

A Survey on 5G Radio Access Network Energy Efficiency: Massive MIMO, Lean Carrier Design, Sleep Modes, and Machine Learning

David López-Pérez, Antonio De Domenico, Nicola Piovesan, Geng Xinli, Harvey Bao, Song Qitao and Mérouane Debbah

Abstract

Cellular networks have changed the world we are living in, and the fifth generation (5G) of radio technology is expected to further revolutionise our everyday lives, by enabling a high degree of automation, through its larger capacity, massive connectivity, and ultra-reliable low-latency communications. In addition, the third generation partnership project (3GPP) new radio (NR) specification also provides tools to significantly decrease the energy consumption and the green house emissions of next generations networks, thus contributing towards information and communication technology (ICT) sustainability targets. In this survey paper, we thoroughly review the state-of-the-art on current energy efficiency research. We first categorise and carefully analyse the different power consumption models and energy efficiency metrics, which have helped to make progress on the understanding of green networks. Then, as a main contribution, we survey in detail —from a theoretical and a practical viewpoint— the main energy efficiency enabling technologies that 3GPP NR provides, together with their main benefits and challenges. Special attention is paid to four key enabling technologies, i.e., massive multiple-input multiple-output (MIMO), lean carrier design, and advanced idle modes, together with the role of artificial intelligence capabilities. We dive into their implementation and operational details, and thoroughly discuss their optimal operation points and theoretical-trade-offs from an energy consumption perspective. This will help the reader to grasp the fundamentals of —and the status on— green networking. Finally, the areas of research where more effort is needed to make future networks greener are also discussed.

I. INTRODUCTION

The industrial revolution and the automation of labour greatly accelerated the natural pace of evolution, and since then, we have significantly transformed the world we are living in [1].

After the Second World War, technology has developed faster than ever before. Automation is making our lives easier, and until very recently, it felt like nothing could stop us. However, we are starting to see the consequences of our unsustainable progress now [2].

From 1937 to 2019, the world population has grown from 2.3 to 7.7 billions [3], and our modern agricultural, manufacturing, transport and life styles have sharply increased energy consumption. As a result, the amount of carbon in the atmosphere raised from 280 to 409 parts per millions in such time period [4], and the levels of greenhouse gas (GHG) emissions reached the historical record of 37.5 gigatonnes of carbon dioxide equivalent (CO_{2e}) in 2018, —a 1.5% increase with respect to 2008 [5].

Importantly, the consequences of such high levels of GHG emissions are already tangible today. The increase of GHG emissions, combined with current trends on deforestation, have contributed to global warming —the rise of the average Earth surface temperature [6]. Between 1906 and 2005, the planet temperature rose 0.6 to 0.9 degrees Celsius due to global warming, and the rate of temperature increase has nearly doubled in the last 50 years [7] [8].

If our societies do not significantly change the manner in which we consume energy, the consequences may be catastrophic [9]. To put the above numbers in perspective, it should be noted that an increase in the Earth's temperature of 1.5 to 2.0 degrees Celsius, above pre-industrial temperatures, has been estimated to be *the limit* —a threat to most natural ecosystems on Earth today, and thus to our planet and everyday lives [10]–[13].

To address this challenge, international policymakers are targeting a dramatic increase in energy efficiency, and a sharp shift from fossil fuels to renewable sources of energy, such as solar, wind and water. This will entail a completely new approach to the generation and use of energy, which must be adopted by every government, industry, business, and individual.

In the following subsections, we first motivate the need for the 5G of mobile technology to enable a green environment, reviewing the important role that 5G will play to revert the global unsustainable energy consumption trend. Then, in the rest of the manuscript, we survey the technical innovations that 5G radio access networks (RANs) bring in terms of energy efficiency with respect to previous technology generations —from a theoretical and practical viewpoint— to enable a greener cellular operation, and in turn, support such environmental change.

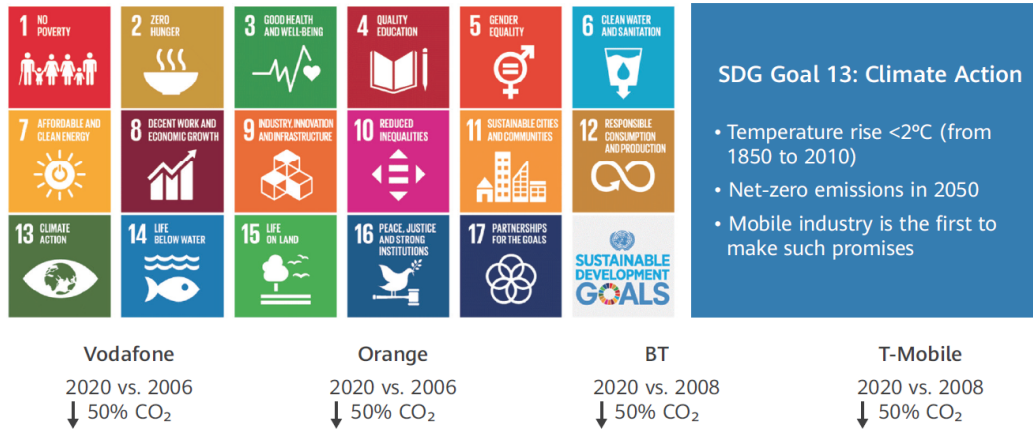


Figure 1. Wireless industry client action goals [14].

A. The Enabling Role of fifth generation (5G)

Governments and industries have set—or are setting—ambitious targets to reduce their GHG emissions and deal with global warming. To date, 77 major economies have already established a net-zero GHG emission target by 2050, and their industries are accordingly (re)defining their energy efficiency and consumption road maps. In this regard, the telecommunication sector has taken the lead—and an set exemplary role—, setting stringent requirements for both the energy efficiency and consumption of their networks, together with clear plans to meet them [14] (see Fig. 1).

Enhancing the energy efficiency and reducing the energy consumption of 5G networks will help reducing GHG emissions—no question about it. Their enabling effect, however, will be—without a doubt—the most important contribution of the mobile industry to address the current climate change [15].

At a macro-scale level, the new 5G technology enables a new type of networking capability able to connect for the first time both everyone and everything together, including machines, objects, and devices, thanks to its higher capacity, lower latency, improved reliability, and larger number of supported connections. Please refer to Fig. 2 for a comparison between international mobile telecommunications (IMT)-Advanced (fourth generation (4G)) and IMT-2020 (5G) specification capabilities. Importantly, these new 5G communication capabilities are already helping governments, current industries, and new forms of businesses to implement novel processes with improved effectiveness, by supporting a more flexible, tailored, and efficient use of resources.

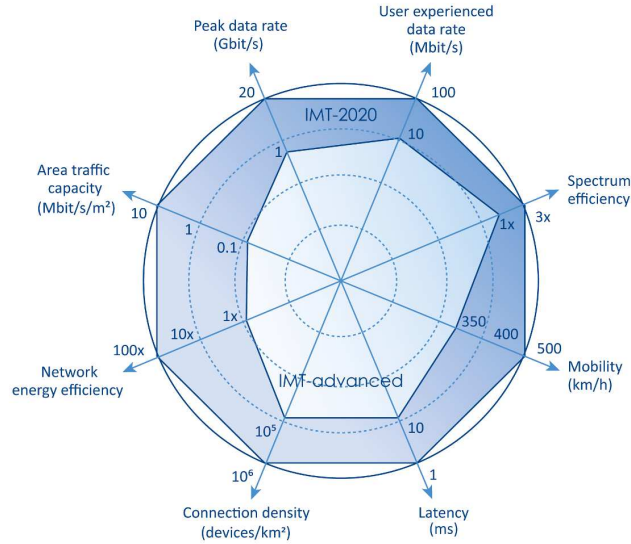


Figure 2. Comparison of key capabilities of IMT-Advanced (4G) with IMT-2020 (5G) according to [17].

As a result, 5G has already become an integral part of governmental and industrial energy efficiency and consumption programs, as it is envisioned that an intelligent exploitation of resources will enable a significant decrease of GHG emissions through different avenues, for example *i*) an improved support for smart city and building energy management, *ii*) reduced requirements for office space and business travel, and *iii*) efficient just-in-time supply chains, enabled by predictive analytics, to cite a few [16].

To give some idea of the magnitude and importance of this 5G enabling effect, it should be noted that, according to the international telecommunication union (ITU) SMART 2020 report [18], the enabling effect of mobile communications alone was estimated to be around 2,135 million tones of CO_{2e} in 2018, and that according to [19], the scale of it will increase in the 5G era, where the enabling effect across all the information and communication technology (ICT) sector was predicted to be equivalent to 15 % of all global emissions by the end of 2020.

Recent reports indicate that industries, such as transportation, health care, and manufacturing, are already significantly benefiting from such 5G enabling effect. For example, through smart city programs and 5G-related innovations, London, Berlin, and Madrid have already reduced GHG emissions of motor vehicles by 30 % each from their peak rates, and Copenhagen by 61 % [20].

To further assess the of breath and depth of the enabling effect of 5G networks, interested

readers are referred to [15], and references there in.

B. The 5G Energy Efficiency Challenge

Unfortunately, the 5G enabling effect is no free lunch. In fact, it comes at the expense of a tremendous challenge for the telecommunication sector in terms of both carried data and energy consumption.

Allowing governments, industries, businesses, and individuals in general to increase their energy efficiencies and reduce their energy consumption through more flexible, tailored, and efficient operations via a telecommunication network entails

- a dramatic growth of data usage in some scenarios, and
- the need for more sophisticated networking to meet the required low-latency, high-reliability, and/or large volume of data connections in some others.

Quantifying such challenge, recent studies already indicate that, by 2030, the number of connected devices is expected to grow to 100 billion [21], and that 5G networks may be supporting up to 1,000 times more data than 4G ones did in 2018 [16].

Importantly,

- to reach their established energy efficiency and consumption targets, and
- reduce their energy bills, up to 40 % by 2030 [16], to make their businesses profitable,

mobile network operators (MNOs) will need to meet the aforementioned more challenging traffic demands and requirements with significantly reduced GHG emissions with respect to those of today's 4G networks.

Considering the aforementioned two aspects, the third generation partnership project (3GPP) stakeholders have already called for a 90 % reduction in energy consumption of 3GPP new radio (NR) compared to 3GPP long term evolution (LTE) [22]. However, whether these gains can be realised or not in practical networks will not only depend on what the new specification can do and/or the energy performance of a single site, but also on how the actual network is deployed and operated as a whole.

In fact, to support the growing use of 5G connectivity and its more stringent requirements, while reducing energy consumption on a per-bit basis through an intelligent use of the network, changes are needed at all levels of it to achieve the maximum holistic effect. MNOs must thus embrace new approaches to network planning, deployment, management, and optimisation that have energy efficiency at heart, and are implemented end-to-end. Without energy efficiency

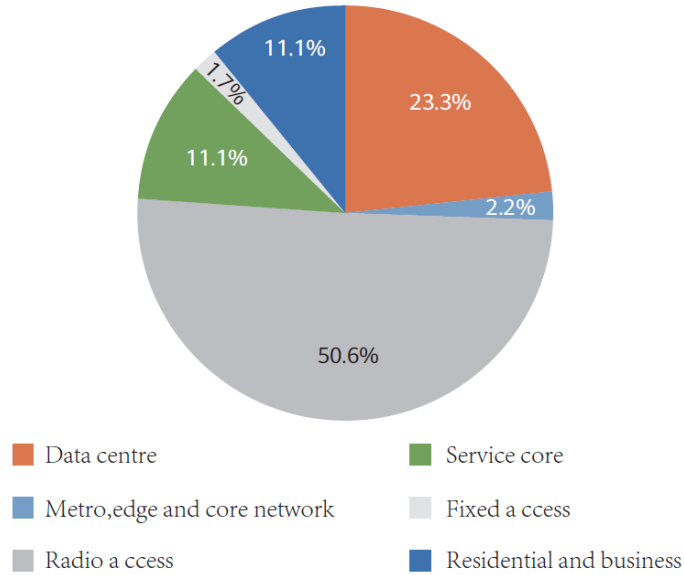


Figure 3. Energy consumption breakdown by network element in 2025 [25].

driving future deployments, the study in [23] indicates that a 5G network, despite of its enhanced energy efficiency in bits per Joule due to its larger bandwidth and better spatial multiplexing capabilities, could typically consume over 140% more energy than a 4G one, with a similar coverage area. This unwanted energy consumption arises from 5G's greater density of base stations (BSs), antennas, cloud infrastructure, and user equipment (UE), among others.

To address this challenge, and better understand where the energy consumption could be meaningfully reduced in a 5G network, thus helping MNOs to take educated decisions, the authors of [16] put together an interesting analysis, reporting that the most consuming part is the RAN, and in particular, the BSs. They currently account for about 57% of the total network energy consumption [24]. By 2025, that figure should be lower, as 5G becomes more prevalent, but the RAN will still be the biggest consumer of energy in the network, a 50.6% according to [25] (see Fig. 3).

Within a BS itself, the radio frequency (RF) equipment, i.e., the power amplifier plus the transceivers and cables, have been identified as the largest energy consumer, typically using about 65% of the total BS energy. The cooling system, the digital signal and base band processing as well as the alternating current (AC)-direct current (DC) converters follow with an energy consumption of around 17.5%, 10%, and 7.5%, respectively.

In this line, 5G can be —and has been— improved for a better energy efficiency. More

efficient power amplifiers have been developed, renewable energy sources for powering on-grid and off-grid sites, including solar power, are starting to be widely adopted. Moreover, smart lithium batteries are becoming an integral part of any 5G site to enhance energy management, and liquid cooling is being implemented to reduce the need for air conditioning [16].

Importantly, since the BSs —and their RF equipment— consume most of the energy at a 5G network, a judicious networking is of imperative importance [26]. BSs and/or its most relevant components must thus only be active —and consume energy— when handling actual data. Plainly speaking, the Joules consumed in a 5G RAN should 'follow' the bits transmitted/received. As a result, and in contrast to 4G, the amount of always-on signalling in next generation RANs must be greatly minimised at all cost. It is also equally necessary that data —and its related signalling— are transmitted to/received from the intended UEs, while consuming the minimum possible energy to meet the end-users' quality of service (QoS) demands [27]. Avoiding the resource waste occasioned by the uncontrolled over-provisioning of end-users' QoS is essential to significantly reduce the GHG emissions of 5G RANs.

To realise such network-wide energy efficient operation, 5G RANs in general and the 3GPP NR specification [28] in particular *i)* have been redesigned and developed using a new and more flexible user-centric principle [29], and *ii)* provide support for new —or further enhance existing— enabling technologies that help reducing energy consumption at the RAN level. In this line, three technologies stand out, i.e. massive multiple-input multiple-output (mMIMO), the lean carrier design, and 5G sleep modes. Moreover, it is important to note that 5G RAN optimization is already benefiting from the latest advancements in the machine learning (ML) field, which in terms of energy efficiency enable a more accurate RAN modeling and the use of more sophisticated optimization techniques. Such over-the-top ML-based optimization frameworks are also supported by new 3GPP NR enhancements, which facilitate, among others, statistics collection and prediction capabilities.

The good understanding of these 5G technologies, how they enable energy efficiency and their optimal operation points in terms of energy efficiency are the main focus of this survey.

C. Comparison to previous 5G Energy Efficiency Surveys

Several surveys have been already published on energy efficiency and related aspects during the last 10 years. In the following, we describe the most relevant ones.

In [30], the authors have surveyed the goals set by the United Nations (UN) for sustainable development, where among others, urgent actions to combat climate change are called upon. Importantly, the wide variety of opportunities in the ICT sector for enabling energy efficiency have been highlighted, but although of relevance to understand the current scene, this survey did not touch on the particularities of wireless networks. The survey in [31], instead, has focused on green wireless networks, and has explained at a high level —and an introductory form— how different metrics such as energy efficiency, spectral efficiency, throughput, and delay relate to each other. The discussion has been organised around the open systems interconnection (OSI) layers, and has presented an exhaustive summary of the different frameworks and techniques that can be used to achieve optimal networking trade-offs in practical RANs. Descriptions, however, are only qualitative, and lack of detail. For example, the paper has considered a simple power consumption model in most explanations, which only accounts for the transmit power, and has neglected the power consumption of equipment hardware. With a similar scope, the authors in [32] have surveyed the literature around the understanding of fundamental energy performance trade-offs, i.e., energy and deployment efficiency, energy and spectral efficiency, bandwidth and power as well as delay and power, this time, however, from a more focused cellular perspective, covering both 3GPP universal mobile telecommunication system (UMTS) and LTE aspects. Importantly, this paper has brought the attention to the importance of accurate BS power consumption models. Based on this work, the authors in [33] have further surveyed techniques to optimise the above mentioned four trade-offs, while putting the spotlight on the potential benefits of two enabling technologies, i.e., multiple-input multiple-output (MIMO) and relay. With a more 3GPP LTE standard focus, the survey in [34] has provided an in-depth review on green aspects, discussing advancements in power amplifier technology, 3GPP LTE protocols for symbol and carrier shutdown, as well as the potential benefits of small cell RANs. On a forward looking note, the authors have also put forward cognitive radio and cooperative relaying as energy saving enabling technologies, and share the latest developments in these areas. Complementarily, the authors in [35] have surveyed the efforts of different academic and industry projects on energy efficiency, emphasizing how to make use of daily network traffic variations to save energy at the RAN level. Several deployment strategies, such as cell-size, heterogeneous networks, cooperative communications and network coding, have been also discussed, together with their merits and challenges. On a similar note, in [36], a complete survey on 3GPP metrics and power consumption models for energy efficiency analysis can be found, with a detailed formulation

of the most relevant energy efficiency trade-offs. An exhaustive —and critical discussion— on standardised enhancements for energy-aware management in 3GPP LTE cellular RANs from a broadband networking point of view has been also provided. In [37], the authors have turned the focus on the survey of efficient resource management schemes, which are capable of controlling how much of the RAN infrastructure is actually needed in a given space-time and which parts can be temporarily powered off to save energy. The survey includes both cellular- and Wi-Fi-based algorithms. Focusing on more modern applications, the work in [26] has overviewed the challenges brought by the big data era, describing issues and solutions around energy efficient data acquisition, communication, storage, computation, and analytics. Discussion on the necessity of avoiding radio resource waste to reduce energy consumption in RANs has been at the core of the survey, and four types of schemes have been surveyed in this line, i.e., power control, time-domain scheduling, spatial resource allocation, and spectrum sharing.

Unfortunately, it should be highlighted that none of the above mentioned surveys has gone into the specifics of 3GPP NR. They are either generic or 3GPP LTE focus, and as a result, they did not survey the assets of this new technology generation to harvest energy savings.

Giving a more related 5G perspective, the authors in [38] have provided guidelines for the development of energy efficient enabling technologies in 3GPP NR, and surveyed the literature around the user-centric concept of *no more cells*. This networking paradigm would enable energy efficiency through a flexible network comprised of heterogeneous cells, decoupled signal- and data-planes as well as downlink and uplink. A dynamic cloud radio access network (C-RAN)-based configuration has also been proposed to handle spatial and temporal mobile traffic variations without energy over-provisioning. In this work, the potential of both mMIMO and full duplex in 5G for power savings has been also surveyed, while considering hardware issues. In the same line, the work in [39] has provided an overview of the latest research on green 5G techniques. The authors have explored ultra-dense sub-6GHz and millimetre wave networks, unlicensed spectrum as well as device to device (D2D) and mMIMO communications, while analysing their potential energy efficiency improvements, as well as circuit power consumption issues. As a main contribution, energy harvesting has been presented as fundamental to meet 5G green requirements, and the different lines of work in this have been discussed (e.g., renewable, RF energy harvesting). Taking a more theoretical and systematic approach, the authors in [40] have extended their work in [32], surveying energy efficient solutions for 5G RANs using the aforementioned fundamental green trade-offs as a driver. This overview has been around three

main pillars, i.e., non-orthogonal access, mMIMO, and heterogeneous networks. Importantly, the paper concludes that mMIMO is the most effective approach to enable high energy efficiencies, provided that issues around channel state information (CSI) acquisition, transceiver hardware impairments, and power inefficient components are addressed. A large number of references has also been provided around carrier and channel (antenna) shutdown techniques for dense small cell networks, and the implications of centralised versus distributed RAN architectures have been discussed. With a more practical—but still vanilla—5G perspective, the authors in [41] provide a comprehensive survey on how ML can be used to address the energy efficiency challenges encountered in generic 5G RANs. Finally, the recent survey in [42] has provided the most up-to-date overview on power saving techniques supported by the 3GPP NR standard, covering developments in Release 15 and 16, and the potential upcoming ones in Release 17. Such overview, however, has mainly focused on—and evaluated—UE power saving mechanisms, such as bandwidth parts, radio resource control (RRC) inactive state, discontinuous reception (DRX) mechanism, wake up signaling, cross-slot scheduling, and MIMO layer adaptation. At the RAN side, the lean carrier and discontinuous transmission (DTX) concepts together with that of dormant cells, have been only briefly touched upon.

As it can be derived from the previous summary, most of the existing energy efficiency surveys were written in a pre-5G era before the 3GPP NR existed/matured, or have a strong focus on the UE—and not on the RAN—side.

For completeness, Table I provides a summary of the contributions and gaps observed in the aforementioned surveys.

D. Objective and Structure of this Survey

In this survey, contrary to the previous ones, which speculatively discussed what 5G RANs could be in terms of energy efficiency, we provide for the first time a detailed, up-to-date overview of the most relevant technologies that a 5G RAN can currently use—or can further leverage—to increase its energy efficiency in sub-6 GHz deployments, from both a theoretical and practical perspective. These technologies are mMIMO, the lean carrier design, 5G sleep modes, and ML. The structure of this survey is built around these technologies. Note that 5G core networks aspects are not covered in this survey.

Importantly, and for motivation purposes, we should highlight that the potential of mMIMO to significantly enhance energy efficiency and/or decrease transmit power has been already shown

Table I
SUMMARY AND COMPARISON OF THE MOST RELEVANT ENERGY EFFICIENCY SURVEYS.

Ref.	Year	Tech. era	Contribution/techniques	Gaps
[30]	2018	ICT	Surveys UN frameworks for sustainable development	No cellular or ML related.
[31]	2016	General wireless	Surveys green metrics and performance trade-offs.	No 5G or ML related.
[32]	2011	3G/4G	Surveys green metrics and performance trade-off.	No 5G or ML related.
[33]	2011	4G	Surveys green performance trade-off optimization. Explores MIMO and relay technologies.	No 5G or ML related.
[34]	2011	4G	Surveys 3GPP LTE green protocols. Focuses on symbol and carrier shutdown technologies. Explores cognitive radio, and cooperative relaying.	No 5G or ML related.
[35]	2013	4G	Surveys green academic and industry projects, and energy minimization under end-users' QoS demands. Explores heterogeneous networks, cooperative communications, and network coding technologies.	No 5G or ML related.
[36]	2014	4G	Surveys 3GPP green metrics, power consumption models, and 3GPP LTE protocols. Detailed formulation of energy efficiency trade-offs.	No 5G or ML related.
[37]	2014	4G and Wi-Fi	Surveys energy efficient resource management schemes (with focus on Network and MAC layers).	No 5G or ML related.
[26]	2018	Big data era and 4G	Surveys the challenges brought by the big data era. With respect to cellular, surveys MAC layer power control, time-domain scheduling, spatial resource allocation, and spectrum sharing green protocols.	No 5G or ML related.
[38]	2014	Pre-5G	Surveys potential 5G energy efficiency enabling technologies. Explores the user-centric concept, downlink and uplink split C-RAN, mMIMO, and full duplex technologies.	Guesses what 5G could be. No 3GPP NR or ML related.
[39]	2017	Pre-5G	Surveys research on green 5G techniques. Explores ultra-dense sub-6GHz and millimetre wave networks, unlicensed spectrum as well as D2D, mMIMO, and energy harvesting technologies.	Guesses what 5G could be. No 3GPP NR or ML related.
[40]	2017	Pre-5G	Surveys green metrics and performance trade-offs. Explores non-orthogonal access, mMIMO, and heterogeneous networks.	Guesses what 5G could be. No 3GPP NR or ML related.
[41]	2020	5G generic	Surveys how ML can be used to address 5G energy efficiency challenges.	No 3GPP NR focus. Lack of detail.
[42]	2020	5G	Surveys power saving techniques supported by the 3GPP NR standard (Release 15/16) with focus on UE	Does not cover network aspects or ML techniques.

in 5G deployments. Moreover, the lean carrier design together with the 5G sleep modes have also taken significant 3GPP specification work, and are critical to allow the adaptation of the RAN capabilities to varying traffic loads at both small and large time scales with symbol as well as carrier and channel (antenna) shutdown mechanisms, respectively. ML frameworks are also having a major impact now on the design and optimization of the aforementioned energy efficiency techniques in practical 5G RANs.

