterface.

In parallel, WG 7 is defining the characteristics that white box hardware needs to satisfy to implement an O-RAN-compliant piece of equipment, e.g., indoor picocells, outdoor microcells and macrocells (all at sub-6 GHz and mmWaves), integrated access and backhaul nodes, and fronthaul gateways. These cover different architectural elements from Fig. 2, including the RAN nodes (CU, DU, RU) and enablers of the fronthaul interface. The specifications clarify the functional parameters corresponding to the scenarios of interest (e.g., frequency bands, bandwidth, inter-site distance, MIMO configurations), and the hardware characteristics (e.g., accelerators, compute, connectivity) of the nodes.

The virtualization for the RAN components and of the O-RAN compute elements is expected to introduce savings and optimization of the power consumption related to the RAN. Virtualization makes it possible to easily and dynamically scale up or scale down the compute resources required to support user requirements, thus limiting the power consumption to the actual network functions that are needed [58], [61]. In this sense, the closed-loop control capabilities described above, together with the virtualization in the RAN, also enable more refined and dynamic sleep cycles for the base stations and the RF components [62], [63], which generally are the cause of most of the power consumption in cellular networks [64].

D. Open Interfaces

Finally, the O-RAN Alliance has introduced technical specifications that describe open interfaces connecting a number of different components of the O-RAN architecture. Figure 4 reports the new, open interfaces defined by O-RAN, as well as the intra-RAN interfaces from the 3GPP specifications. The latter is a partial enabler of the gNB disaggregated architecture, which, however, is complemented by the O-RAN Open Fronthaul between the DU and the RU. The O-RAN interfaces, instead, help overcoming the traditional RAN black box approach, as they expose data analytics and telemetry to the RICs, and enable different kinds of control and automation actions, from RAN control to virtualization and deployment optimization.

Without O-RAN, radio resource management and virtual/physical network functions optimization would be closed and inflexible, i.e., the operators would not have the same level of access to the equipment in their RAN, or it would be performed through a custom, piecemeal approach. Standardization of these interfaces is thus a key step toward breaking the vendor lock-in in the RAN, e.g., allowing a near-RT RIC of one vendor to interact with the base stations of another vendor, or again enabling the interoperability of CUs, DUs and RUs from different manufacturers. This also fosters market competitiveness, innovation, faster update/upgrade cycles, and eases the design and introduction of new softwarized components in the RAN ecosystem [12]. The O-RAN Alliance will need to liason with an Standard-Development Organization (SDO), e.g., the 3GPP or the International Telecommunication Union (ITU) Working Group 5 to properly standardize the interfaces currently being developed, as we discuss in Section IX.

Among the O-RAN-specific interfaces, the E2 interface connects the near-RT RIC to the RAN nodes. E2 enables the near-real-time loops shown in Fig. 3 through the streaming of telemetry from the RAN and the feedback with control from the near-RT RIC. The near-RT RIC is connected to the non-RT RIC through the A1 interface, which enables a non-real-time control loop and the deployment of policy, guidance, and intelligent models in the near-RT RIC. The non-RT RIC also terminates the O1 interface, which connects to every other RAN component for management and orchestration of network functionalities. Finally, the non-RT RIC and the SMO also connect to the O-RAN O-Cloud through the O2 interface, and the O-RAN Fronthaul interface connects DUs and RUs. The O-RAN Alliance has also defined a set of precise and well-defined tests to promote interoperability across different interface implementations, with an initial focus on the fronthaul interface and E2. We will provide details on each interface in Section V.

Thanks to the open interfaces, the O-RAN architecture described in Fig. 4 can be deployed by selecting different network locations (cloud, edge, cell sites) for different pieces of equipment, with multiple configurations described in [14]. An example of deployment (i.e., *Scenario B* [14]) is shown in Fig. 2, with the RICs deployed in the cloud, the CUs and DUs at the edge, and the RU cell sites. Other deployment strategies, with the RICs and the RAN nodes co-located are also possible, for example to support local and private 5G networks.

III. NEAR-RT RIC, XAPPS, AND CONTROL OF THE RAN

The near-RT RIC is the core of the control and optimization of the RAN, thanks to the capabilities offered by the E2 interface. In this section, we will discuss the functionalities of the near-RT RIC and the near-RT RIC implementations available in the open source domain.

As discussed in Section II, the near-RT RIC platform hosts the terminations of three interfaces (O1, A1, and E2), the xApps, and the components required to execute and manage the xApps. The O-RAN specifications describe the requirements and functionalities of the different components of the

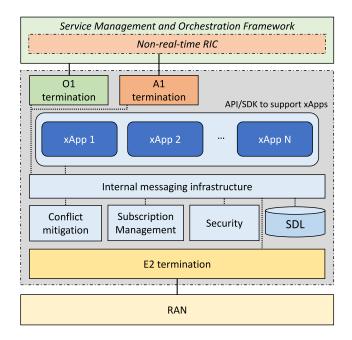


Fig. 5. Near-RT RIC architecture. The near-RT RIC connects to the RAN through the E2 interface, at the bottom of the figure (yellow), and to the non-RT RIC/SMO through the A1 and O1 interfaces, at the top of the figures (orange and green, respectively). The communication among the RIC components (in light blue) is mediated by an internal messaging infrastructure. The near-RT RIC can onboard custom logic as xApps (dark blue).

RIC, so that different specification-compliant implementations can be expected to provide the same set of services, but they do not introduce implementation requirements [50]. However, the O-RAN Alliance, through the O-RAN Software Community (OSC), also provides a reference implementation of a functional near-RT RIC that follows the specifications in [50] and can be used to prototype O-RAN solutions. The OSC near-RT RIC is based on multiple components running as microservices on a Kubernetes cluster [65].

A. Near-RT RIC Internal Components

Figure 5 provides an overview of the architecture of a typical near-RT RIC. The main platform components include:

- Internal messaging infrastructure: The internal messaging infrastructure connects xApps, platform services, and interface terminations to each other. The specifications do not mandate any specific technology (the OSC uses a custom library called RIC Message Router (RMR) [66]), but they list the requirements and functionalities this sub-system needs to provide. The internal messaging infrastructure needs to support registration, discovery, and deletion of endpoints (i.e., internal RIC components and xApps), and provides APIs for sending and receiving messages, either through point-to-point communications or publish/subscribe mechanisms. It also provides routing and robustness to avoid internal data loss;
- Conflict mitigation: This component addresses possible conflicts emerging among different xApps. This is required because different, independent xApps may apply conflicting configurations while trying to achieve independent optimization goals, eventually resulting in performance degradation. The domain of the conflict may be a user, a bearer, or a cell, and can be related to

any control action performed by the RIC. The O-RAN specifications highlight three different classes of conflicts. Direct conflicts can be directly detected by this internal component. For example, multiple xApps can apply different settings for the same parameters in the same control target, or xApps can request more resources than those available. This can be solved by the conflict mitigation component that decides which xApp prevails or limits the scope of a control action (pre-action resolution). Indirect and implicit conflicts, instead, cannot be observed directly and may or may not depend on the relationships among different xApps. For example, configurations that optimize the performance of certain classes of users may degrade others in non-obvious ways. These conflicts may be detected and mitigated through post-action verification, i.e., by monitoring the performance of the system after the application of different control policies. Overall, conflict mitigation is a key component of the RIC but, at the time of this writing, it is not included in the OSC near-RT RIC;

- Subscription manager: The subscription management functionality allows xApps to connect to functions exposed over the E2 interface. It also controls the access that individual xApps can have to E2 messages, and can merge multiple, identical subscription requests to the same E2 node into a single one;
- Security sub-system: According to [50], this component
 has the high-level goal to prevent malicious xApps from
 leaking sensitive RAN data or from affecting the RAN
 performance. The details of this component are still left
 for further studies:
- Network Information Base (NIB) Database and Shared Data Layer API: The RAN NIB (R-NIB) database stores information on the E2 nodes, and the UE-NIB contains entries for the UEs and their identity. The UE identity (i.e., the UE-ID) is a key and sensitive piece of information in the RIC, as it allows UE-specific control, but at the same time it can expose sensitive information on the users. The UE-NIB makes it possible to track and correlate the identity of the same user in different E2 nodes. The database can be queried by the different components of the RIC platform (including the xApps) through the Shared Data Layer (SDL) APIs. The OSC RIC provides an implementation of a SDL library that can be compiled inside xApps, as well as a Redis-based database [67];
- xApp management: The near-RT RIC features services and APIs for the automated life-cycle management of the xApps, from onboarding, to deployment and termination (triggered by the SMO), as well as tracing and logging for Fault, Configuration, Accounting, Performance, Security (FCAPS). In the OSC RIC, this is done through wrappers on the Kubernetes infrastructure.

B. Near-RT RIC xApps

The main components of the near-RT RIC are the xApps. As previously discussed, an xApp is a plug-and-play component that implements custom logic, for example for RAN

data analysis and RAN control. xApps can receive data and telemetry from the RAN and send back control using the E2 interface, as we described in Section V-A.

According to the O-RAN specifications [50], an xApp is defined by a descriptor and by the xApp software image (i.e., the set of files needed to deploy the fully-functional xApp). The xApp descriptor (e.g., a YAML or JSON file) includes information on parameters needed to manage the xApp, such as, for example, autoscaling policies, deployment, deletion, and upgrade information. Additionally, it can describe the data types consumed by the xApp as well as its control capabilities. Specifically, in the OSC RIC, the xApp is defined by a Docker image that can be deployed on a Kubernetes infrastructure by applying the descriptor schema, which is a file that specifies the attributes of the container.

At the time of this writing, the O-RAN specifications only mandate a limited set of APIs that the near-RT RIC platform needs to provide to xApps (including the SDL APIs and the registration/discovery/subscription APIs). The definition of a broader set of APIs into a software development kit, however, would foster the development of xApps that can be seamlessly ported across different near-RT RIC implementations. Efforts in this direction have been promoted by the Telecom Infra Project (TIP) RAN Intelligence and Automation (RIA) subgroup [31], [32].

C. Open Source Near-RT RIC Implementations

Besides the one provided by the OSC, the open source community includes third-party RIC implementations that enrich the Open RAN ecosystem. An example is ColO-RAN, which is an implementation focused on O-RAN experimentation based on the OSC near-RT RIC [68], described in Section X. The SD-RAN project by the Open Networking Foundation (ONF) is developing an open source and cloud-native implementation of O-RAN near-RT RIC, together with xApps to control the RAN, and an Software Development Kit (SDK) to facilitate the design of new xApps [69], [70]. These xApps leverage both standard-compliant E2 Service Models (E2SMs), and custom service models developed by the SD-RAN community. The microservices of this RIC, which is based on the Open Networking Operating System (ONOS) controller, include xApp subscription services, network- and user-based information services, distributed data store services for high availability and operator services. Because of this open implementation, the role of the ONF was key in developing RIC specifications, as well as a host of use cases for the O-RAN ecosystem (e.g., for handover and load balancing), and plans to test them on equipment from telecom vendors. FlexRIC, instead, provides a monolithic near-RT RIC and a RAN agent to interface the OpenAirInterface radio stack with the RIC [71]. It includes Service Models (SMs) for monitoring and slicing programmability use cases, and an SDK to build specialized service-oriented controllers. Finally, 5G-EmPOWER is a near-RT RIC for heterogeneous RANs [72]. It includes non-standard-compliant functionalities like mobility management for Wi-Fi and cellular networks, multi-tenant support, and deployment of custom resource allocation schema within network slices.

IV. NON-RT RIC AND ORCHESTRATION FRAMEWORK

The second key element of the O-RAN architecture is the SMO framework. This component is in charge of handling all orchestration, management and automation procedures to monitor and control RAN components. Primarily, the SMO hosts the non-RT RIC and provides a set of interfaces (described in detail in Section V) that support the interaction between the different network components as well as data collection capabilities to facilitate network monitoring and control via AI/ML [46], [49].

The high-level architecture of the SMO is illustrated in Fig. 6. Its building blocks and their main functionalities will be detailed in the remainder of this section. It is worth mentioning that at the time of writing, the O-RAN specifications do not provide strict guidelines regarding the split between SMO and non-RT RIC functionalities. However, the specifications group such functionalities into three distinct sets [46]. The first set (orange-shaded blocks in Fig. 6) identifies those functionalities and interfaces that are anchored to the non-RT RIC. A second set (green-shaded blocks) identifies functionalities anchored outside the non-RT RIC, while the functionalities from the remaining set (yellow-shaded blocks) are either not yet anchored to any specific SMO component or they span multiple components. Similarly to Fig. 5, the non-RT RIC architecture embeds implementation-specific interfaces that interconnect and regulate the interactions between functionalities and components within the non-RT RIC and the SMO domains. This infrastructure is depicted as the *internal* messaging infrastructure in Fig. 6.

The goal of the next sections is to describe these functionalities and interfaces, as well as to highlight their relevance to O-RAN systems and operations.

A. Non-Real-Time RIC

The non-RT RIC is one of the core components of the O-RAN architecture. Similarly to the near-RT RIC, it enables closed-loop control of the RAN with timescales larger than 1 s. Moreover, it also supports the execution of third-party applications, i.e., the *rApps*, which are used to provide value-added services to support and facilitate RAN optimization and operations, including policy guidance, enrichment information, configuration management and data analytics.

As shown in Fig. 6, the non-RT RIC hosts the R1 termination, which interfaces rApps with the non-RT RIC. This allows them to obtain access to data management and exposure services, AI/ML functionalities, as well as A1, O1 and O2 interfaces through the internal messaging infrastructure. It is worth mentioning that although rApps can support the same control functionalities provided by xApps (e.g., traffic steering, scheduling control, handover management) at larger timescales, they have been designed to derive control policies that operate at a higher level and affect a larger number of users and network nodes. Relevant examples of rApps

¹It is worth mentioning that although SMO functionalities are usually referred to as network-wide orchestration and management procedures (e.g., spanning both core and RAN portions of the network), the O-RAN specifications describe SMO operations and functionalities pertaining to RAN components only.

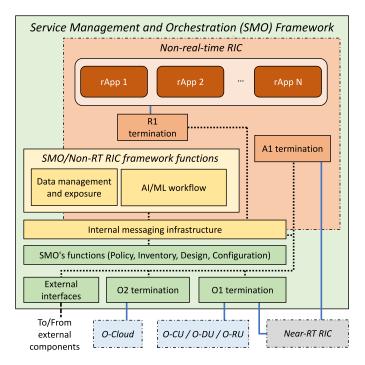


Fig. 6. Non-RT RIC and SMO architecture. The SMO functionalities (in green) enable connectivity to the O-Cloud (through the O2 interface) and the other RAN components (through O1) for management and orchestration. The non-RT RIC features custom logic (rApps, in red), and a termination of the A1 interface to the near-RT RIC (orange). Shared functionalities between the non-RT RIC and the SMO are in yellow.

for non-RT RAN control applications include frequency and interference management, RAN sharing, performance diagnostics, end-to-end Service Level Agreement (SLA) assurance and network slicing [48].

To provide a flexible architecture in which the behavior of each and every network component and functionality can be adjusted in real time to meet the intents and goals of the operators, the non-RT RIC offers the following two high-level management and orchestration services [47]: (i) intent-based network management; and (ii) intelligent orchestration.

Intent-based Network Management: This functionality allows operators to specify their intent via a high-level language through a human-machine interface (it is expected that intents will follow practices used in currently available SMOs, and formalized as YAML or XML configuration files). Intents are then automatically parsed by the non-RT RIC, which determines the policies and the set of rApps and xApps that need to be deployed and executed to satisfy them. In this context, it is worth mentioning that the OSC has already included a human-machine interface (i.e., the Intent Interface [73]) within the design of the non-RT RIC that enables operators to specify intents at a high level.

Intelligence Orchestration: Indeed, the O-RAN architecture enables and facilitates the development, deployment, execution and maintenance of network intelligence. However, this inevitably makes network control more complex due to the increasing number of xApps and rApps that execute at different RICs and locations of the network. This calls for solutions that are capable of coordinating and orchestrating

these applications. Specifically, the non-RT RIC is in charge of orchestrating network intelligence [49] to make sure that selected xApps and rApps are: (i) well-suited to satisfy operator intents and meet their requirements; (ii) instantiated at the appropriate RIC location to ensure control over the specific RAN elements specified in the intent; (iii) fed with relevant data; and (iv) robust enough not to generate conflicts due to multiple applications controlling the same functionalities and/or parameters simultaneously. For example, if the operator has instantiated multiple network slices and wants to control and optimize scheduling policies for each slice in near real time for a selected set of base stations close to a landmark of interest, the non-RT RIC must be able to determine automatically that only xApps executing at a specific near-RT RIC can satisfy the timing requirement. Moreover, the non-RT RIC must select only those xApps that are able to control scheduling decisions, and eventually dispatch them to the near-RT RICs that control the base stations deployed in the area of interest, and are thus capable of generating data to be fed to the xApps.

B. Other SMO/Non-RT RIC Functionalities

In this section, we describe the functionalities and architectural components that can reside both at the SMO and at the non-RT RIC [46], [47].

The internal messaging infrastructure is a composite of several SMO functions that allow all components within the SMO (even those included in the non-RT RIC) to access and utilize interfaces, data and functionalities offered by both the SMO and the non-RT RIC. For example, all interface terminations are tied to interface-specific functions included in the internal messaging infrastructure that are designed to facilitate the exchange of messages between terminations. In this way, policies computed by rApps can reach the non-RT RIC through the R1 termination, and eventually reach the near-RT RIC through the A1 interface.

Data Management and Exposure Services: The O-RAN specifications also include data management and exposure services pertaining to the SMO/non-RT RIC framework. To this purpose, O-RAN follows a consumer/producer protocol in which data producers in the SMO/non-RT RIC can advertise and publish data (e.g., performance reports or AI-based prediction of KPMs and network load). On the other hand, data consumers (e.g., rApps that determine high-level control policies) can discover, subscribe, receive and consume relevant data types from a selected number of nodes in the SMO/non-RT RIC domain. In order to fully support AI/ML solutions, the SMO/non-RT RIC can also perform collection of all data being produced, as well as relevant AI/ML pre-processing operations involving data analytics (e.g., correlation analysis), labeling, and normalization.

AI/ML Workflow: Another important capability offered by the non-RT RIC is the possibility to oversee the entire AI life cycle, and cover all aspects of AI/ML development, including data collection, training, validation, deployment and execution. This AI/ML workflow, which will be detailed in Section VI, is illustrated in Fig. 13.

C. SMO Framework and Open Source SMO Implementations

Besides hosting the non-RT RIC, the SMO also offers additional functionalities and interfaces, summarized in Fig. 6. These include management and interactions with the O-Cloud via the O2 interface (see Section V-E) as well as other O-RAN components via the O1 interface. The SMO takes care of FCAPS management procedures, as well as service and resource inventory, topology and network configuration, in addition to policy management for network orchestration services.

Although several members of the O-RAN Alliance have announced the development and availability of proprietary SMOs compliant with the latest O-RAN releases, these SMOs are closed solutions not open to the general public and whose implementation details and offered functionalities are not generally available. For this reason, we focus on two open-source solutions—with publicly-available code and functionalities—that are currently being integrated with the O-RAN architecture [14]: Open Network Automation Platform (ONAP) [74], [75] and Open Source Management and Orchestration (OSM) [76].

Both ONAP and OSM are comprehensive platforms that enable automation and orchestration in virtualized and softwarized networks. ONAP is one of the main projects being developed and maintained by the Linux Foundation, while OSM is hosted by European Telecommunications Standards Institute (ETSI) and follows ETSI Virtual Network Function (VNF) standard specifications [76]. Being maintained by the Linux Foundation, ONAP provides native integration with other major projects such as Kubernetes, Akraino, Acumos and OpenDaylight [77]. This makes ONAP quite a complex environment compared to OSM which, instead, offers similar services in a more lightweight framework. Despite its complexity, ONAP has been already used by the OSC (which is a joint effort between the Linux Foundation and the O-RAN Alliance) to release an open-source SMO platform [73] that integrates messaging and functionalities necessary to monitor the state of the different O-RAN components as well as to capture and handle messages and events (e.g., fault and heartbeat). It is also worth mentioning that in May 2021 ETSI signed a cooperation agreement with the O-RAN Alliance [78] which is currently facilitating and advancing the integration between OSM and the O-RAN architecture [79].

V. THE O-RAN OPEN INTERFACES

As discussed in Section II, the control loops supported by near-RT and non-RT RICs are enabled by a set of interfaces defined by O-RAN Alliance specifications. Each interface enables services (e.g., reporting of telemetry from the RAN) through the combination of well-defined procedures (e.g., the subscription and indication procedures for E2). Procedures involve the exchange of messages between the endpoints of an interface (e.g., the indication message for E2). This section reviews the logical abstractions and procedures that define such interfaces, providing insights on their role in the Open RAN ecosystem. Specifically, the E2 interface is described in Section V-A, the O1 interface in Section V-B, the A1 interface

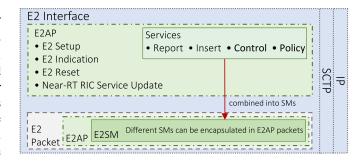


Fig. 7. Representation of an O-RAN E2AP packet (bottom left), which includes an E2SM payload (top left). The E2 payload is then encapsulated in SCTP and IP headers (right). The top part of the figure also summarizes the services provided by the E2 interface.

in Section V-C, and the fronthaul interface in Section V-D. Finally, Section V-E reviews the remaining O-RAN and 3GPP interfaces.

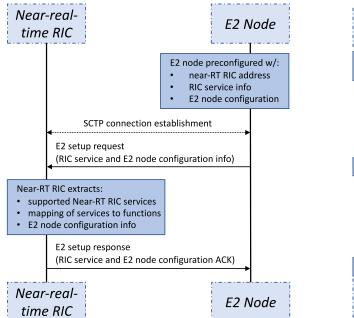
A. E2 Interface

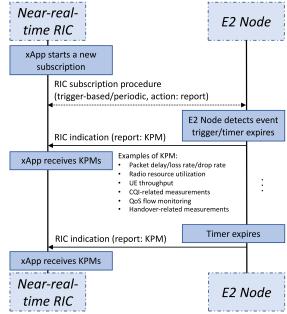
The E2 interface is an open interface between two endpoints, i.e., the near-RT RIC and the so-called E2 nodes, i.e., DUs, CUs, and O-RAN-compliant LTE eNBs [80]. The E2 allows the RIC to control radio resource management and other functionalities of the E2 nodes. Moreover, this interface also enables the collection of metrics from the RAN to the near-RT RIC, either periodically or after pre-defined trigger events. Both control and data collection procedures can pertain to one or more cells, slices, QoS classes, or specific UEs.