It should also be noted that, differently than other surveys, this one provides much more detailed descriptions, including, for example, the formulation of the most relevant BS power consumption models and energy efficiency metrics currently used for energy efficiency optimisation. It also presents the most important energy efficiency bounds, trade-offs, and optimal operation points derived in the literature. This makes this survey a self-contained manuscript. This survey also clearly distinguishes main practical concepts around the analysed enabling technologies already existing in previous technology generations from the latest 3GPP NR specified (supporting) energy efficiency enhancements.

Given that much improvements in terms of energy efficiency are still required in practical 5G RANs, this survey also points out and provides a fresh overview on such challenges, and the potential of other new technologies.

With this in mind, the rest of this survey is organised as follows (see Fig. 4):

- In Section II, we introduce mMIMO, the lean carrier design, 5G sleep modes, and ML as energy efficiency enabling technologies in 5G RANs;
- In Section III and Section IV, we present and discuss in detail existing BS power consumption models and metrics used for energy efficiency optimisation, respectively;
- In Section V¹, we overview the current and well-established theoretical understanding of mMIMO in terms of energy efficiency from a single and a multi-cell perspective and from an uplink and a downlink viewpoint. Important bounds and trade-offs with other key performance indicators are formulated and explained, and the most relevant optimisation frameworks to enhance energy efficiency via mMIMO tuning are presented.

¹Since 3GPP NR mMIMO specification work mostly relates to control signalling and protocols, and not to energy efficiency *per se*, and because the amount of valuable practical work on the energy efficiency of realistic mMIMO-based RANs is limited, with regard to both detailed system-level simulations as well as measurements campaigns —the latter owing to 3GPP NR mMIMO deployments being quite new—, in this section, we focus on well-established theoretical aspects.



Figure 4. Outline and structure of this survey.

- In Sections VI, VII, and VIII² we dive into the details of the lean carrier design, and highlight the importance of sleeping modes and their optimisation at different levels (i.e., symbol, carrier, and channel (antenna) shutdown);
- In Section IX, we highlight the potential of spatio-temporal traffic predictions and ML approaches to maximize energy efficiency, and overview the research in this area;
- Finally, in Section X and XI, we discuss future research directions, and draw the conclusions, respectively.

²Since the body of work on the theoretical understanding of symbol, carrier, and channel (antenna) shutdown methods is limited, given the complexity of their theoretical modelling, as the performance of these schemes highly depends on the details of the 3GPP physical layer specification as well as scenario, deployment, traffic as well as system characteristics and dynamics, in this section, we focus on the 3GPP specification work related these techniques, and the valuable work around their performance analysis and optimization.

II. 3GPP NR ENERGY EFFICIENCY RELATED ENABLING TECHNOLOGIES

To reach the ambitious targets of 5G networks in terms of capacity, latency, reliability, number of supported connections, and energy efficiency, the 3GPP NR specification presents a paradigm shift with respect to any preceding cellular technology [29].

In comparison with 3GPP LTE and with regard to energy efficiency, 3GPP NR introduces a new beam-centric —or mMIMO-centric— design, which enables both

- an extensive use of beamforming through a massive number of antenna elements, not only for data transmissions, but also for control-plane procedures, such as the initial access [43], and
- a larger spatial multiplexing of information on a given time-frequency resource.

Such beamforming and multiplexing gains allow for a reduced transmit power to reach a given targeted distance and/or meet a given end-user QoS demand, with the potential consequent benefits in terms of interference mitigation and energy savings. Similarly, they also allow for larger data rates given a fixed transmit power, which enable a larger BS deactivation time, and thus further energy savings [44].

Moreover, 3GPP NR follows a new ultra-lean design principle, in which control signals are not consistently transmitted in every radio frame, but on demand and more sparsely, based on traffic requirements [45]. This ultra-lean design allows for a more efficient operation of the mMIMO-centric design, and facilitates 5G sleep modes support, i.e., BS (de)activation, including sophisticated symbol, carrier, or channel (antenna) shutdown mechanisms, which can significantly reduce energy consumption [46].

3GPP NR enhancements also allow for a more distributed network architecture, which facilitates the usage of ML to derive optimum network-wide energy efficient operation policies through a centralised data gathering and processing [47][48]. In this line, the current 3GPP NR architecture introduces new functions in the core and the management domains, i.e., the network data analytics function (NWDAF) and the management data analytics function (MDAF), which can either run analytics on collected data or enhance the already supported network functions with statistics collection and prediction capabilities [49][50].

In the following, we introduce the main concept behind these 5G enabling technologies, and the possibilities they enable, whose proper optimisation will be key to minimise the energy

consumption of 5G networks. We will further elaborate on their details and potential energy saving abilities in the rest of this survey.

A. *Massive MIMO*

Full dimension MIMO —generally referred to as mMIMO— is probably one of the most important developments in 5G [43]. Plainly speaking, mMIMO refers to a technology where BSs are equipped with antenna arrays comprised of a large number of antenna elements [51]. At higher frequency bands, due to their more challenging propagation conditions, the large number of antenna elements are primarily used for beamforming to extend coverage. At lower frequency bands, the focus of this survey, in addition to leveraging beamforming gains, the large number of antenna elements is readily used to enable extensive spatial multiplexing and interference mitigation by spatial separation [52] [53]. In more detail, the large mMIMO antenna array can excite a plurality of channel sub-spaces to support multiple simultaneous transmissions to —or receptions from— several UEs. In this way, the network capacity can potentially linearly grow with the number of spatial streams multiplexed.

Despite its significant benefits, mMIMO also presents a number of optimization challenges. Importantly, mMIMO requires CSI at the BS to realise the necessary multi-user uplink signal detection operations and downlink precoding ones. Most practical mMIMO implementations take advantage of the channel reciprocity of time division duplexing (TDD) systems to acquire such CSI with the minimum possible overhead [54]. However, the number of orthogonal uplink pilots in a given time-frequency resource is finite, which limits the number of UEs over which a BS can perform CSI acquisition at once. Thus, such uplink (UL) pilots have to be reused across neighbouring BSs, and UEs allocated to the same UL pilot in neighbouring cells will interfere with each other in channel estimation phase. This effect is known as pilot contamination, and can significantly affect the performance of a mMIMO system [55]. Given that pilot contamination imposes a fundamental limit on what can be achieved with a non-cooperative mMIMO system, coordinated uplink pilot reuse among neighboring cells has become an important area of research and optimization.

Embracing the advantages and challenges of mMIMO, 3GPP NR provides extensive support to realise such mMIMO operation, and the basis for its optimization. Channels and signals, specially those used for control and synchronization, have been re-designed with respect to 3GPP LTE to natively use beamforming. This has been a major specification work. The acquisition of CSI

for the large number of antenna elements in a mMIMO BS is also now supported through *i*) new UE CSI reports estimated over channel state information-reference signals (CSI-RS) in the downlink for frequency division duplexing (FDD) systems, and *ii*) new channel measurements over sounding reference signals (SRSs) in the uplink, exploiting channel reciprocity, for TDD systems. Among others, 3GPP NR is also providing new functionalities to support analog beamforming as well as digital precoding [56].

In Section V, we provide a detailed survey and analysis of the energy efficiency benefits and trade-offs of this important 5G technology —mMIMO— in both single- and multi-cell setups, using well established theoretical results.

B. The Lean Carrier Design

In previous technology generations, signals for BS detection, broadcast of system information, and channel estimation were always-active or transmitted very frequently over the air, regardless of whether the BS was serving UEs or not [57]. It is important to note that, while these always-active signals facilitate UE operations —as UEs always have signals to rely on—, they also

- result in a large overhead in dense deployments,
- introduce inter-cell interference to other cells, thus reducing the achievable throughput,
- reduce the battery lifetime of the UE, and
- worsen the energy efficiency,

thus becoming a burden to efficient network operation.

In 3GPP NR, the transmission of such control signalling and the related procedures have been revisited, following a new lean carrier design [45], to enable larger sleep ratios and longer sleep duration.

The logic behind is as follows. When the traffic load of a cell —or group of them— is low, or the mobility conditions of the UEs allow, larger signaling cycles could be selected to make the carrier leaner, and thus reduce overhead, mitigate interference, and more importantly save energy, if the BS hardware is accordingly shutdown.

With respect to symbol switch-off in 3GPP LTE, 3GPP NR allows for longer and deeper sleep periods up to 160 ms, where more and more BS hardware is progressively shutdown. To accommodate for such sparser signalling cycles, cell search, (re)selection as well as CSI related signalling and procedures, such as the transmission of synchronization Signal/PBCH block (SSB),

system information block 1 (SIB1) and paging, have been accordingly redesigned in the new specification [58].

In Section VI, we survey and discuss how the 5G lean carrier design benefits energy efficiency.

C. 5G Sleep Modes

Cellular networks are usually planned and deployed to meet certain peak-hour requirements, which leads to an over-dimensioning of the network for the less challenging traffic loads during the day time [59]. As the traffic demands fluctuates over both time and space, underutilized BS resources could be dynamically switched off to save energy. The more network components that are shutdown and the longer the time that they are shutdown, the more energy can be saved [60].

Importantly, as mentioned earlier, the less required always-active signalling resulting from the lean carrier design allows for longer micro-sleeps of up to up to 160 ms in the presence of bursty traffic [58]. In addition, when coupled with traffic shaping techniques and an efficient hardware at the BS able to power up and down in fractions of a millisecond, these times with no transmission can be leveraged by advanced sleep modes (ASMs) to allow deeper sleep modes, which can progressively switch off more circuitry depending on such time length [61].

Most of the improvements that a network can achieve by appropriate resource management, however, may not lay on such micro-sleep space, as these sleep periods are usually only opportunistic and generally short. To enable longer BS resource deactivation times, —macro-sleeps— an appropriate management of the carriers and channels (antennas) of the BS according to traffic distributions and demands at a macro-time, e.g. minutes or even hours, can avoid the resource waste emanating from the over-dimensioning of the network to meet peak hour requirements [62] [63].

In 3GPP NR, new functionality to allow large macrocells to (de)activate smaller overlapping ones in an ad-hoc manner has been specified. This signalling, for example, includes messaging for cell activation/deactivation request over $X_n/X_2/F_1$ interface, as well as for confirming activation/deactivation actions [64].

In Sections VI, VII, and VIII we survey and discuss how the 5G sleep modes³ can operate

³It is important to note that, even if most of these 5G sleep mode concepts, i.e. symbol, carrier, and channel (antenna) shutdown, already existed in 4G, the 5G advancements on mMIMO and the lean carrier design, together with other new 5G complexities, such as multi-radio access technology (RAT) 4G and 5G deployments, demand for novel, more complex sleep modes in 5G networks, and thus new specification and/or algorithmic work.

carriers and channels (antennas) to enhance energy efficiency.

D. Machine Learning

The heterogeneous and stringent service requirements of 5G networks, together with their increasing complexity—a pinch of which has been depicted in previous sections—are making traditional approaches to network operation and optimization no longer adequate. Such methods use a significant level of expert knowledge and theoretical assumptions to characterize real environments. Thus, they do not scale well, and cannot handle the complexity of real scenarios with their many parameters and imperfections as well as stochastic and non-linear processes. To bridge this gap, and provide 5G networks with the intelligence required to strike optimum operation points, equipment vendors and MNOs have started to equip their products with ML-based functionalities [26] [48].

Fed by network measurements, supervised and unsupervised learning tools [65], two different branches of ML, are being extensively used nowadays to model 5G network behaviour first, and subsequently, take educated decisions and/or make predictions on complex scenarios [66]. This is particularly relevant to energy efficiency. As one can infer from the previous discussions, minimising 5G energy consumption is a large-scale network problem, which highly depends on complex BS and UE distributions, varying traffic demands and wireless channels as well as hidden network trade-offs. Thus, understanding and predicting UE behaviours and requirements, as well as their evolution in time and space, is critical to tailor the 5G network configuration—mMIMO, lean carrier, and 5G sleep modes—, and address UE specific communications needs with the minimum possible energy consumption. Specifically, the current trend is to replace rule-based heuristics and associated thresholds with e.g., optimal parameters configured through the knowledge acquired by machine learning models.

Additionally, due to the dynamic nature of wireless networks, and the lack of network measurements data for all network procedures and on all the possible configurations they can adopt, reinforcement learning (RL) [68] is also being widely explored to optimise 5G network performance in general, and energy efficiency in particular. For example, shutting down network elements is a combinatorial problem with a large number of variables. RL agents may be used to let the network interact with the environment, and learn optimum resource shutdown policies to minimise the total network energy consumption. In addition, it is expected that current promising machine learning results will be improved with time, while the 5G ecosystem collects and makes

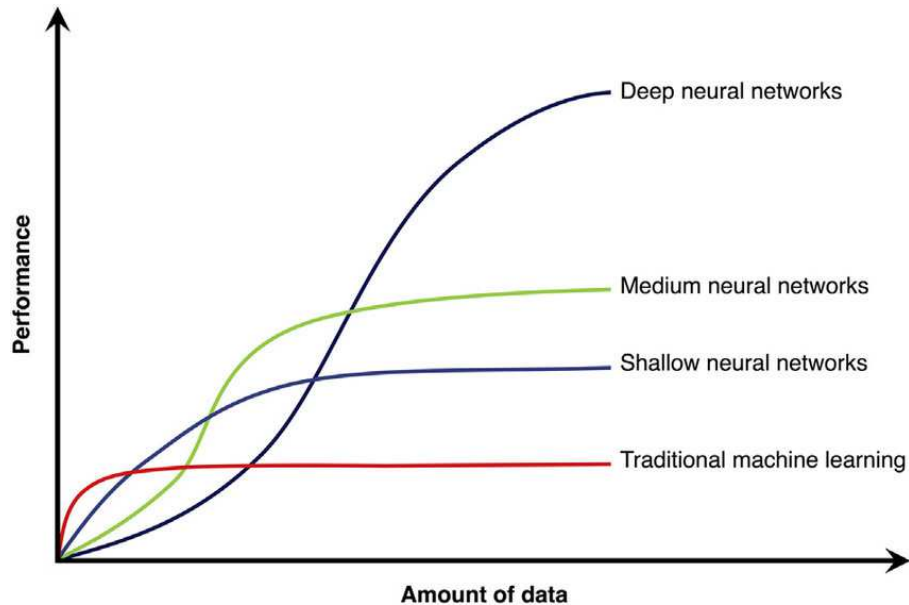


Figure 5. Relation between available training data and ML model performance [67].

available larger data sets related to the different problems. In fact, the amount of training data has a notable impact on the performance of ML algorithms, i.e. adding data generally improves performance, as shown in Fig. 5, particularly for the most recent models.

However, such learning may come at the expense of both *i*) an undesirably long exploration phase, where 5G network performance may be highly suboptimal, and *ii*) large computing powers and storage capabilities [69]. Moreover, continuously adapting an optimal policy, derived from and for a limited set of specific system configurations, to variations of network settings is a challenge, which may drive the need for data-driven model-based approaches.

In Section IX, we review how ML is being used to tackle the energy efficiency problem in 5G networks, and discuss both main its benefits and challenges.

III. POWER CONSUMPTION MODELS

To assess the impact of the different enabling technologies presented earlier on the energy efficiency of a 5G network, it is necessary to define models that provide a good estimation of their energy consumption. Importantly, such energy consumption models need to offer the right balance between accuracy and tractability, while embracing different network and BS architectures, to empower 5G system performance characterisation and optimisation.

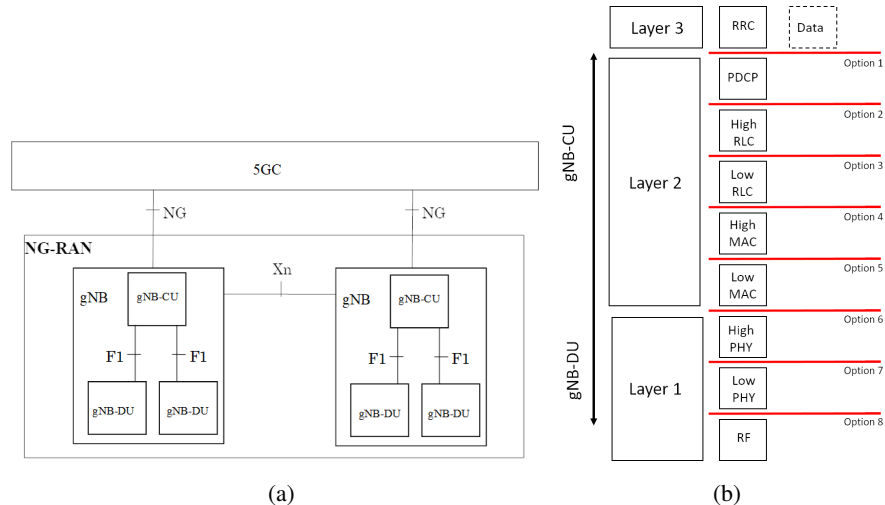


Figure 6. (a) NG-RAN overall architecture [71]; (b) 3GPP options for the function split between gNB-DUs and gNB-CU [70].

Fig. 6a shows the 3GPP NR RAN logical architecture. This architecture, denoted as next generation radio access network (NG-RAN) architecture, consists of a set of next generation NodeBs (gNBs) —the 3GPP NR BSs in the 3GPP terminology— connected *i*) amongst them through the X_n interface and *ii*) to the 5G core network (5GC) through the next generation (NG) interface. To take advantage of virtualization technologies and provide more implementation flexibility, a gNB may also consist of a central unit (CU) and multiple distributed units (DUs), connected to each other through the F1 interface. The 3GPP has studied eight functional split options between CU and DU (see Fig. 6b), and current RAN implementations are focusing on option 2 [70], in which RRC and packet data convergence protocol (PDCP) are in the CU, while radio link control (RLC), medium access control (MAC), physical layer (PHY), and RF are in the DU.

From an implementation perspective, each functional split corresponds to a distinct deployment option, which needs an appropriate power consumption characterisation. Specifically, it is necessary to take into consideration both the power consumption of the sites where the DUs and CUs are deployed, as well as the transport network, also referred to as fronthaul, that connects the DUs and the CU. In addition, with the advent of network functions virtualization (NFV), DUs and CUs functions can be implemented either through standard dedicated hardware or virtual network functions (VNFs) in a network cloud. Therefore, the overall RAN power consumption model may need to include the contribution of the cloud server, where the RAN VNFs are

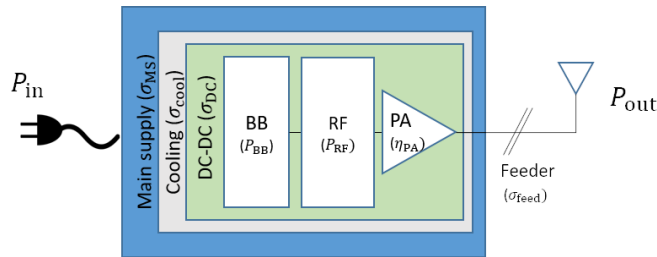


Figure 7. Power consumption block diagram of a BS RF transceiver.

deployed. Accordingly, we can generally express the aggregated RAN power consumption of a 5G network as follows:

$$P_{\text{RAN}} = \sum_i P_{\text{BS}_i} + \sum_j P_{\text{FH}_j} + \sum_k P_{\text{VBBU}_k}, \quad (1)$$

where P_{BS_i} , P_{FH_j} , and P_{VBBU_k} are the power consumption of the i -th gNB, the power consumption of the j -th fronthaul, and the power consumption of the k -th virtualized baseband unit (BBU), respectively. Depending on the specific RAN architecture, distributed or centralised, some of these components may not be considered.

In the following, we survey the most relevant power consumption models for both distributed and centralised RAN, while considering their most relevant characteristics.

A. Power Consumption Model for Distributed RAN

In case of a fully distributed radio access network (DRAN), the RAN power consumption can be modelled by taking into account the contributions of all BSs as in this architecture there are no fronthaul links and virtualized BBUs, whose power consumption contributions correspond to the second and third terms in (1), respectively.

Fig. 7 shows the block diagram of the widely used BS power consumption model defined in [72], where the power consumption of a non-mMIMO BS is computed as a function of the power consumption of all its antennas, each one including an RF transceiver module with its power amplifier (PA), plus that of the BBU associated to them, the DC-DC power supply, the active cooling system and the AC-DCC unit for connection to the electrical power grid. Such model is formulated as follows:

$$P_{\text{BS}} = N_{\text{TRX}} \frac{\frac{P_{\text{out}}}{\eta_{\text{PA}}} (1 - \sigma_{\text{feed}}) + P_{\text{RF}} + P_{\text{BB}}}{(1 - \sigma_{\text{DC}}) (1 - \sigma_{\text{MS}}) (1 - \sigma_{\text{cool}})}, \quad (2)$$

where N_{TRX} is the overall number of RF transceiver modules in a BS, and P_{in} is the input power of each transceiver. P_{out} is the transmit power, P_{RF} and P_{BB} are the RF transceiver module and the BBU power consumption, η_{PA} is the PA power efficiency, and σ_{feed} , σ_{DC} , σ_{MS} , and σ_{cool} are the power losses in the feeder, DC-DC power supply, mains supply, and active cooling, respectively.

It should be noted that the model in eq. (2) is widely represented in the literature by a simplified version of it, which explicitly shows the linear relation between the BS power consumption, P_{BS} , and the transmit power, P_{out} , as follows:

$$P_{\text{BS}} = N_{\text{TRX}} \cdot (P_0 + \delta_p P_{\text{out}}), \quad (3)$$

where P_0 and δ_p are cell-type dependent parameters, which indicate the power consumption at the minimum non-zero output power and the slope of the load-dependent power consumption, respectively. Note that this model is general, and accommodates to macro, micro and small cells. For example, the parameters of eq. (3) for different types of small cells are provided in [72].

Importantly, in the last decade, in addition to appearance of the aforementioned smaller cells [73], two other main solutions have emerged as key enablers to boost the mobile network capacity, i.e., carrier aggregation (CA) (see Section VII) and mMIMO (see Section V). Accordingly, a number of works have evolved the previous presented model to capture the impact of these technologies on the BS power consumption.

1) Carrier Aggregation Power Consumption Model: CA is a 3GPP flagship feature primarily introduced to increase the cell throughput in 3GPP long term evolution advanced (LTE-A). The first version of CA allowed to aggregate up to 5 component carriers (CCs) of up to 20 MHz, and currently, in 3GPP NR, it has been extended to support up to 16 CCs and 1 GHz of bandwidth. The CA framework is flexible by design. It enables to use continuous or discontinuous intra-band CCs as well discontinuous inter-band CCs, which can be characterized by different bandwidth or coverage. It is also a key technology to enable licensed Assisted Access (LAA) [74], heterogeneous network (HetNet) deployments [75] and dual connectivity [76]. For each UE, a CC is defined as its primary cell (PCell) [77], which acts as the anchor CC, and is thus used for basic functionalities, including mobility support and radio link failure (RLF) monitoring. Additionally configured CCs are denoted as secondary cells (SCells), and they can be added, changed, or removed to optimise the BS performance.

When modelling the power consumption of a system using CA, it is necessary to take into account how the power consumption scales with the number of active CCs, N_{CC} . In this line, the work in [78] has presented the following power consumption model for CA:

$$P_{\text{BS}} = \sum_{j=1}^{N_{\text{CC}}} \left(P_{\text{TX}_j} + B_j P_{\text{CP}_j}^{\text{CA}} \right) + P_{\text{CP}}^{\text{CAi}},$$

where $P_{\text{TX}_j} = \frac{P_{\text{out}_j}}{\eta_{\text{PA}}}$, B_j , and $P_{\text{CP}_j}^{\text{CA}}$ are the effective transmit power used by CC j , the bandwidth of CC j , and the variable circuit power consumption, which scales linearly with both the number of active CCs, N_{CC} , and their bandwidth, B_j , respectively, while as in the general model, $P_{\text{CP}}^{\text{CAi}}$ is the load independent circuit power consumption of the CA system, i.e., of the hardware components shared by the distinct CCs.

Embracing the complexity of CA, the variables, $P_{\text{CP}}^{\text{CA}}$ and $P_{\text{CP}}^{\text{CAi}}$, may take different values, depending on the specific CA implementation. For instance, contiguous CA can be realized with a single fast Fourier transform (FFT) and a single RF transceiver module for all CCs, while non-contiguous CA usually requires multiple of them [77]. In the worst case, each CC would require a fully dedicated hardware, and thus the load independent circuit power consumption of the CA system, $P_{\text{CP}}^{\text{CAi}}$, would scale linearly with the number of active CCs, N_{CC} [42].

2) *mMIMO Power Consumption Model:* With regard to mMIMO, and similarly as for the CA case, the linear power consumption model in eq. (3) has also been extended to take into consideration the large number of antenna elements and the new architecture of a mMIMO BS. In this line, the work in [79] proposed the following power consumption model for mMIMO:

$$P_{\text{BS}} = P_{\text{TX}} + N_{\text{TRX}}^{\text{A}} P_{\text{CP}}^{\text{A}} + P_{\text{CP}}^{\text{li}}, \quad (4)$$

where $N_{\text{TRX}}^{\text{A}}$ is the number of RF mMIMO transceiver modules, which do not need to be necessarily equal and may be smaller than the number of antenna elements, P_{CP}^{A} is the power consumption of the RF and digital processing needed to support each RF mMIMO transceiver module, and $P_{\text{CP}}^{\text{li}}$ is the load-independent circuit power consumption.

Importantly, it should be noted that $N_{\text{TRX}}^{\text{A}}$ depends on the type of beamforming architecture implemented at the BS [80]. Digital beamforming provides high flexibility, but requires a large number of RF mMIMO transceiver modules. In contrast, analog beamforming decreases the power consumption by significantly reducing the number of RF mMIMO transceiver modules at the cost of a lower spatial resolution capability of the beamforming and a larger latency to

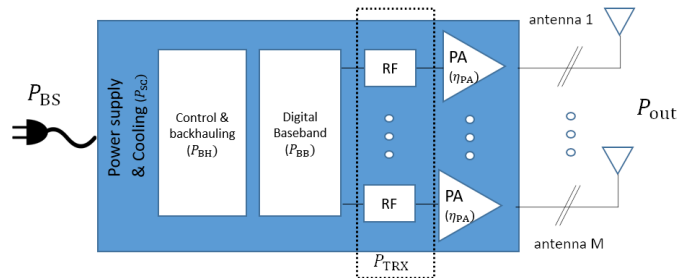


Figure 8. Power consumption block diagram of a mMIMO BS.

select the proper beam weights/configuration. Hybrid beamforming combines the advantages of the two architectures.

The linear mMIMO power consumption model in eq. (4) provides a simple description of the relation between the number of RF mMIMO transceivers and the power consumption in a mMIMO system. However, more advanced works have highlighted that it is of paramount importance to also take into account the impact of multi-user scheduling. Specifically, the research in [81] has described the steps to derive a more complete model for mMIMO BSs, which accounts for both downlink and uplink communications, under the assumption of zero forcing (ZF) processing⁴.

As described in Fig. 8, the main blocks of a mMIMO BS, from a power consumption perspective, are the PA, the analog front-end, the digital baseband, the control and network backhaul platform, and the power system (which includes also the cooling) [82]. Accordingly, the authors of [53] has described the mMIMO BS power consumption as the sum of the effective transmit power, $P_{TX} = \frac{P_{out}}{\eta_{PA}}$, and the circuit power, P_{CP} , as follows:

$$P_{BS} = P_{TX} + P_{CP}, \quad (5)$$

where P_{CP} accounts for the power contributions due to the analog front-end, P_{TRX} , the digital baseband, P_{BB} , the control and network backhaul platform, P_{BH} , and the power system, P_{SC} , and it can be computed as:

$$P_{CP} = P_{TRX} + P_{BB} + P_{BH} + P_{SC} = P_{CP}^{li} + P_{CP}^{ld}. \quad (6)$$

For clarity, we will denote by P_{CP}^{li} the sum of the fixed and load-independent power consumption required for the site power supply, control signalling, backhaul, signal processors, and RF

⁴Models for other specific precoding and combining schemes are discussed in [53].

transceivers, and by $P_{\text{CP}}^{\text{ld}}$ the term that includes the variable load-dependent part of the circuit power consumption, P_{CP} . In addition, it is worth to highlight that the functions that are the main contributors to the power consumption of the digital processing unit are the channel estimation process, the channel coding and decoding units, and the beamforming processing.

Accordingly, the mMIMO BS power consumption model in eq. (6) can be expressed as follows [83] [84]:

$$P_{\text{BS}} = \frac{K \cdot P_{\text{UE}}}{\eta_{\text{PA}}} + P_{\text{CP}}^{\text{li}} + \mathcal{C}_3 K^3 + \mathcal{D}_0 M + \mathcal{D}_1 M \cdot K + \mathcal{D}_2 M \cdot K^2 + \mathcal{A} K \cdot R_{\text{UE}}, \quad (7)$$

where K is the number of simultaneously multiplexed UEs at the BS, P_{UE} is the downlink output power per UE (i.e., $P_{\text{out}} = K \cdot P_{\text{UE}}$), \mathcal{C}_3 is the part of the beamforming processing, P_{SP} , which scales linearly with K^3 , \mathcal{D}_0 is the power consumed by the transceiver module attached to each antenna, \mathcal{D}_1 is the part of the beamforming processing, P_{SP} , which scales linearly with $M \cdot K$, \mathcal{D}_2 is the sum of the contributions of the channel estimation process, P_{CE} , and the beamforming processing, P_{SP} , which scale linearly with $M \cdot K^2$, R_{UE} is the UE throughput, and \mathcal{A} is the aggregated power consumption per bit of information required by the coding/decoding operations and by the load-dependent part of the backhaul. Table II describes typical values of the parameters in eq. (7).

Table II
TYPICAL VALUES FOR mMIMO POWER CONSUMPTION MODEL PARAMETERS [81] [83] [84].

Parameter	Value	Parameter	Value
η_{PA}	0.39	$P_{\text{CP}}^{\text{li}}$	20 [W]
\mathcal{C}_3	10^{-7} [W]	\mathcal{D}_0	1 [W]
\mathcal{D}_1	$3 \cdot 10^{-3}$ [W]	\mathcal{D}_2	$9.4 \cdot 10^{-7}$ [W]
\mathcal{A}	1.15 [W/Gbps]		

To conclude this section, Table III summarizes the main contributions of the presented works on the power consumption models for DRAN architectures.

B. Power Consumption Model for Centralised RAN

Considering a more sophisticated RAN architecture, the work in [85] has studied the power consumption modelling of a centralized radio access network (CRAN) considering different

Table III
SUMMARY OF POWER CONSUMPTION MODELS FOR DRAN.

Paper	Model	Main parameters
[72]	non-mMIMO BSs	Number of RF transceivers, load-dependent power consumption and static power consumption
[78]	non-mMIMO BSs with CA	Number of CCs, transmit power per CC, bandwidth per CC, variable circuit power, and load independent circuit power
[79]	mMIMO BSs	Number of RF transceivers, transmit power consumption, power consumption of the RF and digital processing per transceiver, and static power consumption
[83] [84]	mMIMO BSs	Number of multiplexed UEs, number of antennas, transmit power consumption, beamforming power consumption

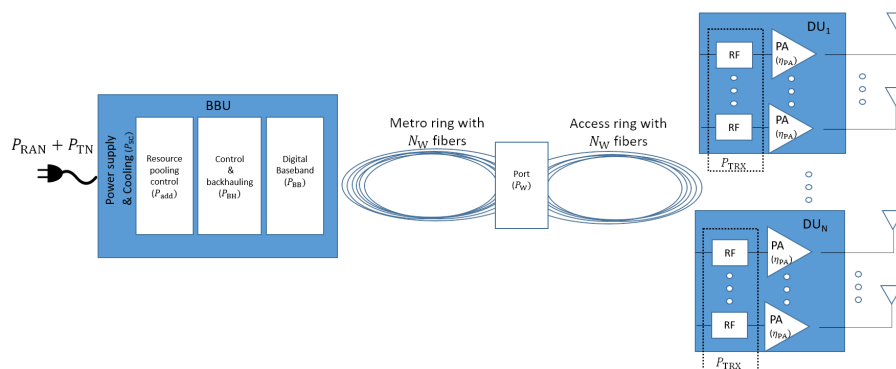


Figure 9. Power consumption block diagram of a CRAN system with optical transport network.

functional splits between the gNB-DUs and the gNB-CU. In terms of energy efficiency, centralisation enables three main benefits with respect to a decentralised architecture: stacking gain, pooling gain, and cooling gain [85]. The stacking gain refers to the capability of deploying less processing units (in the central node) to serve the same amount of cells in a given area, while the pooling gain refers to the capability of using a limited amount of centralised resources to operate a large amount of cells, by exploiting the load variations in the network. The cooling gain appears due to the reduced amount of energy required to cool the cell site and the more advanced cooling solutions that can be implemented at the central node [86].

In case of a CRAN, the network power consumption shall be modelled by taking into account the contribution of both the access network, P_{RAN} , including all CUs and DUs, and the transport network, P_{TN} , (see Fig. 9).

In this centralized architecture, the aggregated power consumption at the RAN, P_{RAN} , can be

computed as [85]:

$$\sum_i P_{BS_i} = \sum_i P_{DU_i} \left(1 - \frac{P_{co_i} + P_{NF_i}}{100} \right) + P_{BBU}, \quad (8)$$

where P_{DU_i} is the power consumption of the site where the i -th DU is deployed, which can be computed as, e.g., eq. (7), P_{BBU} is the power consumption of the BBU host where the CU is located⁵, P_{NF_i} is the fraction of power consumption that corresponds to the network functions (NFs) moved to the CU, and P_{co_i} is the fraction of power consumption that corresponds to the cooling related to the i -th DU. Importantly, note that the higher the number of NFs moved to the CU, the higher the values of both parameters, P_{co_i} and P_{NF_i} .

Moreover, the power consumption of the BBU host, P_{BBU} , can be computed as:

$$P_{BBU} = \sum_i P_{DU_i} \left(\frac{1/G_{co}}{100} + \frac{P_{NF_i}/N_{BS}}{100} \left[\frac{N_{BS}}{G_{st}G_{po}} \right] P_{add} \right), \quad (9)$$

where G_{st} , G_{po} , G_{co} , and P_{add} are the stacking gain, the pooling gain, the cooling gain, and the additional power consumption in the BBU needed to enable resource pooling, respectively. The work in [85] has provided a preliminary analysis to estimate the values of these parameters. The authors in [87] have used tele-traffic theory and simulations to evaluate the resource pooling gains at the BBU and the fronthaul for different functional splits, and then investigated how to optimize the cost for BBU pool, fronthaul capacity, and radio resource utilization. More recently, the authors in [88] have further extended this work, and provided an analysis of the power consumption of the different RAN functional splits, considering the multiplexing gain arising in each split option.

It is important to note that eq. (9) assumes that the CU is deployed on a dedicated hardware platform, which has a similar architecture to the one used for a single DU but with larger computational capacity. In the case of a virtualized radio access network (VRAN), the goal is to move the processing to general purpose processors (GPPs), which are capable to provide real-time processing to maintain the timing in the RAN protocols, while equipped with efficient and elastic resources —CPU, memory, and networking— to perform intensive digital processing. This approach promises further improved energy savings, which depend on the specifically used architecture [89]. One option is to use dedicated hardware (e.g., system on chip) for managing the layer 1 functions on dedicated hardware, while higher-layer functions are implemented in a software-based architecture. This solution may enable, however, limited additional energy

⁵This model can be easily generalised to the case where a larger network with multiple CUs is considered.

savings with respect to CRAN. To reduce the power consumption and increase flexibility, an alternative option is to deploy only the most computationally-intensive functions on dedicated hardware, such as turbo decoding and encryption/decryption. As an extreme option, in the full GPP architecture, all the RAN functions are implemented in a virtual environment, which may enable large power savings at the expense, however, of performance if the GPPs and infrastructure around cannot cope with the workload. Considering the GPP architecture, the work in [90] proposes a power consumption model for virtualized BBUs in a cloud node, which takes into account the impact of the cooling system, the workload dispatcher switch, and the GPPs as follows:

$$P_{\text{VBBU}} = P_{\text{Vco}} + P_{\text{Dis}} + \sum_i P_{\text{GPP}_i},$$

where P_{Vco} , P_{Dis} and P_{GPP_i} are the power consumption contributions in the virtualized BBU due to the cooling system, the switch, and the i -th GPP, respectively.

To characterize the GPP power consumption, a linear power consumption model can be used [71] [91], i.e.,

$$P_{\text{GPP}} = P_{\text{GPP}_0} + \Delta_{\text{GPP}} P_{\text{GPP}_m} \rho_{\text{GPP}},$$

where P_{GPP_0} and P_{GPP_m} are the power consumption of the GPP when it is in idle mode and at maximum load, respectively, and Δ_{GPP} and ρ_{GPP} are the slope of the power consumption model, which is related to the specific GPP architecture, and the load of the GPP, respectively, where the latter depends on the GPP computational capacity and the computational resources required by the hosted VNFs.

To model the GPP computational load, the work in [91] has used an experimental platform to infer the relationship between the downlink throughput and the percentage of central processing unit (CPU) usage at the BBU. Instead, the authors in [92][93] have proposed an analytical characterisation of the GPP computational load, which jointly depends on the functional split, the number of used transceiver modules/antennas, the bandwidth, the rate and the number of spatial MIMO-layers used at the virtualized BS. Importantly, the authors in [94] have further analytically investigated the computational complexity of lower layer functions in 3GPP NR, and this research has shown that, today, due to the form factor and power consumption of GPP, dedicated hardware is the only feasible option for deploying full 3GPP NR capabilities.

Complementing the above, the authors in [95] have also recently developed an empirical model to describe the computational requirements of the RAN NFs, focusing on the PHY layer, which

includes the most computationally expensive functionalities. Their results highlight that NFs can be classified according to their complexity into three classes:

- PHY NFs, whose computational complexity only depends on the system configuration, and does not change with time (e.g., FFT),
- PHY NFs, whose computational complexity depends on both throughput requirements and channel quality (e.g., encoding/decoding), and
- higher layer NFs, whose computational complexity only depends on throughput requirements (e.g., packet scheduling).

Finally, to connect each DU to the CU, a fronthaul link, whose capacity fits the functional split throughput and latency requirements is needed. Therefore, the total power consumption of the fronthaul is a function of the transport network technology, its topology, the capacity required by each BS, and the number of BSs deployed in the RAN. For instance, when the fronthaul is based on optical dense wavelength division multiplexing with a ring architecture and optical switching, the power consumption of the transport network, P_{TN} , is due to all the fiber links connecting the DUs and the associated CUs. The contribution of a link connecting the i -th DU with its CU, P_{FH_i} , can be modelled as [85]:

$$P_{\text{FH}_i} = \left\lceil \frac{C_{\text{FS}_i}}{R_{\text{tr}}} \right\rceil \left(2 \cdot P_{\text{tr}} + \frac{P_{\text{W}}}{N_{\text{W}}} \right),$$

where C_{FS_i} is the transport capacity required by the functional split at the i -th gNB, P_{W} is the power consumption of a port used to interconnect access and metro rings in the transport network, N_{W} is the number of wavelengths per fiber, and R_{tr} and P_{tr} are the rate and the associate power consumption of the transport network nodes, respectively. It is important to highlight that a centralized architecture leads to multiplexing gains not only at the BBU side, but also in the communication network. Specifically, the work in [96] has proposed an analytical model to derive an upper bound of the fronthaul statistical multiplexing gain, and shown that even pooling a moderate number of BSs can result in notable radio resource savings.

Table IV provides a summary of the main contributions presented in this section, which discuss the power consumption models for CRAN architecture. To provide a general view, and highlight the impact of the transport network on the overall system power consumption, we show in Fig. 10a the network power consumption with respect to the BS deployment density for a classic DRAN and different CRAN architectural options, using the parameters indicated in Table V. Moreover, Fig. 10b shows the impact of each contributor, i.e., DUs, fronthaul nodes, and BBU,

Table IV
SUMMARY OF POWER CONSUMPTION MODELS FOR CRAN.

Paper	Model	Main parameters
[85]	CRAN architecture	BBU power consumption, DU power consumption, and stacking, pooling, and cooling gains
[90]	CRAN architecture with GPPs	Power consumption due to the cooling system, the switch, and GPPs
[71] [91]	GPP power consumption	GPP power consumption in idle and maximum load, load-dependent slope of the power model
[92][93]	GPP power consumption	Functional split, number of antennas, bandwidth, rate, MIMO-layers
[85]	FH power consumption	Transport network capacity, number of wavelengths per fiber, bandwidth, rate and power of the transport network nodes

with respect to the aggregated CRAN power consumption. Fig. 10a highlights that only a CRAN with functional split option 6 (see Fig. 6b) has similar power consumption to the classic DRAN. CRAN architectures with lower functional splits provide larger centralisation gains, but lead to higher power consumptions. In addition, Fig. 10b shows that the contribution from the transport network cannot be neglected in the overall power consumption characterisation, particularly when some architectures are adopted. Specifically, while for split option 6, it only amounts for around 2% of the network power consumption, in split option 7 and 8, it contributes for around 30% and 60% of the network power consumption, respectively.

To conclude this section, let us also highlight that the research and industry communities have spent notable efforts to characterize the RAN power consumption —whether distributed or centralised— while considering the most relevant architectures. However, this is still an open research field. In particular, we argue that due to the large ecosystem of equipment vendors and their different implementation of solutions, there is a need for further data and experimentally driven research to validate the proposed models and find their appropriate parameters on real 3GPP NR equipment. In addition, further research is also needed to characterise CA and mMIMO BS power consumption in the presence of multiple configurations and more complex technologies, such as coordinated multi-point (CoMP), and the characterisation of the power consumption of sites where multiple RATs co-exist.

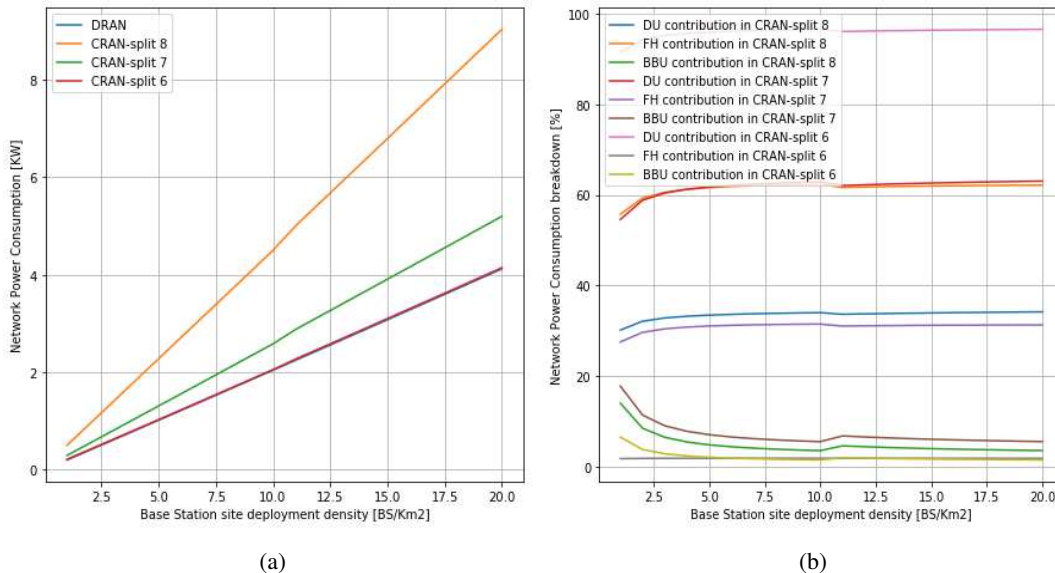


Figure 10. Network power consumption (a) and its breakdown (b) as a function of the BS deployment density for classic DRAN and different CRAN architectural options.

Table V

TYPICAL VALUES FOR CRAN POWER CONSUMPTION MODEL PARAMETERS [85] [97].

CRAN split 8	Value	CRAN split 7	Value	CRAN split 6	Value
P_{co}	10	P_{co}	10	P_{co}	0
P_{NF}	15	P_{NF}	10.5	P_{NF}	3
G_{po}	5	G_{po}	5	G_{po}	5
G_{st}	2	G_{st}	2	G_{st}	2
P_{add}	2	P_{add}	2	P_{add}	2
G_{co}	0.2	G_{co}	0.2	G_{co}	1
C_{FS}	3044 Gbps	C_{FS}	882 Gbps	C_{FS}	10 Gbps
R_{tr}	100 Gbps	R_{tr}	100 Gbps	R_{tr}	10 Gbps
N_W	40	N_W	40	N_W	80
P_W	2.2 W	P_W	2.2 W	P_W	2.2 W
P_{tr}	4.5 W	P_{tr}	4.5 W	P_{tr}	2 W

IV. ENERGY EFFICIENCY METRICS

To evaluate the impact of energy efficiency mechanisms on energy savings, energy efficiency metrics are as important as the used power consumption models. These metrics must be com-

prehensive, reliable and widely accepted to allow comparisons. In addition, they have to capture both the energy consumed by the system under investigation as well as the performance measured at network level (such as coverage, capacity, and delay).

To achieve these goals, the European telecommunications standards institute (ETSI) Environmental Engineering technical committee, the ITU-T Study Group 5 and the 3GPP Technical Specification Group RAN have specified metrics to assess mobile network energy efficiency under different operating conditions. The contribution of these standard development organizations (SDOs) has mainly focused on global system performance, while considering different demographic areas, load scenarios, and radio access technologies. In contrast, the academic research community has contributed to this effort by proposing energy efficiency link-specific metrics, which enable a more tailored energy efficiency network optimisation [98].

In the following, we will first overview the contributions from different SDOs, considering both interference- and noise-limited scenarios, and then touch on energy-delay related metrics. Subsequently, we also describe alternative metrics proposed by the academic research community to drive energy efficiency optimisation problems, indicating the merits and the drawbacks of each one of them.

A. Energy Efficiency Metrics for Interference-Limited Networks

As one of the main targets of 5G networks to enhance energy efficiency is to adapt the system capacity —and the associated power consumption— to the network load, load-aware metrics are key for the next generation of green communication networks. In this context, ETSI has defined the mobile network data energy efficiency metric, EE_{DV} [bit/J] [99], which is the ratio between the data volume, DV , delivered in the network and the network energy consumption, EC , observed during the time period required to deliver such data, i.e.,

$$EE_{DV} = \frac{DV}{EC}, \quad (10)$$

where EC should be computed by integrating eq. (1) over an observation period that includes distinct load levels⁶. Note that the data volume, DV , includes both downlink and uplink traffic for both circuit switched services and packet switched services, and that this metric can be used to characterise a single or multiple BSs, operating in urban and/or dense-urban areas, in which the

⁶the 3GPP recommends that the performance should be evaluated considering at least 3 load levels [100].

network experiences highly variable loads, i.e. interference-limited scenarios. However, it is not appropriate for scenarios where the traffic load is low. Specifically, since the energy efficiency metric, EE_{DV} , is not weighted according to the global reference, a small energy saving in the low energy consumption region in a low load scenario may comparatively lead to apparently large energy efficiency gains, while a large energy saving in the high energy consumption region may comparatively result into limited energy efficiency gains.