To support the above operations, the O-RAN Alliance uses a variety of unique identifiers. Specifically, O-RAN uses identifiers based on 3GPP specifications for the gNB, slice, and QoS class [81]. Regarding specific UEs, the O-RAN Alliance has instead introduced a common user identifier (i.e., the UE-ID) across its specifications. This provides a consistent and uniform user identity across the system without exposing sensitive information related to the user.

The E2 interface has been logically organized in two protocols: E2 Application Protocol (AP) and E2 SM. The E2 AP [82] is a basic procedural protocol that coordinates how the near-RT RIC and the E2 nodes communicate with each other, and provides a basic set of services, as shown in Fig. 7. E2AP messages can embed different E2 SMs [83], which implement specific functionalities (i.e., the reporting of RAN metrics or the control of RAN parameters). The E2 interface runs on top of the SCTP protocol [84].

Each E2 node exposes a number of *RAN functions*, i.e., the services and capabilities it supports. For example, DUs from different vendors may expose different control knobs depending on which parameters and functionalities can be tuned, as well as their capability in collecting and reporting different performance metrics. By using publish-subscribe mechanics, E2 nodes can publish their data and the xApps on the near-RT RIC can subscribe to one or more of these RAN functions through the E2 interface. This makes it possible to clearly separate the capabilities of each node and to define how the xApps interact with the RAN.





- the near-RT RIC.
- (a) Procedure for the setup of an E2 session between the near-RT RIC and (b) Procedures related to the streaming of KPMs from the E2 node to the an E2 node. The procedure is initiated by the E2 node which interacts with near-RT RIC. The subscription procedure is started by an xApp on the near-RT RIC, which then receives the reports.

Fig. 8. Procedures for E2 setup and E2SM KPM. The vertical lines represent the temporal evolution of the process, while horizontal lines are the messages exchanged by the near-RT RIC and the E2 node.

At the lowest level, the E2AP handles interface management (setup, reset, reporting of errors for the E2 interface itself) and near-RT RIC service updates (i.e., the exchange of the list of the RAN functions supported by the E2 node). As an example, the E2 setup procedure is shown in Fig. 8(a). At first, the SCTP connection is established between the near-RT RIC and the E2 node (which is aware of the IP address and port of the E2 termination of the near-RT RIC). Then, the E2 node transmits an E2 setup request, in which it lists the RAN functions and configuration it supports, together with the identifiers for the node. The near-RT RIC processes this information and replies with an E2 setup response.

After the connection is established, the E2AP provides four services which can be combined in different ways to implement an E2SM [83]. These services, also shown in Fig. 7, are [82]:

- E2 Report: The report service involves E2 RIC Indication messages that contain data and telemetry from an E2 node. The E2 report service is activated upon subscription from an xApp to a function offered by the E2 node. During the subscription negotiation, the xApp in the near-RT RIC can specify trigger events or the periodicity with which the E2 node should send report messages. Based on this periodicity, a timer is set in the E2 node and a report is sent whenever the timer expires. The RIC Indication message is of type report.
- E2 Insert: Similarly, the insert service involves messages sent from an E2 node to an xApp in the near-RT RIC to notify the xApp about a specific event in the E2 node (e.g., a UE signaling the possibility to perform a handover). It is activated upon subscription from an xApp

- and involves a RIC Indication message (of type *insert*). In this case, the trigger is associated to a RAN radio resource management procedure which is suspended when the insert message is sent. A wait timer is also started, and, if the RIC does not reply before the timer expires, the procedure in the E2 node can be resumed or definitely halted.
- E2 Control: The control service can be autonomously initiated by the RIC, or it can be the consequence of the reception of an insert message at the near-RT RIC. This service is based on a procedure with two messages, a RIC Control Request from the RIC to the E2 node, and a RIC Control Acknowledge in the opposite direction. The control services can influence parameters exposed by the RAN functions of the E2 node.
- E2 Policy: This service involves a subscription procedure that specifies (i) an event trigger; and (ii) a policy that the E2 node should autonomously follow to perform radio resource management.

These services are then combined to create a service model. The service model message is inserted as payload in one of the E2AP messages, as shown in Fig. 7. The actual content is encoded using ASN.1 notation, i.e., through well-defined types for numbers and key-value pairs [85]. We provide examples of an E2 Subscription Request message and of an E2 Indication message (of type report) in Appendices A and B, respectively.

E2 Service Models: At the time of this writing, the O-RAN Alliance WG3 has developed four service models: (i) the E2SM KPM [86]; (ii) the E2SM Network Interfaces (NI) [87]; (iii) the E2SM Cell Configuration and Control (CCC) [88]; and (iv) the E2SM RAN Control (RC) [89].

The E2SM KPM [86] reports performance metrics from the RAN, using the E2 report service. The procedures associated to the KPM service model are shown in Fig. 8(b). During the E2 setup procedures, the E2 node advertises the metrics it can expose. An xApp in the near-RT RIC can then send a subscription message specifying which KPMs are of interest, and whether the reporting is periodic or trigger-based. Finally, the E2 node uses E2 Indication messages of type report to stream the selected KPMs. Different KPM messages are generated from different E2 nodes, i.e., the specifications defines performance metrics containers with different fields to be populated for DU, CU-UP, and CU-CP. Additionally, the recent Version 2 of the KPM specifications [86] has introduced cellspecific and user-specific performance containers, which are based on the metrics defined in 3GPP technical specifications for LTE and NR [90], [91].

The E2SM NI [87] is used to take the messages received by the E2 node on specific network interfaces and forward them to the near-RT RIC domain via E2 *report* messages over the E2 interface. The E2 node advertises which interfaces it supports during the subscription procedure, and they include X2 (which connects LTE eNBs), Xn (which connects different NR gNBs), and F1, which connects DUs and CUs).

The E2SM CCC [88], introduced in October 2022, performs control and re-configuration of the E2 nodes at a cell or node level, e.g., for the bandwidth part configuration. It primarily leverages E2 report and control messages. The current version of the technical specification for CCC [88] includes control targets such as the selection of X2 and Xn neighbors, RAN slicing, and parameters related to the bandwidth part and the synchronization signals of a cell.

Closed-loop control with E2SM RAN Control (RC): One of the goals of the xApps is to optimize the radio resource management in the E2 nodes. The E2SM RC, introduced in July 2021, implements control functionalities through E2 control services. Compared to E2SM CCC, it focuses on more granular control (up to the UE or bearer level). The first release of the specifications for E2SM includes a broad set of control domains:

- radio bearer control, to modify parameters for the bearers of an E2 node or of a specific UE, for example related to QoS parameters, bearer admission control, split bearer and PDCP duplication control;
- radio resource allocation control, to modify, among others, discontinuous reception, semi-persistent scheduling, for slice-level Physical Resource Blocks (PRBs) for an E2 node, a cell, a slice, a UE, or QoS classes;
- connected mode mobility control, to initiate a mobility procedure for UEs in RRC connected state, i.e., a handover to a specific cell, or a conditional handover² to a set of cells;
- radio access control, to set parameters for random access backoff, UE admission to a cell, among others;
- Dual Connectivity (DC) control, to configure and trigger the handover of a UE to selected target secondary cells, secondary cell updates, or DC release;

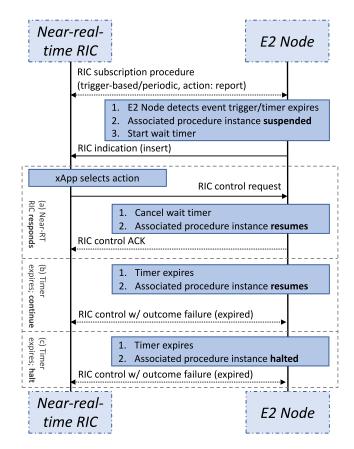


Fig. 9. E2 insert service with subsequent E2 control service response. The vertical lines represent the temporal evolution of the process, while horizontal lines are the messages exchanged by the near-RT RIC and the E2 node.

- Carrier Aggregation (CA) control, to initiate CA and modify the component carriers for a specific UE;
- idle mobility control, to modify functions for mobility procedures of UEs in a RRC idle state, including cell re-selection priorities and idle timers.

E2SM RC also provides capabilities for UE identification and UE information reporting. The control actions and policies that E2SM specifies relate to specific parameters standardized by the 3GPP, i.e., the control action will usually carry the value for a 3GPP Information Element (IE) defined in the E2SM specifications [89]. For example, the handover control IE includes (among other things) a target cell ID for the handover encoded as a 3GPP *Target Cell Global ID* IE in [92, Sec. 9.2.3.25].

To effectively implement control actions and/or enforce policies, the E2SM leverages a combination of the E2 services described in Section V-A. Figure 9 shows an example of the E2 messages that an E2 node and an xApp in the near-RT RIC can exchange during a typical control procedure. First, the xApp subscribes to the E2 node, which needs to expose a specific RAN function called *RAN control*. During the subscription phase, the xApp can specify a triggering event or a timer. Then, if the timer expires or the condition defined in the trigger is verified, the E2 node sends an E2 Insert message to the RIC. For example, this may happen when a metric that the gNB monitors exceeds a certain threshold, e.g., when the Reference Signal Received Power (RSRP) for a neighboring cell (or target cell) becomes better than that of the current cell

²The conditional handover is a feature in 3GPP NR Release 16 that mandates conditions for which the UE should handover to a target cell, but does trigger an immediate handover [43].



Fig. 10. O-RAN O1 interface and Management Services payload (left), which is encapsulated in HTTPS/TCP/IP packets (right).

plus an offset (also called *A3 event* in the 3GPP specifications [43]). In this case, the E2 node sends a handover control request insert message (in which it can also specify the target cell defined by the E2 node itself), suspends the radio resource management procedure, and starts a timer. Three different outcomes can follow this event. The near-RT RIC can reply with an E2 control message that either denies the control request, or accepts it and replies with a control action (for example, a target cell ID selected by the xApp). Alternatively, if the E2 node timer expires before the reception of a message from the RIC, the procedure may continue autonomously, or it may be halted, according to the specific procedure and configuration of the E2 node.

The control action sent by the xApp can also be asynchronous, i.e., it does not depend on the reception of an insert message from the E2 node. Additionally, the E2SM can also be used to specify policies, i.e., to alter pre-defined behaviors in the E2 nodes. Policies can be of two different types: (i) control policies that allow the E2 node to perform radio resource management actions without the interaction with the RIC, when certain conditions are satisfied; and (ii) offset policies, which change 3GPP- or vendor-defined thresholds by adding or removing offsets and thus modifying how the E2 node performs specific functions.

B. O1 Interface

Besides E2, the other interface that connects O-RAN specific components with RAN nodes is the O1 interface [93]. In general, O-RAN-managed elements (including the near-RT RIC, RAN nodes) are connected via O1 to the SMO and the non-RT RIC. The O1, thus, is an open interface which adopts and extends standardized practices for *operations and maintenance*.

The O1 interface supports *Management Services (MnS)*, which include the management of the life-cycle of O-RAN components (from startup and configuration to fault tolerance and heartbeat services [94]), performance assurance and trace collection through Key Performance Indicators (KPIs) reports, and software and file management (see Fig. 10). The O1 interface generally connects one MnS provider (i.e., generally the node managed by the SMO) to one MnS consumer (i.e., the SMO).

The *Provisioning Management Services* allow the SMO to push configurations to the managed nodes, and the reporting of external configuration updates from managed nodes to the SMO. For this, O1 uses a combination of REST/HTTPS APIs and NETCONF [95], which is a protocol standardized by the

Internet Engineering Task Force (IETF) for the life-cycle management of networked functions. The supported provisioning services include those defined in 3GPP technical specifications [96], [97]. An additional *Fault Supervision MnS* is used to report errors and events to the SMO. It is also based on 3GPP-defined fault events [94], [97], [98], [99], and can be used by the RAN nodes to report errors (through standardized JSON payloads) using REST APIs. For each node, the SMO can also query a list of alarms (i.e., probes that monitor the status of specific elements and components in the node), and, in case, acknowledge or clear them. Finally, the SMO can provision heartbeats on the devices it manages, through the *Heartbeat MnS*, and manage not only virtual but also Physical Network Functions (PNFs). Heartbeat messages are used to monitor the status and availability of services and nodes.

The Performance Assurance MnS can be used to stream (in real time) or report in bulk (through file transfer) performance data to the SMO, to enable, for example, data analytics and data collection for AI/ML. The SMO can select the KPIs (also referred to as counters) to be reported and the frequency of reporting. It relies on use cases and formats for KPI reporting defined by the 3GPP [94], [97], [99], [100] or by the VNF Event Stream (VES) project. The performance metrics are also either based on 3GPP documents [90], vendor specific, or specified by the different WGs of the O-RAN Alliance. For bulk transfer, a file-ready notification is first sent from the MnS provider (e.g., the specific node of the RAN) to the SMO through HTTP APIs. Then, a file transfer through SFTP is performed. The SMO can also download files for which the file-ready notification was received in previous instants. A WebSocket is instead used for real-time streaming, following a handshake. The SMO can also monitor trace-based events through the Trace MnS, e.g., to profile calls, RRC connection establishment or radio link failures.

The O1 interface can also be used to push and/or download files on the nodes managed by the SMO. This enables, for example, software updates, new beamforming configuration files for the RUs, and deployment of ML models and security certificates.

C. Al Interface

The A1 interface connects two O-RAN-specific components, i.e., the non-RT RIC (or SMO) and the near-RT RIC, as shown in Fig. 4 [101]. It allows the non-RT RIC to deploy policy-based guidance for the near-RT RIC (e.g., to set highlevel optimization goals), to manage ML models used, for example, in xApps, and to negotiate and orchestrate the transfer of enrichment information for the near-RT RIC. This is done through well-defined mechanisms that are based on a specific syntax (based on JSON schema) which can express policies and high-level intent. The policies, ML models, and enrichment information can refer to a group of UEs, or even to a specific UE. Notice that, at the time of writing, the A1-based ML model management is still considered for further study [49], [101].

The A1 interface, illustrated in Fig. 11, relies on the A1AP application protocol, whose functionalities are then further

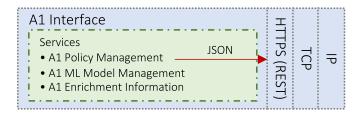


Fig. 11. O-RAN A1 JSON payload (left), which is encapsulated in HTTPS/TCP/IP packets (right). The figure also summarizes the services that can be provided over A1.

specified for each service it supports [102]. The A1AP is based on a 3GPP framework for policy deployment for network functions [81], which combines REST APIs over HTTP for the transfer of JSON objects. For each service, both the non-RT and the near-RT RICs feature a pair of HTTP clients and servers, which are used alternatively for service management and for the actual data transfer and/or notifications.

The A1 Policy management is used by the non-RT RIC to drive the functionalities of the near-RT RIC to achieve highlevel intent for the RAN, as we will discuss in Section IV. This intent is generally defined through QoS or KPI goals for all users or subsets of users (e.g., a slice) and monitored using the reporting functionalities of the O1 interface and feedback over A1.³ The policies are defined by the non-RT RIC and then deployed over A1. The non-RT RIC is also tasked with monitoring and managing the life cycle of the policies, thanks to APIs for deleting, updating, and querying policies in the near-RT RIC.

Each policy is based on specific JSON schema which are grouped according to different policy types [103]. All JSON schema have in common a policy identifier, which is unique for the non-RT RIC, a scope identifier, and one or more policy statements. The scope can be a single UE, a group of UEs, slices, cells, bearers, and application classes. The policy itself is then expressed through a sequence of policy statements, which can cover policy resources (i.e., the conditions for resource usage for a policy) and policy objectives (i.e., the goal of the policy in terms of QoS or KPI targets). The O-RAN technical specification [103] lists several types for policy statements, which depend on specific use cases (e.g., throughput maximization, traffic steering preferences, QoS targets, among others).

Finally, the A1 Enrichment Information (EI) service aims at improving RAN performance by providing information that is generally not available to the RAN itself (e.g., capacity forecasts, information elements from sources outside the RAN, aggregate analytics). The non-RT RIC and the SMO have indeed a global perspective on the network and access to external sources of information, and can relay it to the xApps in the near-RT RIC using A1 EI. The flow of information can also bypass the non-RT RIC, which can instruct the near-RT RIC (over A1) to connect directly to sources of EI.

D. O-RAN Fronthaul

The O-RAN Fronthaul (FH) interface connects a DU to one or multiple RUs inside the same gNB [104], [105]. The

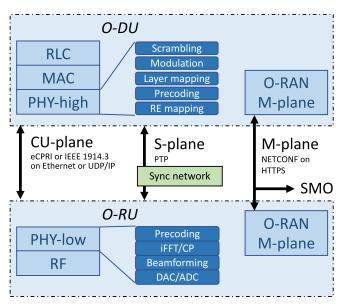


Fig. 12. O-RAN fronthaul interface and fronthaul planes. This interface enables the 7.2x split of the physical layer functionalities across DU and RU, with data and control between the PHY-high and PHY-low transported over the Control and User- (CU-) planes. The S-plane provides synchronization, and the M-plane management and orchestration functionalities.

O-RAN FH interface makes it possible to distribute the physical layer functionalities between the RU and the DU, and to control RU operations from the DU. As discussed in Section II, the O-RAN Alliance has selected a specific configuration (split 7.2x) for the splitting of the physical layer among those proposed by the 3GPP [13]. As shown in Fig. 12, the lower part of the physical layer (low PHY) resides in the RU and performs Orthogonal Frequency Division Multiplexing (OFDM) phase compensation [106], the inverse FFT and Cyclic Prefix (CP) insertion for frequency-to-time conversion in downlink, and FFT and CP removal in uplink. More capable RUs (i.e., category B RUs) can also perform precoding. This functionality is implemented at the DU for less capable RUs (i.e., category A). To complement the DU capabilities, category A RUs need to support low PHY processing for at least 8 data streams. The physical layer in the DU (high PHY) performs scrambling, modulation, layer mapping, and resource element mapping.

According to O-RAN specifications [104], the 7.2x split strikes a balanced trade-off among the simplicity of the interface and of the RU design, the potential for interoperability (fewer parameters to configure than higher layer splits), and the data rate required for fronthaul transport (lower with respect to configurations that split the physical layer even further). The latter can be based on Ethernet or UPD/IP encapsulation, carrying either an eCPRI [107] or an IEEE 1914.3 [108] payload. Note that one DU can support more than one RU, e.g., to serve carriers of the same cells from different RUs, or to process multiple cells with one DU and multiple RUs. To do this, the O-RAN FH specifications foresee an additional component which multiplexes a fronthaul stream to multiple RUs, or the daisy-chaining of RUs. In addition, the O-RAN FH interface has been designed to support reliable, low-latency communications between DUs and RUs with

³This can only notify if the policy is enforced or not.

timing that matches the requirements of URLLC flows. For example, the fronthaul interface includes different modulation compression techniques, to reduce the load on the fronthaul network [109], and fronthaul networks can be designed to support URLLC flows with minimal jitter [110].

The O-RAN FH protocol includes four different functionalities (or *planes*). Besides the user (U-) and control (C-) planes (for transport of data and PHY-layer control commands) [104], the O-RAN FH also features a synchronization plane (S-plane), for timing management among DUs and RUs [104], and a management plane (M-plane), for the configuration of the RU functionalities from the DU itself [105], [111]. In the following, we describe the functionalities of the different FH planes.

C-plane: The C-plane takes care of transferring commands from the high-PHY in the DU to the low-PHY of the RU, including scheduling and beamforming configurations, management of different NR numerologies in different subframes, downlink precoding configuration, and spectrum sharing control. For the latter, the specifications include Licensed-Assisted Access (LAA) procedures, such as the possibility to perform listen-before-talk in the RU and DU pair, and dynamic spectrum sharing operations, with the possibility to specify which PRBs are dedicated to spectrum sharing between LTE and NR.

The C-plane messages are encapsulated in eCPRI or IEEE 1914.3 headers and payloads, with specific fields and commands for the different control procedures. The O-RAN FH specifications also provide details on how specific C-plane directives (e.g., related to the usage of a specific beamforming vector) can be coupled with specific U-plane packets (and thus symbols to be transmitted).

The combination of C-plane and M-plane can be used to configure and manage beamforming capabilities of the RU, a key feature in 5G networks, especially at FR2 [112]. In particular, each RU can support multiple antenna panels, each with multiple Transmitter (TX) or Receiver (RX) arrays [105]. Each array can be mapped to one or more data flows. The fronthaul interface allows the control of amplitude and phase of the radiating elements in the phased arrays of the RU (for beamforming in the time domain), or the selection of digital precoding weights (for beamforming in the frequency domain), with four different beamforming options. With predefinedbeam beamforming, the DU dynamically selects (through the C-plane) time and/or frequency beamforming vectors⁴ that the RU advertises as available at startup (through the M-plane). Alternatively, with attributed-based beamforming, the DU can select beams based on specific attributes, e.g., azimuth and elevation angles. With weight-based beamforming, the DU also specifies the weights for generic time and/or frequency domain beamforming vectors. The last option is channelinformation-based beamforming, in which the DU provides the RU with Channel State Information (CSI) for a specific user and the RU computes the beamforming weights. The O-RAN FH supports a well-defined model for the RU antennas, so that the DU can unambiguously identify antenna

elements, their polarization, position, and orientation of the panel.

U-plane: The main functionality of the U-plane is transferring I/Q samples in the frequency domain between the RU and the DU. Typically, a C-plane message specifies scheduling and beamforming configuration, and is followed by one or more U-plane messages with the I/Q samples to be transmitted in the corresponding transmission opportunities. The U-plane also takes care of timing the transmission of its messages so that they are received at the RU with enough time for processing before transmission. Additionally, the U-plane specifies the digital gain of the I/Q samples, and can compress them for more efficient data transfer.