To characterise the energy efficiency, while considering distinct deployment scenarios, e.g., dense-urban, urban, sub-urban, rural, or deep rural areas, ETSI has extended the previous metric to the total energy efficiency metric, EE_{Total} [bit/J] [99], which can be defined as the weighted sum of the energy efficiency in each deployment scenario, i.e.,

$$EE_{Total} = \frac{\sum_m PoPP_m \cdot EE_{DV_m}}{\sum_m PoPP_m},$$

where m is the index of the m -th scenario, EE_{DV_m} is the energy efficiency of the m -th scenario, and $PoPP_m$ is the weight (percentage) representing the typicality of the m -th scenario in the network under test. This metric shall be used to characterise the energy efficiency performance of a large scale network, e.g., the network of an operator in a specific country.

To jointly consider deployment scenarios and traffic loads, the 3GPP has further complemented the work from ETSI, introducing the network energy efficiency metric, EE_{global} [bit/J] [100], which can be defined as the sum of the energy efficiencies in multiple deployment scenarios and under different traffic loads, i.e.,

$$EE_{global} = \sum_m b_m \cdot EE_{DS_m} = \sum_m \sum_l b_m \cdot a_l \cdot EE_{DV_{m,l}},$$

where EE_{DV_l} is the energy efficiency of the network for an observed deployment scenario, m , and traffic load level, l , while b_m and a_l are the weights for the corresponding deployment scenario, m , and traffic load level, l , respectively. The 3GPP recommends to compute b_m considering the proportion amongst the different deployment scenarios in terms of *i*) power consumption, *ii*) traffic load, or *iii*) connection density [101]. With respect to a_l , and giving an example, if 10%, 30%, and 50% traffic load levels are investigated, the corresponding weights based on a daily traffic model could be 6/24, 10/24 and 8/24, respectively [101]. EE_{global} shall be used to characterise the energy efficiency performance of a large scale network, considering the multiple hours of the day, during which the spatial distribution of the load in the network may significantly change due to the daily human habits.

Although the above presented metrics provide a useful indication on how the energy consumption scales with the increase of the data rate requirements, it does not give any information about actual economic costs. To address this challenge, the work in [102] has introduced the economical energy efficiency metric, E^3 [bit/J], which is defined as the ratio of effective system throughput to the associated energy consumption, weighted by a cost coefficient, i.e.,

$$E^3 = \frac{\sum_{k \in \mathcal{K}} \alpha_k R_k}{\sum_{n \in \mathcal{N}} P_{Tn} + P_{0n} C_n},$$

where \mathcal{K} and \mathcal{N} are the sets of UEs and serving BSs, respectively, R_k is the effective throughput perceived by the k -th UE, α_k is the priority weight related to the k -th UE, P_{0n} and P_{Tn} are the static and the load-dependent power consumption, respectively, and C_n is the cost coefficient of the n -th BS. Note that the cost coefficient, C_n , is calculated as the ratio of the corresponding device cost to a predefined benchmark cost, such that the metric, E^3 , has the same unit as the energy efficiency. In scenarios where multiple access, fronthaul, computing, and other mechanisms coexist to provide efficient services to the end-users, this economical energy efficiency metric, E^3 , is able to discriminate amongst solutions, not only in terms of provided network throughput and energy consumption, but also in terms of additional (deployment and operational) costs. This allows the analysis of advanced network deployment and resource management schemes.

B. Energy Efficiency Metrics for Noise-Limited Networks

To complement the energy efficiency metrics presented in the previous section and deal with scenarios with sustained low data volumes, in particular in rural or in deep rural areas, ETSI also introduced the mobile network coverage energy efficiency metric, $EE_{MN,CoA}$ [m^2/J] [99], which is the ratio between the area covered by the network, CoA, and the network energy consumption observed during one year, EC, i.e.,

$$EE_{CoA} = \frac{CoA}{EC}. \quad (11)$$

This metric shall be used to characterise the energy efficiency performance in rural areas, where coverage is the main objective, or in the context of low-power wide-area networks.

In addition to the above, the 3GPP has extended the coverage energy efficiency metric in eq. (11) to consider a global metric, the mobile network data energy efficiency metric,

$EE_{\text{global, CoA}}$ [71], which can be used to analyse the energy efficiency over distinct areas, each one having distinct size or relevance for the operator, i.e.,

$$EE_{\text{global, CoA}} = \sum_i C_i \cdot EE_{\text{CoA},i},$$

where C_i is the weight for each such distinct area/deployment scenarios, which is computed by taking into account the area covered per deployment scenario and the relevance of the scenario itself to the total power consumption, e.g., the percentage of power consumption in dense urban areas to that in rural areas. The energy computation in the actual area under coverage—and how it is affected by a network technology—, however, can be complex to estimate, and it may require the collection of a large number of UE measurement reports.

In this context, ETSI has defined a coverage quality factor, CoAQ [%] [99], to estimate the quality of the coverage and measure the amount of connection failures due to coverage issues, load congestion or significant interference effects, i.e.,

$$\text{CoAQ} = (1 - \text{FR}_{\text{RRC}})(1 - \text{FR}_{\text{RAB}_S})(1 - \text{FR}_{\text{RAB}_R}),$$

where FR_{RRC} , FR_{RAB_S} , and FR_{RAB_R} are the radio resource control setup failure ratio, the radio access bearer setup failure ratio, and radio access bearer release failure ratio, respectively.

C. Energy Efficiency of E2E Network Slicing

Network slicing is a promising 5G technology to simultaneously support multiple services with diverse characteristics and requirements in a single 5G network. The 3GPP is currently focusing on three main families of services, i.e., enhanced mobile broadband (eMBB), ultra-reliable low-latency communication (URLLC), and massive machine type communication (mMTC). In Release 17, the 3GPP has worked towards an approach to evaluate the energy efficiencies for these families of slices [103]. Importantly, a slice is typically defined end-to-end, including RAN, 5GC, and transport network.⁷

For eMBB slices, the 3GPP has specified the usage of the metric, EE_{DV} , which we have defined in eq. (10) [103]. More interestingly, for URLLC and mMTC slices, data volume is not the main key performance indicator (KPI), and the 3GPP has introduced more appropriate metrics to characterize energy efficiency in these cases.

⁷However, these metrics can be used also for evaluating the performance of a slice defined only over one of these (sub)networks, by considering only the corresponding measurements.

For URLLC slices, the 3GPP, in line with ETSI efforts [99], has defined the metric, $EE_{\text{URLLC, Lat}}$ [s^{-1}/J] [103], which is the inverse of the average end-to-end latency of the network slice divided by the energy consumption of the network slice, i.e.,

$$EE_{\text{URLLC, Lat}} = \frac{1}{T_{e2e} \cdot EC},$$

where T_{e2e} is the overall system end-to-end latency.

In cases where latency and throughput are both important KPIs of the URLLC slice, or when the operator wants to evaluate the slice energy efficiency over different periods of time with distinct loads, the 3GPP has specified the metric, $EE_{\text{URLLC, DV, Lat}}$ [$\text{bit}/\text{s}/\text{J}$] [103]. This metric is defined as the slice data volume divided by the product between its average end-to-end latency and its energy consumption, i.e.,

$$EE_{\text{URLLC, DV, Lat}} = \frac{DV}{T_{e2e} \cdot EC}.$$

Finally, for mMTC slices, where the system main target is to provide service to a large number of mobile devices, the 3GPP has specified the metric, EE_{mMTC} [103], which is defined as the ratio of the number of UEs⁸ in the slice divided to the associated energy consumption:

$$EE_{\text{mMTC}} = \frac{N_{\text{UE}}}{EC}.$$

D. Link-aware Energy Efficiency Metrics

The energy efficiency metrics presented in Sections IV-A, IV-B, and IV-C should be used to evaluate the impact of network-wide solutions specifically designed to increase the system energy efficiency (as those presented in Sections VI, VII, and VIII), and to compare these solutions with the current baseline systems. In contrast, the metrics presented in this section are more appropriate to make *greener* the classic RRM schemes, such as beamforming or scheduling, which mainly target spectral efficiency. Hence, they should be used to optimize specific mechanisms rather than to evaluate a system. The way in which RRM schemes are optimized in terms of energy efficiency is not discussed in this survey, as this is well reviewed in previous manuscripts (see for instance [98]).

⁸The 3GPP has currently defined two variants of this metric, where the first one considers the maximum number of registered UEs, while the second one considers the average number of active UEs.

Specifically, to aid a finer grain performance optimization, the research community has extended the network energy efficiency metric, EE_{DV} , presented in eq. (10) by taking into consideration the different requirements of its distinct links. In fact, since the energy efficiency metric, EE_{DV} , can be seen as the aggregated sum of the energy efficiencies of each network link, its optimisation through radio resource management may tend to favour the links that can provide the largest throughput, which may limit, e.g., the cell-edge UEs performance, and thus the network fairness.

To address this issue, and denoting by \mathcal{L} the set of links in the network, the weighted sum of the energy efficiencies (WSEE) metric, WSEE [bit/J] [98], is defined as the weighted mean of the different energy efficiencies measured at each link $i \in \mathcal{L}$, i.e.,

$$\text{WSEE} = \sum_{i \in \mathcal{L}} w_i \cdot EE_{DV,i}, \quad (12)$$

where $EE_{DV,i}$ is the link energy efficiency, defined as in eq. (10), and w_i is the weight of i -th link. Importantly, note that the use of different weights for different links enables assigning them different priorities during the resource allocation to increase the system fairness. Accordingly, the weights are defined a-priori, for instance based on the data plan of the specific user.

To enable an even more fair resource allocation, researchers have also proposed the weighted product of the energy efficiencies (WPEE) metric, WPEE [bit/J] [98], which is defined as the exponentially weighted product of the different energy efficiencies measured at each link $i \in \mathcal{L}$, i.e.,

$$\text{WPEE} = \prod_{i \in \mathcal{L}} (EE_{DV,i})^{w_i}.$$

In particular, the WPEE metric maximisation ensures that no link experiences a zero throughput, and has been shown to converge to the Nash Bargaining solution [104]. Following this approach, however, it is not possible to improve the energy efficiency of the i -th link, $EE_{DV,i}$, without decreasing the energy efficiency of other link, e.g. the j -th link, $EE_{DV,j}$.

To achieve a better trade-off among overall system performance and fairness, a max-min fair resource allocation policy can be adopted. An energy efficiency resource allocation scheme that is max-min fair can be designed by maximising the weighted minimum of the energy efficiencies (WMEE) metric, WMEE [bit/J], i.e.,

$$\text{WMEE} = \min_{\{i \in \mathcal{L}\}} (w_i \cdot EE_{DV,i}).$$

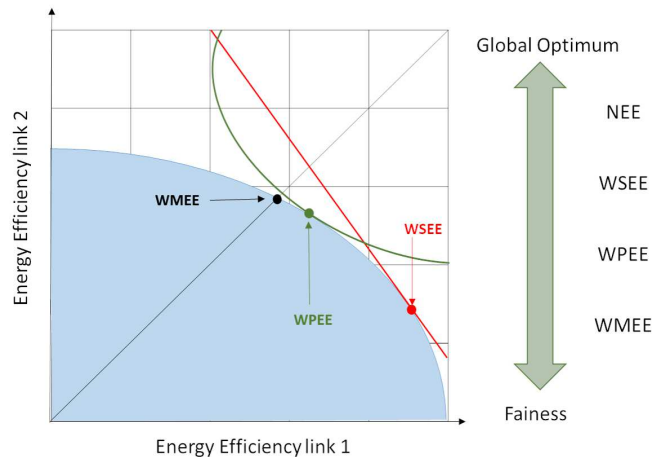


Figure 11. Operating regions of the energy efficient optimisation based on the NEE, the WMEE, the WPEE, and the WSEE metrics [98], where $NEE = \sum_{i \in \mathcal{L}} EE_{DV,i}$.

Note that when the WMEE metric is optimised, the resource allocation achieves the same product, $w_i \cdot EE_{DV,i}$, for all links $i \in \mathcal{L}$. Thus, if the weights are set equal, this max-min fair resource allocation will provide the same energy efficiency for each link.

To facilitate the reader's understanding, Fig. 11 provides a qualitative comparison of the performance achieved by a network with two links when optimising its energy efficiency using the WMEE, the WSEE and the WPEE metrics with respect to the Pareto boundaries. Each of the three approaches leads to a solution that belongs to the energy efficient Pareto region. However, the WSEE and the WPEE metrics allocate more resources to the link 1, which is characterized by a better energy efficiency in order to get closer to the global optimum. In contrast, the WMEE metric shares the resources between the two links such that they are characterized by the same energy efficiency.

To conclude this section, Table VI provides a summary of the most relevant energy efficiency metrics for 5G system optimisation, where it is important highlighting that these metrics have been mainly designed for assessing 4G networks, focusing on eMBB services, and considering only data rate, latency, and coverage requirement. In contrast, with respect to other 5G use cases, e.g. mMTC and URLLC, there is a lack of specific and well understood metrics, and further research is needed in this area. mMTC applications call for energy efficiency metrics that takes into account the number of connection handled by the network in combination with the area

covered by the network, e.g., an extension of the mobile network coverage energy efficiency metric, $EE_{MN,CoA}$. Similarly, URLLC applications require a metric able to capture the system reliability in combination with the end-to-end delay, e.g., an extension of the latency metric, EE_L .

Table VI
SUMMARY OF EE METRICS.

Metric	Unit	Calculation	KPI	Application
EE_{DV} [99]	[bit/J]	$\frac{DV}{EC}$	Load-aware	Evaluate network energy efficiency of eMBB systems, characterized by large load variations
EE_{Total} [99]	[bit/J]	$\frac{\sum_m PoPP_m \cdot EE_{DV,m}}{\sum_m PoPP_m}$	Load-aware	Evaluate network energy efficiency of eMBB systems spanning across multiple areas
EE_{global} [100]	[bit/J]	$\sum_m \sum_l b_m \cdot a_l \cdot EE_{DV_l}$	Load-aware	Evaluate network energy efficiency of eMBB systems considering scenarios with distinct loads
EE_{CoA} [99]	[m ² /J]	$\frac{CoA}{EC}$	Coverage-aware	Evaluate network energy efficiency of rural areas
$EE_{global, CoA}$ [71]	[m ² /J]	$\sum_i C_i \cdot EE_{CoA,i}$	Coverage-aware	Evaluate network energy efficiency of rural areas with distinct sizes or priorities
E^3 [102]	[bit/J]	$\frac{\sum_{k \in \mathcal{K}} \alpha_k R_k}{\sum_{n \in \mathcal{N}} P_{Tn} + P_{0n} C_n}$	Load & cost-aware	Evaluate network energy efficiency with consideration on deployment costs
$EE_{URLLC, Lat}$ [103]	[s ⁻¹ /J]	$\frac{1}{T_{e2e} \cdot EC}$	Latency	Evaluate network energy efficiency for URLLC systems
$EE_{URLLC,DV,Lat}$ [103]	[bit/s/J]	$\frac{DV}{T_{e2e} \cdot EC}$	Latency & load	Evaluate network energy efficiency for URLLC systems with consideration on distinct load scenarios
EE_{mMTC} [103]	[UE/J]	$\frac{N_{UE}}{EC}$	Connectivity	Evaluate network energy efficiency of mMTC systems
WSEE [98]	[bit/J]	$\sum_{i \in \mathcal{L}} w_i \cdot EE_{DV,i}$	Link & Load-aware	Optimize energy efficient resource management considering link priorities
WPEE [98]	[bit/J]	$\prod_{i \in \mathcal{L}} (EE_{DV,i})^{w_i}$	Link & Load-aware	Optimize energy efficient resource management preventing links with zero throughput
WMEE [98]	[bit/J]	$\min_{\{i \in \mathcal{L}\}} (w_i \cdot EE_{DV,i})$	Load-aware	Optimize energy efficient resource management with max-min fairness across radio links

V. THEORETICAL UNDERSTANDING OF ENERGY EFFICIENCY: MASSIVE MIMO

5G networks are targeted at a 100 times higher energy efficiency with respect to previous generations (see Fig. 2), and mMIMO has been identified as a key technology to reach such target. As discussed in previous sections, by leveraging its extensive beamforming and spatial multiplexing capabilities, mMIMO can significantly network capacity, but also reduce the transmit power required at the BS to achieve a targeted rate, given a frequency band of operation and

coverage area. However, running such larger number of antennas at the BS, together with the more signal processing required to handle the larger capacity in a mMIMO cell, also increases the energy consumption of the BS.

Multiple studies have set out to fundamentally understand the challenges of mMIMO networks in general, and the above presented energy efficiency trade-off in particular. Mainly concentrating on providing insights for network design, a large body of research has focused on deriving theoretical bounds on the capacity and the energy consumption of mMIMO systems, as well as the interplay between different network parameters. More practical research, on the other hand, has tackled mMIMO-based network design, BS deployment, as well as radio resource management and optimisation problems, while taking into account end-users' QoS demands. Extensive specification work has also been carried out in 3GPP NR to accommodate mMIMO capabilities in both the control- and user-planes.

One of the main challenges faced in mMIMO is around the accuracy and the tractability of the models used to characterise performance at the network-level. mMIMO capacity and rates, for instance, significantly degrade in the presence of channel correlation⁹ and/or pilot contamination, which is a function, among others, of the BS and the UE distributions, the scenario topology and the wireless channel, as well as of independent BS scheduling decisions and the resulting interference. However, this is hard to capture in tractable models. Same issues also revolve around the accuracy and complexity of the power consumption models, as discussed in Section III.

In light of these tractability issues, most theoretical findings on energy efficiency related to mMIMO networks, which are well grounded now, are limited to single-cell scenarios, where model simplifications can be more easily justified. In contrast, the estimation of energy efficiency metrics and trade-offs in multi-cell scenarios tends to be addressed through computer aided numerical evaluations.

3GPP NR mMIMO specification work and enhancements can be framed within the broader context of improved mobile broadband, and that such 3GPP work does not necessarily relate to energy efficiency *per se*, but mostly to control signalling definition and related procedures. In contrast, in the following, we concentrate on providing an overview of the fundamental understanding of energy efficiency in the context of mMIMO systems, in both single- and multi-

⁹Spatial correlation means that there is a correlation between the received average signal gain and the angle of arrival of such signal. Spatial correlation generally degrades the performance of multi-antenna systems, as it decreases the number of independent channels that can be created by precoding.

cell scenarios, and focus on green large-scale mMIMO network deployment and optimization, where energy efficiency is defined as the ratio of the achievable data rate to the related power consumption [bit/Joule].¹⁰

A. Single-cell scenario

In this subsection, we focus —and survey— energy efficiency bounds and trade-offs when considering a single-cell mMIMO scenario, which are fundamentally well understood, and already helping MNOs to design their networks. Explicit closed-form expressions are provided, and the impact of different mMIMO parameters into energy efficiency and power consumption are carefully analyzed.

1) *Bounds:* The pioneering work in [106]–[108] provided a first analysis of the energy efficiency in a single-cell mMIMO system, mostly based on the assumption of perfect CSI being available at the BS. The research in [106] showed that the performance of a mMIMO system with M antennas at the BS and a BS transmit power, P_{TX}/M , is equal to the performance of a single input single output (SISO) system with a BS transmit power, P_{TX} , without any intra-cell interference. This result indicated that, by using a large number, M , of BS antennas the BS transmit power, P_{TX} , can be proportionally scaled down by a factor, $1/M$. This work also suggested that the spectral efficiency in a mMIMO system can be increased by a factor, K , when serving K UEs in the same time-frequency resource. The findings in [107], [108] also resonated with these conclusions, reporting that a power reduction proportional to $1/M$ can be achieved in TDD systems¹¹, while maintaining non-zero rates, as the number, M , of antennas grows to infinity. This was a promising result, indicating that the energy efficiency of the system under study could be monotonically increased with the number, M , of antennas, without any trade-off. However, it is important to mention that such conclusions were obtained with significant assumptions in the BS power consumption model.

Soon after, however, further developments in this area of research showed that the energy efficiency bounds in mMIMO systems are hidden behind more complex BS power consumption models, and that inaccuracies and/or oversimplifications, not reflecting the essence of mMIMO hardware implementations, can lead to misleading practical insights. For example, when not

¹⁰For more details on 3GPP NR related mMIMO specification, the reader is referred to [105], and references therein.

¹¹The power reduction is proportional to $1/\sqrt{M}$ in the case of imperfect CSI.

considering the circuit power consumption related to running a larger number of antennas, as it was the case in [106]–[108], one can be let to believe that an unbounded energy efficiency can be achieved by adding more and more antennas. Deploying more and more antennas, however, requires additional circuitry, incurring a larger power consumption, which needs to be accurately captured by the analysis to draw the appropriate conclusions. In this line, the pioneering work in [81] stands out, which considered a more sophisticated —and realistic— mMIMO BS power consumption model, and in turn, found significantly different conclusions to the previously stated ones.

In the following, we survey in detail this work, which represent the most advanced —and settled— understanding of energy efficiency in single-cell mMIMO systems.

In more detail, the authors in [81], using the more detailed BS power consumption model already presented in eq. (6), analyzed for the first time in a systematic manner the total energy efficiency of the UL and downlink (DL) in a single-cell mMIMO network with respect to

- the total transmit power, $P_{\text{TX}}^{\text{tot}}$, accounting for the downlink and the uplink transmit powers,
- the number, K , of simultaneously multiplexed UEs, and
- the number, M , of antennas, where $M \geq K + 1$.

Remarkably, closed-form expressions for the energy efficiency optimal operating point with respect to each of these parameters, when the others are fixed, were provided for the case where both ZF processing and perfect CSI knowledge are considered. In particular, the authors considered a TDD system, in which the pilot signaling occupies $\tau^{(\text{ul})}K$ and $\tau^{(\text{dl})}K$ symbols in the uplink and downlink, respectively, where the inequality, $\tau^{(\text{dl})}, \tau^{(\text{ul})} \geq 1$, must be true to enable orthogonal pilot sequences among UEs, and with such system model, they provided a formulation of the gross downlink rate, expressed in bit per second, as follows:

$$\bar{R} = B \log(1 + \rho(M - K)), \quad (13)$$

where ρ is a design parameter proportional to the received signal to noise ratio (SNR). Following this formulation, the total transmit power, $P_{\text{TX}}^{\text{tot}}$, required to serve each UE with a gross rate, \bar{R} , was shown to be:

$$P_{\text{TX}}^{\text{tot}} = \frac{B\sigma^2\rho\mathcal{S}_{\mathbf{x}}}{\eta}K, \quad (14)$$

where σ^2 is the channel noise power, $\mathcal{S}_{\mathbf{x}}$ is a term that accounts for the UE distribution and propagation environment, and η is a term that accounts for the efficiency of the PAs.

From the above formulation, it is important to restate that the value of the design parameter, ρ , is proportional to the UE SNR, which in turn, is directly proportional to the total transmit power, $P_{\text{TX}}^{\text{tot}}$, when considering ZF processing. Thus, finding the optimal total transmit power, $P_{\text{TX}}^{\text{tot}*}$, which maximizes the joint UL and DL energy efficiency involves deriving the optimal design parameter, ρ^* . After some manipulations, the authors of in [81] found a lower bound for the optimal design parameter, ρ^* ,

which shows that the optimal total transmit power, $P_{\text{TX}}^{\text{tot}*}$, increases with the load-independent circuit power consumption, $P_{\text{CP}}^{\text{li}}$, the power required by all the circuit components of each single-antenna UE, \mathcal{C}_1 , and the power consumed by the transceiver module attached to each antenna, \mathcal{D}_0 . This result also shows that the optimal strategy to improve the energy efficiency is to increase the total transmit power, $P_{\text{TX}}^{\text{tot}}$, with the number, M , of antennas, not in an arbitrary manner, but while considering the effect of the circuit power consumption, P_{CP} .

These conclusions are in stark contrast with those in [106]–[108], which as depicted earlier, concluded instead that the BS transmit power, P_{TX} , —the downlink part of the total transmit power, $P_{\text{TX}}^{\text{tot}}$ — monotonically decreases with the increase of the number, M , of antennas, while the system remains energy efficient. This implies a linear relationship between the energy efficiency and the number, M , of antennas. This optimistic finding only applies, however, in the idealistic case where the circuit power consumption, P_{CP} , plays a negligible role in the overall power consumption.