S-plane: The S-plane takes care of time, frequency, and phase synchronization between the clocks of the DUs and of the RUs. This is key to a correct functioning of a time-and frequency-slotted system distributed across multiple units. Thanks to the shared clock reference, the DU and RU can properly align time and frequency resources for the transmission and reception of the different LTE and NR data and control channels.

The O-RAN FH S-plane can be deployed with different topologies, specified in [104], which differ according to whether a direct DU-RU link exists, or if the two elements are connected through a fabric of Ethernet switches. Additionally, the FH specifications include different synchronization profiles, based on different protocols, such as Physical Layer Frequency Signals (PLFS) or Precision Time Protocol (PTP), which can achieve sub-microsecond time accuracy [114].

M-plane: The O-RAN FH M-plane is a protocol that runs in parallel to the C-, U-, and S-planes, with dedicated endpoints in the DU and RU that establish an IPv4 or IPv6 tunnel [105]. It enables the initialization and the management of the connection between the DU and the RU, and the configuration of the RU. In this context, the specifications foresee two architectural options, i.e., hierarchical, in which the SMO manages the DU and the DU manages the RU, and hybrid, in which the SMO can also interact directly with the RU. The M-plane of the O-RAN FH can thus function as the O1 interface of the RU. As for O1, the management directives are based on NETCONF. Finally, contrary to the C-, U-, and S-planes, the M-plane is end-to-end encrypted through SSH and/or TLS.

The M-plane takes care of several operations related to the life cycle of the RU. First, it manages the start-up, during which the RU establishes the management with the DU and/or the SMO thanks to pre-defined IP addresses or DHCP configurations. Then, it enables software updates, configuration management, performance and fault monitoring, and file management for bulk transfer of data. Among others, the M-plane manages the registration of the RU as PNF, the parameters of the RU-to-DU link (including timing), and the update of beamforming vectors (from the deployment of new codebooks, to the tilting of existing ones, and calibration of the antennas).

Besides specifying the FH interface, the O-RAN Alliance is also developing a set of specifications to characterize the transport and synchronization capabilities of an open fronthaul or crosshaul network that supports the connectivity between DUs and RUs. For example, [115] reviews network-enabled

⁴The combination of time and frequency beamforming enables hybrid beamforming strategies [113].

synchronization, by discussing PTP profiles, support required by the Ethernet substrate, points of failures, among others. Other areas of interest are related to the management of the open fronthaul network, wave-division-multiplexing-based networks, and packet-switched architectures with modern features such as, for example, slicing. In this context, it is worth mentioning that ETSI has adopted the O-RAN Alliance FH user, control, and synchronization plane document [104] as an ETSI standard [116].

E. Other Interfaces

The O2 interface connects the SMO to the O-RAN O-Cloud, enabling the management and provisioning of network functions in a programmatic manner [117]. It allows the definition of an inventory of the facilities controlled by the O-Cloud, monitoring, provisioning, fault tolerance and updates. Being these functionalities common to several Network Function Virtualization (NFV)-related problems, infrastructures and standards, the O-RAN Alliance WG6 considers plausible candidates for implementing this interface a set of well-known standards and open-source solutions, e.g., relevant ETSI NFV standards, 3GPP service-based interfaces, and the Kubernets, Open Stack, and ONAP/OSM projects.

Finally, as shown in Fig. 4, the O-RAN disaggregated architecture also leverages additional interfaces defined by the 3GPP. Notably, the E1 interface connects the CU control and user functions [118]. The F1 interface connects the CU to the DU, with dedicated sub-interfaces for user and control planes [119]. The Xn (X2) interface connects different gNBs (eNBs), for example to perform handovers and to enable dual connectivity [92], [120]. The Uu interface exists between an UE and the gNB [38], and the NG interface connects the gNB to the 5G core, i.e., to the User Plane Function (UPF) for the user plane and the Access and Mobility Management Function (AMF) for the control plane [38].

VI. AI/ML WORKFLOWS

The goal of this section is to provide a detailed overview of the procedures and operational steps that regulate the AI/ML workflow in the O-RAN architecture. In the following, we provide a step-by-step guide on the life cycle of this workflow—from data collection to the actual deployment and execution of network intelligence in real time—that relies on the architectural components described in Sections II, III, and IV, and on the interfaces discussed in Section V.

The AI/ML workflow is being investigated by O-RAN WG2, with its specifications described in [49]. However, not all the procedures, features and functionalities have been finalized yet, with some of them left for further studies. This workflow is composed of six main steps, which we will describe next: (i) data collection and processing, (ii) training; (iii) validation and publishing; (iv) deployment; (v) AI/ML execution and inference; and (vi) continuous operations. The different steps of the end-to-end AI/ML workflow help network designers to avoid common issues in AI/ML for networks—including, for example, limited availability or generality of data for training, improper training, action space definition,

and agent design, limited testing and validation, and lack of infrastructure support for resilient AI/ML workloads.

For the sake of illustration and clarity, in the following we will describe the procedures involved in the AI/ML lifecycle within O-RAN systems. We take as an example the case where an operator aims at developing, training and deploying an xApp that controls RAN slicing policies by adapting them in near real time according to current network load and traffic demand via AI-based algorithms. In this example, each base station hosts three slices, namely a URLLC slice for ultralow latency services, an Enhanced Mobile Broadband (eMBB) slice for high-throughput traffic (e.g., video streaming and file transfer), and a Massive Machine-Type Communications (mMTC) slice to handle traffic generated, for example, by small sensors and Internet of Things (IoT) devices. The goal of the xApp is to control RAN slicing policies by assigning the available PRBs to each slice so that the diverse performance requirements of each slice are satisfied.

Data Collection and Processing: First, data is collected over the O1, A1 and E2 interfaces and stored in large datasets (e.g., *data lake* centralized repositories [121]) where it can be extracted upon request. The comprehensive design of data reporting solutions for these interfaces promotes the collection of large-scale datasets, which can help AI/ML solutions generalize better [122].

Additionally, since different AI/ML solutions might use different KPM types collected over different time periods and with different granularity (e.g., throughput, latency, Modulation and Coding Scheme (MCS), Channel Quality Information (CQI), data demand, jitter, to name a few), the O-RAN specifications consider also a preliminary data preprocessing (or preparation) step. Here, data for both training and online inference is shaped and formatted according to the input size of the specific AI/ML model being considered [123]. This step might include the use of autoencoders for dimensionality reduction [12], [124], as well as well-established AI data processing procedures such as normalization, scaling and reshaping. Such pre-processing avoids dimensionality issues or problems with poorly formatted data fed to the AI/ML agents [123].

Example: with respect to the xApp controlling RAN slicing policies, this step involves collecting data and performance metrics over the O1 interface to generate a training dataset to be used in the next step (i.e., training phase). For example, since the xApp must be able to adapt RAN slicing policies for the different slices according to current data demand and required minimum performance levels, the collected data must include how many PRBs are necessary to transmit the data requested by each user of the three slices, as well as throughput (eMBB), number of transmitted packets (mMTC) and latency (URLLC) measurements.

Training: The O-RAN specifications do not allow the deployment of any untrained data-driven solution [49]: all the AI/ML models are required to go through an offline training phase. This ensures the reliability of the intelligence and avoids typical issues of poorly trained or untrained AI/ML models that might result in outages or inefficiencies in the network, e.g., inaccurate predictions, classifications, and/or

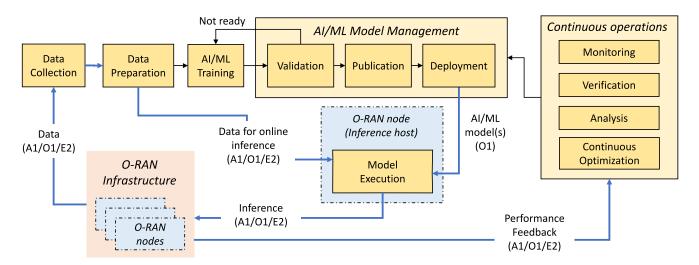


Fig. 13. AI/ML workflow in the O-RAN architecture. The RAN infrastructure provides data through O-RAN interfaces to the data collection and preparation logical blocks. The AI/ML models are then trained, validated, and deployed, and execute on O-RAN nodes (e.g., the RICs). Models can be further refined through monitoring and continuous optimization based on performance feedback from the RAN.

actions [125], [126], [127]. However, this does not preclude online training, which is still supported by O-RAN provided that it is only used to fine-tune and update a model previously trained offline [49], [128].

Example: In our example, the operator can train a variety of AI algorithms all controlling the number of PRBs allocated to each slice but differing one from another with respect to their implementation details. For example, the operator can train a set of Deep Reinforcement Learning (DRL) agents and decision trees and explore different combination and input formats (e.g., the specific subset of KPMs and their amount), different architectures (e.g., depth and width of a DRL agent, number of neurons, among others). The goal of this procedure is to train a large number of AI algorithms and identify which ones are the most suitable to accomplish a specific task.

Validation and Publishing: Once models are trained, they go through a validation phase to make sure they are reliable, robust and effective in performing classification, prediction or control tasks. If the validation is successful—and the models are deemed ready for deployment—they are published and stored in an AI/ML catalog on the SMO/non-RT RIC. Otherwise, they are required to go through additional redesign and re-training phases until the validation tests are successful [129]. Typically, agents do not meet validation requirements when the input and output spaces are not properly designed, or when the modeling performed by the agent's neural network(s) fails at capturing key relationships between the input and output of the agent. For an extended discussion on this, we refer the interested readers to [130].

Example: Once training has been completed, the different AI algorithms are compared one another and against diverse validation datasets including previously unseen data to identify which models are the most effective in controlling RAN slicing policies. For example, a typical validation test includes evaluating how well diverse AI solutions perform under diverse traffic patterns and demand, number and distribution of users, available bandwidth and operational frequencies. This procedure can either point out AI solutions that are not performing

well and need to be retrained, as well as determine the subset of AI algorithms that can be published to the AI/ML catalog as well as provide side information on the ideal network conditions (e.g., network load, mobility pattern, size of deployment) under which the specific AI solution delivers the best performance so that the operator can deploy the AI solution that is best suited to a specific deployment.

Deployment: Models stored in the AI/ML catalog can be downloaded, deployed and executed following two different options, namely *image-based* and *file-based deployments*. In both cases, the deployment of the model is performed by using the O1 interface, and the node where the model executes is referred to as *inference host*.

In the image-based deployment, the AI/ML model executes as a containerized image in the form of an O-RAN application (e.g., xApps and rApps) deployed at the O-RAN nodes, where it is executed to perform online inference. At the time of this writing, these nodes are limited to the RICs and the execution of AI at the CUs/DUs is left for further studies. The file-based deployment, instead, considers the case where the AI/ML model is downloaded as a standalone file that executes within an inference environment—outside the O-RAN application domain—that forwards the inference output of the model to one or more O-RAN applications. The main challenge associated with AI/ML deployment is the selection of the model that satisfies the operator's needs and/or intent, as discussed in Section IV [131].

Example: In our case, the operator will select the pre-trained AI-based RAN slicing models from the AI/ML catalog and deploy them as xApps that will be executed in the near-RT RIC.

AI/ML Execution and Inference: Once models are deployed on the inference host, they are fed with data to perform diverse online inference tasks. These include classification and prediction tasks, deriving policies at both RICs (transmitted over the A1 and E2 interfaces), and taking management and control actions (over the O1 and E2 interfaces, respectively). At this step, it becomes extremely important to

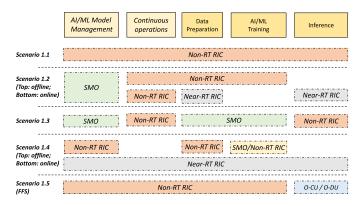


Fig. 14. O-RAN AI/ML deployment scenarios, adapted from [49]. Different scenarios embed different components of the O-RAN AI/ML workflows in different components of the RAN, from the SMO/non-RT RIC to the near-RT RIC or the RAN nodes.

execute inference algorithms that comply with the timelines of each control loop. For example, at the near-RT RIC an xApp with an AI/ML component needs to complete the computations in hundreds of milliseconds at most, possibly less, to guarantee that the control loop closes in less than 1 s. While inference is a less complex task than training (as we discuss in [68]), the loop execution time needs to be tested and profiled to ensure that it delivers a satisfactory performance also at the expected scale (see Fig. 3).

Example: Once the xApp has been deployed on the near-RT RIC, RAN slicing control is performed by executing the operations described in Fig. 9 where the xApp (i) is fed with KPMs (e.g., requested PRBs, latency and throughput measurements) collected over the E2 interface; and (ii) computes control actions that are used to pilot the DU and assign the available PRBs to the different slices in near-RT.

Continuous Operations: An important aspect of the AI/ML workflow is the ability to monitor and analyze the intelligence deployed throughout the network to verify that the inference outputs of AI/ML models are effective, accurate and do not negatively affect the performance of the network. Continuous operations ensure that models that perform poorly online can be refined and re-trained to improve their functionalities [132], [133], [134], avoid issues related to data and/or service unavailability, and help updating the AI/ML models without service interruption.

Example: In our case, the operator can monitor constantly the performance of the RAN slicing xApp and, if any anomalies or inefficiencies are detected, it can decide to re-train the AI/ML model embedded in the xApp over new data collected through the O1 and E2 interfaces.

A. AI/ML Deployment Scenarios

One of the main features of 5G networks is the ability to support a large variety of use-cases and applications. Given the diversity of the 5G ecosystem, it is clear that a one-size-fits-all solution in deploying and controlling network intelligence is unlikely to exist. For this reason, the O-RAN Alliance has specified five different deployment scenarios that define the location where the different components of the AI/ML work-flow of Fig. 13 are instantiated and executed [49]. Although

these deployment scenarios, which are shown in Fig. 14, cover a large set of real-world use cases, practical deployments might deviate from them to accommodate operator- and application-specific requirements.⁵

As mentioned before, O-RAN specifications are specifically designed for the RAN portion of the network and its functionalities. However, it is worth mentioning that O-RAN and its RICs can influence decisions regarding the core network and the Multi-access Edge Computing (MEC) infrastructure. Indeed, the SMO can act as a gateway between the RAN and other network components, as it has the capability of orchestrating functionalities across the whole network. In this way, xApps and rApps executing within the O-RAN environment can be leveraged to gather information on the RAN (e.g., traffic load forecast, mobility prediction, network state dynamics) that can be used by the SMO (e.g., ONAP or OSM) to take informed decisions about MEC service instantiation and delivery as well as network slicing policies in the core network [135].

B. Gathering Inputs for Online Inference

Data for online inference can be collected from multiple data producers over the O1, A1 and E2 interfaces, which are designed to support O-RAN control loops operating at different time scales. The O1 interface allows components in the SMO/non-RT RIC domain to gather data from any O-RAN management component and perform non-real-time optimization. The A1 interface can be used by the non-RT RIC to send enrichment information from the SMO/non-RT RIC domain to the near-RT RIC and its applications. For example, an rApp in the non-RT RIC can send enrichment information to the near-RT RIC on the predicted KPI evolution over the next few seconds. Finally, the E2 interface allows the near-RT RIC and its xApps to collect data from E2 nodes (e.g., KPMs) for near-RT control of the RAN. It is worth mentioning that the O-RAN specifications consider the case where input data might be generated within the same inference host by different data producers [49]. This is the case of chained AI/ML models, where one control task consists of the execution of several sub-tasks, each involving a different AI/ML model and requiring diverse input data and types.

To regulate data production and consumption between applications hosted in the same node (e.g., the near-RT RIC), the O-RAN specifications lay the basis for a *data access plat-form* that regulates data production, sharing and access. This platform acts as a middleware between the applications and a common data repository where data is stored and shared. To provide a high-level example of this, in Fig. 15 we show how data produced by different sources can be consumed by O-RAN applications performing heterogeneous tasks. We consider the case of an rApp X forecasting traffic load (data type A). Two xApps (xApp Y and xApp Z) leverage AI to control network slicing and scheduling strategies, respectively. xApp Y consumes data type A received from the data access platform and produces a slicing profile (data type C), while

⁵Although the O-RAN specifications consider the case in which AI is executed at the CUs/DUs (Scenario 1.5 in Fig. 14), in practice this is left for further studies.

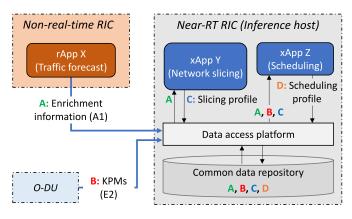


Fig. 15. An example of chained AI/ML models with diverse input types and data producers. Letter A indicates the traffic forecast generated by rApp X. Letter B represents the KPMs from the RAN. xApp Y uses data A to generate control action C (i.e., a slicing profile). xApp Z, instead, using data A, B and control C as input to generated a new control action D (i.e., a scheduling profile).

xApp Z consumes data types A, B and C coming from the data access platform (with type B consisting of KPMs sent from a DU over the E2 interface) and selects a scheduling profile D to be used by the controlled DU.

C. Using AI/ML to Design O-RAN Systems

One particular application that is worth mentioning is that of leveraging the very same AI/ML infrastructure offered by O-RAN to make its design and operations more efficient and effective.

Although in principle one could use AI/ML to improve any of the operations or functionalities executed within the O-RAN architecture, the most relevant applications include the management of the software infrastructure, the open interfaces and security.

For example, a quite relevant application where AI/ML plays an important role is that of taking intelligent decisions on regulating the lifecycle of all software components within the O-RAN infrastructure. Indeed, AI/ML can be used to automate and optimize the instantiation, placement and management of RICs to dynamically adapt to network conditions and traffic load, so as to minimize energy consumption while reducing the burden of each RIC. Similarly, AI/ML can be used to determine the optimal instantiation of CUs and DUs [136], [137], [138].

Another interesting application is that of using AI/ML to help designing of the open interfaces. AI/ML can be used to determine the optimal data retrieval rate on each interface so as to adapt and capture current network dynamics. For example, in slowly changing environments, AI/ML can reduce the rate at which data is collected over the open interfaces (e.g., E2 and O1), thus reducing overhead. Similarly, one could increase the data retrieval rate in the case of rapidly chaining conditions that need frequent and up-to-date data to let the RICs take optimal and timely decisions at the price of increased overhead.

Finally, security is another possible use-case where AI/ML can be used to fine-tune security mechanisms to secure the O-RAN architecture. For example, tools such as DRL

can be used to adaptively regulate the level of security in the network [139] by controlling encryption rates, detecting intruders and moving services and functionalities to reduce the likelihood they can be targeted by adversarial nodes [140].

VII. OPEN RAN USE CASES AND RESEARCH

The capabilities introduced by the RICs, the open interfaces, and the AI/ML workflow described in Sections III-VI make it possible to support advanced use cases and scenarios for the RAN control and deployment self-optimization.

Recent years have seen an increase of O-RAN-driven research on applications and use cases, e.g., on the design of xApps and rApps or on the optimal configuration of O-RAN networks. This section describes the areas where the application of the Open RAN paradigm is the most promising, and recent research results that show how O-RAN-based solutions can be used to optimize the RAN.

The O-RAN Alliance has collected an extensive list of 19 exemplary use cases for Open RAN deployments in [141], [142], and the literature further discusses at a high level some of these in [10], [15], [16], [17], [18], [21], [22], [23], [24], [25], [26], [27], [28], [29], [30], [31], [32]. At a high level, the scenarios and use cases can be classified in different ways, e.g., by considering the control knob or inference target, or the domain that is being controlled or optimized (e.g., a UE, a slice, a RAN node, or the whole network). In terms of control knobs, we discussed E2-specific targets in Section V-A. More generally, it is possible to identify several areas of interest, as follows.

Mobility: Open RAN networks can influence the mobility management or the performance of mobile users by tuning handover, load balancing, multi-connectivity, access barring, and beamforming parameters in the RAN. Differently than in traditional 3GPP networks, this can be done in a closedloop fashion by exploiting knowledge on the state of multiple base stations, and predictions on the user mobility based on information from the RAN or external enrichment information. For example, [143] presents an Open-RAN-based mechanism for the prediction of the load of multiple base stations in a cellular network, with application to the dynamic routing of autonomous vehicles to avoid introducing congestion. The O-RAN documents [141], [142] also include context-based handover for vehicular scenarios, in which xApps exploiting enrichment information from the non-RT RIC and inference on AI/ML RAN data to manage handovers, and applications related to Unmanned Aerial Vehicles (UAVs), with configuration of RAN parameters based on the expected trajectory of the UAV.

Resource allocation: Control in this area spans network slicing, scheduling, and provisioning of services and network functions. As for mobility control, the advantages of Open RAN compared to traditional cellular networks lie in the possibility of adapting to dynamic, evolving contexts, to new user requirements (e.g., for different slices), and to external events that alter the state and configurations of the RAN.

In this sense, the research on intelligent O-RAN control has adopted network slicing as one of the most interesting and promising areas for ML-based optimization. Bonati et al.

analyze data-driven approaches in O-RAN, and provide the first demonstration of closed-loop control of a softwarized cellular network instantiated on a large-scale experimental platform [12]. The performance of the RAN—implemented on the Colosseum testbed (see Section X)—is optimized through xApps that control the scheduling policies of various network slices on the base stations. The slices have different optimization targets, e.g., for the URLLC slice the reward minimizes the buffer occupancy as a proxy for the endto-end latency. Polese et al. propose ColO-RAN, a pipeline for the design, training, testing and experimental evaluation of DRL-based control loops in O-RAN [68]. The capabilities of ColO-RAN—prototyped on the Colosseum and Arena testbeds—are showcased through xApps to control the RAN slicing allocation and scheduling policies, and for the online training of ML models. Johnson et al. propose NexRAN, an xApp to perform the closed-loop control of the slicing resources of softwarized base stations, and demonstrate it on the POWDER platform of the U.S. PAWR program [144], [145]. Sarikaya and Onur consider the placement of RAN slices in a multi-tier 5G Open RAN and formulate it mathematically, showing the benefits of flexible functional splits compared to fixed split options [146]. Niknam et al. analyze the principles and requirements of the O-RAN specifications, and propose an intelligent scheme for the management of radio resources and traffic congestion. The effectiveness of this solution has been proved on real-world operator data [147]. Similarly, Mungari assesses the performance of ML-driven radio resource management in O-RAN-managed networks through an xApp deployed on the near-RT RIC, and evaluates it in a small laboratory setup [148]. The authors of [149] focus on the cell selection process, showing how O-RAN can help improve the allocation of users to specific cells based on forecasted throughput metrics rather than simple signal level metrics.

Lien et al. [150] and Filali et al. [151] explicitly consider xApp-based optimization of radio resources for URLLC users. In particular, the first paper shows how a reinforcement learning agent running in the near-RT RIC can effectively control the instantiation of new URLLC guaranteed bitrate sessions and configure the session-level parameters to increase the probability of successfully onboarding new URLLC users [150]. The paper by Filali et al., instead, studies an O-RAN-based slicing mechanism that controls the resource allocation for URLLC users and manages to achieve quality of service targets [151].

In [152], Iturria-Rivera et al. present a framework for coordinating multiple xApps and rApps via multi-agent team learning, so that resources can be allocated via an AI/ML infrastructure that leverages proximity and hierarchy relationships between multiple xApps/rApps.

Additional examples can be found in [141], [142], e.g., resource allocation for mobile users (including UAVs) with anticipatory mobility prediction, and QoS-based resource allocation. The latter aims at dynamically provisioning the set of resources that selected users require to satisfy their QoS, for example through ad hoc slicing and subsequent allocation of

PRBs to the QoS-driven slices. Finally, the O-RAN infrastructure can also be used to predict the emergence of congestion and apply appropriate remediations by adding more resources (e.g., carriers, MIMO layers, cells).

QoE/QoS-based control: The optimization of RAN resources to meet specific QoS and Quality of Experience (QoE) requirements also extends beyond resource allocation. For example, Bertizzolo et al. consider drone-enabled video streaming applications and propose a control system for the Open RAN for the joint optimization of transmission directionality and the location of the drone [153]. This solution is evaluated both experimentally, in a multi-cell outdoor RAN testbed, and through numerical simulations.

The O-RAN Alliance use cases also include a QoE optimization scenario, where inputs from external systems, monitoring of application performance, and multi-dimensional data are combined to identify the best RAN configuration for the optimization of the QoE on a user basis [141]. Here, O-RAN interfaces play a unique role, as they make it possible to combine and fuse data input from different, heterogeneous sources (RAN, external) in a way that is not usually possible in closed, traditional RAN deployments.

RAN sharing: This use case covers different scenarios, from spectrum sharing to infrastructure sharing on a neutral host architecture, which are expected to increase the spectral and energy efficiencies, and to reduce operational and capital expenditures. This scenario combines the flexibility of O-RAN softwarization and virtualization with the dynamicity and reconfigurability of closed-loop control (including with external information). The O-RAN specifications include multiple mechanisms for spectrum sharing, including slicing and control of dynamic spectrum access at the DU/RU.

In this regard, [154] studies a sharing scenario between a 5G RAN and a low-Earth orbit satellite constellation (in particular, the uplink between a ground station and a satellite), managed by a RIC. The study includes input from a RAN-independent spectrum-sensing framework and a sharing mechanism applied on a UE-basis by the RIC. Paper [155] also proposes an inter-technology spectrum sharing solution enabled by O-RAN, but, in this case, the sensing is performed inside the DU and RU themselves, with I/Q-based deep learning analysis of the received signals. The proposed sharing scheme detects Wi-Fi users (or unknown users occupying the spectrum of interest) and adapts the configuration of the RAN to avoid interference. Reference [156] proposes a framework based on the aggregation of contextual information from multiple sources to create spectrum maps that are then used by xApps for dynamic spectrum access control. An outlook to spectrum sharing enabled by O-RAN controllers in future 6G applications is provided in [157], where a backhaul link carrier frequency is changed by a centralized controller based on external information on incumbents that may suffer from interference from the communication links.

Blockchain-based approaches are also considered in [158], [159]. Notably, [158] embeds a blockchain framework on top of the O-RAN infrastructure to enable secure and trusted exchange of RAN resources (as for example DUs, RUs, etc.)

POLESE et al : UNDERSTANDING O-RAN

among multiple operators. Reference [159] combines Open-RAN-enabled neutral infrastructure and dynamic allocation based on blockchain as a practical way to bridge the digital divide gap. Finally, [160] explores the usage of smart contracts to activate and deactivate carrier aggregation across different operators on a shared O-RAN infrastructure.

Massive MIMO: This technology represents a key enabler of 5G networks. Through an Open RAN architecture, it is possible to embed dynamic control and adaptability to the configuration of the MIMO codebook (or group of beams)⁶ or of the beam selection process, and to make mobile experience more reliable and robust.

From the signal processing point of view, there is an extensive body on research that benefits from C-RAN and virtualized, centralized CUs and DUs [36], [161]. Paper [162] studies how channel state information available at the RU or DU may need to be shared across the two nodes. The authors consider specifically the capabilities provided by the O-RAN FH. In this context, [109], [163] discuss different compressions schemes for the O-RAN FH interface, also considering multi-stream capabilities typical of MIMO setups. Finally, the authors of [164] analyze different beamforming options based on the O-RAN 7.2x split of the physical layer.

When it comes to the RICs introduced by O-RAN, the control and optimization is generally related to beam parameters and to the codebook or group of beams in the DU and RU. For example, the O-RAN technical document on use cases investigates two different data-driven solutions [141]. The first adapts the group of beams based on telemetry collected at the non-RT RIC, e.g., user activity and reports, measurement reports, GPS coordinates, and reconfiguration through the O1 and O-RAN FH interfaces. The second is an optimization on the near-RT RIC of the mobility configuration, e.g., the beam-specific offsets that will determine whether a user should change beam or not. Another scenario of interest is the grouping of the users into multi- and single-user MIMO groups, which would then benefit from capacity enhancement or diversity. This can be done through policy guidance and enrichment information from the non-RT RIC, and with the actual control being performed by the near-RT RIC. The dynamic, data-driven reconfiguration of beamforming with the RICs is relatively unexplored in the Open RAN literature.

Security: We will discuss security in details in Section VIII.

New applications: The next generations of wireless cellular networks also embed and expand to other use cases. An example is the support for UAVs and vehicular communications, which we discussed as part of the mobility and resource allocation use cases.

Another example is represented by industrial IoT scenarios, which require high reliability and precise timing and synchronization achieved with data duplication, multi-connectivity, dedicated QoS and packet compression techniques [165]. Considering the number of parameters that can be tuned in this context, the closed-loop control enabled by the O-RAN RICs

can provide elastic configurations that adapt to the evolving conditions on the factory floor.

Moving away from communication scenarios, 5G networks from Release 16 also support positioning through dedicated signals on the air interfaces (both LTE and NR) and a location management function in the core network [166]. Location services will be used, for example, in indoor scenarios (e.g., factories, malls), to provide value-added services, or to implement location-based safety information broadcasting. However, relaying the location information to the core for analysis and processing may be affected by delays or jitter, making precise and timely estimation more difficult. In this case, the O-RAN specifications [141] see the near-RT RIC with a dedicated positioning xApp as an alternative to the 5G core location management function, leveraging the deployment of the near-RT RIC at the edge of the network.

O-RAN deployments optimization: Besides optimization of the RAN through O-RAN, several papers also study how to optimize the deployment of the Open RAN components themselves (e.g., RICs, xApps and rApps, RAN nodes). This leverages the management functionalities of the SMO, as well as its centralized point of view and the telemetry and statics from the RAN.

D'Oro et al. design a zero-touch orchestration framework that optimizes network intelligent placement in O-RANmanaged networks, and prototype it at scale on Colosseum using open source RIC and RAN components [131]. In this context, the paper [134] introduces the concept of RLOps, or reinforcement learning operations, i.e., a framework to manage the life-cycle of intelligence (specifically, reinforcement learning) for Open RAN. RLOps cover the whole end-to-end workflow of AI/ML for the RAN (Section IV), from the design and development to the operations (i.e., deployment, updates), and the management of safety and security during the overall intelligence life-cycle. Huff et al., instead, develop a library, namely RFT, to make xApps fault-tolerant while preserving high scalability [167]. This is achieved through techniques such as state partitioning, partial replication and fast re-route with role awareness.

Other papers focus on the virtualized components and of the disaggregated base stations. Tamim et al. maximize the network availability by proposing deployment strategies for the virtualized O-RAN units in the O-Cloud [168]. Pamuklu et al. propose a function split technique for green Open RANs [62]. The proposed solution, which is based on DRL, is evaluated on a real-world dataset. The authors of [169] develop a matching scheme between DUs and RUs with a 2D bin packing problem. Finally, the O-RAN specifications consider a similar use case, with data-driven pooling of the CUs and DUs on shared, virtualized resources [141], orchestrated through the non-RT RIC.

O-RAN white papers and surveys: Finally, overviews of O-RAN and of its components are given by Lee et al. in [15], which implements AI/ML workflows through open-source software frameworks; by Abdalla et al. in [16], which reviews O-RAN capabilities and shortcomings; by Garcia-Saavedra and Costa-Pérez in [17], which gives a succinct overview of O-RAN building blocks, interfaces and services; and by

⁶The group of beam is a set of beams on which the base stations transmits synchronization signals for initial access, to improve the initial access performance especially at higher frequencies [112], [142].

Brik et al. in [18] and Arnaz et al. in [19], which discuss deep learning and artificial intelligence applications for the Open RAN. Reference [14] discusses open source software that can be used to deploy 5G and Open RAN networks. Several white papers [10], [21], [22], [23], [24], [25], [26], [27], [28], [29], [30], [31], [32] provide high-level overviews of the Open RAN vision and of the O-RAN architecture. Differently from the above works, in this paper, we present a comprehensive overview of the O-RAN specifications, deep-diving into the open interfaces, protocols and services, and discussing in detail use cases, AI/ML workflows, deployment scenarios, and open platforms for O-RAN-enabled experimental research.

VIII. SECURITY IN THE OPEN RAN

It is undeniable that the introduction of new architectural components, open interfaces, disaggregation, and the integration of custom and possibly data-driven control logic will make next-generation cellular networks more efficient and flexible. On the other hand, however, this revolution comes with unprecedented security challenges that primarily stem from the fact that the distributed and disaggregated nature of the O-RAN infrastructure effectively extends the attack surface for malicious users, thus posing severe threats to the network and its users [28]. For example, the authors of [170] have highlighted that partial and limited support for authentication can lead to vulnerabilities in O-RAN interfaces. At the same time, the unparalleled monitoring capabilities, the intelligence and the cloud-native deployment that characterize O-RAN architectures truly add insights on the state of the network and provide the necessary tools to implement advanced solutions to monitor, detect, prevent and counteract threats [27]. In this regard, the O-RAN Alliance has created a dedicated working group (WG 11) for the analysis and definition of threat models for O-RAN networks, as well as for the definition of security measures and policies for the components of the O-RAN architecture toward a zero-trust model [171], [172], [173]. Therefore, this section presents a breakdown of the security stakeholders of Open RAN deployments (Section VIII-A), of the threat surface (Section VIII-B), and of the related countermeasures that are discussed by the O-RAN Alliance documents and by the literature on Open RAN (Section VIII-C).

A. Security Stakeholders

The O-RAN Security Focus Group (SFG) (recently promoted to a full WG, i.e., WG11) has defined a list of stakeholders that need to proactively secure the RAN. This extends the interested parties beyond those considered in traditional 4G and 5G networks, e.g., vendors, operators, and system integrators. As also discussed in [27], operators will assume a more predominant role in securing the infrastructure, as the openness of the platform and the usage of multivendor components allows them to customize the build (and thus the security) of the infrastructure. This also means that operators can assess and vet the security level of the open components introduced in the network, which is often not possible in close architectures that are fully vendor-driven. Reference [27] also

identifies network functions and virtualization platform vendors as new stakeholders (e.g., third-party xApp and rApps developers, O-Cloud providers), along with administrator profiles that manage virtualized and disaggregated components. In addition, the orchestrator (e.g., the entity that manages the SMO) also has a role in securing the operations of the network.

B. Extended Threat Surface

According to the threat analyses in [171], [174], soft-warization and disaggregation will inevitably extend the threat surface of cellular systems. So far, WG11 has identified seven threat categories that include a total of 84 possible threats:

• Threats against the O-RAN system: The Open RAN architecture introduces new architectural elements and interfaces (from the fronthaul to control and management interfaces), which become part of an extended threat surface. These components can be subject to different attacks, which may compromise (i) the availability of the infrastructure (e.g., unauthorized access to disaggregated RAN components aiming at deteriorating the performance of the network, or the malicious deployment of xApps that intentionally introduce conflicts with other xApps); (ii) data and infrastructure integrity (e.g., compromised software trust chains or the misconfiguration of interfaces); and (iii) data confidentiality (e.g., through attacks that disable over-the-air encryption, or facilitate unregulated access of user data from xApps and rApps). As an example, an attacker could exploit insecure design of the O-RAN components (e.g., interfaces, functions, etc.), their complexity, or their misconfiguration, to gain unauthorized access to the system. Or again, weak authorization/access control mechanisms could allow an attacker to penetrate the network boundaries and access O-RAN services exposed through Web servers. Datarelated threats encompass (i) information transported over O-RAN interfaces for control, management, and configuration of the RAN; (ii) data used for training and testing of the ML models; (iii) sensitive data on the users, e.g., their identities; and (iv) the cryptographic keys deployed on the elements of the network. The threat surface also expands to the new logical components of the architecture, i.e., the RICs, the SMO, and the software frameworks (e.g., xApps, rApps) that execute on the RICs. At the near-RT RIC and xApps attacks include malicious xApps that leverage unauthenticated RIC APIs and allow attackers to request services ("theft of services") or sensitive information ("data leakage"), track UEs and change their priority, access resources. Attacks could also compromise the isolation among xApps, or send them false policies to degrade their performance. Concerning the non-RT RIC and rApps an attacker could cause denials of service or degrade performance, track UEs, data corruption and modification. Finally, attacks on the Open Fronthaul 7.2x split interface are also often mentioned as a potential vulnerability [28]. As of today, the 7.2x interface is not encrypted on the control plane, because of the challenging timing

POLESE et al.: UNDERSTANDING O-RAN

requirements that encryption would introduce. This introduces man-in-the-middle attacks, in which the attacker impersonates the DU (or RU), and compromises user data or configurations in either of the two endpoints, or gains access to the DU and beyond through attacks on the RU or fronthaul interface, or compromises the availability/confidentiality of such interface. Another attack can be carried out against the S-Plane, with a malicious actor compromising the synchronization infrastructure and thus causing performance degradation, e.g., by spoofing the master clock, removal of PTP packets, or through a rogue instance wanting to become a PTP grand master.

- Threats against the O-Cloud: The O-Cloud provides a virtualization environment that encompasses RAN elements as well as O-RAN components. To this end, threats identified for the O-Cloud relate to attacks in virtualized environments. Possible attacks include (i) compromising virtual network functions, either being executed or their snapshots or images (e.g., to leak embedded cryptographic secrets); (ii) exploitation of the O2 interface between the O-Cloud and the SMO to gain access and escape isolation; (iii) misuse of containers or virtual machines for network functions to attack other entities in the network; (iv) spoofing or compromising the underlying networking or auxiliary services; and (v) abuse of privileged Virtual Machines (VMs)/containers, or escape from their isolated environments, data theft, modification of the resources allocated to the VMs/containers, image tampering, and secred disclosure in virtualized environments.
- Threats against open-source code: The softwarization
 of the RAN and of O-RAN components opens new
 vulnerabilities because of backdoors in the O-RAN
 code by (i) trusted developers, which intentionally
 compromise software components, or by (ii) upstream
 libraries that are not under the control of the O-RAN
 developers.
- Physical threats: The additional hardware introduced to support the gNB split and the O-RAN infrastructure can be compromised by attackers that gain physical access to the infrastructure. The attacks can range from power availability attacks, to cabling reconfigurations, or addition of hardware backdoors, to the damage of sensitive data.
- Threats against the wireless functionalities: Attacks on the RU or the Open Fronthaul interface between RUs and DUs can lead to performance degradation on the air interface, with typical attacks related to jamming of data or synchronization signals, or denials of service. Additionally, combining equipment from different vendors may reduce the performance of the network in case of mismatch in configurations or supported functionalities.
- Threats against the protocol stack: Attacks can also be carried out toward the protocol stack layers. Attacks concerning injection, cross-site scripting, denial of service, exposure of objects identifiers without proper authorization, and Web tokens through REST APIs, JSON or HTTP exploits can be performed.

• Threats against the AI/ML components: Finally, the O-RAN specification [171] and the literature [27], [28] also describe a new class of threats, i.e., attacks against AI/ML models used for inference and control in xApps and rApps. A practical instance of such an attack is that of poisoning attacks, in which an attacker exploits unregulated access to the data stored in the SMO/non-RT RIC to inject altered and misleading data into the datasets used for offline training of AI/ML algorithms. Another example is that of an adversary gaining unrestricted control over one or more O-RAN nodes to produce synthetic data fed in real time to AI/ML solutions being fine-tuned online, or being used to perform online inference. These attacks are extremely relevant as they might result in AI/ML solutions that output wrong predictions, or make wrong control decisions that result in performance degradation or-even worse-outages. Similar attacks can also target the ML model directly (e.g., by modifying the weights or configurations of the model) [175].

Based on this security analysis, the O-RAN WG11 has so far identified 32 O-RAN-specific critical assets related to interfaces and data, and 16 related to logical components, also partially discussed in [27], [28], [176]. Additionally, other issues discussed in [174] include low product quality for Open RAN disaggregated components, immature and in-progress technical specifications, supply chain tampering, and support infrastructure (e.g., power grid) failure.

C. O-RAN Security Principles and Opportunities

While the new architecture and its interfaces introduce new threats and opportunities for attackers, they also come with the opportunity to re-think the security principles and best practices for designing, deploying and operating cellular networks and align them to the best practices of cloud-native deployments [27]. In general, openness is associated with increased visibility into the processes and operations of the RAN, which puts operators in control of their network. The O-Cloud and the virtualized nature of the O-RAN platforms enable quick deployment of security patches and updates, automated testing and deployment, with full control over the entire end-to-end process including information on vendors and software components being used at any moment. The disaggregation also makes it possible to deploy simpler network functions, with more atomic components that are easier to test and profile. Finally, the virtualized CU is generally deployed in a centralized data center, which makes it easier to physically secure the RAN cryptographic keys. These phenomena are not completely new to the networking industry, as some of the very same challenges were faced, and practically solved, already in the data center and cloud networking computing domains where softwarization and disaggregation principles brought the data and control plane disaggregation together with the logical split of network functions into atomic VNFs. Therefore, the O-RAN Alliance will be able to build on top of the existing security measures (and coutermeasures) for data centers and cloud networking systems, thus facilitating the transition from physical to virtual appliances.

In this sense, the O-RAN WG11 has published a number of technical specifications that mandate authentication and encryption procedures across the different elements of the O-RAN architecture [172], [173]. These documents complement the 3GPP security requirements and specifications [177], [178], [179], [180], [181] to address security issues specific to the O-RAN architecture. They define: (i) the supported version of the Secure Shell (SSH) protocol and cryptographic ciphers to use with it, as well as procedures for keys and key exchange, symmetric algorithms for encryption of the transferred data, and message authentication codes; (ii) the supported version of the Transport Layer Security (TLS) protocol, how to configure it, and ciphers that should be used with it; (iii) the support of Network Configuration Protocol (NETCONF) over secure transport operations; (iv) the requirements for interfaces using Datagram Transport Layer Security (DTLS) for mutual authentication, and integrity, replay, and confidentiality protection; (v) the requirements for interfaces using IPSec for authentication, confidentiality, and integrity functionalities; (vi) the use of the Certificate Management Protocol (CMP) for the base stations to obtain certificates signed by operators; and (vii) the use of the OAuth 2.0 framework for service producers to authorize requests from service consumers, and for the latter to obtain token-based authorization credentials (mandatory for operators and cloud providers). These functionalities are mandatory for vendors, with optional elements that—at the time of writing—refer to the compilation of a software bill of materials and spoofing and resiliency mechanisms for the S-plane of the fronthaul interface. The document [173] also identifies the deployment of port-based network access control for devices that connect to the fronthaul interface (RUs and DUs) with IEEE 802.1x-2020 as optional, with mandatory deployment left for further study.

Finally, the availability of data and the insights on the RAN that the different interfaces (E2, O1) provide can also be leveraged to increase the security of the RAN itself. This is due to the intelligent, data-driven self-monitoring of the RAN performance, which can automatically trigger warnings and alarms in case unintended behaviors are detected. In this sense, Open RAN deployments can leverage and extend techniques discussed in the significant body of literature on data-driven anomaly detection, for which a survey can be found in [182]. Some examples applied to cellular networks are [183], where the authors use machine learning to detect RAN components that are not following standard patterns and may be classified as rogue, and [184], which uses deep learning to detect and classify anomalies in cellular IoT applications.

IX. O-RAN DEVELOPMENT AND STANDARDIZATION

The O-RAN Alliance is a consortium of operators, vendors, research institutions, and industry partners that focuses on reshaping the RAN ecosystem toward an intelligent, open, virtualized and interoperable architecture. To this end, the efforts of the Alliance cover three macro-areas [185]: (i) *specification*, aimed at extending RAN standards from SDOs such as 3GPP, ETSI and ITU to include openness and intelligence; (ii) *software development*, focused on developing and contributing

TABLE I O-RAN WORKING GROUPS

Working Group	Focus
WG1	Architecture and use cases
WG2	Non-RT RIC and A1 interface
WG3	Near-RT RIC and E2 interface
WG4	Open Fronthaul interfaces
WG5	Open F1/W1/E1/X2/Xn interfaces
WG6	Cloudification and orchestration
WG7	White-box hardware
WG8	Stack reference design
WG9	Open X-haul transport
WG10	OAM and O1 interface
WG11	Security

open source software for the RAN components of the O-RAN architecture; and (iii) *testing and integration*, in which it provides guidance to members of the Alliance willing to perform testing and integration of the O-RAN-compliant solutions they develop. Note that the O-RAN Alliance is currently not an SDO, thus any standard-related activity will require liasons with SDOs, e.g., 3GPP, ITU working group 5, or ETSI, as discussed in Section V-D.