Moreover, the study on the optimal number, K^* , of multiplexed UEs, provided in [81], highlights that, if the power consumption for the beamforming processing, P_{SP} , and channel estimation, P_{CE} , are negligible, then the optimal number, K^* , of multiplexed UEs can be approximated as follows:

$$K^* \approx \left\lceil \left(\mu \sqrt{1 + \frac{U}{(\tau^{(\text{ul})} + \tau^{(\text{dl})}) \mu}} - 1 \right) \right\rceil, \quad (15)$$

where U is the channel coherence time (in symbols), and

$$\mu = \frac{P_{\text{CP}}^{\text{li}} + \frac{B\sigma^2 S_x}{\eta} \rho K}{\mathcal{C}_1 + \bar{\beta} \mathcal{D}_0}, \quad (16)$$

where $\bar{\beta}$ is the ratio of the number, M , of antennas to the number, K , of multiplexed UEs, (i.e., $\bar{\beta} = M/K$),

This bound indicates that the optimal number, K^* , of multiplexed UEs decreases with the power required by all the circuit components of each single-antenna UE, \mathcal{C}_1 , and the power

consumed by the transceiver module attached to each antenna, \mathcal{D}_0 , whereas it increases with the load-independent circuit power consumption, $P_{\text{CP}}^{\text{li}}$, the total transmit power $P_{\text{TX}}^{\text{tot}}$ (proportional to ρ), the noise power, σ^2 , and the parameter, \mathcal{S}_x . Note that this last term, \mathcal{S}_x , increases proportionally with the coverage radius of the cell, meaning that a larger number, K , of multiplexed UEs must be served when the coverage area increases in order to maximize the energy efficiency.

In addition, the study on the optimal number, M^* , of antennas, provided in [81], also indicates that such number, M^* , is lower bounded by:

$$M^* \geq \frac{\frac{B\sigma^2\mathcal{S}_x}{\eta\mathcal{D}'}\rho + \frac{C'}{\mathcal{D}'} + K - \frac{1}{\rho}}{\ln(\rho) + \ln\left(\frac{B\sigma^2\mathcal{S}_x}{\eta\mathcal{D}'}\rho + \frac{C'}{\mathcal{D}'} + K - \frac{1}{\rho}\right) - 1} - \frac{1}{\rho}. \quad (17)$$

This bound indicates that the optimal number, M^* , of antennas increases with the load-independent circuit power consumption, $P_{\text{CP}}^{\text{li}}$, and the power required by all the circuit components of each single-antenna UE, \mathcal{C}_1 , whereas it decreases with the power consumed by the transceiver module attached to each antenna, \mathcal{D}_0 .

Importantly, when the design parameter, ρ , becomes large, the lower bound given in eq. (17) can be approximated as:

$$M^* \approx \frac{B\sigma^2\mathcal{S}_x}{2\eta\mathcal{D}'} \frac{\rho}{\ln(\rho)}, \quad (18)$$

which shows that there is an almost linear scaling law of the optimal number, M^* , of antennas with respect to the design parameter, ρ , in its high value regime, and thus with respect to the total transmit power, $P_{\text{TX}}^{\text{tot}}$.

It should also be noted the linear dependence of the optimal number, M^* , of antennas with the UE distribution and propagation environment, captured by the variable, \mathcal{S}_x , which implies that a larger optimal number, M^* , of antennas is needed, as the size of the coverage area increases, since the variable, \mathcal{S}_x , increases with the cell radius.

For the sake of clarity, in Fig. 12, we illustrate the achievable energy efficiency for four different values of the number, M , of antennas, when varying the number, K , of multiplexed UEs. Note that this energy efficiency curves are characterized by a concave shape in all the considered antenna configurations, and that the maximum achievable values are highlighted in red. This plot highlights that the optimal number, K^* , of multiplexed UEs increases sub-linearly with the number, M , of antennas.

In a similar way, in Fig. 13, we illustrate the achievable energy efficiency for different values of the number, K , of multiplexed UEs, when varying the number, M , of antennas. The energy

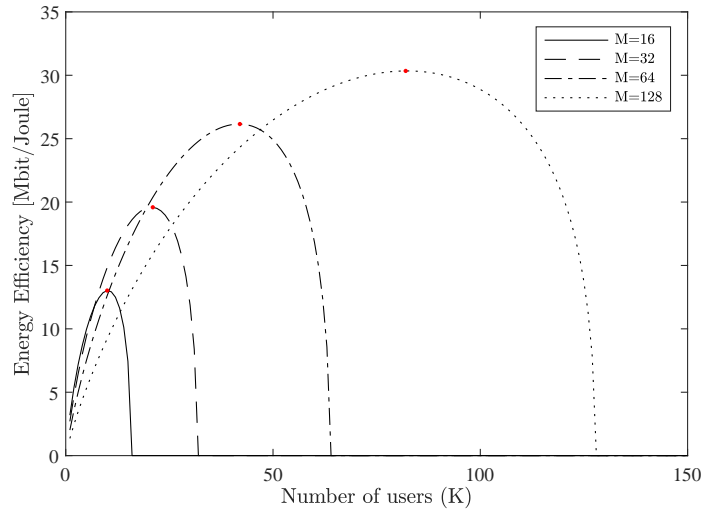


Figure 12. Energy efficiency with respect to the number, K , of multiplexed UEs for a given number, M , of antennas. The maximum energy efficiency values are highlighted in red.

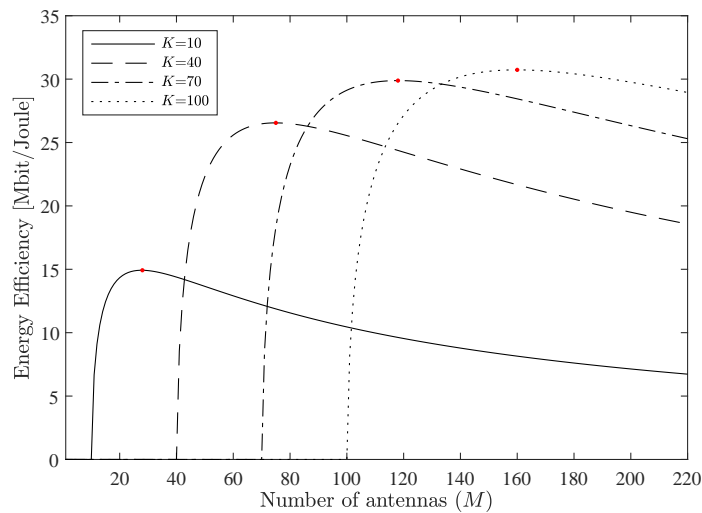


Figure 13. Energy efficiency with respect to the number, M , of antennas for a given number, K of multiplexed UEs. The maximum energy efficiency values are highlighted in red.

efficiency curves have, in this case, a quasi-concave shape, and the maximum achievable values are also highlighted in red. These results confirm that, for a given value of the number, K , of multiplexed UEs, augmenting the number, M , of antennas increases the energy efficiency up to a maximum energy efficiency value, where the UE rate gain due to the further increasing the number, M , of antennas is not sufficient anymore to counterbalance the cost incurred by their associated power consumption.

Finally, it should be noted from Fig. 12 and Fig. 13 that deploying hundreds of antennas to serve a large number of UEs is the optimal solution from an energy efficiency perspective, confirming the energy efficiency-enabler role of mMIMO. These results are well established and at the forefront of the state-of-the-art.

2) *Trade-offs*: The spectral efficiency, defined as the system throughput per unit of bandwidth [bit/s/Hz], has been historically adopted as the key optimisation metric to optimise mobile network deployments. As energy efficiency has become a new important KPI to the operation of 5G networks, the relation between these two metrics is thus of fundamental importance. How much spectral efficiency is traded to realize an energy efficiency gain?

Early works in the literature, not accounting for the circuit power consumption, P_{CP} , in the BS power consumption model, concluded that a 10x spectral efficiency improvement could be achieved for a given BS transmit power, P_{TX} , when adopting hundreds of antennas at the BS [51]. In this idealized scenario, since the energy efficiency grows proportionally to the number, M , of antennas, the corresponding energy-spectral efficiency curve is pushed outwards, meaning that both metrics can be simultaneously maximized.

As shown earlier, however, the circuit power consumption, P_{CP} , linearly grows with the number, M , of antennas, and when this number is large, it dominates the BS overall power consumption, implying that the energy efficiency cannot be enhanced, unless the efficiency of the BS hardware is improved accordingly [81]. In this way, the circuit power consumption, P_{CP} , breaks the monotonic relation between the energy and spectral efficiency, and makes these two metrics not consistent and conflicting with each other. For this reason, their trade-off must be carefully analyzed, and proper network optimisation techniques need to be developed to strike the right balance between these two metrics.

Studies on such energy and spectral trade-off are often carried out based on optimisation problems, aiming, e.g., at maximizing the energy efficiency given a spectral efficiency requirement. However, in some more sophisticated approaches, the balance between the energy and spectral efficiency is achieved by maximizing the resource efficiency (RE) [bit/Joule] [109], which is defined as the weighted sum of the energy and spectral efficiency, and the weights assigned to the two terms allows to explore the trade-off. It is important to highlight, however, that all such frameworks do not generally provide explicit equations for the energy-spectral efficiency trade-off. Instead, they mostly build on the top of results obtained from tailored optimisation algorithms aimed at solving such problem, which usually turns out to be intricate [110]–[113].

In this line, the fundamental trade-off between the energy and spectral efficiency has been carefully studied in [113], when considering a single-cell mMIMO system with linear precoding, i.e. ZF and maximum ratio transmission (MRT), and transmit antenna selection. In particular, the BS transmit power, P_{TX} , and the number, $M_a < M$, of active antennas are jointly considered as a resource to balance the energy and spectral efficiency, and the adopted BS power consumption model is equivalent to the one introduced in eq. (6). Note, however, that this research did not take into account the impact of the number, K , of multiplexed UEs, which may have a significant influence to the trade-off. Nonetheless, this study has importantly shown that different numbers, M_a , of active antennas lead to different energy efficiency-spectral efficiency curves. In more details, for a given number, M_a , of active antennas, the energy efficiency is a quasi-concave function with respect to the spectral efficiency, which confirms the existence of a clear trade-off between these two metrics. Their numerical results in this work have also shown that, in the low SNR region, which corresponds to a low spectral efficiency, MRT achieves a higher energy efficiency than ZF due to its lower complexity. In contrast, in the high spectral efficiency regime, ZF outperforms MRT in terms of spectral efficiency, for a given energy efficiency, owing to its ability of canceling intra-cell interference.

As discussed earlier, most theoretical mMIMO works in literature assume spatially uncorrelated channels and perfect CSI knowledge due to tractability reasons, as in the aforementioned study. However, CSI acquisition is challenging, especially when dealing with the large antenna arrays in the mMIMO case. As explained in Section II-A, in TDD systems, the acquisition of downlink CSI can be facilitated via uplink training by taking advantage of the channel reciprocity. However, even in this case, the CSI may still be inaccurate due to practical hardware limitations, such as calibration errors in the transceivers [114], and in high mobility scenarios, it can quickly become outdated.

Tackling this challenge, the authors in [111], [112] have explored this problem, and provided analyses of the energy and spectral efficiency trade-off, considering statistical CSI knowledge, instead of an instantaneous one. This type of feedback, at the expense of a lower network capacity, has the advantage of being stable during longer time periods and to be more easily obtainable by the BS through long-term feedback or covariance extrapolation.

Particularly, the work in [112] has analysed the single-cell mMIMO downlink case, where one BS with M antennas simultaneously transmits signals to K UEs, the channel spatial correlation is captured using a jointly correlated Rayleigh fading model [115], statistical CSI is available at

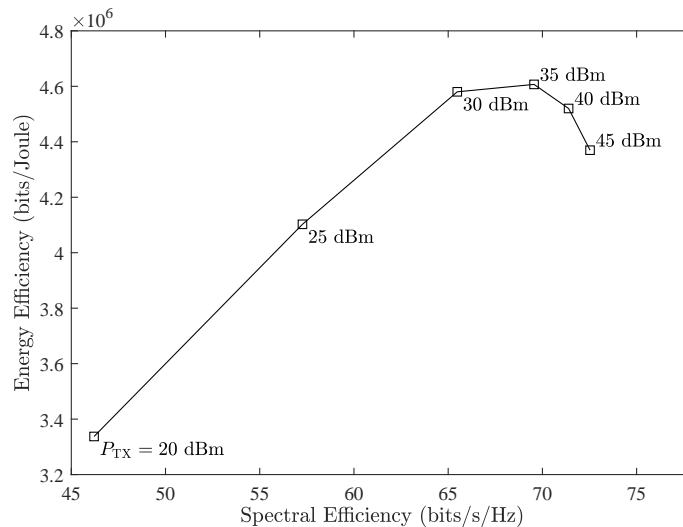


Figure 14. Trade-off between energy and spectral efficiency for different values of the BS transmit power, P_{TX}

the transmitter, and the considered BS power consumption model accounts for the BS transmit power, P_{TX} , and the static circuit power consumption, P_{CP} . The energy-spectral efficiency trade-off has been investigated by maximizing the RE metric [109], which strikes for a energy-spectral efficiency balance.

Fig. 14 shows the derived trade-off when considering different values of the BS transmit power, P_{TX} . Note that the energy-spectral efficiency curve is characterized by a concave shape. In particular, the energy and spectral efficiency can be jointly augmented by increasing the BS transmit power, P_{TX} , until reaching, in this case, the optimal energy-spectral efficiency point with a BS transmit power, $P_{TX} = 35$ dBm. After reaching such optimal point, which corresponds to the maximum achievable values of both energy and spectral efficiency, the BS transmit power, P_{TX} , becomes the dominating consuming factor. Thus, increasing the BS transmit power, P_{TX} , allows to increase the spectral efficiency only at the expense of a reduced energy efficiency.

Finally, before concluding this section, a summary of the main works on single-cell energy efficiency bound and trade-offs is presented in Table VII.

As concluding remark, it should also be noted that, to enhance the energy efficiency of mMIMO systems, large research efforts have also being spent on the understanding and optimization of both practical precoding and UE scheduling techniques, while considering single-cell mMIMO scenarios. The optimal precoding to achieve optimal energy efficiency under ideal

Table VII

SUMMARY OF SINGLE-CELL ENERGY EFFICIENCY BOUND AND TRADE-OFFS IN THE LITERATURE

Paper	Type	KPI	Parameters
[116]	bound	EE	bandwidth, BS transmit power
[81]	bound	EE	transmitting antennas, multiplexed UEs, BS transmit power
[113]	trade-off	EE-SE	transmitting antennas, BS transmit power
[112]	trade-off	RE	BS transmit power
[111]	trade-off	RE	BS transmit power

and known channel conditions has been derived in [117]. Optimal energy efficient precoding schemes, while considering imperfect CSI at the transmitter, have been studied in [118]. UE scheduling algorithms that take channel orthogonality into account, avoiding to schedule two *nearby* UEs in the same time-frequency resource, have been proposed in [119]–[122]. Other more experimental approaches considered to increase energy efficiency have been modulation diversity [123], cognitive radios [124], and spatial modulation [125].

B. Multi-cell scenario

In this subsection, we focus —and survey— the latest most representative developments on the understating of the energy efficiency in large-scale multi-cell mMIMO networks, introducing the most relevant tools used to carry the analyses out and differentiating among the uplink and downlink case.

In contrast to the single-cell mMIMO scenario, it should be noted, however, that the multi-cell mMIMO one is not so well theoretically understood as of today, due to its much higher complexity. Indeed, large-scale networks are in general hard to model and analyse (see Fig. 15), due to their

- complex topology,
- evolving UE distributions,
- dynamic traffic demands,
- fluctuating wireless channels,
- sophisticated protocols and algorithms, and
- large number of parameters to tune.

This poses a challenge on their theoretical comprehension. In fact, there are no available explicit —and general— closed-form expressions that describe the energy efficiency bounds and trade-

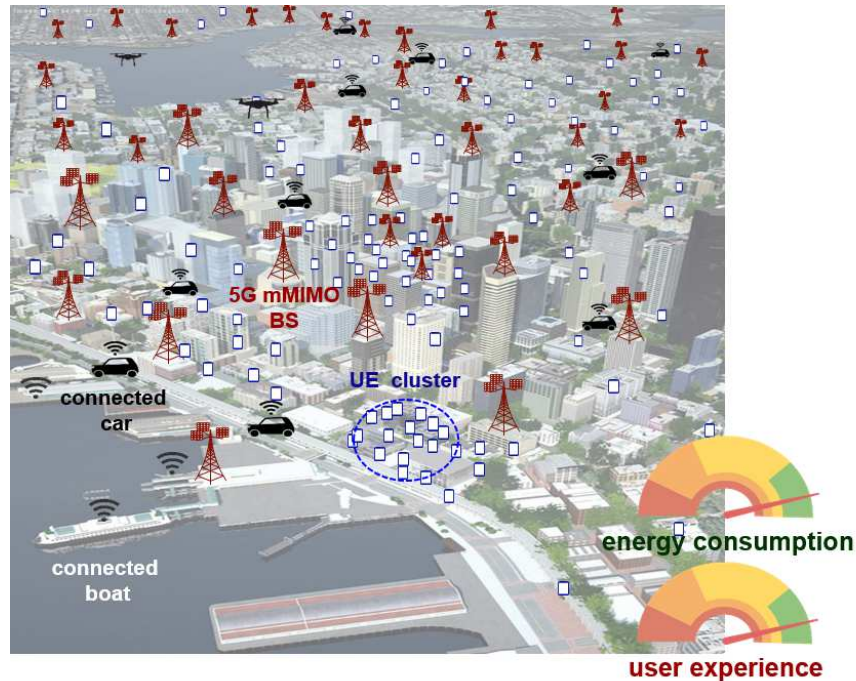


Figure 15. Example of a large-scale multi-cell mMIMO network.

offs of a multi-cell mMIMO network in a holistic manner. Instead, the research in this area is scattered and focused on particular aspects of the system, which together with some of the assumptions taken, have helped to increase the tractability of the problem, and derive some initial insights on the energy efficiency problem. Further research on the topic is needed.

1) *Available Mathematical Tools*: Given that the comprehension of the power consumption of a mMIMO BS has reached some degree of maturity (see Sections III and V-A), one of the main difficulties to unlock a fundamental understanding of the energy efficiency in multi-cell mMIMO networks is the derivation of its capacity. Stochastic geometry [126], the state-of-the-art tool for the theoretical analysis of multi-cell networks, particularly small cell ones [127]–[130], has been recently started to be used to address this issue.

Embracing the randomness of today’s deployments, the work in [131] has analysed the uplink mMIMO performance, considering a large-scale multi-cell network and the pilot contamination problem. This paper has presented a mMIMO network capacity scaling law as a function of the number of antennas and multiplexed UEs per mMIMO BS, and investigated the performance

gains attainable through a practical fractional uplink pilot reuse¹², used to mitigate pilot contamination. It is important to note, however, that significant assumptions were taken for tractability reasons. The study assumed that there was a large number of UEs in each BS at all times, and thus that all BSs in the network were active and that all uplink pilots available were in use in all BSs. This makes it difficult to study off-peak hours where energy savings are more likely to be harvested. Moreover, a pure non-line-of-sight (NLoS) path loss model with Rayleigh multi-path was considered, which cannot capture the important effect of channel correlation on mMIMO performance. Unfortunately, even though these important assumptions and simplifications were made, the resulting expressions are still not tractable, requiring few folds of integrals, making difficult to infer relationships among system parameters. This is particularly true when capacity expressions need to be coupled with complex BS power consumption models, as those presented in Section III, to derive the energy efficiency.

Using similar stochastic geometry tools, the authors in [132] studied the downlink mMIMO performance, while considering a heterogeneous network comprised of mMIMO macrocells and single-antenna small cells, together with the effect of line-of-sight (LoS) and NLoS transmissions, pilot contamination, and cross-tier interference (i.e., interference from macro to small cells and vice versa). However, similar as in the previous case, a fully loaded network is considered with all uplink pilots at use, which significantly affects the pilot contamination, and does not allow to analyse low and medium traffic load scenarios. In addition, important assumptions also revolve around the use of Rayleigh fading for the LoS transmissions. Some of these issues have been recently addressed by the research in [133], where pilot allocation schemes are considered with the objective of reducing the pilot contamination. However, the expressions are still too complex to infer parameter relationships without a numerical evaluation.

2) *Uplink mMIMO network deployment perspectives:* Aware of such complexities, the authors in [134] proposed a tractable stochastic geometry-based analysis of the energy efficiency of a TDD multi-cell mMIMO network subject to end-users' QoS demands, while using a simplified but yet complete system model. In more detail, the authors targeted at maximizing the uplink energy efficiency of such multi-cell mMIMO network, while ensuring a minimum uplink spectral efficiency to the average UEs, and used this formulation to obtain some practical insights into the

¹²Fractional uplink pilot reuse refers to a class of schemes where cell-centre UEs of every cell reuse the available set of uplink pilots for CSI estimation, while cell-edge ones can only reuse a subset of them, with this subset being orthogonal in neighboring or more coupled cells.

mMIMO energy efficiency problem. The adopted homogeneous Poisson point process (HPPP)-based system model¹³ accounted for

- the BS density, λ ,
- the number, M , of antennas per BS,
- the number, K , of multiplexed UEs per transmission time interval (TTI) at each BS,
- an idealised uplink fractional power control,
- imperfect channel estimation through pilot contamination¹⁴,
- hardware impairments, modelled as a reduction of the signal power by a factor, $1 - \epsilon^2$,
- maximal ratio combining (MRC) received filters at the BS,
- single-antenna UEs, and
- an uplink pilot reuse scheme.

Importantly, the complexity of calculating the average spectral efficiency of the network as the sum of the spectral efficiencies of all UEs through this framework was acknowledged, indicating the need for heavy numerical evaluations of integrals. To address this issue, and obtain explicit expressions, the authors focused on the performance of the average UE instead, and derived the following tractable—but yet tight—lower bound for its uplink signal to interference plus noise ratio (SINR):

$$\text{SINR}_{\text{UL}} = \frac{M(1 - \epsilon^2)^2}{\left(K + \frac{\sigma^2}{\zeta}\right)\left(1 + \frac{2}{\beta(\alpha-2)} + \frac{\sigma^2}{\zeta}\right) + \frac{2K}{\alpha-2}\left(1 + \frac{\sigma^2}{\zeta}\right) + \frac{K}{\beta}\left(\frac{4}{(\alpha-2)^2} + \frac{1}{\alpha-1}\right) + M(1 - \epsilon^2)\left(\frac{1}{\beta(\alpha-1)} + \epsilon^2\right)}, \quad (19)$$

where α is the path loss exponent, β is the pilot reuse factor, ζ is the path loss compensation power control coefficient, and σ^2 is the noise power.

Leveraging this expression, the authors formulated the uplink average UE spectral efficiency and the resulting uplink area spectral efficiency (ASE), and with that, they derived the uplink area power consumption (APC) of the multi-cell mMIMO network, using a linear version of the BS power consumption model presented in (7), i.e.,

$$\text{APC}_{\text{UL}} = \lambda \left(\left(1 - \frac{\beta K - 1}{S}\right) \frac{\zeta \omega \Gamma\left(\frac{\alpha}{2} + 1\right)}{\eta (\pi \lambda)^{\left(\frac{\alpha}{2}\right)}} K + P_{\text{CP}}^{\text{li}} + \mathcal{C}_1 K + \mathcal{D}_0 M + \mathcal{D}_1 M K \right) + \mathcal{A} \cdot \text{ASE}, \quad (20)$$

where S is coherent block length, which is related to the channel coherence time, U , in eq. (15), η is the PA power efficiency, $\Gamma(\cdot)$ is the Gamma function, and ASE is the area spectral efficiency.

¹³It should be also noted that this general system model allows to compare mMIMO setups, with few BSs and many antennas per BS, to small cell ones, with many BSs and few antennas per BSs, thus providing guidance on the design of future green wireless networks, from the uplink perspective.

¹⁴Rayleigh fading was assumed, and thus the effect of spatial correlation was ignored in this analysis.

With these formulations, the authors defined an optimisation problem to find the most uplink energy efficient network deployment, while provisioning the average UE with a minimum SINR. Through some mathematical manipulations detailed in [134], the authors derived

- the uplink spectral efficiency feasibility region,
- the optimal uplink pilot reuse factor, β^* ,
- the optimal BS density, λ^* ,
- the optimal number, K^* , of multiplexed UEs per TTI at each BS, subject to a given ratio, $\frac{M}{K}$, of the number, M , of antennas to the number, K , of multiplexed UEs, and
- the optimal number, M^* , of antennas per BS, subject to a given number, K , of multiplexed UEs,

and demonstrated that reducing the cell size —increasing the BS density— is beneficial for energy efficiency, but that such positive effect saturates when the circuit power consumption dominates over the transmission power. Their results also showed that adding more antennas in a controlled manner to the BS to bring it to the mMIMO regime also enhances the energy efficiency¹⁵. In more detail, their numerical examples on a typical scenario resulted in the maximum energy efficiency when having 91 antennas and 10 UEs multiplexed per BS, which resembles a mMIMO —and not a small cell— setup. It should be noted that the energy efficiency gains, in this case, mostly came from the intra-cell interference suppression provided by mMIMO, and by sharing the circuit power costs among the various multiplexed UEs in the same time-frequency resource. Moreover, the analysis showed that a large pilot reuse factor can be used to protect the network against inter-cell interference, and that it can be tailored to guarantee a certain average UE spectral efficiency.

Using a different modelling tool than stochastic geometry, based on numerical evaluations, but also focusing on the uplink energy efficiency, the authors in [135] provided an analysis of the area energy efficiency (AEE) and the ASE trade-off for different system parameters, such as the pilot reuse ratio and number of antennas and multiplexed UEs per BS. Importantly, even if this work does not focus on minimising power consumption, but maximising the AEE, it reaches similar general conclusions than those of [134]. The multi-cell mMIMO network always

¹⁵These conclusion are inline with that presented in the analysis of the single-cell mMIMO case in the previous section, indicating that the optimal strategy to improve the energy efficiency is to increase the total transmit power, P_{TX}^{tot} , with the number, M , of antennas, not in an arbitrary manner, but while considering the circuit power consumption, P_{CP} .

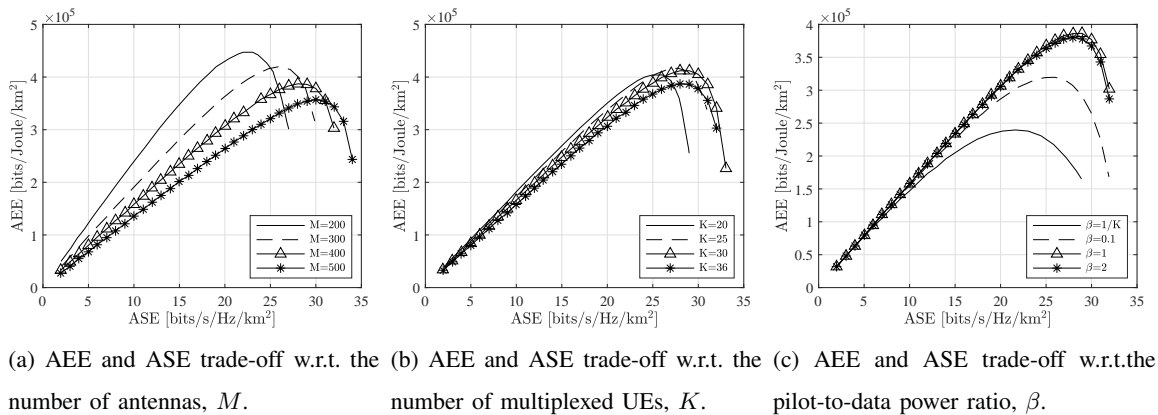


Figure 16. Multi-cell AEE and ASE trade-offs

performs better in terms of ASE with an increasing number of antennas. However, adding more antennas or multiplexed UEs might not achieve the optimal AEE. The reason is that the power consumed by the transceiver module attached to each antenna and the related signal detection and processing have a non-negligible effect on the total BS power consumption. In more detail, their results have shown that, when considering smaller ASE targets, a lower number of antennas (see Fig. 16a) and multiplexed UEs (see Fig. 16b) at the BS suffices to achieve the optimal AEE. However, more antennas and multiplexed UEs are required to satisfy higher ASE requirements, which makes the network less energy efficient, resulting in a smaller AEE, as soon as the transmit power dominates the circuit and processing power consumption. Finally, in this work, the authors have also derived the pilot-to-data power ratio that maximizes the AEE, and studied how such parameters affects the AEE and the ASE trade-off (see Fig. 16c). Moreover, they show that in relevant scenarios, the optimal number of multiplexed UEs for maximizing the AEE is much smaller than half of the coherent clock length, as it is for maximizing the ASE [51].

3) *Downlink mMIMO network deployment perspectives:* Using a similar stochastic geometry approach and BS power consumption model as in [134], the authors in [136] extended the previous work to the downlink case, aiming at optimising the downlink energy efficiency of a TDD multi-cell mMIMO network, while providing a minimum spectral efficiency to the average UE. It is important to note that perfect CSI was assumed in this case due to the complexity of modelling in the same framework both the channel estimation phase in the uplink and the

data transmission phase in the downlink¹⁶. As a result, the effect of pilot contamination was neglected, decreasing the accuracy of the model. ZF precoders at the BS and single-antenna UEs were considered. Importantly, it should be noted that embracing the same methodology as in [134], the authors also focused their analysis on the performance of the average UE, and derived the following tractable lower bound for its downlink SINR:

$$\text{SINR}_{\text{DL}} = \frac{(1 - \epsilon^2)(M - K)}{\frac{2K}{(\alpha-2)} + \epsilon^2(M - K) + \frac{\Gamma(\alpha/2+1) \omega \sigma^2}{(\pi \lambda^2) \rho}}, \quad (21)$$

where ρ now represents the downlink transmit power allocated by the BS to the average UE.

With such expression, and using the corresponding downlink APC model, shown in the following for the sake of clarity:

$$\text{APC}_{\text{DL}} = \lambda \left(\frac{K\rho}{\eta} + P_{\text{CP}}^{\text{li}} + \mathcal{C}_1 K + \mathcal{D}_0 M + \mathcal{D}_1 MK \right) + \mathcal{A} \cdot \text{ASE}, \quad (22)$$

the authors derived, following similar steps as in the uplink counterpart work, the following optimal operation points:

- the optimal downlink transmit power per UE, ρ^*
- the optimal BS density, λ^* ,
- the optimal number, K^* , of multiplexed UEs per TTI at each BS, subject to a given ratio, $\frac{M}{K}$, of the number, M , of antennas to the number, K , of multiplexed UEs, and
- the optimal number, M^* , of antennas per BS, subject to a given number, K , of multiplexed UEs.

The analysis of the obtained closed-form expressions showed that the same conclusions obtained from the uplink analysis apply to the downlink one. The optimal energy efficiency is achieved by a mMIMO-like deployment, where in their particular example, the optimum number of antennas and multiplexed UEs per BS are 193 and 21, respectively. Note that the results for the downlink case in this section and those for the uplink one in the previous section differ due to the different UE requirements and the different power consumption of a BS and a UE while transmitting.

Finally, before concluding this section, a summary of the main works on multi-cell energy efficiency modelling, bound derivation, optimization and trade-offs is presented in Table VIII.

¹⁶Such downlink dependency on the uplink is the main reason why the theoretical downlink energy efficiency of mMIMO networks has been less rigorously studied.

Table VIII
SUMMARY OF MULTI-CELL ENERGY EFFICIENCY BOUND AND TRADE-OFFS IN THE LITERATURE

Paper	Type	KPI	Model features
[131]	uplink/stochastic geometry (SG)	coverage probability/ASE	macrocells, pilot contamination
[132]	downlink/SG	coverage probability/ASE	HetNet, pilot contamination, LoS/NLoS
[134]	uplink/SG and optimization	APC/EE	macrocells, pilot contamination, hardware impairments
[135]	downlink/analytical model/trade-off	AEE-ASE	macrocells, pilot contamination, pilot reuse
[136]	downlink/SG and optimization	APC/EE	macrocells, hardware impairments

As a concluding remark, it is also important to highlight that, the presented uplink and downlink results in this section advocate for mMIMO deployments —and their dimensioning optimisation—, according to end-users’ QoS demands, as an important tool to increase energy efficiency in 5G networks. Using an inadequate number of BS or antennas per BS to meet a given end-users’ QoS demands can result in highlight suboptimal energy efficiency performances. This is an important area of research were more work is needed.

Moreover, it should me mentioned that, despite the lack of work on a holistic mMIMO deployment dimensioning and operation optimization, there is a large body of work around mMIMO power control optimization for energy efficiency maximization in multi-cell mMIMO networks. In this line, the following research stands out. In [98], systematic approaches to solve energy efficiency maximization problems are extensively discussed. In this regard, the framework presented in [137] has provided network- and UE-centric downlink power control algorithms, where minimum rate constraints are imposed and the SINR takes a general form, able to deal with complex mMIMO systems/configurations. Centralised algorithms are also developed, which are guaranteed to converge, with affordable computational complexity, to a Karush–Kuhn–Tucker point of the considered non-convex optimisation problem. Building on such framework, the work in [138] has proposed a framework to compute suboptimal power control strategies with even more affordable complexity. This is achieved by jointly using fractional programming and sequential optimisation. Numerical evidence has shown that such sequential fractional programming framework achieves global optimality in several practical communication scenarios.

VI. TIME-DOMAIN (SYMBOL SHUTDOWN) ENERGY SAVING-BASED SOLUTIONS

As motivated in Section II, future wireless communication systems require efficient hardware and mechanisms that enable the adaptation of the network functionalities and parameters to the load variations in an on-line manner in order to avoid excessive energy consumption, while ensuring the end-users' QoS demands. These enabling technologies can be classified, for example, by observing the time-scale and the domain in which they operate, e.g., time, frequency, or spatial (antenna) domains.

In this section, we concentrate on time-domain mechanisms, also known as symbol shutdown. These schemes refer to those solutions that are targeted at harvesting energy savings by considering the short-term traffic load variations in a cell. They operate in the time scale of hundreds of microseconds to hundreds of milliseconds, taking advantage of the 3GPP NR lean carrier design, to rapidly (de)activate hardware components of the BS according to the presence or absence of traffic. Note that such short-term traffic load variations are a function of several factors, such as the number of active UEs, the traffic types, and the interference, and that to realise such energy savings, fast reacting BS hardware and radio resource management mechanisms, which operate in the 3GPP NR orthogonal frequency division multiplexing (OFDM) or slot time scale, are thus required to dynamically minimum energy consumption, while avoiding any UE performance degradation. Given that such radio resource management mechanisms may operate at an OFDM symbol time scale, and that they need to guarantee cell service continuity, it is important to note that they must be designed and operated embracing the 3GPP NR signalling framework provided by the lean carrier design, to make sure they do not violate any control-plane protocol.

In the following, we overview the 3GPP NR lean carrier design capabilities, and present state of the art symbol shutdown mechanisms and frameworks for their optimization.

In 3GPP LTE, cell DTX can be used for deactivating BS hardware components, such as the PA, when transmissions are absent in a given frame [139]. The benefits of cell DTX are, however, limited by the control signalling required by 3GPP LTE to drive UE cell camping procedures and others, even in the absence of traffic. In more detail, the 3GPP LTE frame lasts 10 ms, and it is composed by 10 sub-frames, each one including 14 OFDM symbols with a duration of $71.4 \mu\text{s}$. In the unicast mode, cell-specific reference symbols (CRSSs) are transmitted in every sub-frame, primary synchronisation channels (PSSs) and secondary synchronisation channels (SSSs) are transmitted every 5 ms, and broadcast channels (BCHs) are repeated in every first

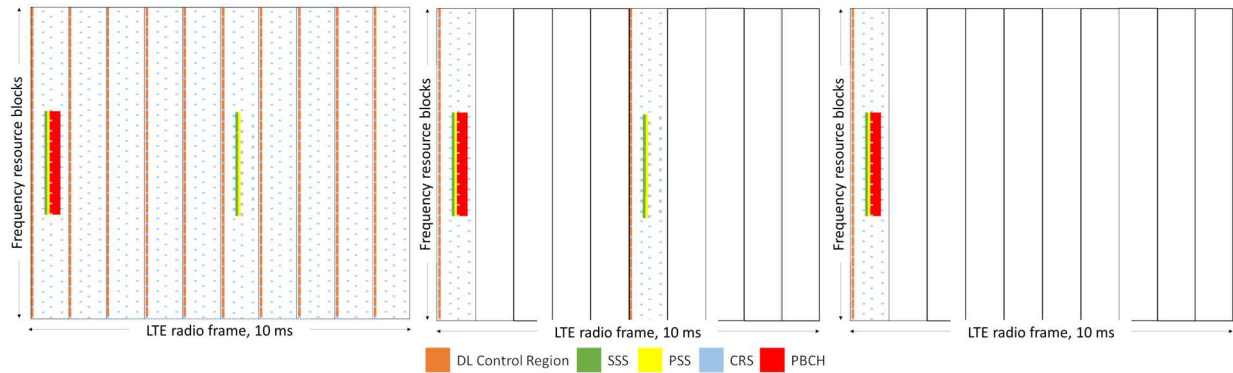


Figure 17. LTE frame for (left) unicast cell; (middle) Release 14 MBSFN-dedicated cell; (right) Release 14 FeMBMS/unicast-mixed cell.

sub-frame of a frame (see the left side of Fig. 17). Therefore, the BS is only able to sleep for a few OFDM symbols, and needs to wake-up one OFDM symbol before the transmission of any control signal. In this case, cell DTX can achieve at most a 33% sleep ratio at zero load [61].

To allow longer sleep periods in 3GPP LTE, the use of the multicast-broadcast single-frequency network (MBSFN) frame was proposed, which is a technology introduced to enable mobile television broadcasting, characterized by the need of less frequent signalling. To standardise such technology, the 3GPP evaluated in [58] both the sleep ratio (at symbol level and sub-frame level) and the sleep duration, when using

- 3GPP LTE Release 14 FeMBMS/unicast-mixed cells, in which two out of ten sub-frames, i.e. sub-frames 0 and 5, contain unicast control signaling (see the central plot in Fig. 17), and
- 3GPP LTE Release 14 MBSFN-dedicated cells, where only one non-MBSFN sub-frame with standard unicast signaling is transmitted every 40 ms (see the right side of Fig. 17).

The results from the 3GPP studies showed that, in the FeMBMS/unicast-mixed mode, a BS could stay in energy saving mode for up to 4 ms, which leads to an 80% sleep ratio. Importantly, with the MBSFN-dedicated mode, a BS could sleep even further up to 39 ms, which results into a 93.75% sleep ratio.

As discussed in Section II-B, in contrast to 3GPP LTE, 3GPP NR is characterized by a user/data-specific signalling instead of a cell-specific one —the lean carrier (see Fig. 18). Specifically:

- The CRS is not used anymore in 3GPP NR, and a synchronization signal (SS) burst set,

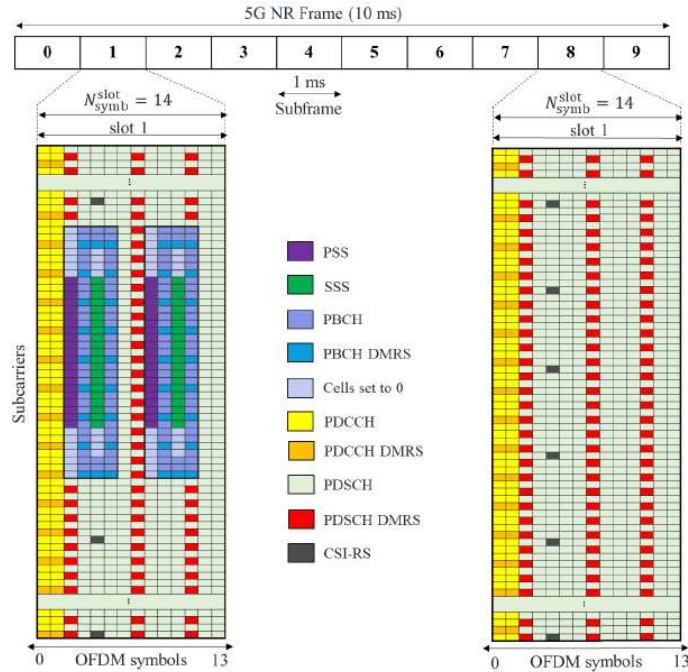


Figure 18. NR frame structure example [46].

including one or multiple SSBs —each one of them in turn comprised of PSS, SSS, and BCH— is transmitted to support UE cell (re)selection and handover procedures with a larger periodicity, i.e., 5, 10, 20, 40, 80, and 160 ms [29].

- Since the CRS is not used, CSI acquisition procedures have also been redesigned, and on demand CSI-RSs are reused —and further extended— to provide support for beam and mobility management as a complement to the SSB [29].
- The minimum required system information (SI) broadcast in 3GPP NR has also been reduced with respect to 3GPP LTE, and the part not strictly necessary for network entry is transmitted on-demand now [140].
- Moreover, in 3GPP NR, a single-antenna port can be used to transmit the mandatory control signals, while, in 3GPP LTE, all the antenna ports are used to transmit such mandatory control signals [61].

This lean carrier design enables both larger sleep ratios and longer sleeping duration, whose particular values depend on the specific system numerology, i.e., subcarrier spacing, the number of SSBs per SS burst set, and the SS burst set periodicity. For example, when considering a subcarrier spacing of 15 kHz and two SS blocks per SS burst set, a 3GPP NR cell can stay in

sleep mode up to 19 ms, which leads to a sleeping ratio of 95% for a SS burst set periodicity of 20 ms. Importantly, when the SS burst set periodicity is maximised to 160 ms, a 3GPP NR cell can stay in sleep mode up to 159 ms, which leads to a sleeping ratio of 99.38% [58].

As introduced earlier, the longer sleeping duration enabled by the lean carrier design also allows for deeper sleeps, where more BS hardware components can be switched off. In this line, multiple BS sleep states have been defined, where each sleep mode is associated to a given sleeping duration [82] [141]. As a rule of thumb, all the BS hardware components that can enter and exit a sleep mode fast enough with respect to the sleep mode duration can be easily deactivated in such period, while the components having a longer latency need to remain active. In more detail, and to give an example according to 3GPP discussions [141]:

- In sleep mode 1 (i.e., cell DTX [139]), characterized by a duration (deactivation plus reactivation time) of 71 μ s, the PA and some components of the digital BBU and the analog front-end (FE) can be deactivated.
- In sleep mode 2, which has a minimum duration of 1 ms, additional components of the analog FE can be switched off.
- In sleep mode 3, which has a minimum duration of 10 ms, the BS can additionally deactivate all the digital BBU processing, and almost all the analog FE (except the clock generator).
- Finally, in sleep mode 4, which has a minimum duration of 1 s, only the wake-up functionalities are maintained.

Interested readers should note that the complete list of BS elements that can be switched off for each sleep mode can be found in [142].

In periods with absence of traffic, a BS can go through the presented sleep modes subsequently to reduce its energy consumption. This process is often referred as ASM (see Fig. 19). When the cell load rises and UE traffic appears, such UE traffic is buffered, and the BS has to immediately switch on its functionalities, and serve the required data to satisfy the end-users' QoS demands. However, it should be noted that, since hardware activation and deactivation times are not negligible, and their lengths increase with the number of involved BS hardware components [141], ASMs may increase the UE perceived latency, and this can accordingly affect the overall UE performance [143]. Therefore, there is a need for optimising the 'path' through the different sleep modes, and proactively activate the BS hardware components before traffic arrives in order to limit performance losses.

To assess the positive impact of different ASMs and their optimization on energy efficiency,

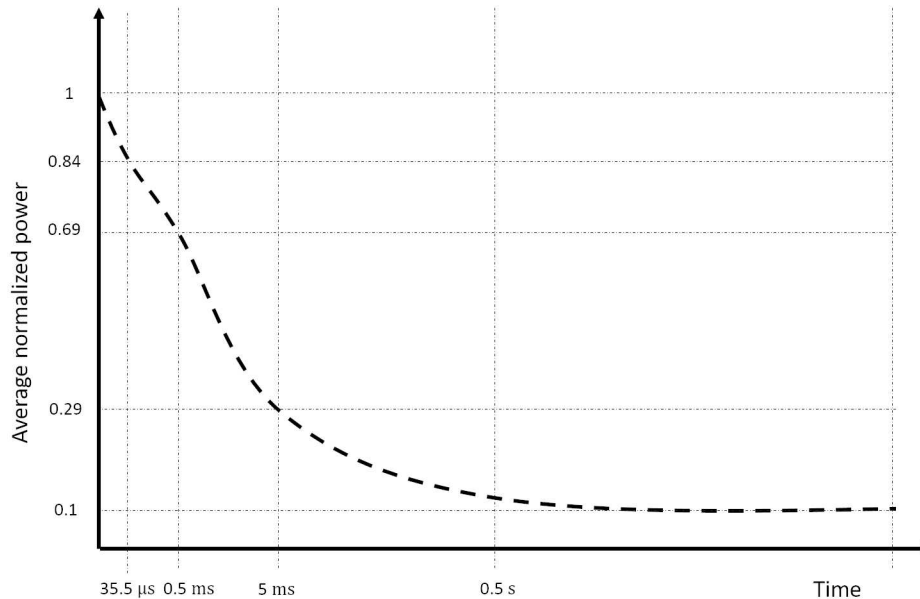


Figure 19. Power consumption trend as the BS enters in successive sleep modes [82]. Activation time is assumed to be equal to the half of the minimum sleep duration [141].

a new BS power consumption model was proposed in [144], i.e.,

$$P_{BS} = \begin{cases} P_{TX} + P_{CP}, & \text{if } P_{TX} > 0, \\ \delta_s P_{CP}^{li}, & \text{if } P_{TX} = 0 \text{ and sleep mode is active,} \end{cases} \quad (23)$$

where P_{TX} , P_{CP} , and P_{CP}^{li} are the transmit power, the circuit power, and the load independent part of the circuit power consumption, respectively, and δ_s is the fraction of the load-independent circuit power consumption, P_{CP}^{li} , required by the cell in sleep mode. In this case, the authors assume that the fraction, δ_s , is equal to 0.84, 0.69, and 0.29 for sleep mode 1, sleep mode 2, and sleep mode 3, respectively, while for sleep mode 4, the fraction, δ_s , can be lower than 0.1. It is important to note that current 3GPP NR implementation cannot realise this last mode of operation, as it requires 1 s of continuous sleeping period, and 3GPP NR can only do up to 160 ms.

Embracing this framework, the work in [143] has proposed the use of dedicated timers to control when to deactivate components and go into deeper sleep modes. The authors have highlighted that these timers should depend on the type of traffic carried out by each single BS, and to make a more flexible usage of the sleep modes, they have designed a RL model, based on Q-Learning, to optimise the duration of each sleep state. Energy savings and experienced

delay are balanced using this technique, using as enabler the average packet inter-arrival time. Importantly, their results have shown that it is possible to implement sleep modes and achieve significant energy savings, even with stringent delay constraints, for low to medium traffic load scenarios. In more detail, up to 80% energy savings can be obtained when replacing 3GPP LTE with 3GPP NR technology and using the proposed sleep modes. Nevertheless, it should be noted that this type of scheme relies on a continuous exchange of network signalling, which may impact, e.g., the performance of cell (re)selection processes [145]. More work is thus needed in this direction to understand how sleep modes impact the network load distribution, the resulting inter-cell interference, the related radio resource management (RRM) procedures, and finally the end-users' QoS.

To finalise this section, let us highlight that to further enhance the performance of sleep modes, avoiding the large access delay for UEs due to the time incurred by some BS hardware components to power up, the latest proposals in 3GPP NR Release 18 suggest the use of some potential assisting information sent from the UE to the BS to setup ASM and transmission parameters and maximize the power saving without greatly degrading the end-users' QoS. Such information may include —but it is not limited to— latency requirements. Multi-carrier operation schemes where one carrier can transmit the control signalling —or a lighter version of it— of another sleeping carrier due to sleep mode operation have also been proposed to avoid UE performance degradation due to transients [146].