The specification tasks of the O-RAN Alliance are divided among 10 WGs, each responsible for specific parts of the O-RAN architecture. The main focus of each of these WGs is summarized in Table I:

- WG1 (Use Cases and Overall Architecture): This WG identifies key O-RAN use cases, deployment scenarios and development tasks of the overall O-RAN architecture. It includes three task groups: (i) Architecture Task Group, focused on specifying the overall O-RAN architecture, on describing its functions and interfaces, on illustrating relevant implementation options, and on facilitating cross-WG architectural discussions; (ii) Network Slicing Task Group, focused on studying network slicing in O-RAN and on defining its use cases, requirements and extensions to the O-RAN interfaces; and (iii) Use Case Task Group, focused on identifying, defining and disseminating use cases enabled by O-RAN. We discussed WG1 activities primarily in Sections II and VII.
- WG2 (Non-RT RIC and A1 Interface): This WG specifies
 the architecture and functionalities of the non-RT RIC and
 of the A1 interface, which are discussed in Sections IV
 and V-C of this paper.
- WG3 (Near-RT RIC and E2 Interface): This WG specifies the architecture and functionalities of the near-RT RIC and of the E2 interface. It also takes care of providing support to AI/ML and data analytics model design to train models and enhance radio resource management and allocation. WG3 content is discussed in Sections III and V-A of this paper.
- WG4 (Open Fronthaul Interfaces): This WG focuses on defining an Open Fronthaul interface that supports interoperability of DUs and RUs manufactured by different vendors, as discussed in Section V-D.
- WG5 (Open F1/W1/E1/X2/Xn Interface): This WG provides multi vendor profile specifications for the

TABLE II O-RAN FOCUS GROUPS

Focus Group	Focus	State
OSFG SDFG TIFG	Open source issues, establishment of the OSC Standardization strategies Testing and integration, PlugFests	Dormant Active Active

F1/W1/E1/X2/Xn interfaces that are compliant with the 3GPP specifications. These interfaces are briefly discussed in Sections II and V of this paper.

- WG6 (Cloudification and Orchestration): This WG identifies use cases to demonstrate the benefits of the software/hardware decoupling of the O-RAN elements (e.g., RICs, CU, DU, RU) and deployment scenarios, and develops requirements and reference designs for the cloud platform. This WG also develops life-cycle flows and commonalities of O2 interface APIs between the SMO and the O-Cloud. This is primarily described in Section II.
- WG7 (White-box Hardware): This WG specifies and releases a reference design toward a decoupled software/hardware platform, as discussed in Section II.
- WG8 (Stack Reference Design): This WG develops software architecture, design, and plans for the CU and DU compliant with the 3GPP NR specifications, and a set of tests that promote interoperability across different implementations of the O-RAN interfaces. We reviewed this in Section II.
- WG9 (Open X-haul Transport): This WG focuses on the network transport, including transport equipment, physical media and protocols, as discussed in Section V-D.
- WG10 (OAM): This WG focuses on the O1 interface Operations, Administration and Maintenance (OAM) specifications (e.g., unified O1 operation and notification) and on creating OAM architecture and requirements for the O-RAN architecture and use cases identified by WG1. O1 and related topics are presented in Section V-B.
- WG11 (Security Work Group (SWG)): This WG focuses on the security aspects of the O-RAN ecosystem, as discussed in Section VIII.

Besides WGs, O-RAN also includes groups that focus on topics that are relevant to the whole Alliance. These are named Focus Groups (FGs) and they are summarized in Table II:

- Open Source Focus Group (OSFG): This FG used to deal with open source-related issues in O-RAN. These included the planning, preparation, and establishment of the OSC, and the coordination with other open source communities. Since the launch of the OSC in 2019, most of the activities of the OSFG have been taken care directly by the OSC. As a result, this FG is currently dormant [188].
- Standard Development Focus Group (SDFG): This FG
 programs the development strategies of the O-RAN
 Alliance and interfaces with other standard development
 organizations (e.g., 3GPP Service and System Aspects,
 3GPP RAN, ETSI, ITU-T, Small Cell Forum, IEEE 1914,
 NGMN Alliance). The SDFG also collects requirements

- and suggestions from the other entities of the Alliance (e.g., the WGs), and provides them guidance on use cases.
- Test & Integration Focus Group (TIFG): This FG defines the overall approach of O-RAN for testing and integration, including the coordination of tests across the WGs. Examples are the test and integration of specifications, the creation of profiles to facilitate O-RAN commercialization, and the specification of processes for O-RAN integration and solution verification. The TIFG also plans and coordinates the O-RAN PlugFests, and sets guidelines for third-party Open Testing & Integration Centres (OTICs). These are vendor-independent entities in charge of performing interoperability and conformity tests to verify that: (i) individual software and hardware components from vendors comply with O-RAN specifications; (ii) components from multiple vendors can interoperate seamlessly; and (iii) end-to-end systems (possibly multivendor) are functional and able to deliver the required minimum performance levels without any degradation due interoperability and integration issues. At the time of writing, there are already eight established OTICs spread across Europe, Asia and USA that grant (and have already granted) conformance certificates, as well as interoperability and end-to-end system integration badges [189]. OTICs will play a major role in ensuring interoperability and compliance with O-RAN specifications, and will be key in offering operators with the necessary tools to monitor and determine which components are really open and can used to deploy multi-vendor end-to-end cellular networks.

The O-RAN Alliance also features a Technical Steering Committee (TSC), which guides the activities of the Alliance, and features four sub-committees. These aim at (i) defining a minimum viable plan for a fully-compliant O-RAN network; (ii) define the procedures for the O-RAN Alliance; (iii) promote industry engagement; and (iv) develop research activities toward next generation or 6G networks.

Finally, the O-RAN alliance is also working together with external organization that are contributing to the development of the Open RAN. Examples of this are the ONF, that contributed to the specifications of the near-RT RIC, or the Telecom Infra Project (TIP), that entered a liason agreement with the Alliance to ensure the development of interoperable and disaggregated solutions for the RAN [190].

X. EXPERIMENTAL WIRELESS PLATFORMS FOR O-RAN

The recent adoption of the softwarization paradigm in next generation cellular networking, and the emergence of open frameworks and interfaces for the Open RAN, such as the ones described in this work, bring unprecedented opportunities to the field of experimental research on cellular networks. In this regard, experimental research will be fundamental in developing and validating the use cases described in Section VII. Once hard to experiment upon—due to the complex, hardware-based and closed implementations of RAN nodes, e.g., the base stations—in recent years, the cellular networking ecosystem

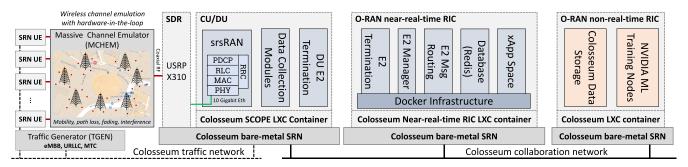


Fig. 16. OpenRAN Gym components [186] deployed on the Colosseum experimental testbed [187], adapted from [68]. The non-RT RIC, near-RT RIC, and RAN nodes (either base stations or users) are deployed as containers on different Colosseum Standard Radio Nodes (SRNs). Different Colosseum networks can be used to provide application traffic and support for the O-RAN interfaces, while the wireless communications among RAN components is emulated by the Massive Channel Emulator (MCHEM).

TABLE III
EXPERIMENTAL WIRELESS PLATFORMS AND FRAMEWORKS FOR O-RAN

Platform	O-RAN-compatible	Deployment
Arena [193]	yes (container ready to use)	sub-6 GHz indoors
Colosseum [187]	yes (container ready to use)	sub-6 GHz network emulator
AERPAW [194]	yes	aerial outdoors
COSMOS [195]	yes	mmWave outdoors, sub-6 GHZ indoors
POWDER [196]	yes (container ready to use)	sub-6 GHz outdoors & indoors
5GENESIS [197]	will be (platform under development)	5G outdoors
Framework	O-RAN-compatible	Deployment
OpenRAN Gym [186]	yes	end-to-end pipeline for AI/ML O-RAN research
Open AI Cellular [198, 199]	yes	AI-enabled framework

has seen this task facilitated by open and open-source protocol stacks (e.g., srsRAN [191] and OpenAirInterface [192]). These software-based implementations enable the instantiation of 3GPP-compliant network elements on general-purpose, off-the-shelf devices, allowing virtually anyone to instantiate a complete and operational cellular network with multiple nodes, and to experimentally validate solutions for cellular applications.

Thanks to the open protocol stacks mentioned above, in the last few years, the wireless community has seen the creation and broader adoption of publicly-available platforms and frameworks open and available to the research community. The prominent ones are summarized in Table III. These platforms play a vital role as they provide the means—and scale—to virtualize cellular stacks and controllers on their publicly-available infrastructure, and to design and prototype solutions in deployments as close as possible to those of commercial networks. Moreover, when it comes to AI/ML solutions—which require large amounts of data for the training and testing processes—they operate as wireless data factories, providing users with the tools to perform data collection at scale in controlled—yet realistic—wireless environments.

Open RAN experimentation relies on a combination of (i) compute resources, to run the virtualized O-RAN

components (e.g., the RICs); and (ii) radio resources, to host the over-the-air component of the RAN. For example, the base stations can be implemented through open and softwarized protocol stacks, such as srsRAN [191] and OpenAirInterface [192]. These stacks leverage Software-defined Radios (SDRs) (e.g., NI USRPs) as radio front-ends, and serve users implemented through analogous protocol stacks, or via commercial smartphones. The O-RAN software can be either developed ad hoc, e.g., as FlexRIC [71], or be based on the OSC, ONF, Linux foundation, or ETSI open source frameworks, as discussed in Sections III and IV. For example, when focusing on RAN control, the research pipeline generally involves a near-RT RIC, the RAN, the E2 interface, and custom xApps implementing the desired control on the RIC.

All the platforms in Table III provide these two ingredients, at different extents, and are thus fit for Open RAN research. Some of them, as discussed next, are already equipped with software and pipelines for this, while others can be used in combination with other frameworks.

A. Experimental Open RAN Research With OpenRAN Gym

The Open RAN experimental workflow is enabled, for example, by OpenRAN Gym [186]. OpenRAN Gym combines software-defined cellular stacks with a lightweight RIC, which can be deployed on multiple experimental platforms from Table III, E2 termination, and an end-to-end AI/ML pipeline for O-RAN.

OpenRAN Gym offers a ready-to-use Linux Container (LXC)-based implementation with the main components of the OSC near-RT RIC, as shown in Fig. 16 [68]. This can be instantiated on top of any bare metal compute and includes RIC services implemented as Docker containers. They include the O-RAN E2 termination, manager and message routing containers, which are used to communicate with the E2 nodes, and Redis database container, used to keep a record of the E2 nodes associated with the near-RT RIC. Additionally, external xApps can be instantiated in what is shown in the figure as xApp space [200]. As part of OpenRAN Gym, we provide a sample xApp that manages the connectivity to and from the RIC platform.

Through the E2 termination, the near-RT RIC can connect to the E2 nodes (e.g., CUs/DUs) to implement softwarized

control of the RAN (see Fig. 16, left). In OpenRAN Gym, the latter can be implemented through the SCOPE framework that extends srsRAN [191] with additional slicing, MAC- and PHY-layer functionalities, control APIs, and data collection capabilities [201]. This can be paired with the OSC RAN-side E2 termination to interface with the E2 termination on the near-RT RIC.

B. Colosseum and Arena

Two of the prominent platforms that allow users to instantiate an O-RAN-compliant network (e.g., with OpenRAN Gym) and components on a white-box infrastructure are Colosseum and Arena [187], [193]. Colosseum is the world's largest wireless network emulator with hardware-in-the-loop. Through a first-of-its-kind FPGA fabric, Colosseum empowers researchers with the tools to capture and reproduce different conditions of the wireless channel, and to experiment at scale through 256 USRP X310 SDRs [187]. Colosseum provides researchers with access to 128 Standard Radio Nodes (SRNs), i.e., a combination of a server and of an USRP X310 SDR acting as a RF front-end. SRNs can be used to instantiate the RICs (without considering the radio component) or softwarized base stations and users. Colosseum also provides data storage and NVIDIA Graphics Processing Units (GPUs) for the training of ML models (see Fig. 16, right). These resources can be used, for example, as a component of the SMO/non-RT RIC.

Finally, the SDRs are connected through coaxial RF cables to Colosseum Massive Channel Emulator—which takes care of emulating in FPGA different conditions of the wireless environment (Fig. 16, left)—and through the traffic network to Colosseum Traffic Generator—which leverages Multi-Generator [202] to generate and stream IP traffic flows to the SRNs.

After prototyping O-RAN-powered solutions on Colosseum with the setup of Fig. 16, users can port them to other experimental testbeds [203], such as Arena and the platforms of the U.S. PAWR program [193], [204]. Arena is an indoor testbed with 24 SDRs (16 USRPs N210 and 8 USRPs X310, synchronized in phase and frequency through 5 OctoClock clock distributors) connected to a ceiling grid with 64 antennas deployed across a 2240 ft² office space and controlled by a set of 12 high-performance servers. The combination of servers, SDRs, and antenna layout offers the ideal setup for testing of MEC capabilities [135] and private indoor cellular deployments [205], [206]. An example of O-RAN-related research that combines Colosseum and Arena is described in [68], [203].

C. Other Experimental Research Platforms

Other publicly-available testbeds include the city-scale platforms of the U.S. National Science Foundation PAWR program [204]. These consist of POWDER [196], focused on sub-6 GHz cellular deployments, COSMOS [207], on mmWave communications, and AERPAW [194], on aerial cellular deployments.⁷ POWDER, deployed in Salt Lake City,

UT allows users to perform experimentation in the sub-6 GHz spectrum in both outdoor and indoor environments. The former is achieved through fixed ground-level and rooftop nodes, also supporting massive MIMO applications, as well as mobile nodes deployed on university shuttles. The latter, instead, enables controlled experiments either in an indoor over-the-air laboratory, or through SDRs wired through a RF attenuation matrix. COSMOS, in New York City, NY focuses on mmWave experimentation through outdoor nodes deployed in a densely populated area, optical experimentation through an xHaul network developed by the National Science Foundation Center for Integrated Access Networks, and indoor sub-6 GHz experimentation through an array of nodes deployed as part of the ORBiT testbed. Finally, AERPAW, in the Research Triangle of North Carolina, allows users to perform experiments with in the sub-6 GHz spectrum through SDRs-either mounted on UAVs or on fixed nodes—as well as in the mmWave spectrum through Terragraph radios mounted on light poles and rooftops, LoRA devices, and phased arrays for radar applications.

Namely, all of these platforms are compatible with the O-RAN paradigm, as they allow users to instantiate white-box base stations managed by the O-RAN RICs. However, at the time of this writing, the only testbed that offers a ready-to-use O-RAN implementation (in the form of a pre-compiled container image) is POWDER [144], [209]. OpenRAN Gym has been tested on POWDER and COSMOS [203].

In Europe, the 5GENESIS consortium is working on the implementation of various 5G components, and the validation of different use cases across its several testbeds [197]. These include edge-computing NFV-enabled heterogeneous radio infrastructure, orchestration and management frameworks, terrestrial and satellite communications, and ultra-dense network deployments. Upon completion, these testbeds will be compatible with the O-RAN ecosystem. Similarly, the Scientific LargeScale Infrastructure for Computing/Communication Experimental Studies (SLICES) aims at creating a scientific instrument for research in networking—equipped with cuttingedge Open RAN, NFV, orchestration, and cloud computing solutions and well-defined research data and experiments management practices and techniques [210]. In addition, the SLICES effort aims at creating an inter-connect to share infrastructure across multiple European testbeds, to reduce fragmentation and increase the scale of networking experiments.

Among the notable open-source initiatives, Open AI Cellular (OAIC) proposes a framework that is integrated with the O-RAN ecosystem. This framework allows users to manage cellular networks through AI-enabled controllers, and to interact with systems that locate implementation, system-level, and security flaws in the network itself [199], [211].

XI. CHALLENGES AND FUTURE RESEARCH AND DEVELOPMENT DIRECTIONS FOR OPEN RAN

While the foundational principles and the main specifications for O-RAN have been drafted, partially enabling the use

⁷A fourth platform, ARA, which will focus on rural broadband connectivity, has been announced [208]. However, this platform is not yet operational.

cases described in Section VII, there are still several open challenges for standardization, development, and research. While some issues are associated to the identification of business and market opportunities for Open RAN [20], we focus here on the technical obstacles. Specifically, we identified some of them as follows.

A. Use Cases

- Identification of key O-RAN use cases: While O-RAN provides the infrastructure to implement RAN closedloop control, the identification of a key set of use cases that leverage these extended capabilities is still ongoing. The O-RAN Alliance provides a list of relevant use cases for the RIC-enabled control, which include classic radio resource management optimization related to handover optimization, resource allocation, QoE optimization, traffic steering, among others [141]. Nonetheless, as the capabilities of the 3GPP RAN evolve toward, for example, non-terrestrial networks and support for Augmented Reality (AR)/Virtual Reality (VR) in the metaverse, it becomes necessary to further refine and evaluate future O-RAN use cases and the role intelligent, data-driven closed-loop control can have in future domains.
- Open RAN beyond public cellular networks—the private cellular use case: Private cellular networks are quickly emerging as a key 5G deployment scenario, with applications in industrial automation, warehouses, healthcare industry, education, and entertainment. Greenfield private 5G deployments can easily embed Open RAN solutions for network control and optimization, as well as to reduce ownership and operation costs thanks to disaggregated and virtualized nodes. The design of O-RAN-enabled private networks introduces challenges in terms of domain-specific optimization, integration with edge systems and local breakouts, and support for connectivity in constrained environments.
- Spectrum sharing solutions enabled by Open RAN:
 Network controllers and programmable RAN nodes open new opportunities for the development of spectrum sharing systems [155]. The O-RAN specifications already include capabilities for LTE/NR dynamic spectrum sharing, but the design of algorithms to enable this is still an open challenge. Future research can investigate how to practically enable O-RAN-based sensing, reactive and proactive spectrum adaptation solutions, considering 3GPP and non-3GPP systems, as well as sharing-related extensions of the O-RAN architecture.

B. Interfaces

 Interoperability and testing: The introduction of new interfaces in the architecture defined by the O-RAN Alliance brings further need for interoperability testing in the telecom ecosystem. This is required to verify that multi-vendor deployments feature implementations that comply with the specifications and are interoperable. The O-RAN Alliance has defined interoperability and

- testing requirements for the fronthaul interface [212], but additional steps need to be taken by vendors to fully commit to implementing specification-compliant interfaces across the board. In addition, a truly interoperable ecosystem will foster the development of xApps and rApps that can be ported across multiple near-RT and non-RT RICs. In this sense, the definition of APIs or of a SDK for the RICs is a key interoperability enabler. Finally, the fronthaul interface and the deployment of efficient, scalable fronthaul networks is one of the key challenges in the design and scaling of Open RAN deployments.
- Service models development and implementation: As discussed in Section V, the E2 service models play a key role in the definition of what O-RAN and, specifically, the near-RT RIC can control in a 3GPP-defined RAN. A key challenge is to design service models that are comprehensive and track new use cases developed for the Open RAN, and to identify profiles and basic set of functionalities that RAN equipment vendors need to implement to be O-RAN compliant. Indeed, the near-RT RIC effectiveness for RAN analysis and control ultimately depends on the E2 service models implemented in the RAN.

C. Architecture and Platform

- O-RAN Architecture and its evolution: The foundational elements of the O-RAN architecture include the disaggregated RAN nodes (CUs, DUs, RUs) and the near-RT and non-RT RICs hosting xApps and rApps, respectively. There are several open questions as of how this architecture can be effectively deployed, e.g., in terms of distribution of networking elements across the edge and cloud network, or ratio among RAN nodes and RIC elements. In addition, further research can help designing extensions of the O-RAN architecture toward 6G networks. For example, what we define as dApps in [213] can enable real-time control in the RAN nodes. These elements can work together with xApps to leverage data that cannot be transferred for analysis from the RAN to the RIC (e.g., I/Q samples, or fine-grained channel estimation information). Another extension can combine the centralized control of the O-RAN architecture with cell-free cellular networks (i.e., a version of massive MIMO with distributed antennas and centralized processing) [214], providing support for coordination across the different antenna endpoints.
- Energy efficiency with Open RAN: As discussed in Section II, virtualization and closed-loop control provide useful primitives for the dynamic network function allocation and thus for the energy efficiency maximization. Further research is required to develop orchestration routines at the non-RT RIC/SMO that embed energy efficiency in the optimization goal, as well as xApps and rApps that adopt control actions or policies that include energy efficiency targets.
- Security in O-RAN: As discussed in Section VIII, the openness of the RAN increases the threat surface but also

POLESE et al : UNDERSTANDING O-RAN

enables new approaches to network security. For example, the improved visibility into the RAN performance and telemetry and the possibility of deploying plug-in xApps and rApps for security analysis and threat identification make it possible to explore novel approaches for securing wireless networks and make them more robust and resilient. The research and development of security approaches that leverage O-RAN capabilities and improve the integrity, resiliency, and availability of its deployments is a key step toward making Open RAN approaches a viable and future-proof alternative to traditional RAN deployments.

D. Intelligent Control

- Multi-time-scale control: When considering the full O-RAN architecture (and possible extensions as discussed above), different control loops will operate and have visibility on the system at different time scales. This opens challenges in terms of multi-scale control. Further research is required on the design of the multi-scale algorithms, on the identification of instability in the system as well as conflicts across the different control loops, and on the automated selection of the optimal control loops that can be used to reach specific high-level intents [131].
- Effective AI/ML algorithm design, testing, and deployment: The AI/ML workflow described in Section VI positions O-RAN to be a framework for the practical deployment of ML solutions in the RAN. While this workflow is being specified, several challenges still remain. They are related to how to (i) collect training and testing datasets that are heterogeneous and representative of large-scale deployments; (ii) test and/or refine data-driven solutions through online training without compromising the production RAN performance; and (iii) design AI/ML algorithms that work with real, unreliable input, and can easily generalize to different deployment conditions [68].

XII. CONCLUSION

This paper presented a comprehensive overview of the O-RAN specifications, architectures and operations. We first introduced the main architectural building blocks and the principles behind the design of O-RAN networks. Then, we described the components of the near-RT and non-RT RICs and of the SMO, and discussed the O-RAN interfaces, including E2, O1, A1, the fronthaul interface, and O2. The second part of this paper focused on topics spanning multiple components and interfaces in the O-RAN architecture. We provided details on the AI and ML workflow that O-RAN enables, on O-RAN use cases and research, and on the O-RAN security challenges and potential. Finally, we reviewed the structure and development efforts of the O-RAN Alliance, and discussed research platforms, and future research directions.

We believe that these insights, together with the deep dive on the O-RAN specifications, architecture, and interfaces, will foster and promote further efforts toward more open, programmable, virtualized, and efficient wireless networks.

APPENDIX

1403

In the following appendices, we provide examples for messages exchanged over the E2 interface (the subscription message, in Appendix A, and the Indication message of type report, in Appendix B), and a list of acronyms used throughout this work (Appendix C).

A. Example of E2 Subscription Request Message

Listing 1 shows an example of the fields of an E2AP message for a subscription request, which is generated in the near-RT RIC and sent to the E2 termination of an E2 node. The XML format is used to describe the set of fields and their entries (i.e., the IEs), then the actual message is encoded in ASN.1 format (i.e., a sequence of bytes) before being encapsulated and transmitted on the SCTP socket. The first field is the message type, which contains a procedure code (8 for the subscription) and the actual type of message (an initiating message, as it begins a procedure on E2). Then, the message contains two IDs, one that uniquely identifies the RIC request, and the other the RAN function to which the RIC wants to subscribe. The core of the message is the set of IEs with details on the subscription request. In particular, the message defines a list of actions (with only one entry in this example). Each action has a type (in this case, report), a definition (which is optional, and depends on what the specific RAN function supports), and, possibly, a subsequent action to perform once the first is completed (in this case, continue, after waiting for a timer of 10 ms to expire).

B. Example of E2 Indication (Report) Message

Listing 2 features the XML for an E2AP Indication message, which is generated by the E2 node and sent to the near-RT RIC upon triggering of an event defined through the related subscription procedure. As for the E2 Subscription Request message, also this message is encoded using ASN.1, and features the message type and procedure code at the beginning of the message. Then, it lists IEs related to the RIC request, the function identifier, the corresponding action identifier (which is a unique value for each RIC request), a sequence number (which is optional), and the type of indication, which in this case is report (the possible values are insert or report). The last three fields (indication header, indication message, and call process identifier) are encoded by the E2AP message but their semantics are defined by the specific E2SM used to populate the message (e.g., E2SM KPM).

C. Acronyms

3GPP	3rd Generation Partnership Project
AAL	Acceleration Abstraction Layer
AI	Artificial Intelligence
AMF	Access and Mobility Management Function
AP	Application Protocol
API	Application Programming Interface
AR	Augmented Reality
ASIC	Application-specific Integrated Circuit
C-RAN	Cloud RAN

```
<E2AP-PDU>
       <initiatingMessage>
           cedureCode>8
           <criticality><ignore/></criticality>
           <value>
               <RICsubscriptionRequest>
                   otocolIEs>
                       <RICsubscriptionRequest-IEs>
                           <id>29</id>
                           <criticality><reject/></criticality>
10
                           <value>
11
                                <RICrequestID>
12
                                    <ricRequestorID>123</ricRequestorID>
13
                                    <ricInstanceID>34</ricInstanceID>
14
15
                                </RICrequestID>
                           </value>
16
17
                       </RICsubscriptionRequest -IEs>
18
                       <RICsubscriptionRequest-IEs>
19
                           <id>5</id>
                           <criticality><reject/></criticality>
20
21
                           <value>
22
                                <RANfunctionID>1</RANfunctionID>
                           </value>
23
                       </RICsubscriptionRequest -IEs>
24
                       <RICsubscriptionRequest-IEs>
25
                           <id>30</id>
26
27
                           <criticality><reject/></criticality>
28
                           <value>
                                <RICsubscriptionDetails>
29
30
                                    <ricEventTriggerDefinition>31 32 33 34</ricEventTriggerDefinition>
                                    <ri><ricAction-ToBeSetup-List>
31
32
                                        <ProtocolIE - SingleContainer>
33
                                            <id>19</id>
34
                                            <criticality><ignore/></criticality>
35
                                            <value>
                                                <RICaction-ToBeSetup-Item>
36
                                                     <ricActionID>1</ricActionID>
37
                                                     <ricActionType><report/></ricActionType>
38
                                                     <ricActionDefinition>35 36 37 38</ricActionDefinition>
39
                                                     <ricSubsequentAction>
40
                                                         <ricSubsequentActionType><continue/></ricSubsequentActionType>
41
42
                                                         <ricTimeToWait><w10ms/></ricTimeToWait>
43
                                                     </ricSubsequentAction>
                                                </RICaction-ToBeSetup-Item>
44
45
                                            </value>
                                        </ProtocolIE - SingleContainer>
47
                                    </ricaction-ToBeSetup-List>
48
                                </RICsubscriptionDetails>
49
                           </value>
                       </RICsubscriptionRequest -IEs>
50
51
                   </RICsubscriptionRequest>
52
           </value>
53
       </initiatingMessage>
54
  </E2AP-PDU>
55
```

Listing 1. Example of E2 Subscription Request message, compliant with E2AP V2.0 [82]. Generated using the E2 simulator library from [215].

CA	Carrier Aggregation	FCAPS	Fault, Configuration, Accounting, Performance,		
CCC	Cell Configuration and Control		Security		
CMP	Certificate Management Protocol	FEC	Forward Error Correction		
CP	Control Plane	FFT	Fast Fourier Transform		
CP	Cyclic Prefix	FG	Focus Group		
CQI	Channel Quality Information	FH	Fronthaul		
CSI	Channel State Information	FPGA	Field Programmable Gate Array		
CU	Central Unit	gNB	Next Generation Node Base		
DC	Dual Connectivity	GPU	Graphics Processing Unit		
DRL	Deep Reinforcement Learning	ΙE	Information Element		
DTLS	Datagram Transport Layer Security	IETF	Internet Engineering Task Force		
DU	Distributed Unit	IoT	Internet of Things		
E2SM	E2 Service Model	ITU	International Telecommunication Union		
EI	Enrichment Information	KPI	Key Performance Indicator		
eMBB	Enhanced Mobile Broadband	KPM	Key Performance Measurement		
eNB	evolved Node Base	LAA	Licensed-Assisted Access		
ETSI	European Telecommunications Standards	LTE	Long Term Evolution		
	Institute	LXC	Linux Container		

```
<E2AP-PDU>
       <initiatingMessage>
           code > 5 
           <criticality><ignore/></criticality>
           <value>
               <RICindication>
                   otocolIEs>
                       <RICindication - IEs>
                            <id>29</id>
10
                            <criticality><reject/></criticality>
11
                            <value>
                                <RICrequestID>
12
13
                                    <ricRequestorID>123</ricRequestorID>
                                    <ricInstanceID>26</ricInstanceID>
                                </RICrequestID>
                            </value>
                        </RICindication-IEs>
17
18
                        <RICindication-IEs>
19
                            <id>5</id>
                            <criticality><reject/></criticality>
20
21
                            <value>
                                <RANfunctionID>0</RANfunctionID>
22
                            </value>
23
                        </RICindication-IEs>
24
25
                        <RICindication-IEs>
                            <id>15</id>
26
27
                            <criticality><reject/></criticality>
                            <value>
29
                                <RICactionID>1</RICactionID>
30
                            </value>
                        </RICindication-IEs>
31
32
                        <RICindication - IEs>
33
                            <id>27</id>
                            <criticality><reject/></criticality>
34
35
                            <value>
                                <RICindicationSN>24</RICindicationSN>
36
                            </value>
37
                        </RICindication-IEs>
38
                        <RICindication-IEs>
39
                            <id>28</id>
40
                            <criticality><reject/></criticality>
41
43
                                <RICindicationType><report/></RICindicationType>
44
                            </value>
45
                        </RICindication-IEs>
46
                        <RICindication-IEs>
47
                            <id>25</id>
                            <criticality><reject/></criticality>
48
                            <value>
49
                                <RICindicationHeader>
50
                                    <--- E2SM Header --->
51
                                </RICindicationHeader>
52
                            </value>
53
                        </RICindication-IEs>
54
                        <RICindication-IEs>
55
                            <id>26</id>
57
                            <criticality><reject/></criticality>
58
                            <value>
59
                                <RICindicationMessage>
60
                                    <--- E2SM Payload --->
                                </RICindicationMessage>
61
                            </value>
62
                        </RICindication-IEs>
63
                        <RICindication - IEs>
64
                            <id>20</id>
65
                            <criticality><reject/></criticality>
66
67
                            <value>
                                <RICcallProcessID><--- E2SM process identifier ---></RICcallProcessID>
68
                            </value>
                        </RICindication-IEs>
                   </protocolIEs>
               </RICindication>
73
           </value>
74
       </initiatingMessage>
75 </E2AP-PDU>
```

Listing 2. Example of E2 Indication message of type report, compliant with E2AP V2.0 [82]. Generated using the E2 simulator library from [215].

MAC	Medium Access Control	MIMO	Multiple Input, Multiple Output
MCHEM	Massive Channel Emulator	ML	Machine Learning
MCS	Modulation and Coding Scheme	mMTC	Massive Machine-Type Communications
MEC	Multi-access Edge Computing	MnS	Management Services
MGEN	Multi-Generator	NETCONF	Network Configuration Protocol

UP

UPF

User Plane

User Plane Function

NFV	Network Function Virtualization
nGRG	next Generation Research Group
NI	Network Interfaces
NIB	Network Information Base
OAM	Operations, Administration and Maintenance
OFDM	Orthogonal Frequency Division Multiplexing
ONAP	Open Network Automation Platform
ONF	Open Networking Foundation
ONOS	Open Networking Operating System
OSC	O-RAN Software Community
OSFG	Open Source Focus Group
OSM	Open Source Management and Orchestration
OTIC	Open Testing & Integration Centre
PAWR	Platforms for Advanced Wireless Research
PDCP	Packet Data Convergence Protocol
PLFS	Physical Layer Frequency Signals
PNF	Physical Network Function
PRB	Physical Resource Block
PTP	Precision Time Protocol
QoE	Quality of Experience
QoS	Quality of Service
R-NIB	RAN NIB
RAN	Radio Access Network
RC	RAN Control
RF	Radio Frequency
RIC	RAN Intelligent Controller
RLC	Radio Link Control
RMR	RIC Message Router
RRC	Radio Resource Control
RSRP	Reference Signal Received Power
RU	Radio Unit
RX	Receiver
SCTP	Stream Control Transmission Protocol
SDAP	Service Data Adaptation Protocol
SDFG	Standard Development Focus Group
SDK	Software Development Kit
SDL	Shared Data Layer
SDN	Software-defined Networking
SDO	Standard-Development Organization
SDR	Software-defined Radio
SFG	Security Focus Group
SLA	Service Level Agreement
SM	Service Model
SMO	Service Management and Orchestration
SRN	Standard Radio Node
SSH	Secure Shell
SWG	Security Work Group
TB	Transport Block
TGEN	Traffic Generator
TIFG	
TIP	Test & Integration Focus Group Telecom Infra Project
TLS	Transport Layer Security
TSC	
TX	Technical Steering Committee Transmitter
UAV	Unmanned Aerial Vehicle
UE	User Equipment

URLLC	Ultra	Reliable	and	Low	Latency
	Commi	unications			
USRP	Univers	sal Software	Radio F	Peripheral	
VES	VNF E	vent Stream			
VM	Virtual	Machine			
VNF	Virtual	Network Fu	nction		
VR	Virtual	Reality			
WG	Workin	g Group.			

REFERENCES

- [1] M. Giordani, M. Polese, M. Mezzavilla, S. Rangan, and M. Zorzi, "Toward 6G networks: Use cases and technologies," *IEEE Commun. Mag.*, vol. 58, no. 3, pp. 55–61, Mar. 2020.
- [2] T. L. Marzetta, "Noncooperative cellular wireless with unlimited numbers of base station antennas," *IEEE Trans. Wireless Commun.*, vol. 9, no. 11, pp. 3590–3600, Nov. 2010.
- [3] I. F. Akyildiz, J. M. Jornet, and C. Han, "Terahertz band: Next frontier for wireless communications," *Phys. Commun.*, vol. 12, pp. 16–32, Sep. 2014. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S1874490714000238
- [4] A. Bourdoux et al., "6G white paper on localization and sensing," 2020, arXiv:2006.01779.
- [5] S. D'Oro, L. Bonati, F. Restuccia, and T. Melodia, "Coordinated 5G network slicing: How constructive interference can boost network throughput," *IEEE/ACM Trans. Netw.*, vol. 29, no. 4, pp. 1881–1894, Aug. 2021.
- [6] S. D'Oro, F. Restuccia, and T. Melodia, "Toward operator-to-waveform 5G radio access network slicing," *IEEE Commun. Mag.*, vol. 58, no. 4, pp. 18–23, Apr. 2020.
- [7] S. D'Oro, F. Restuccia, A. Talamonti, and T. Melodia, "The slice is served: Enforcing radio access network slicing in virtualized 5G systems," in *Proc. IEEE INFOCOM*, Paris, France, May 2019, pp. 442–450.
- [8] S. D'Oro, F. Restuccia, T. Melodia, and S. Palazzo, "Low-complexity distributed radio access network slicing: Algorithms and experimental results," *IEEE/ACM Trans. Netw.*, vol. 26, no. 6, pp. 2815–2828, Dec. 2018.
- [9] T. O'Shea and J. Hoydis, "An introduction to deep learning for the physical layer," *IEEE Trans. Cogn. Commun. Netw.*, vol. 3, no. 4, pp. 563–575, Dec. 2017.
- [10] "The open future of radio access networks." Deloitte—Telecom Engineering Centre of Excellence (TEE). 2021. [Online]. Available: https://www2.deloitte.com/content/dam/Deloitte/pt/Documents/ technology-media-telecommunications/TEE/The-Open-Future-of-Radio-Access-Networks.pdf
- [11] U. Challita, H. Ryden, and H. Tullberg, "When machine learning meets wireless cellular networks: Deployment, challenges, and applications," *IEEE Commun. Mag.*, vol. 58, no. 6, pp. 12–18, Jun. 2020.
- [12] L. Bonati, S. D'Oro, M. Polese, S. Basagni, and T. Melodia, "Intelligence and learning in O-RAN for data-driven NextG cellular networks," *IEEE Commun. Mag.*, vol. 59, no. 10, pp. 21–27, Oct. 2021.
- [13] "Study on new radio access technology: Radio access architecture and interfaces, version 14.0.0," 3rd Gener. Partnership Project (3GPP), Sophia Antipolis, France, Rep. (TR) 38.801, Apr. 2017. [Online]. Available: http://www.3gpp.org/DynaReport/38801.htm
- [14] L. Bonati, M. Polese, S. D'Oro, S. Basagni, and T. Melodia, "Open, programmable, and virtualized 5G networks: State-of-the-art and the road ahead," *Comput. Netw.*, vol. 182, pp. 1–28, Dec. 2020.
- [15] H. Lee, J. Cha, D. Kwon, M. Jeong, and I. Park, "Hosting AI/ML workflows on O-RAN RIC platform," in *Proc. IEEE GLOBECOM Workshops*, Dec. 2020, pp. 1–6.
- [16] A. S. Abdalla, P. S. Upadhyaya, V. K. Shah, and V. Marojevic, "Toward next generation open radio access network—what O-RAN can and cannot do!" Nov. 2021, arXiv:2111.13754.
- [17] A. Garcia-Saavedra and X. Costa-Pérez, "O-RAN: Disrupting the virtualized RAN ecosystem," *IEEE Commun. Standards Mag.*, vol. 5, no. 4, pp. 96–103, Dec. 2021.
- [18] B. Brik, K. Boutiba, and A. Ksentini, "Deep learning for B5G open radio access network: Evolution, survey, case studies, and challenges," *IEEE Open J. Commun. Soc.*, vol. 3, pp. 228–250, 2022.

- [19] A. Arnaz, J. Lipman, M. Abolhasan, and M. Hiltunen, "Toward integrating intelligence and programmability in open radio access networks: A comprehensive survey," *IEEE Access*, vol. 10, pp. 67747–67770, 2022.
- [20] D. Wypiór, M. Klinkowski, and I. Michalski, "Open RAN—Radio access network evolution, benefits and market trends," *Appl. Sci.*, vol. 12, no. 1, p. 408, 2022. [Online]. Available: https://www.mdpi. com/2076-3417/12/1/408
- [21] M. Dryjanski and R. Lundberg, "The O-RAN whitepaper," Rimedo Labs, Poznań, Poland, White Paper, 2021. [Online]. Available: https:// www.rimedolabs.com/blog/the-o-ran-whitepaper
- [22] T. Nolle, "Transformation and 5G O-RAN," VMWare, Palo Alto, CA, USA, White Paper, 2021. [Online]. Available: https://telco.vmware.com/content/dam/digitalmarketing/vmware/en/pdf/microsites/telco/vmware-transformation-and-5g-o-ran.pdf
- [23] "5G Open RAN ecosystem whitepaper," NTT Docomo, Tokyo, Japan, White Paper, 2021. [Online]. Available: https://ssw.web.docomo.ne.jp/ orec/5g_open_ran_ecosystem/whitepaper/OREC_WP.pdf
- [24] "Accelerating 5G virtual RAN deployment," Comcores, Lyngby, Denmark, White Paper, 2020. [Online]. Available: https://www. comcores.com/wp-content/uploads/2020/07/O-RAN-Intro-whitepaper-1.pdf
- [25] "Dell Technologies, VMware, and Mavenir 5G O-RAN reference architecture guide." 2021. [Online]. Available: https://tinyurl.com/39bj9z4n
- [26] "O-RAN: An open ecosystem to power 5G applications," VIAVI Solutiuons, Scottsdale, AZ, USA, White Paper, 2021. [Online]. Available: https://www.viavisolutions.com/en-us/literature/o-ran-open-ecosystem-power-5g-applications-white-papers-books-en.pdf
- [27] "Security in Open RAN," Altiostar, Tewksbury, MA, USA, White Paper, Jan. 2021. [Online]. Available: http://altiostar.com/wp-content/ uploads/2021/02/Open-RAN-Security-White-Paper-January-2021.pdf
- [28] S. P. Jason S. Boswell, "Security considerations of Open RAN," Ericsson, Stockholm, Sweden, White Paper, Aug. 2021. [Online]. Available: https://www.ericsson.com/4a67b7/assets/local/reports-papers/further-insights/doc/02092021-12911-security-considerations-for-cloud-ran.pdf
- [29] "Open RAN security white paper," Deutsche Telekom, Bonn, Germany, Orange, Paris, France, Telefónica, Madrid, Spain, TIM, Rome, Italy, and Vodafone, Berkshire, U.K., White Paper, 2022. [Online]. Available: https://cdn.brandfolder.io/D8DI15S7/at/45zqtkzjp4n9ncn77mkqbx8/Open_RAN_MoU_Security_White_Paper_-FV.pdf
- [30] S. Teral, "RIC as the next generation SON for Open RAN and more," LightCounting, Eugene, OR, USA, White Paper, 2021. [Online]. Available: https://gsma.force.com/mwcoem/servlet/servlet. FileDownload?file=00P6900002qUSTHEA4
- [31] "OpenRAN RAN intelligence and automation." Telecom Infra Project (TIP). 2021. Accessed: Dec. 2021. [Online]. Available: https://cdn.brandfolder.io/D8DI15S7/at/xq2qrcwgszxpb49zt93bwt/ RIA_OpenRAN_ataglance_Glossy_v08_2021_06_16.pdf
- [32] G. Brown. "TIP OpenRAN: Toward disaggregated mobile networking." May 2020. Accessed: Jul. 2020. [Online]. Available: https://cdn.brandfolder.io/D8DI15S7/as/qc19tk-54bsw-305pae/TIP_OpenRAN_-Heavy_Reading_May_2020-_White_Paper.pdf
- [33] N. McKeown et al., "OpenFlow: Enabling innovation in campus networks," ACM SIGCOMM Comput. Commun. Rev., vol. 38, no. 2, pp. 69–74, Mar. 2008.
- [34] "xRAN forum merges with C-RAN alliance to form ORAN alliance." xRAN Forum. 2018. [Online]. Available: https://www.businesswire. com/news/home/20180227005673/en/
- [35] C.-L. I, J. Huang, R. Duan, C. Cui, J. Jiang, and L. Li, "Recent progress on C-RAN centralization and cloudification," *IEEE Access*, vol. 2, pp. 1030–1039, 2014.
- [36] A. Checko et al., "Cloud RAN for mobile networks—A technology overview," *IEEE Commun. Surveys Tuts.*, vol. 17, no. 1, pp. 405–426, 1st Quart., 2015.
- [37] O-RAN Working Group 1, "O-RAN architecture description 5.00," O-RAN, Alfter, Germany, document O-RAN.WG1.O-RAN-Architecture-Description-v05.00 Technical Specification, Jul. 2021.
- [38] NG-RAN; Architecture Description, Version 17.0.0, 3GPP Standard (TS) 38.401, Apr. 2022. [Online]. Available: http://www.3gpp.org/ DynaReport/38401.htm
- [39] F. W. Murti, J. A. Ayala-Romero, A. Garcia-Saavedra, X. Costa-Pérez, and G. Iosifidis, "An optimal deployment framework for multi-cloud virtualized radio access networks," *IEEE Trans. Wireless Commun.*, vol. 20, no. 4, pp. 2251–2265, Apr. 2021.