Table IX provides a summary of the main contributions of the time-domain energy saving schemes discussed in this section.

Table IX
SUMMARY OF MAIN TIME-DOMAIN ENERGY SAVING SOLUTIONS.

Paper	Focus	Main contribution/idea
[139]	Cell DTX	Enable cell shutdown if there is no data in a given frame
[58]	Compare EE in NR and LTE	Evaluate sleep ratio and sleep duration in NR and LTE
[82] [141]	ASM	Leveraging lean carrier design, introduce deeper sleep modes
[144]	Sleep mode EE evaluation	Introduce a power model for sleep modes
[143]	Study and improve sleep modes	Introduce timers to control sleep modes

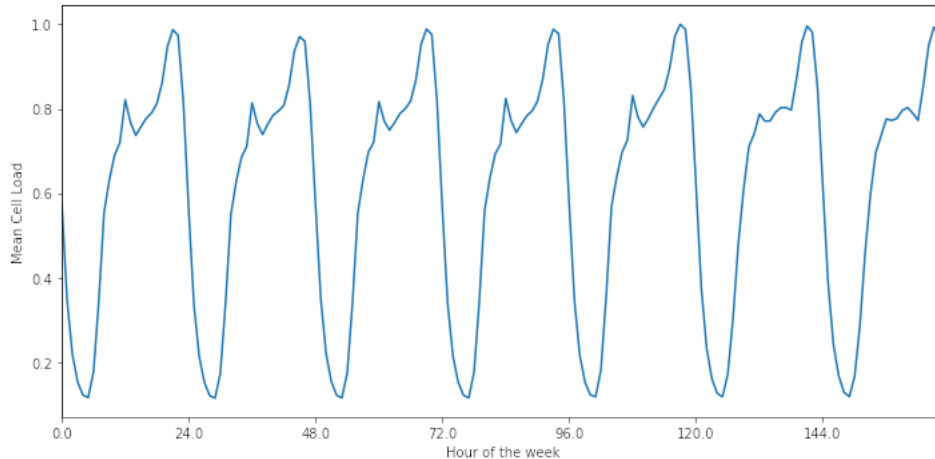


Figure 20. Normalized weekly load for a typical cell in a dense urban scenario.

VII. CARRIER-DOMAIN (CARRIER SHUTDOWN) ENERGY SAVING-BASED SOLUTIONS

In this section, we concentrate on carrier-domain mechanisms, also known as carrier shutdown. These schemes refer to those solutions that take advantage of deep dormancy to deactivate hardware components of the BS for long periods of time, e.g. time scales of minutes, even hours. In more detail, we focus on the mechanisms used to harvest energy savings by considering the long-term traffic load variations of a cell. To give an example, in most of the cells, the traffic daily profile shows a regular trend, with low load periods early in the morning, medium loads during the work-hours and high data rate in the late evening (see Fig. 20). Weekends may be characterized by lower traffic demands with respect to workdays. As a result, since mobile networks are sized to satisfy peak time traffic, energy may be wasted during low and medium load periods, and thus energy saving mechanisms able to adapt the network configuration to this long-term traffic load variations by fully (de)activating cells are necessary. Note that these schemes allow deeper sleeps than the symbol shutdown ones.

It is important to note that, since these solutions operate at large time scales, they are not usually designed for providing cell service continuity, but for switching off carriers. As a result, carrier-domain mechanisms do not need to be designed and integrated with the 3GPP NR signalling framework provided by the lean carrier design, and can be operated over the top. However, 3GPP NR specification work was still required to coordinate carrier (de)activation among BSs.

According to the degree of coordination among BSs, and the nature of such BSs, three main

approaches to carrier-domain mechanisms have been considered by the 3GPP, i.e. intra-BS, inter-BS, and inter-RAT energy saving mechanism [147], which are further discussed in the following.

A. *Intra-BS energy saving mechanisms*

In the first case, intra-BS energy saving, a BS may activate energy saving mechanisms to locally adapt its capacity to traffic requirements with the main aim of reducing RF amplifier power consumption. At this level, the tasks of each BS is to independently control its number of active cells.

In this line, to improve the system energy efficiency, besides switching off one or more sectors of a BS, which may be risky in terms of coverage, it is practically more important and feasible to dynamically control the number of active CCs, when a multi-carrier system or CA is deployed. As discussed in Section III-A1, CA is a 3GPP flagship technology introduced in 3GPP LTE-A, which allows a BS to simultaneously operate on different bands. In 3GPP NR Release 15, the dormant state was introduced in CA, such that (de)activation delay for SCells could be reduced, and the set of active CCs could be rapidly adapted to match UEs requirements [42]. In addition, the concept of bandwidth part (BWP) was also introduced in 3GPP NR Release 15 [148]. With this mechanism, when the cell load is low, the BS can configure only a part of a given CC for actual transmission/reception, which is referred to as a BWP¹⁷. Importantly, control and data signalling only occur within this part of the spectrum, enabling thus a reduced power consumption at both the BS and UE sides, as they need to handle/monitor smaller bandwidths. BWP can be (de)activated by a timer, downlink control information (DCI), and/or RRC signalling, which can enable faster bandwidth adaptation with respect to the original CA framework. In this context, the work in [78] stands out, which investigated the optimal transmit power and CA configuration to optimise the WSEE metric (see eq. (12) in Section III), while satisfying the downlink and the uplink UEs requirements. Important energy savings were reported.

It should be noted, however, that these intra-BS energy saving mechanisms may have unwanted side effects on the overall network performance, as hinted earlier, as they only take a BS-centric viewpoint. For example, while reducing its energy consumption, a BS may significantly impact the network layout (i.e., coverage) as well as the load and interference distributions across the

¹⁷For each CC, at most four BWPs can be defined for the downlink and the uplink communications.

network [149]. The turning off of such BS may leave some UE without coverage or proper service, and increase the traffic loads of multiple neighbouring BSs. If these BSs operate on the same spectrum band, BS (de)activation will also impact the inter-cell interference pattern, which will also have important consequences on the rank and the modulation and coding scheme (MCS) selection, and thus on packet success rate. To provide a proper network-wide energy consumption optimisation, while ensuring end-users' QoS, it is thus necessary to take these network effects into account, when adjusting a BS configuration.

Table X presents the main contribution on intra-BS energy saving mechanism discussed in this section.

Table X
SUMMARY OF MAIN INTRA-BS ENERGY SAVING MECHANISMS.

Paper	Focus	Main contribution/idea
[42]	CA dormant state	Enable fast bandwidth adaptation in CA framework
[148]	BWP concept	When the cell load is low, only a part of a given CC is used
[78]	CA optimization	Study optimal transmit power and CA configuration
[149]	BS ON-OFF switching	Discuss challenges in Intra-BS energy saving

B. Inter-BS energy saving mechanisms

Inter-BS mechanisms are suited to address the above mentioned challenge, operating over multiple neighbouring BSs, in a centralised or a distributed manner, to jointly optimise the use of radio resources and provide energy savings without significantly affecting the end-users' QoS. This resource management problem, however, is typically challenging, as it involves multiple network functionalities and complex mathematical formulations.

Taking a global perspective, without too many practical constraints, the authors of [150] investigated how to optimise the network energy efficiency by jointly managing long-term BS activation, UE association, and power control in a HetNet with mMIMO capabilities. They modelled this problem using mixed-integer programming, and to address complexity, proposed a sub-optimal centralised scheme, where *i*) the integer variables (i.e., the cell (de)activation and the UE association variables) are relaxed, and *ii*) the BS activation, UE association, and power

control problems are iteratively solved. Since their centralised solution was still characterized by a large complexity, they also proposed a distributed solution based on game theory, which provided lower performance in terms of energy savings and UE performance, but it is proven to converge to the Nash equilibrium.

From a more practical view point, the 3GPP NR has defined in Release 15 and beyond functionalities to support the (de)activation of 3GPP NR capacity booster cells, which are overlaid on a larger 3GPP NR coverage cell [64], [151]. This solution builds upon the possibility for the capacity booster cell to autonomously decide to shutdown. This switch off decision is typically based on cell load information, and can also be taken by the operation, administration and maintenance (OAM), if needed be. The capacity booster cell may initiate handover actions in order to off-load the cell being switched off, and may indicate the reason for handover with an appropriate cause value to support the target node in taking subsequent actions. Dual connectivity may also be used to support the offloading procedure. During the shutdown time period of capacity booster cell, idle mode UEs are prevented from camping on this cell or handover to it. All neighbouring cells are informed by the coverage cell owning such capacity booster cell about its switch-off actions over the X_n interface. Moreover, the coverage cell owning such capacity booster cell can request an inter-BS cell re-activation over the X_n interface to wake up the capacity booster cell, if the former cell overloads. The switch on decision may also be taken by the OAM. All neighbouring cells are informed by the coverage cell owning such capacity booster cell about the re-activation by an indication on the X_n interface.

In addition to their impact on network performance, i.e. coverage, rate, latency, energy inter-BS saving mechanisms can have a detrimental effects on the BS hardware life-time, e.g. due to ping-pong effects. Specifically, deep and frequent transients between different status lead to large temperature gradients in the involved hardware components, which increase their failure rates, thus augmenting the maintenance costs to fix or replace the BS. The impact due to these long-term energy saving mechanisms can be measured by the acceleration factor, which is defined as the ratio between the failure rate observed by applying such mechanisms over time and the one experience when keeping the BS always active. In this line, the work in [152] has modelled hardware failure rate due to cell (de)activation, and proposed a heuristic to control the statuses of the BSs of a network, which minimises the acceleration factor growth over time, while satisfying the end-users' QoS demands. In contrast to baseline solutions, which maximize the energy savings at the cost of increasing the BS failure rate over time, the proposed approach

achieves around a 30% of power savings in a 3GPP LTE scenario, while keeping the acceleration factor close to one, which can be translated in further energy savings as BSs are maintained or replaced less often.

To further avoid excessively frequent status changes of the BSs and associated network performance losses, the BS control policy in charge of (de)activation decisions could consider the load distribution and the manner in which it varies in time and space. To this end, load can either be characterized statistically or using data-driven approaches. Embracing this concept, the authors of [153] have considered a dense HetNet, where small cells and UEs are randomly deployed, following a HPPP distribution, and packets arrive to the transmission buffers according to an exponential distribution. Using this model, they have characterized the probability density function of the cell load using a gamma distribution, and used this information to elaborate multiple BS (de)activation strategies, comparing them in terms of complexity, blocking rate probability, throughput, and energy efficiency.

Following the same line of thinking, the authors of [154] have proposed a stochastic game, where distinct BS instances in a CRAN platform take advantage of spatial correlation, and jointly estimate the network traffic, exploiting past observations. The CRAN then uses these estimates to decide the status of the remote radio heads (RRHs) in the network, distribute the UEs among the cells, control the RRHs transmit power, and setup cooperative transmission schemes to guarantee that coverage requirements are satisfied. The authors also demonstrated, through a CRAN experimental platform, that the proposed solution using traffic estimations leads to large energy savings with respect to a dynamic BS switching solution [63], which is not aware of the traffic evolution.

As a different approach to inter-BS energy saving, cell zooming has also become a popular mechanism, which is not based on carrier shutdown, but consists in reducing the cell coverage of lightly loaded cells, while simultaneously increasing the area covered by neighbouring compensating ones [155] (see Fig. 21). When using this mechanism, topology changes should be smoothly implemented to limit service outage [156]. To face this challenge, the work in [157] has proposed a data-driven approach to optimise the cell zooming mechanism. Specifically, this framework has represented the network through a graph, and used a BS connectivity metric related to user-level data to construct the graph adjacency matrix. The authors have run a Markov process on the graph to identify network communities on which implementing the cell zooming. This process is realized through two steps. The first step consists on using a polynomial model to

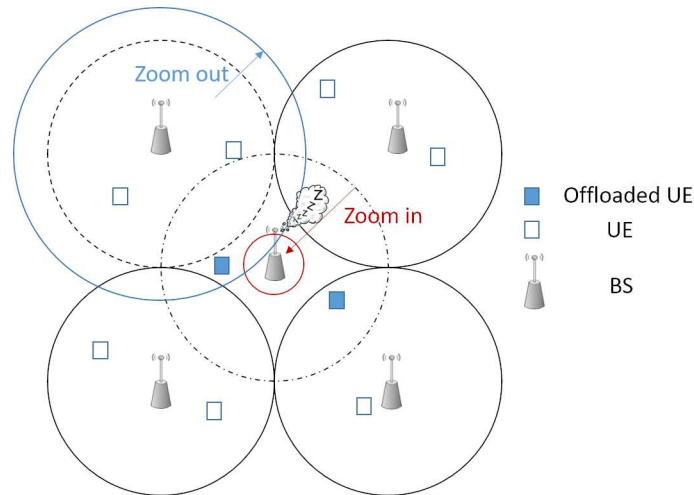


Figure 21. Cell zooming procedure [157].

predict the expected community traffic load in the next hour, and the second step operates the cell zooming to identify the BSs to deactivate and how to distribute the load among the active BSs. With respect to baseline solutions based on the knowledge of instantaneous traffic [63], [155], the authors have highlighted that the proposed solution leads to 20% energy saving gains at the cost of increasing the blocking rate only by 0.1%. They have also shown that the achieved blocking rate is greatly affected by the large prediction noise, which results in inaccurate forecasts.

Table XI summarizes the main contributions of the inter-BS energy saving schemes discussed in this section.

Table XI
SUMMARY OF MAIN INTER-BS ENERGY SAVING MECHANISMS.

Paper	Focus	Main contribution/idea
[150]	Network coordination in HetNets	Joint control of BS activation, UE association, and power control
[64]	Network coordination in 5G NR	Defines message to support the (de)activation of capacity booster cells
[152]	Impact of BS switching on HW failure	Propose a heuristic that limits the HW failure rate with BS (de)activation
[155][156][157]	Cell zooming	Control network layout to prevent coverage losses due to energy saving schemes

C. Inter-RAT energy saving mechanisms

It should also be noted that, in the early stage of 5G deployments, 3GPP NR BSs are not uniformly distributed across a city area, and thus there is a need for a tight inter-working between the 3GPP NR network and the underlying 3GPP LTE one. As previously mentioned, 3GPP NR BSs are characterized by a larger power consumption with respect to 3GPP LTE ones, as they integrate more complex hardware to operate on a wider bandwidth and use a larger number of antennas and transceiver modules (due to mMIMO). Therefore, the wireless community is also currently developing inter-RAT energy saving solutions to switch off 3GPP NR booster capacity cells, when their traffic demand is low [158].

From a general and conceptual perspective, many of the inter-BS energy saving optimization strategies presented earlier apply to this inter-RAT, provided that such frameworks can capture the different characteristics of 3GPP NR and LTE BSs. Thus, we do not provide further details in this space.

From a more practical view point, however, it is important to highlight that the 3GPP NR Release 16 has done specification work in this field. Similarly as described in the inter-BS energy saving case, the 3GPP NR capacity booster cell can autonomously decide to switch off based on its own load in the inter-RAT energy saving one. The new enhancements, however, allow now to declare a 3GPP LTE cell as the coverage cell of the 3GPP NR capacity booster cell, and thus such 3GPP LTE coverage cell can request an inter-system cell re-activation over the S1/NG interface based on its own cell load information or neighbour cell load information. Further details can be found in [64] [159].

VIII. ANTENNA-DOMAIN (CHANNEL SHUTDOWN) ENERGY SAVING-BASED SOLUTIONS

While the mMIMO frameworks presented in Section V provided general important insights on the deployment and optimization of energy-efficient mMIMO networks in both the uplink and the downlink, it should be noted that, once the network is deployed and running, different approaches can be used to minimise energy consumption. For instance, when the traffic load is low, the energy consumption of a mMIMO system may be reduced by using only a subset of the available BS antennas and/or transceiver modules¹⁸, according to traffic requirements and

¹⁸These two terms, antennas and transceivers, although they represent different concepts, abusing notation, they are interchangeably used in this section.

avoiding resource waste. This type of approaches are referred to as antenna selection or channel shutdown [160], [161], and are the focus of this section.

A number of frameworks have investigated channel shutdown subject to end-users' QoS demands on short time scales on the basis of multi-path fast fading variations, i.e., activating those BS antennas with favorable channel conditions at each —or a small number of— TTIs [162]–[167]. However, in wideband systems such as 5G with many subcarriers per carrier, it is unlikely that a BS antenna is simultaneously not selected on all such subcarriers. Moreover, antenna selection based on multi-path fast fading also requires all BS antennas to be activated at least for channel estimation, thus limiting their sleeping time.

Taking a more practical approach, the authors in [168] investigated from a generic view point the antenna selection problem in the downlink considering larger time scales, basing the decision-making in large-scale fading parameters. This work is targeted at finding both the optimal number of BS antennas and their transmit powers to minimise the downlink power consumption of a mMIMO network, while meeting end-users' QoS demands in the form of a minimum SINR per UE. Both single-cell and multi-cell scenarios were analysed, where the system model in the latter accounted for multiple BSs, a number of antennas and simultaneously multiplexed UE per BS, and imperfect channel estimation through pilot contamination. It should be noted, however, that the authors adopted a basic BS power consumption model, which depends on the PA efficiency, and is only linear with the number of BS antennas. Signal processing power consumption, for example, as a function of the number of simultaneously multiplexed UEs in each TTI is not considered. For the single-cell case, the authors derived the optimal number of BS antennas and their transmit powers in closed-form. Importantly, these expressions prove that, only when the circuit power consumption per BS antenna is small, the minimum BS power consumption can be attained by activating all the BS antennas. Otherwise, the BS can save energy by deactivating a subset of them. For the multi-cell case, and contrary to the single-cell one, since pilot contamination was considered, a coherent interference term appeared in the SINR formulation, which scales with the number of BS antennas in the pilot-sharing cells, thus limiting the achievable SINRs of UEs. As a consequence, the results indicated that increasing the number of BS antennas still leads to lower transmit powers, but this is not necessarily the optimal to minimise the BS power consumption, as there is a cost associated with using such BS antennas. Unfortunately, the authors concluded that it is hard to obtain closed-form expressions for the optimal number of BS antennas and their transmit powers for this more complex multi-cell

case. However, they showed that the joint optimisation problem can be relaxed as a geometric programming problem that can be solved efficiently, and suggested that their algorithm can be used to optimally (de)activate BS antennas to deliver the requested traffic with minimal power consumption. However, a detailed system-level analysis is missing in this paper.

The work in [83] complements the above study investigating how to solve the antenna selection problem, while considering daily load variations in multi-cell mMIMO networks. To do so, the authors computed the distribution of the active UEs in a cell for different network loads using queuing theory. They then modelled the distributed daily energy efficiency maximization problem, using the active number of BS antennas in each cell as a variable, and solved it through a game theory framework. As a result, a best response algorithm was proposed, where each BS iteratively selects the strategy that produces its most favorable outcome given other BS strategies. This *selfish* approach did not achieved the global optimum, but lead to a Nash equilibrium without the need for coordination across cells, thus being adequate for an intra-BS operation without any network coordination. Noting that this work did not considered specific UE rate requirements in the optimization, results showed that, at very low loads, the proposed adaptive antenna selection scheme could achieve significant energy efficiency gains as high as 250%, when compared with a baseline system that does not adapt the number of antennas to the traffic profile, at the expense of around a 50% reduction in the average UE data rate.

It should be considered, however, that even if most practical antenna selection schemes deployed in current BSs aim at turning off BS antennas for a long period of time to achieve deeper sleeps and more power savings. This may lead in some cases to a large activation delays, and thus to a degraded UE experience. For example, if parts of the transceiver chains are turned off semi-statically, the rate of cell-edge UEs could be significantly reduced, and latency may become unacceptable for some URLLC services.

Considering both rate and latency requirements, the work in [169] has also recently investigated the power consumption minimisation problem in multi-cell mMIMO networks at shorter time scales, by implementing cell DTX (see Section VI) in conjunction with precoding and antenna selection. Note that differently than the references presented earlier in this section, this work does not base its decision-making on the basis of multi-path fast fading variations. The authors, instead, have considered a multiple frame optimisation window, and proposed a strategy to select the right precoding technique for each transmission frame such that the total transmission time and latency are minimised. Moreover, this work have proposed a technique to trade the UE

latency for additional energy savings, by reducing the number of active BS antennas used in each frame. Numerical results have highlighted that the proposed adaptive antenna selection scheme provides large energy efficiency gains in lightly loaded scenarios without impacting the end-users' QoS. In more detail, for rate and latency requirements between 1 and 6 Mbps per UE and 0.001 and 5 duty cycles per frame, respectively, the authors claim that the proposed technique can provide energy efficiency improvements between 125 % and 1124 % in the suburban scenario, and between 196 % and 952 % in the rural scenario, without compromising QoS.

To support this type of schemes, and avoid the mentioned performance losses due to transients, the latest proposals in 3GPP NR Release 18 suggest the development of new specification enhancements to support a fine-grained (de)activation of BS antennas, e.g., in the unit of an OFDM symbol or a slot. In practice, considering the buffer sizes of active UEs, their data rate requirements and the expected transmission capacity, the number of active BS antennas can be adjusted to match the end-users' QoS demands and maximize the potential power savings without incurring large performance loss [146].

Table XII summarizes the main contribution presented in this section on antenna-domain energy saving schemes.

Table XII
SUMMARY OF MAIN ANTENNA-DOMAIN ENERGY SAVING MECHANISMS.

Paper	Focus	Main contribution/idea
[168]	Antenna selection problem	Optimize the BS active antennas and transmit powers, based on slow fading
[83]	Antenna selection in multi-cell mMIMO	Use game theory to adjust the active antenna number to the network daily load variations
[169]	Cell DTX+precoding and antenna selection in multi-cell mMIMO	Propose a multi-stage optimization scheme that trades-off user latency and energy efficiency

IX. MACHINE LEARNING AND DATA-DRIVEN ENERGY EFFICIENCY OPTIMIZATION

As previously discussed in Section II-D, energy efficiency optimisation highly depends on the accuracy of the embraced models, and unfortunately, many of the current models are rigid, mostly the theoretical ones, unable to adapt to specific channel characteristics, enabling technologies, or environment changes. This may yield a considerable theory to practice gap. Instead, data-driven optimisation may be able to close this gap, learning the practical state of the network and

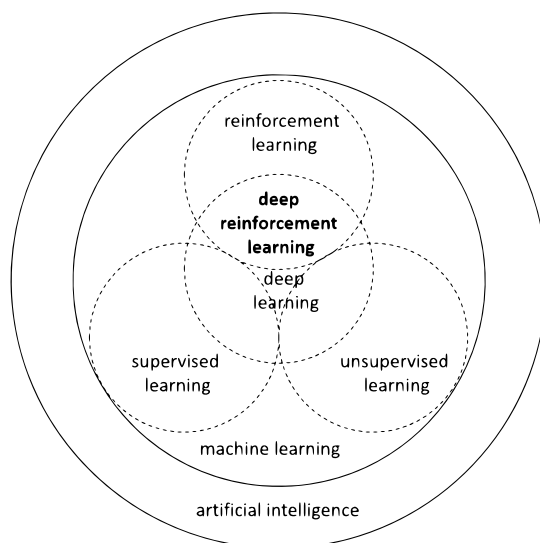


Figure 22. Relationships among deep reinforcement learning, deep learning, reinforcement learning, supervised learning, unsupervised learning, machine learning, and AI.

inferring optimum network operation policies by means of artificial intelligence (AI) to enhance the energy efficiency of mobile networks [26]. In this line, recently, several works in the literature are aiming at improving energy efficiency by exploiting state-of-the-art ML algorithms.

In Fig. 22, we illustrate the relationships between the main concepts developed in the framework of AI and ML [170]. Specifically, today, within the AI world, ML comprises the family of algorithms which uses data to develop intelligent systems. We can identify three main models: supervised learning, unsupervised learning, and RL.