- [40] NR; Radio Link Control (RLC) Protocol Specification, Version 15.0.0, 3GPP Standard (TS) 38.322, Jan. 2018. [Online]. Available: http://www.3gpp.org/DynaReport/38322.htm
- [41] NR; Medium Access Control (MAC) Protocol Specification, Version 15.0.0 3GPP Standard (TS) 38.321, Jan. 2018. [Online]. Available: http://www.3gpp.org/DynaReport/38321.htm
- [42] NR; Physical Layer; General Description, Version 15.0.0, 3GPP Standard (TS) 38.201, Jan. 2018. [Online]. Available: http://www.3gpp. org/DynaReport/38201.htm
- [43] NR; Radio Resource Control (RRC); Protocol Specification, Version 15.0.0, 3GPP Standard (TS) 38.331, Jan. 2018. [Online]. Available: http://www.3gpp.org/DynaReport/38331.htm
- [44] Evolved Universal Terrestrial Radio Access (E-UTRA) and NR; Service Data Adaptation Protocol (SDAP) Specification, Version 17.0.0, 3GPP Standard (TS) 37.324, Apr. 2022. [Online]. Available: http://www.3gpp. org/DynaReport/37324.htm
- [45] NR; Packet Data Convergence Protocol (PDCP) Specification, Version 15.0.0, 3GPP Standard (TS) 38.323, Jan. 2018. [Online]. Available: http://www.3gpp.org/DynaReport/38323.htm
- [46] O-RAN Working Group 2, "O-RAN non-RT RIC architecture 1.0," O-RAN, Alfter, Germany, document O-RAN.WG2.Non-RT-RIC-ARCH-TS-v01.00 Technical Specification, Jul. 2021.
- [47] O-RAN Working Group 2, "O-RAN non-RT RIC: Functional architecture 1.01," O-RAN, Alfter, Germany, document O-RAN.WG2.Non-RT-RIC-ARCH-TS-v01.00 Technical Specification, Jul. 2021.
- [48] O-RAN Working Group 2, "O-RAN non-RT RIC & A1 interface: Use cases and requirements 4.00," O-RAN, Alfter, Germany, document O-RAN.WG2.Use-Case-Requirements-v04.00 Technical Specification, Jul. 2021.
- [49] O-RAN Working Group 2, "O-RAN AI/ML workflow description and requirements 1.03," O-RAN, Alfter, Germany, document O-RAN.WG2.AIML-v01.03 Technical Specification, Jul. 2021.
- [50] O-RAN Working Group 3, "O-RAN near-RT RAN intelligent controller near-RT RIC architecture 2.00," O-RAN, Alfter, Germany, document O-RAN.WG3.RICARCH-v02.00, Mar. 2021.
- [51] Evolved Universal Terrestrial Radio Access (E-UTRA) and Evolved Universal Terrestrial Radio Access Network (E-UTRAN); Overall Description; Stage 2, Version 14.2.0, 3GPP Standard (TS) 36.300, Apr. 2017. [Online]. Available: http://www.3gpp.org/DynaReport/ 36300.htm
- [52] O-RAN Working Group 6, "O-RAN cloud architecture and deployment scenarios for O-RAN virtualized RAN 2.02," O-RAN, Alfter, Germany, document O-RAN.WG6.CAD-v02.02 Technical Specification, Jul. 2021.
- [53] R. Jain and S. Paul, "Network virtualization and software defined networking for cloud computing: A survey," *IEEE Commun. Mag.*, vol. 51, no. 11, pp. 24–31, Nov. 2013.
- [54] O-RAN Working Group 6, "O-RAN acceleration abstraction layer FEC profiles 1.0," O-RAN, Alfter, Germany, document O-RAN.WG6.AAL-FEC.0-v01.00 Technical Specification, Jul. 2021.
- [55] O-RAN Working Group 6, "O-RAN acceleration abstraction layer general aspects an principles 1.01," O-RAN, Alfter, Germany, document O-RAN.WG6.AAL-GAnP-v01.01 Technical Specification, Jul. 2021.
- [56] A. Nasrallah et al., "Ultra-low latency (ULL) networks: The IEEE TSN and IETF DetNet standards and related 5G ULL research," *IEEE Commun. Surveys Tuts.*, vol. 21, no. 1, pp. 88–145, 1st Quart., 2019.
- [57] A. Kelkar and C. Dick, "NVIDIA aerial GPU hosted AI-on-5G," in Proc. IEEE 4th 5G World Forum (5GWF), Oct. 2021, pp. 64–69.
- [58] G. Garcia-Aviles, A. Garcia-Saavedra, M. Gramaglia, X. Costa-Perez, P. Serrano, and A. Banchs, "Nuberu: Reliable RAN virtualization in shared platforms," in *Proc. 27th Annu. Int. Conf. Mobile Comput. Netw.*, 2021, pp. 749–761.
- [59] J. S. Panchal, S. Subramanian, and R. Cavatur, "Enabling and scaling of URLLC verticals on 5G vRAN running on COTS hardware," *IEEE Commun. Mag.*, vol. 59, no. 9, pp. 105–111, Sep. 2021.
- [60] E. A. Papatheofanous, D. Reisis, and K. Nikitopoulos, "LDPC hard-ware acceleration in 5G open radio access network platforms," *IEEE Access*, vol. 9, pp. 152960–152971, 2021.
- [61] D. Sabella et al., "Energy efficiency benefits of RAN-as-a-service concept for a cloud-based 5G mobile network infrastructure," *IEEE Access*, vol. 2, pp. 1586–1597, 2014.
- [62] T. Pamuklu, M. Erol-Kantarci, and C. Ersoy, "Reinforcement learning based dynamic function splitting in disaggregated green Open RANs," in *Proc. IEEE ICC*, Montreal, QC, Canada, Jun. 2021, pp. 1–6.
- [63] J. Luo, Q. Chen, and L. Tang, "Reducing power consumption by joint sleeping strategy and power control in delay-aware C-RAN," *IEEE Access*, vol. 6, pp. 14655–14667, 2018.

- [64] D. López-Pérez et al., "A survey on 5G radio access network energy efficiency: Massive MIMO, lean carrier design, sleep modes, and machine learning," *IEEE Commun. Surveys Tuts.*, vol. 24, no. 1, pp. 653–697, 1st Quart., 2022.
- [65] "D release." O-RAN Software Community. Accessed: Jul. 2021. [Online]. Available: https://wiki.o-ran-sc.org/pages/viewpage.action? pageId=20878658
- [66] "RIC message router documentation." O-RAN Software Community. Accessed: Feb. 2021. [Online]. Available: https://docs.o-ran-sc.org/projects/o-ran-sc-ric-plt-lib-rmr
- [67] J. Carlson, Redis in Action. New York, NY, USA: Simon and Schuster, 2013
- [68] M. Polese, L. Bonati, S. D'Oro, S. Basagni, and T. Melodia, "ColO-RAN: Developing machine learning-based xApps for Open RAN closed-loop control on programmable experimental platforms," *IEEE Trans. Mobile Comput.*, early access, Jul. 4, 2022, doi: 10.1109/TMC.2022.3188013.
- [69] "ONF's software-defined RAN platform consistent with the O-RAN architecture," ONF, Palo Alto, CA, USA, White Paper, Feb. 2020. [Online]. Available: https://www.opennetworking.org/wp-content/uploads/2020/03/SD-RAN-White-Paper.pdf
- [70] "SD-RAN." ONF. Accessed: Dec. 2021. [Online]. Available: https:// opennetworking.org/sd-ran
- [71] R. Schmidt, M. Irazabal, and N. Nikaein, "FlexRIC: An SDK for next-generation SD-RANs," in *Proc. ACM CoNEXT*, Dec. 2021, pp. 411–425.
- [72] "5G-EmPOWER." Fondazione Bruno Kessler (FBK). Accessed: Dec. 2021. [Online]. Available: https://5g-empower.io
- [73] "Non-RealTime RIC (NNONRTRIC)." O-RAN Software Community. Accessed: Dec. 2021. [Online]. Available: https://wiki.o-ran-sc.org/pages/viewpage.action?pageId=3604819
- [74] "ONAP 5G blueprint overview." ONAP. 2019. [Online]. Available: https://www.onap.org/wp-content/uploads/sites/20/2019/ 07/ONAP_CaseSolution_5G_062519.pdf
- [75] "Architecture overview." ONAP. 2019. [Online]. Available: https://www.onap.org/wp-content/uploads/sites/20/2019/07/ ONAP_CaseSolution_Architecture_062519.pdf
- [76] "OSM scope, functionality, operation and integration guidelines." Open Source MANO End User Advisory Group. Feb. 2019. [Online]. Available: https://osm.etsi.org/images/ OSM_EUAG_White_Paper_OSM_Scope_and_Functionality.pdf
- [77] "ONAP." ONF. Accessed: Jul. 2020. [Online]. Available: https://onap.org
- [78] "ETSI enjoy magazine—July 2021." ETSI. Accessed: Dec. 2021. [Online]. Available: https://www.etsi.org/images/files/Magazine/ ETSI_Enjoy_MAG_2021_N03_July.pdf
- [79] S. Pal and R. Armada. "OSM for 5G O-RAN." 2020. [Online]. Available: http://osm-download.etsi.org/ftp/osm-8.0-eight/OSM10-hackfest/EcosystemDay/E%2%20O%20f%20%20O-R%20(Altran).pdf
- [80] O-RAN Working Group 3, "O-RAN near-real-time RAN intelligent controller architecture & E2 general aspects and principles 2.00," O-RAN, Alfter, Germany, document O-RAN.WG3.E2GAP-v02.01 Technical Specification, Jul. 2021.
- [81] System Architecture for the 5G System (5GS), Version 16.4.0, 3GPP Standard (TS) 23.501, Mar. 2020. [Online]. Available: http://www. 3gpp.org/DynaReport/23501.htm
- [82] O-RAN Working Group 3, "O-RAN near-real-time RAN intelligent controller, E2 application protocol 2.00," O-RAN, Alfter, Germany, document O-RAN.WG3.E2AP-v02.00 Technical Specification, Jul 2020
- [83] O-RAN Working Group 3, "O-RAN near-real-time RAN intelligent controller E2 service model 2.00," O-RAN, Alfter, Germany, document ORAN-WG3.E2SM-v02.00 Technical Specification, Jul. 2021.
- [84] R. Stewart and C. Metz, "SCTP: New transport protocol for TCP/IP," IEEE Internet Comput., vol. 5, no. 6, pp. 64–69, Nov./Dec. 2001.
- [85] J. Larmouth, ASN. 1 Complete. San Francisco, CA, USA: Morgan Kaufmann, 2000.
- [86] O-RAN Working Group 3, "O-RAN near-real-time RAN intelligent controller E2 service model (E2SM) KPM 2.0," O-RAN, Alfter, Germany, document ORAN-WG3.E2SM-KPM-v02.00 Technical Specification, Jul. 2021.
- [87] O-RAN Working Group 3, "O-RAN near-real-time RAN intelligent controller E2 service model (E2SM), RAN function network interface (NI) 1.0," O-RAN, Alfter, Germany, document ORAN-WG3.E2SM-NI-v01.00.00 Technical Specification, Feb. 2020.

- [88] O-RAN Working Group 3, "O-RAN E2 service model (E2SM), cell configuration and control 1.0," O-RAN, Alfter, Germany, document O-RAN.WG3.E2SM-CCC-v01.00 Technical Specification, Oct. 2022.
- [89] O-RAN Working Group 3, "O-RAN near-real-time RAN intelligent controller E2 service model, RAN control 1.03," O-RAN, Alfter, Germany, document ORAN-WG3.E2SM-RC-v01.03 Technical Specification, Oct. 2022.
- [90] Management and Orchestration; 5G Performance Measurements, Version 17.6.0, 3GPP Standard (TS) 28.552, Mar. 2022. [Online]. Available: http://www.3gpp.org/DynaReport/28552.htm
- [91] Telecommunication Management; Performance Management (PM); Performance Measurements Evolved Universal Terrestrial Radio Access Network (E-UTRAN), Version 17.1.0, 3GPP Standard (TS) 32.425, Jun. 2021. [Online]. Available: http://www.3gpp.org/ DynaReport/32425.htm
- [92] NG-RAN; Xn Application Protocol (XnAP), Version 17.0.0, 3GPP Standard (TS) 38.423, Apr. 2022. [Online]. Available: http://www.3gpp. org/DynaReport/38423.htm
- [93] O-RAN Working Group 1, "O-RAN operations and maintenance interface 4.0," O-RAN, Alfter, Germany, document O-RAN.WG1.O1-Interface.0-v04.00 Technical Specification, Nov. 2020.
- [94] Telecommunication Management; Generic Network Resource Model (NRM) Integration Reference Point (IRP); Information Service (IS), Version 17.1.1, 3GPP Standard (TS) 28.622, Mar. 2022. [Online]. Available: http://www.3gpp.org/DynaReport/28622.htm
- [95] R. Enns, M. Bjorklund, J. Schoenwaelder, and A. Bierman, "Network configuration protocol (NETCONF)," IETF, RFC 6241, Jun. 2011. [Online]. Available: http://www.rfc-editor.org/rfc/rfc6241.txt
- [96] Management and Orchestration; Provisioning, Version 17.3.0, 3GPP Standard (TS) 28.531, Mar. 2022. [Online]. Available: http://www. 3gpp.org/DynaReport/28531.htm
- [97] Management and Orchestration; Generic Management Services, Version 17.0.0, 3GPP Standard (TS) 28.532, Mar. 2022. [Online]. Available: http://www.3gpp.org/DynaReport/28532.htm
- [98] Management and Orchestration; Fault Supervision (FS), Version 17.0.0, 3GPP Standard (TS) 28.545, Jun. 2021. [Online]. Available: http://www.3gpp.org/DynaReport/28545.htm
- [99] Telecommunication Management; Generic Network Resource Model (NRM) Integration Reference Point (IRP); Solution Set (SS) Definitions, Version 17.1.0, 3GPP Standard (TS) 28.623, Mar. 2022. [Online]. Available: http://www.3gpp.org/DynaReport/28623.htm
- [100] Management and Orchestration; Performance Assurance, Version 17.0.0, 3GPP Standard (TS) 28.550, Mar. 2022. [Online]. Available: http://www.3gpp.org/DynaReport/28550.htm
- [101] O-RAN Working Group 2, "A1 interface: General aspects and principles 2.03," O-RAN, Alfter, Germany, document ORAN-WG2.A1.GAPv02.03 Technical Specification, Jul. 2021.
- [102] O-RAN Working Group 2, "A1 interface: Application protocol 3.01," O-RAN, Alfter, Germany, document O-RAN.WG2.A1APv03.01 Technical Specification, Mar. 2021.
- [103] O-RAN Working Group 2, "A1 interface: Type definitions 2.0," O-RAN, Alfter, Germany, document O-RAN.WG2.A1TD-v02.00 Technical Specification, Jul. 2021.
- [104] O-RAN Working Group 4, "O-RAN Fronthaul control, user and synchronization plane specification 7.0," O-RAN, Alfter, Germany, document ORAN-WG4.CUS.0-v07.00 Technical Specification, Jul. 2021.
- [105] O-RAN Working Group 4, "O-RAN management plane specification 7.0," O-RAN, Alfter, Germany, document O-RAN.WG4.MP.0-v07.00 Technical Specification, Jul. 2021.
- [106] D. Petrovic, W. Rave, and G. Fettweis, "Effects of phase noise on OFDM systems with and without PLL: Characterization and compensation," *IEEE Trans. Commun.*, vol. 55, no. 8, pp. 1607–1616, Aug. 2007.
- [107] "Common public radio interface: ECPRI specification V2.0," CPRI Consortium, Bengaluru, India, document O-RAN.SFG.O-RAN-Security-Protocols-Specifications-v02.00 Technical Specification, May 2019.
- [108] IEEE SA, IEEE Standard for Radio Over Ethernet Encapsulations And Mappings, IEEE Standard 1914.3-2018, Oct. 2018.
- [109] S. Lagen, L. Giupponi, A. Hansson, and X. Gelabert, "Modulation compression in next generation RAN: Air interface and fronthaul tradeoffs," *IEEE Commun. Mag.*, vol. 59, no. 1, pp. 89–95, Jan. 2021.
- [110] D. Chitimalla, K. Kondepu, L. Valcarenghi, M. Tornatore, and B. Mukherjee, "5G fronthaul-latency and jitter studies of CPRI over Ethernet," *J. Opt. Commun. Netw.*, vol. 9, no. 2, pp. 172–182, Feb. 2017.

- [111] O-RAN Working Group 4, "O-RAN management plane specification—YANG models 7.0," O-RAN, Alfter, Germany, document O-RAN.WG4.MP-YANGs-v07.00 Technical Specification, Jul. 2021.
- [112] M. Giordani, M. Polese, A. Roy, D. Castor, and M. Zorzi, "A tutorial on beam management for 3GPP NR at mmWave frequencies," *IEEE Commun. Surveys Tuts.*, vol. 21, no. 1, pp. 173–196, 1st Quart., 2019.
- [113] F. Gómez-Cuba, T. Zugno, J. Kim, M. Polese, S. Bahk, and M. Zorzi, "Hybrid beamforming in 5G mmWave networks: A full-stack perspective," *IEEE Trans. Wireless Commun.*, vol. 21, no. 2, pp. 1288–1303, Feb. 2022.
- [114] R. L. Scheiterer, C. Na, D. Obradovic, and G. Steindl, "Synchronization performance of the precision time protocol in industrial automation networks," *IEEE Trans. Instrum. Meas.*, vol. 58, pp. 1849–1857, 2009.
- [115] O-RAN Working Group 9, "Synchronization architecture and solution specification 2.0," O-RAN, Alfter, Germany, document ORAN-WG9.XTRP-SYN.0-v02.0, Mar. 2021.
- [116] "ETSI releases first O-RAN specification." ETSI. Sep. 2022. [Online]. Available: https://www.etsi.org/newsroom/press-releases/2120-2022-09-etsi-releases-first-o-ran-specification
- [117] O-RAN Working Group 6, "O-RAN O2 general aspects and principles specification 1.01," O-RAN, Alfter, Germany, document O-RAN.WG6.O2-GA&P-v01.01 Technical Specification, Jul. 2021.
- [118] NG-RAN; E1 General Aspects and Principles, Version 17.0.0, 3GPP Standard (TS) 38.460, Apr. 2022. [Online]. Available: http://www.3gpp. org/DynaReport/38460.htm
- [119] NG-RAN; F1 Application Protocol (F1AP), Version 17.0.0, 3GPP Standard (TS) 38.473, Apr. 2022. [Online]. Available: http://www.3gpp. org/DynaReport/38473.htm
- [120] Evolved Universal Terrestrial Radio Access Network (E-UTRAN); X2 Application Protocol (X2AP), Version 17.0.0, 3GPP Standard (TS) 36.423, Apr. 2022. [Online]. Available: http://www.3gpp.org/ DynaReport/36423.htm
- [121] N. Miloslavskaya and A. Tolstoy, "Big data, fast data and data lake concepts," *Procedia Comput. Sci.*, vol. 88, pp. 300–305, Oct. 2016.
- [122] M. M. Najafabadi, F. Villanustre, T. M. Khoshgoftaar, N. Seliya, R. Wald, and E. Muharemagic, "Deep learning applications and challenges in big data analytics," *J. Big Data*, vol. 2, no. 1, pp. 1–21, 2015.
- [123] S. García, J. Luengo, and F. Herrera, Data Preprocessing in Data Mining, vol. 72. Cham, Switzerland: Springer, 2015.
- [124] A. Ng, "Sparse autoencoder," Lecture Notes CS294A, Standford Univ., Stanford, CA, USA, 2011.
- [125] S. Levine, A. Kumar, G. Tucker, and J. Fu, "Offline reinforcement learning: Tutorial, review, and perspectives on open problems," 2020, arXiv:2005.01643.
- [126] R. Agarwal, D. Schuurmans, and M. Norouzi, "An optimistic perspective on offline reinforcement learning," in *Proc. 37th Int. Conf. Mach. Learn.*, vol. 119, Jul. 2020, pp. 104–114. [Online]. Available: https://proceedings.mlr.press/v119/agarwal20c.html
- [127] A. Nair, M. Dalal, A. Gupta, and S. Levine, "AWAC: Accelerating online reinforcement learning with offline datasets," 2020, arXiv:2006.09359.
- [128] C. Käding, E. Rodner, A. Freytag, and J. Denzler, "Fine-tuning deep neural networks in continuous learning scenarios," in *Proc. Asian Conf. Comput. Vis.*, 2016, pp. 588–605.
- [129] A. Zappone, M. Di Renzo, and M. Debbah, "Wireless networks design in the era of deep learning: Model-based, AI-based, or both?" *IEEE Trans. Commun.*, vol. 67, no. 10, pp. 7331–7376, Oct. 2019.
- [130] R. S. Sutton and A. G. Barto, Reinforcement Learning: An Introduction. Cambridge, MA, USA: MIT Press, 2018.
- [131] S. D'Oro, L. Bonati, M. Polese, and T. Melodia, "OrchestRAN: Network automation through orchestrated intelligence in the Open RAN," in *Proc. IEEE INFOCOM*, May 2022, pp. 270–279.
- [132] C. Ebert, G. Gallardo, J. Hernantes, and N. Serrano, "DevOps," *IEEE Softw.*, vol. 33, no. 3, pp. 94–100, May/Jun. 2016.
- [133] S. Mäkinen, H. Skogström, E. Laaksonen, and T. Mikkonen, "Who needs MLOps: What data scientists seek to accomplish and how can MLOps help?" in *Proc. IEEE/ACM 1st Workshop AI Eng. Softw. Eng. AI (WAIN)*, May 2021, pp. 109–112.
- [134] P. Li et al., "RLOps: Development life-cycle of reinforcement learning aided Open RAN," 2021, arXiv:2111.06978.
- [135] S. D'Oro, L. Bonati, F. Restuccia, M. Polese, M. Zorzi, and T. Melodia, "SI-EDGE: Network slicing at the edge," in *Proc. ACM Mobihoc*, Oct. 2020, pp. 1–10.

- [136] F. Z. Morais et al., "PlaceRAN: Optimal placement of virtualized network functions in beyond 5G radio access networks," *IEEE Trans. Mobile Comput.*, early access, Apr. 29, 2022, doi: 10.1109/TMC.2022.3171525.
- [137] S. Matoussi, I. Fajjari, S. Costanzo, N. Aitsaadi, and R. Langar, "5G RAN: Functional split orchestration optimization," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 7, pp. 1448–1463, Jul. 2020.
- [138] S. Mollahasani, M. Erol-Kantarci, and R. Wilson, "Dynamic CU-DU selection for resource allocation in O-RAN using actor-critic learning," Oct. 2021, arXiv:2110.00492.
- [139] T. T. Nguyen and V. J. Reddi, "Deep reinforcement learning for cyber security," *IEEE Trans. Neural Netw. Learn. Syst.*, early access, Nov. 1, 2021, doi: 10.1109/TNNLS.2021.3121870.
- [140] W. Soussi, M. Christopoulou, G. Xilouris, and G. Gür, "Moving target defense as a proactive defense element for beyond 5G," *IEEE Commun. Standards Mag.*, vol. 5, no. 3, pp. 72–79, Sep. 2021.
- [141] O-RAN Working Group 1, "O-RAN use cases detailed specification 6.0," O-RAN, Alfter, Germany, document O-RAN.WG1.Use-Cases-Detailed-Specification-v06.00 Technical Specification, Jul. 2021.
- [142] O-RAN Working Group 1, "O-RAN use cases analysis report 6.0," O-RAN, Alfter, Germany, document O-RAN.WG1.Use-Cases-Analysis-Report-v06.00 Technical Specification, Jul. 2021.
- [143] M. Polese, R. Jana, V. Kounev, K. Zhang, S. Deb, and M. Zorzi, "Machine learning at the edge: A data-driven architecture with applications to 5G cellular networks," *IEEE Trans. Mobile Comput.*, vol. 20, no. 12, pp. 3367–3382, Dec. 2021.
- [144] D. Johnson, D. Maas, and J. Van Der Merwe, "NexRAN: Closed-loop RAN slicing in POWDER—A top-to-bottom open-source open-RAN use case," in *Proc. ACM WiNTECH*, New Orleans, LA, USA, Oct. 2021, pp. 17–23.
- [145] D. Johnson, D. Maas, and J. Van Der Merwe, "Open source RAN slicing on POWDER: A top-to-bottom O-RAN use case," in *Proc. ACM MobiSys*, Jun. 2021, pp. 507–508.
- [146] E. Sarikaya and E. Onur, "Placement of 5G RAN slices in multi-tier O-RAN 5G networks with flexible functional splits," in *Proc. IEEE CNSM*, Izmir, Turkey, Oct. 2021, pp. 274–282.
- [147] S. Niknam et al., "Intelligent O-RAN for beyond 5G and 6G wireless networks," May 2020. arXiv:2005.08374.
- [148] F. Mungari, "An RL approach for radio resource management in the O-RAN architecture," in *Proc. IEEE SECON*, Rome, Italy, Jul. 2021, pp. 1–2.
- [149] E. Coronado, S. Siddiqui, and R. Riggio, "Roadrunner: O-RAN-based cell selection in beyond 5G networks," in *Proc. IEEE/IFIP Netw. Oper. Manag. Symp. (NOMS)*, Apr. 2022, pp. 1–7.
- [150] S.-Y. Lien, D.-J. Deng, and B.-C. Chang, "Session management for URLLC in 5G open radio access network: A machine learning approach," in *Proc. Int. Wireless Commun. Mobile Comput. (IWCMC)*, Jun. 2021, pp. 2050–2055.
- [151] A. Filali, B. Nour, S. Cherkaoui, and A. Kobbane, "Communication and computation O-RAN resource slicing for URLLC services using deep reinforcement learning," 2022, arXiv:2202.06439.
- [152] P. E. Iturria-Rivera, H. Zhang, H. Zhou, S. Mollahasani, and M. Erol-Kantarci, "Multi-agent team learning in virtualized open radio access networks (O-RAN)," *Sensors*, vol. 22, no. 14, p. 5375, 2022
- [153] L. Bertizzolo et al., "Streaming from the air: Enabling drone-sourced video streaming applications on 5G Open-RAN architectures," *IEEE Trans. Mobile Comput.*, early access, Nov. 18, 2021, doi: 10.1109/TMC.2021.3129094.
- [154] R. Smith, C. Freeberg, T. Machacek, and V. Ramaswamy, "An O-RAN approach to spectrum sharing between commercial 5G and government satellite systems," in *Proc. IEEE Mil. Commun. Conf. (MILCOM)*, Nov. 2021, pp. 739–744.
- [155] L. Baldesi, F. Restuccia, and T. Melodia, "ChARM: NextG spectrum sharing through data-driven real-time O-RAN dynamic control," in *Proc. IEEE INFOCOM*, May 2022, pp. 240–249.
- [156] Ł. Kułacz and A. Kliks, "Dynamic spectrum allocation using multisource context information in OpenRAN networks," *Sensors*, vol. 22, no. 9, p. 3515, 2022. [Online]. Available: https://www.mdpi.com/1424-8220/22/9/3515
- [157] M. Polese et al., "Dynamic spectrum sharing between active and passive users above 100 GHz," Commun. Eng., vol. 1, no. 1, pp. 1–9, 2022.
- [158] L. Giupponi and F. Wilhelmi, "Blockchain-enabled network sharing for O-RAN in 5G and beyond," 2021, arXiv:2107.02005.

- [159] M. Enayati et al., "Blockchain-based location sharing in 5G Open RAN infrastructure for sustainable communities," in *Intelligent Sustainable Systems*, A. K. Nagar, D. S. Jat, G. Marin-Raventos, and D. K. Mishra, Eds. Singapore: Springer Nat., 2022, pp. 571–585.
- [160] N. Suzuki, T. Yoshioka, A. Hasegawa, H. Yokoyama, and T. Maeyama, "Implementation and evaluation of spectrum sharing technology using smart contracts," in *Proc. IEEE 12th Annu. Comput. Commun. Workshop Conf. (CCWC)*, Jan. 2022, pp. 922–928.
- [161] S. Parsaeefard, R. Dawadi, M. Derakhshani, T. Le-Ngoc, and M. Baghani, "Dynamic resource allocation for virtualized wireless networks in massive-MIMO-aided and fronthaul-limited C-RAN," *IEEE Trans. Veh. Technol.*, vol. 66, no. 10, pp. 9512–9520, Oct. 2017.
- [162] T. Hewavithana, A. Chopra, B. Mondal, S. Wong, A. Davydov, and M. Majmundar, "Overcoming channel aging in massive MIMO basestations with Open RAN fronthaul," in *Proc. IEEE Wireless Commun. Netw. Conf. (WCNC)*, 2022, pp. 2577–2582.
- [163] S. Lagen, X. Gelabert, L. Giupponi, and A. Hansson, "Fronthaul-aware scheduling strategies for dynamic modulation compression in next generation RANs," *IEEE Trans. Mobile Comput.*, early access, Nov. 17, 2021, doi: 10.1109/TMC.2021.3128700.
- [164] M. Mohsin, J. M. Batalla, E. Pallis, G. Mastorakis, E. K. Markakis, and C. X. Mavromoustakis, "On analyzing beamforming implementation in O-RAN 5G," *Electronics*, vol. 10, no. 17, p. 2162, 2021.
- [165] I. Godor et al., "A look inside 5G standards to support time synchronization for smart manufacturing," *IEEE Commun. Standards Mag.*, vol. 4, no. 3, pp. 14–21, Sep. 2020.
- [166] R. Keating, M. Säily, J. Hulkkonen, and J. Karjalainen, "Overview of positioning in 5G new radio," in *Proc. 16th Int. Symp. Wireless Commun. Syst. (ISWCS)*, Aug. 2019, pp. 320–324.
- [167] A. Huff, M. Hiltunen, and E. P. Duarte, "RFT: Scalable and fault-tolerant microservices for the O-RAN control plane," in *Proc. IFIP/IEEE IM*, Bordeaux, France, May 2021, pp. 402–409.
- [168] I. Tamim, A. Saci, M. Jammal, and A. Shami, "Downtime-aware O-RAN VNF deployment strategy for optimized self-healing in the O-Cloud," Oct. 2021, arXiv:2110.06060.
- [169] X. Wang, J. D. Thomas, R. J. Piechocki, S. Kapoor, R. Santos-Rodríguez, and A. Parekh, "Self-play learning strategies for resource assignment in open-RAN networks," *Comput. Netw.*, vol. 206, Apr. 2022, Art. no. 108682.
- [170] C. T. Shen et al., "Security threat analysis and treatment strategy for ORAN," in *Proc. 24th Int. Conf. Adv. Commun. Technol. (ICACT)*, Feb. 2022, pp. 417–422.
- [171] O-RAN Working Group 11, "O-RAN security threat modeling and remediation analysis 4.0," O-RAN, Alfter, Germany, document O-RAN.WG11.O-RAN-Threat-Model-v04.00 Technical Specification, Jul. 2022.
- [172] O-RAN Working Group 11, "O-RAN security protocols specifications 4.0," O-RAN, Alfter, Germany, document O-RAN.WG11.Security-Protocols-Specifications-v04.00 Technical Specification, Jul. 2022.
- [173] O-RAN Working Group 11, "O-RAN security requirements specifications 4.0," O-RAN, Alfter, Germany, document O-RAN.WG11.Security-Requirements-Specifications-v04.00 Technical Specification, Oct. 2022.
- [174] "Cybersecurity of open radio access networks." European Commission. May 2022. [Online]. Available: https://digital-strategy.ec.europa.eu/en/library/cybersecurity-open-radio-access-networks
- [175] G. Mcgraw, R. Bonett, H. Figueroa, and V. Shepardson, "Security engineering for machine learning," *Computer*, vol. 52, no. 8, pp. 54–57, Aug. 2019.
- [176] D. Mimran et al., "Evaluating the security of open radio access networks," 2022, arXiv:2201.06080.
- [177] Catalogue of General Security Assurance Requirements, Version 17.0.0, 3GPP Standard (TS) 33.117, Jun. 2021. [Online]. Available: http://www.3gpp.org/DynaReport/33117.htm
- [178] Security Architecture and Procedures for 5G System, Version 17.5.0, 3GPP Standard (TS) 33.501, Mar. 2022. [Online]. Available: http://www.3gpp.org/DynaReport/33501.htm
- [179] Security Assurance Specification (SCAS) for the Next Generation Node B (gNodeB) Network Product Class, Version 17.1.0, 3GPP Standard (TS) 33.511, Dec. 2021. [Online]. Available: http://www.3gpp.org/ DynaReport/33511.htm
- [180] "Security assurance methodology (SECAM) and security assurance specification (SCAS) for 3GPP virtualized network products, version 17.1.0," 3rd Gener. Partnership Project (3GPP), Sophia Antipolis, France, Rep. (TR) 33.818, Sep. 2021. [Online]. Available: http://www. 3gpp.org/DynaReport/33818.htm

- [181] "Study on security impacts of virtualisation, version 0.12.0," 3rd Gener. Partnership Project (3GPP), Sophia Antipolis, France, Rep. (TR) 33.848, May 2022. [Online]. Available: http://www.3gpp.org/DynaReport/33848.htm
- [182] D. Chou and M. Jiang, "A survey on data-driven network intrusion detection," ACM Comput. Surv., vol. 54, no. 9, pp. 1–36, Oct. 2021. [Online]. Available: https://doi.org/10.1145/3472753
- [183] J. Jin, C. Lian, and M. Xu, "Rogue base station detection using a machine learning approach," in *Proc. 28th Wireless Opt. Commun. Conf. (WOCC)*, May 2019, pp. 1–5.
- [184] M. Savic et al., "Deep learning anomaly detection for cellular IoT with applications in smart logistics," *IEEE Access*, vol. 9, pp. 59406–59419, 2021.
- [185] "O-RAN WhitePaper—Building the next generation RAN," O-RAN Alliance, Alfter, Germany, White Paper, Oct. 2018. [Online]. Available: https://www.o-ran.org/resources
- [186] L. Bonati, M. Polese, S. D'Oro, S. Basagni, and T. Melodia, "OpenRAN gym: An open toolbox for data collection and experimentation with AI in O-RAN," in *Proc. IEEE WCNC Workshop Open RAN Archit.* 5G Evol. 6G, Austin, TX, USA, Apr. 2022, pp. 518–523.
- [187] L. Bonati et al., "Colosseum: Large-scale wireless experimentation through hardware-in-the-loop network emulation," in *Proc. IEEE DySPAN*, Dec. 2021, pp. 105–113.
- [188] Z. Ghadialy. "O-RAN technical steering committee (TSC) & Workgroups." Jul. 2021. [Online]. Available: https://www. parallelwireless.com/blog/o-ran-technical-steering-committee-tsc-workgroups
- [189] "O-RAN testing & integration." O-RAN Alliance. Jan. 2023. [Online]. Available: https://www.o-ran.org/testing-integration
- [190] A. Weissberger. "TIP OpenRAN and O-RAN alliance liaison and collaboration for open radio access networks," Feb. 2020. [Online]. Available: https://tinyurl.com/dnxkdns4
- [191] I. Gomez-Miguelez, A. Garcia-Saavedra, P. D. Sutton, P. Serrano, C. Cano, and D. J. Leith, "srsLTE: An open-source platform for LTE evolution and experimentation," in *Proc. ACM WiNTECH*, New York, NY, USA, Oct. 2016, pp. 25–32.
- [192] F. Kaltenberger, A. P. Silva, A. Gosain, L. Wang, and T.-T. Nguyen, "OpenAirInterface: Democratizing innovation in the 5G era," *Comput. Netw.*, vol. 176, Jul. 2020, Art. no. 107284.
- [193] L. Bertizzolo et al., "Arena: A 64-antenna SDR-based ceiling grid testing platform for sub-6 GHz 5G-and-beyond radio spectrum research," *Comput. Netw.*, vol. 181, pp. 1–17, Nov. 2020.
- [194] A. Panicker et al., "AERPAW emulation overview and preliminary performance evaluation," *Comput. Netw.*, vol. 194, pp. 1–11, Jul. 2021.
- [195] D. Raychaudhuri et al., "Challenge: COSMOS: A city-scale programmable testbed for experimentation with advanced wireless," in *Proc. ACM MobiCom*, London, U.K., Sep. 2020, pp. 1–13.
- [196] J. Breen et al., "POWDER: Platform for open wireless data-driven experimental research," *Comput. Netw.*, vol. 197, pp. 1–18, Oct. 2021.
- [197] H. Koumaras et al., "5GENESIS: The genesis of a flexible 5G facility," in *Proc. IEEE CAMAD.* Barcelona Spain Sep. 2018, pp. 1–6.
- in *Proc. IEEE CAMAD*, Barcelona, Spain, Sep. 2018, pp. 1–6. [198] "OpenAI cellular documentation." OpenAI Cellular. Accessed: Jun. 2022. [Online]. Available: https://openaicellular.github.io/oaic
- [199] P. S. Upadhyaya, A. S. Abdalla, V. Marojevic, J. H. Reed, and V. K. Shah, "Prototyping next-generation O-RAN research testbeds with SDRs," 2022, arXiv:2205.13178.
- [200] L. Bonati, M. Polese, S. D'Oro, S. Basagni, and T. Melodia, "Intelligent closed-loop RAN control with xApps in OpenRAN gym," in *Proc. Eur. Wireless*, Dresden, Germany, Sep. 2022, pp. 1–6.
- [201] L. Bonati, S. D'Oro, S. Basagni, and T. Melodia, "SCOPE: An open and softwarized prototyping platform for NextG systems," in *Proc. ACM Int. Conf. Mobile Syst. Appl. Serv. (MobiSys)*, Jun. 2021, pp. 415–426.
- [202] "MGEN traffic emulator." U.S. Naval Research Laboratory. Sep. 2021.
 [Online]. Available: https://www.nrl.navy.mil/Our-Work/Areas-of-Research/Information-Technology/NCS/MGEN
- [203] L. Bonati, M. Polese, S. D'Oro, S. Basagni, and T. Melodia, "OpenRAN gym: AI/ML development, data collection, and testing for O-RAN on PAWR platforms," Jul. 2022. arXiv:2207.12362.
- [204] "PAWR." Accessed: Dec. 2021. [Online]. Available: https://www.advancedwireless.org
- [205] L. Bonati et al., "CellOS: Zero-touch softwarized open cellular networks," Comput. Netw., vol. 180, pp. 1–13, Oct. 2020.
- [206] L. Bonati, S. D'Oro, F. Restuccia, S. Basagni, and T. Melodia, "SteaLTE: Private 5G cellular connectivity as a service with full-stack wireless steganography," in *Proc. IEEE INFOCOM*, Vancouver, BC, Canada, May 2021, pp. 1–10.

- [207] M. Kohli et al., "Open-access full-duplex wireless in the ORBIT and COSMOS testbeds," *Comput. Netw.*, vol. 199, Nov. 2021, Art. no. 108420.
- [208] H. Zhang et al., "ARA: A wireless living lab vision for smart and connected rural communities," in *Proc. ACM WiNTECH*, New Orleans, LA, USA, Oct. 2021, pp. 9–16.
- [209] "POWDER O-RAN profile." Accessed: Nov. 2021. [Online]. Available: https://gitlab.flux.utah.edu/powder-profiles/oran
- [210] S. Fdida et al., "SLICES, a scientific instrument for the networking community," *Comput. Commun.*, vol. 193, pp. 189–203, Sep. 2022. [Online]. Available: https://www.sciencedirect.com/science/article/pii/ S0140366422002663
- [211] "Open AI cellular." Accessed: Nov. 2021. [Online]. Available: https://sites.google.com/msstate.edu/oaic/home
- [212] O-RAN Working Group 4, "O-RAN fronthaul interoperability test specification (IOT) 8.0," O-RAN, Alfter, Germany, document O-RAN.WG4.IOT.0-v08.00 Technical Specification, Oct. 2022.
- [213] S. D'Oro, M. Polese, L. Bonati, H. Cheng, and T. Melodia, "dApps: Distributed applications for real-time inference and control in O-RAN," *IEEE Commun. Mag.*, vol. 60, no. 11, pp. 52–58, Nov. 2022.
- [214] E. Björnson and L. Sanguinetti, "Scalable cell-free massive MIMO systems," *IEEE Trans. Commun.*, vol. 68, no. 7, pp. 4247–4261, Jul. 2020.
- [215] "sim-e2-interface repository." O-RAN Software Community. 2022. Accessed: Mar. 2022. [Online]. Available: https://github.com/o-ran-sc/sim-e2-interface



Michele Polese (Member, IEEE) received the Ph.D. degree from the Department of Information Engineering, University of Padova in 2020. He has been a Principal Research Scientist with the Institute for the Wireless Internet of Things, Northeastern University, Boston, since March 2020. His research interests are in the analysis and development of protocols and architectures for future generations of cellular networks (5G and beyond), in particular for millimeter-wave and terahertz networks, spectrum sharing and passive/active user coexistence,

open RAN development, and the performance evaluation of end-to-end, complex networks. He has contributed to O-RAN technical specifications and submitted responses to multiple FCC and NTIA notice of inquiry and requests for comments. He was awarded with several best paper awards. He is currently an Associate Technical Editor of the IEEE Communications Magazine and a Guest Editor of the IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS. He is a member of the Committee on Radio Frequency Allocations of the American Meteorological Society from 2022 to 2024. He has served as the TPC Co-Chair for WNS3 2021–2022 and as an Organizer for the Open 5G Forum in Fall 2021.



Leonardo Bonati (Member, IEEE) received the Ph.D. degree in computer engineering from Northeastern University, Boston, MA, USA, in 2022, where he is an Associate Research Scientist with the Institute for the Wireless Internet of Things. His main research focuses on softwarized approaches for the Open Radio Access Network (RAN) of the next generation of cellular networks, on O-RAN-managed networks, and on network automation and orchestration. He served multiple times on the technical program committee of the

ACM Workshop on Wireless Network Testbeds, Experimental Evaluation & Characterization, and as the Guest Editor of the special issue of Elsevier's *Computer Networks* journal on Advances in Experimental Wireless Platforms and Systems.



Salvatore D'oro (Member, IEEE) received the Ph.D. degree from the University of Catania in 2015. He is a Research Assistant Professor with Northeastern University. He is one of the contributors to OpenRAN Gym, the first open-source research platform for AI/ML applications in the Open RAN. His research interests include optimization, artificial intelligence, security, network slicing and their applications to 5G networks and beyond, with specific focus on Open RAN systems. He is an Area Editor of Computer Communications (Elsevier). He serves

on the Technical Program Committee of multiple conferences and workshops, such as IEEE INFOCOM, IEEE CCNC, IEEE ICC, and IFIP Networking.



Stefano Basagni (Senior Member, IEEE) received the first Ph.D. degree in computer science from the University of Milano, Italy, in 1998, and the second Ph.D. degree in electrical engineering from the University of Texas at Dallas in 2001. He is with the Institute for the Wireless Internet of Things and a Professor with the ECE Department, Northeastern University, Boston, MA, USA. His current interests concern research and implementation aspects of mobile networks and wireless communications systems, wireless sensor networking for IoT

(underwater, aerial, and terrestrial), and definition and performance evaluation of network protocols. He has published over 12 dozen of highly cited, refereed technical papers, and book chapters. His h-index is currently 49 (November 2022). He is also co-editor of three books. He served as a guest editor of multiple international ACM/IEEE, Wiley, and Elsevier journals. He has been the TPC co-chair of international conferences. He is a Distinguished Scientist of the ACM and a Member of Council for Undergraduate Education.



Tommaso Melodia (Fellow, IEEE) received the Ph.D. degree in electrical and computer engineering from the Georgia Institute of Technology in 2007. He is the William Lincoln Smith Chair Professor with the Department of Electrical and Computer Engineering, Northeastern University, Boston. His research on modeling, optimization, and experimental evaluation of Internet of Things and wireless networked systems has been funded by the National Science Foundation, the Air Force Research Laboratory, the Office of Naval Research,

DARPA, and the Army Research Laboratory. He is also the Founding Director of the Institute for the Wireless Internet of Things and the Director of Research for the PAWR Project Office. He is a recipient of the National Science Foundation CAREER Award. He has served as an Associate Editor for IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS, IEEE TRANSACTIONS ON MOBILE COMPUTING, and Computer Networks (Elsevier). He has served as the Technical Program Committee Chair for IEEE Infocom 2018, and the General Chair for IEEE SECON 2019, ACM Nanocom 2019, and ACM WUWnet 2014. He is the Director of Research for the Platforms for Advanced Wireless Research (PAWR) Project Office, a \$100M public–private partnership to establish 4 city-scale platforms for wireless research to advance the U.S. wireless ecosystem in years to come. He is a Senior Member of the ACM.