Supervised and unsupervised learning models have been investigated to characterise and forecast network traffic and predict 5G network behaviour, leveraging rich data-sets. These models are, for example, becoming increasingly popular to define new solutions for PHY layer functions [171]. In contrast to supervised and unsupervised learning, the essence of RL concerns learning to make online decisions through interactions with the environment to control network operations. Therefore, ML models are key to enable an intelligent mobile network able to characterize its environment, predict system changes in time and space, and react accordingly in a real time manner.

In the following, we review the literature related to the use of ML techniques for traffic prediction and network optimization in green 5G networks.

A. ML for Traffic Prediction

One of the fundamental challenges along the path to enable full network adaptation to end-users' QoS requirements is the accurate forecasting of the network traffic. Such data forecasting can help driving energy efficient network decisions, e.g., carrier shutdown and others, as it will be shown in the next section.

The forecasting of the network traffic presents, however, important challenges:

- End-users have different QoS demands at different moments of the day and in different places. Therefore, traffic demands change in time and space, making the prediction task difficult.
- The mobility of UEs introduces spatial dependencies among neighboring cells. Moreover, spatial dependencies can occur between distant cell towers, as efficient urban transportation systems easily enable UEs to travel across cities within half an hour.
- The spatial distribution of UEs at the urban scale is further influenced by many factors, including land use, population, holidays, and various social activities. These further complicate the spatio-temporal dependencies among traffic in distinct cell towers.
- The prediction time scale should match the decision periodicity of the energy saving mechanism. For instance, when adjusting every hour the number of active carriers, the forecasting model should provide predictions of the cell load every hour. In contrast, if the mechanism works on a daily basis, a longer prediction, i.e., 24 hours, is required, making the task more challenging.

The studies in the literature aiming at predicting the network traffic can be differentiated into two groups, according to the adopted methods, i.e., statistical-based and ML-based approaches.

1) *Statistical-based methods*: Statistical-based methods rely on capturing the statistics of the network traffic. One of the most popular statistical approaches when predicting network traffic is autoregressive integrated moving average (ARIMA) [172], which originates from three models: the auto-regressive model, the moving average model, and their combination (ARMA). The predictions performed by this model are based on considering the lagged values of a given time-series, while accommodating for non-stationarity. The main limitation of ARIMA is its inability to capture the seasonality—a time series with a repeating cycle—of network traffic. To overcome such limitation, an extension of this algorithm, named seasonal autoregressive integrated moving average (SARIMA), has been proposed [173]. SARIMA adds three new hyper-parameters to

specify the auto-regression, the moving average, and the differencing for the seasonal component of the series, as well as an additional parameter for the period of the seasonality.

Statistical methods like this, however, are not able to capture rapid traffic variations, since they rely on the mean value of the historical data. Moreover, they are mainly linear, and it has become clear that they cannot provide high accuracy when predicting network traffic, especially when considering complex network traffic behaviors observed in real scenarios [174].

2) *ML-based methods*: In contrast to statistical methods, data-driven approaches based on ML have been recently investigated as a solution for network traffic prediction, as they allow to model non-linearities, while taking advantage of the big amounts of data currently being collected by the BSs.

Traditional ML algorithms such as k -nearest neighbours (KNN) [175] and support vector regression (SVR) [176] are able to model non-linear relationships. However, they require well-tuned parameters to achieve accurate prediction results. Moreover, these methods are known to have short memory due to their limited parameter set and inefficient computing, which is detrimental for improving the prediction accuracy.

a) *Recurrent Neural Networks*: Research has then moved into recurrent neural networks (RNNs) to model more complicated nonlinear sequence patterns, which has provided promising results in many fields, such as speech recognition, image caption, and natural language processing. In particular, long short term memory cell (LSTM) has been proposed as a solution to the problem of vanishing gradient in traditional RNNs [177]. This neural network architecture allows to learn long-term dependencies from the time series provided in the input. A LSTM unit is characterised by three gates, i.e., input gate, forget gate, and output gate. These gates control the unit operations by considering three inputs, i.e., the input vector, the memory of the previous time-step, and the output of the previous time-step. The non-linearity is modeled through a sigmoid unit and a hyperbolic tangent unit, which implement the respective functions. As an example, Figure 23 shows a neural network architecture composed of three stacked LSTM units.

Based on such initial definition, the authors in [178] improved the LSTM state-of-the-art with an LSTM architecture using an encoder-decoder model based on gated dilated causal convolution. In the encoder, the long-range memory capacity is enhanced by gated dilated convolutions without increasing the number of model parameters in order to learn a vector representation of the input sequence. Subsequently, different temporal-independent and temporal-dependent features, such

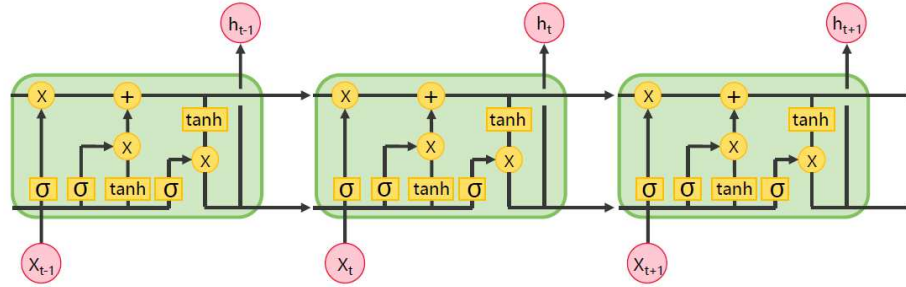


Figure 23. Example of RNN composed of three stacked LSTM units. The terms x_t and h_t are respectively the input and the output at time t . Moreover, the sigmoid and hyperbolic tangent activation functions are represented by the σ and \tanh symbols, respectively.

as the daytime, holidays, weather, are fused with the representation vector. This allows to provide to the model additional relevant information with respect to the network traffic time series. In the decoder, the model applies a RNN with multiple LSTM units to map the fused vector representation back to the variable-length target sequence.

Attention mechanisms are also often used when adopting a LSTM architectures to weight the importance of previous observations [179]. However, experiments have highlighted that simply using lagged inputs (i.e., data points from one year ago, half a year ago, and a quarter before) allows reaching better prediction accuracy than using complex attention mechanisms due to the strong periodicity characterizing the network traffic time series.

It should be noted, however, that the aforementioned LSTM methods do not take into consideration the spatial dependencies between the traffic experienced by different BSs, although it has become apparent that capturing this information may be fundamental to provide accurate forecasting of network traffic [180]. Different extensions of the previous presented approaches have been proposed in this direction.

The authors in [181] have, for example, proposed a novel prediction model to forecast traffic congestion, such that the uplink to downlink resource ratio can be adjusted to improve networking efficiency. The proposed model is composed of a tree-based deep model, followed by an LSTM. This tree-based model uses convolutional layers, which are useful to capture spatial information. Moreover, using a deep model allows to reduce the computational cost, because the convolution operations are performed in parallel in a tree-like structure.

In [182], a hybrid deep learning model for spatio-temporal prediction is proposed instead

to incorporate spatial correlation, in which the temporal dependence is captured by a LSTM, whereas the spatial dependence is encapsulated by auto-encoders. Specifically, an auto-encoder is a neural network architecture used in unsupervised learning to represent a set of data, while reducing its dimensionality [183]. The auto-encoder learns to compress data from the input layer into a short code, i.e., the embedding, and then decompresses that code into a data structure that closely matches the original data i.e., the output layer. With this architecture, the auto-encoders are used to model historical information from the neighboring BSs, and capture spatial dependencies. The proposed model is shown to offer an average Mean Absolute Error (MAE) improvement of 31% and 40% with respect to SVR and ARIMA, respectively.

Even though the aforementioned attempts to capture spatial information, LSTM is not fundamentally adequate for it. In particular, the gates that characterize this model are usually fully connected, and as a consequence, the number of parameters is large, requiring high memory and computation time for training the model. This model is thus highly complex and frequently turns overfitted.

b) Convolutional Neural Networks: An evolution of LSTM, named convolutional LSTM (ConvLSTM) has been proposed to solve this problem, by replacing the inner dense connections with convolution operations [184]. This architecture significantly reduces the number of parameters, and enhances the ability of capturing spatio-temporal information. Indeed, convolutional neural networks (CNNs) are widely adopted now to deal with image classification problems, and capture spatial information. Similarly, when considering the network traffic prediction case, network traffic data is treated as images, where the geographical space is modeled by a matrix, and the traffic distribution in different areas of the city is described by the elements of such matrix.

As a good example, the authors in [185] provide a traffic prediction architecture, in which spatio-temporal dependencies are captured by utilizing densely connected CNNs. In a similar way, the authors in [186] have proposed a prediction algorithm, which can model both temporal and long-distance spatial dependencies. The proposed model follows an encoder-decoder paradigm, where a stack of ConvLSTM and CNN elements are combined. Numerical results show that such model can achieve up to 61% lower normalised root mean square error (NRMSE) as compared to ARIMA and other statistical methods, while requiring up to 600 times shorter ground truth measurement durations.

A main limitation of this type of approaches, however, is that they only work with regular grid-

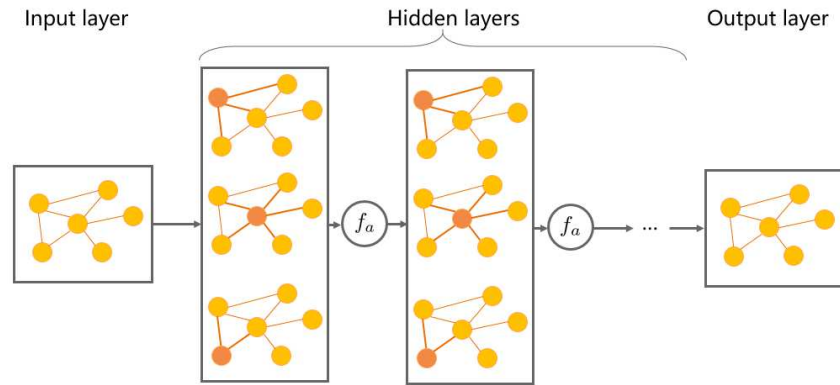


Figure 24. Example of a GNN architecture. The activation function unit is indicated by the f_a symbol.

based region partitions, which are not practical for cellular networks, and limits the prediction performance.

c) Graph Neural Networks: To overcome this limitation, an architecture based on graph neural network (GNN) has recently been proposed to model the network traffic spatio-temporal dependencies using a graph representation. In particular, given a direct graph, each BS is modeled as a vertex, and each edge defines the spatial relations between adjacent BSs.

The authors in [180] have adopted such GNN-based architecture, and decomposed the total data traffic volume into in-tower traffic and inter-tower traffic, which corresponds, respectively, to the traffic serviced to the UEs residing within the coverage of a BS and the traffic serviced to the UEs moving among areas covered by different BSs. In the proposed architecture, each edge has a weight that depends on the total data traffic moving from the corresponding BSs vertices. Importantly, it should be noted that complete directed graphs can contain a huge amount of edges, which hinder the efficient learning of model parameters. Therefore, low weights are treated as noise, and the corresponding edges are pruned by defining a threshold that allows to balance prediction accuracy and computing efficiency. The presented numerical results have shown that this architecture achieves 16% lower MAE compared to a LSTM model. Moreover, the performance analysis have shown that the traffic mobility induced by the roaming of humans plays a large role in the prediction accuracy, and showed that a combination of in-tower and inter-tower traffic patterns can be applied for network or social event forecasting.

d) External inputs: Point of Interests: While the aforementioned research has mainly focused on the network traffic itself, it is also well understood that external factors, such as the

point of interest (POI) distribution, may influence the demand of network traffic. In particular, the analysis provided in [187] reveals that the dynamic urban network traffic usage exhibits five basic time domain patterns, which are correlated to the city functional zones, e.g., residential, transport, office.

In this line, the work in [188] has targeted network traffic prediction by exploiting cross-domain data. Specifically, three types of data sets are considered, i.e., BSs location, POIs distribution, and social activity level. In particular, the latter contains information generated by the end-users when using social networks, such as location and keywords, which may allow to better capture particular social events, such as concerts and football matches. The correlation between these data sets and the network traffic is analyzed and used to improve the prediction accuracy. A novel deep learning based traffic prediction architecture is then proposed. This architecture can fuse the cross-domain data sets into a unified representation. Spatial, temporal, and external factors are then captured and processed by ConvLSTMs. In order to consider the pattern diversity and similarity of the network traffic of different city functional zones, the authors have also proposed an algorithm for grouping city areas into different clusters. Then, inter-cluster transfer learning is proposed to capture regional similarities and differences. The achieved results have shown that cross-domain data sets have high correlation with the network traffic, and thus the introduction of the aforementioned data sets benefits the prediction accuracy. In details, the authors show that, in the considered scenario, cross-domain data sets allowed to improve the MAE by up to 14%

e) Model re-usability: Another important issue related to traffic prediction is that of re-usability. Prediction algorithms generally lack of re-usability, which require them to be re-trained to learn a new representation of the spatio-temporal information, when adopted in a new —or dramatically changing— scenario. The generalization problem of prediction algorithms has been recently discussed in [189], where the authors have proposed a model based on auto-encoders, which learns the embedding of BSs based on raw data. In this framework, the embeddings are vectors, which contain spatio-temporal information of the BSs, and their size is much smaller than the raw data, which allows to improve the generalisation capabilities of the prediction algorithm, while also reducing its computational cost. The architecture is composed of three main modules: an encoder, a spatial adder and a decoder. In more detail, the encoder is designed to extract information from the BS, and infer its embedding. The spatial adder is in charge of building the relation among different BSs, whereas the decoder restores original data from the embeddings.

In this way, the training phase makes the encoder learn how to generate an embedding that conforms to the spatial relation with neighboring BSs. After training the model, the encoder is able to use the raw data from the BS itself to infer its embedding, which contains information about how this BS influence other BSs. Numerical results have shown that this approach helps temporal models to achieve similar performance as spatio-temporal ones, at the cost of a small increase in the training time.

B. ML for 5G Energy Efficiency Optimisation

The wireless network environment is complex and stochastic by nature as already discussed earlier, and in more detail, traffic requirements, user mobility, interference and channel variations in time and space make system-wide optimisation a hard problem. In the past, most of solutions proposed in the literature to configure mobile network parameters have not considered the dynamic nature of wireless networks. More specifically, state-of-the-art algorithms, as many of the once already surveyed, are typically based on perfect —or partial— knowledge of the instantaneous system conditions, which requires to re-compute the solution of a problem whenever a notable change has occurred in the environment. With the complexity carried by such approaches, they may lead to significant computation and signaling overhead. Thus, there is an urgent need for more light-weight, flexible and adaptive solutions with respect to environment dynamics to minimise the energy consumption of practical networks.

In the last decade, RL, and more recently deep reinforcement learning (DRL), have emerged as potential tools to pave the way for artificial intelligence driven optimisation in 5G systems and beyond. For more details on the motivation, refer to [190] and [65] and the references therein.

In RL, the learning process is represented by the interaction between a learner and decision maker, named agent, and the environment. The agent selects actions and the environment responds to those actions by returning new states of the environment. The environment also returns a numerical value named reward that the agent tries to maximize over time. The agent interactions and learning process are represented for clarity in Fig. 25.

This problem is characterised by multiple challenges. For instance, the environment is typically partially observable, i.e., an agent does not have full knowledge of the system state, but only a partial observation. In the context of wireless network optimisation, this is very likely, and it can represent the case where a BS is not aware of the load of other BSs. Moreover, the perception of the reward related to a state-action pair is often delayed, which makes hard to evaluate the

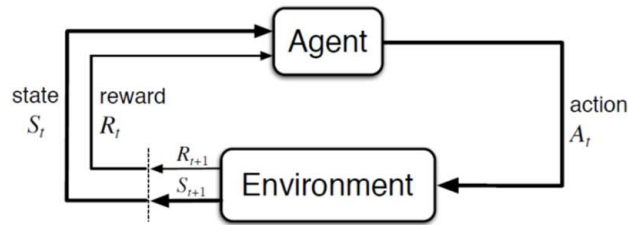


Figure 25. The perception-action-learning loop.

effect of an action, and thus to learn the optimal policy. This effect is known as the temporal credit assignment problem [68]. Similarly, in a multi-agent system, a perceived reward depends on the actions of multiple agents behaving independently from each other, i.e., a cell throughput depends on how each neighbouring BS schedules its own resources.

When the environment can be fully modelled, dynamic programming can be used to solve the learning problem through algorithms whose complexity is polynomial in the size of the set of states [68].

The objective of RL is to learn how to map experienced situation (i.e., the state of the system) into action to take so as to maximise a numerical reward, through continuous interactions with the environment. Through this learning by interaction loop, an exploration-exploitation trade-off arises, i.e., exploiting the information collected so far to benefit the locally optimal decision, or exploring for achieving a better characterisation of the environment and achieve higher long-term gains.

Recently, in the context of energy efficiency, new RL schemes have been proposed to manage sleep modes in 5G BSs. In [141], the authors have mapped the time-domain sleep mode control—see Section VI—as a decision making problem in which a BS sequentially sets the sleeping level length. Specifically, when the cell becomes idle, this approach first puts the BS in the deepest level of sleep, and then gradually switches it on. At each stage, the BS decides the number of slots during which the current sleep mode status will be kept. If a UE request arrives during a sleep period, the associated data is saved in a buffer until the BS wakes up. Accordingly, the state set includes the possible states in which a BS can operate, i.e., idle, active, or one of the sleep modes enabling energy saving. The action set includes the values of the possible number of time slots that can be associated to a given sleep mode. The reward is defined as the weighted sum of the energy saving gain due to the sleep mode and the additional latency

experienced at the UEs due to the buffering of their traffic. In this way, different sleep mode policies can be defined according to whether the operator wants to trade end-users' QoS for energy savings. The authors have used a popular RL scheme, named as Q-Learning, to find the optimal policy. This scheme is a model-free algorithm, which uses state transition experiences to iteratively construct an estimate of the optimal state-action function, also referred as Q-function. A learning parameter allows to control how the estimates are iteratively blended together over time. If each state-action pair is visited infinitely often, and the learning rate is decreased over time, the estimated Q-function converges to the optimum [191]. Note that the optimal state-action pairs are stored into a look-up-table. The results in [141] have shown that, if delay is critical, the sleep mode should not be activated for a cell load larger than 30%. In contrast, for very low loads, up to 55% of energy savings can be achieved, even when prioritizing the end-users' QoS.

A well known problem of RL algorithms is the so-called curse of dimensionality, meaning that their computational requirements grow exponentially with the size of the state and action spaces [68]. To deal with this challenge, function approximation can be used to approximate the state-action value function when the state and/or action spaces are large or continuous. Multiple function approximation methods have been investigated in the literature for RL, e.g., linear functions or artificial neural networks (ANNs).

In [192], the authors proposed a fuzzy Q-learning model to deal with complexity, where a network controller jointly optimises the DTX of the underlying cells and backhaul nodes to minimise the energy consumption and satisfy the end-users' QoS. Specifically, to reduce complexity, the controller maintains a distinct model for each cell. The state space of each cell characterises its buffer state in terms of rate and latency requirements, the BS capacity, and the estimated spectral efficiency loss due to the interference of the nearby BSs, which are expected to be active. Then, in each time slot, the controller observes individually each state space, and decides in a distributed manner which cells (and associated backhaul nodes) to keep in energy saving mode or activate. The authors have associated to each state-action pair a cost function, which models the weighted sum of the BS power consumption and the packet lost, either due to latency constraints or interference. As the state space is composed by continuous variables, this would prevent a classic RL algorithm to converge to an optimal policy in a finite time. Accordingly, the authors have integrated a fuzzy inference system (FIS) to their framework, and reduced the state space by mapping the state representation in fuzzy sets [193]. The authors have shown that the proposed framework is able to coordinate the activation and

deactivation of neighbouring BSs, thus limiting the inter-cell interference. Moreover, this scheme takes advantage of the latency-energy trade-off, and achieves up to 38% of energy savings with respect to a baseline DTX, which does not exploit data buffering.

To manage curse of dimensionality, deep neural networks (DNNs) are currently widely used as a powerful global function approximator, where a neural network is used to compress the Q-table [194]. However, the combination of DNN and RL, i.e., DRL can lead to instability, and even divergence during the training process [195]. To address these issues, recently, the authors of [194] have proposed a deep Q-learning (DQL) framework that leverages two main ideas, i.e., the usage of experience replay, and the introduction of the target network. Experience replay consists in the usage of a buffer, where tuples of experiences (i.e., interactions with the environment) are saved and continuously replayed to break the correlation across subsequent observations during training. Moreover, during training, two distinct deep networks are used. One that is continuously updated, and another one, updated less frequently. These modifications make the algorithm training more stable.

These enhancements have led to continuous research innovations and DRL architectures that can successfully deal with problems that were previously considered intractable. For instance, the authors of [196] have proposed a DQL model to dynamically (de)activate BSs based on traffic requests. This framework has introduced few enhancements with respect to the baseline DQL in [194]. First, they have observed that non-stationary traffic leads to oscillation between waking- and sleeping-dominating regimes. To break this correlated sequence of actions, the authors propose an *action-wise* experience replay, where experiences related to different actions are saved into distinct buffers, which are uniformly sampled during the training process.

In the literature, other mechanisms have been proposed to improve the effectiveness of the experience replay process, e.g., the well known prioritized experience replay [197]. Moreover, and although the reward is clipped to $[-1; 1]$ in classic DQL, to capture the strong variations characterising the wireless environment, this work has also proposed an adapting reward re-scaling scheme, which consists into dividing the instantaneous reward by a positive adapting scaling factor, and summing a saturation penalty to the DQL loss function. In addition, the authors of this work have also used an interrupted Poisson process to model the traffic requests, and generate additional pseudo experiences, which, using the DynaQ framework [68], are periodically stored into the replay memory along with real experiences, and used indiscriminately for training. Empowered by these innovations, their DQL algorithm attempts to learn the optimal policy that

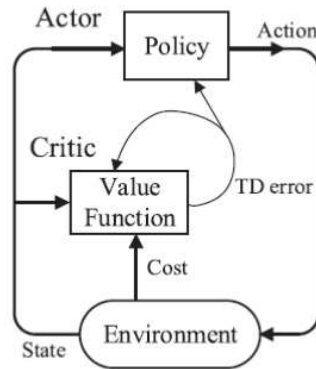


Figure 26. The actor-critic framework.

control the BS status, based on a reward that takes into account the served requests, the queued or re-transmitted requests, and the failed ones. The reward considers the cost to wake up the BS and the one for changing the BS status. Their experiments have shown that modelling the traffic requests and generating pseudo experiences does not lead to large gains. In contrast, action-wise experience replay and adaptive reward scaling improve the stability and adaptability of the proposed framework. Overall, the proposed scheme achieves large gains with respect to a baseline Q-learning approach, in terms of energy saving and end-users' QoS.

Similarly, the authors of [198] have proposed a DRL model to control the small cell (de)activation in dense HetNets. In this framework, the system state comprises the status of each small cell and the estimate of its traffic arrival rate, while the action set includes the (de)activation actions. The cost function provides a qualitative description of the small cell network power consumption, the additional latency experienced due to deactivating BSs, and the switching cost due to the change of status of the small cells (i.e., from off to on and vice-versa). This work has improved the state-of-the-art solutions by considering an actor-critic DRL scheme. In the actor-critic algorithm, the actor selects the action given the state of the environment, and the critic estimates the value function, given the state and the action. Then, it delivers a feedback to the actor (see Fig. 26). Importantly, it should be noted that this type of DRL based on actor-critic has emerged as a powerful solution to deal with continuous action spaces [199]. Conventionally, the actor provides a probability distribution of the possible actions at a given state. In [198], the action space has the size, $2^{N_{SC}}$, where N_{SC} is the number of small cells in the network. Therefore, the output of the actor is defined as a single vector of continuous values. To compensate for the lack of

exploration in the actor's side, this work proposed to add noise to the output action vector. The noisy vector is then converted in a hard decision (i.e., the proto-action), and then, the algorithm explores the set of actions close to the proto-action, and selects the action with the minimum estimated cost. For training the proposed model, the authors have used a deep deterministic policy gradient framework, in which the policy and value function are both approximated by DNNs. The authors have shown for that this approach limits the cumulative network cost over time with respect to baseline RL algorithms, achieving up to 30% of gain with respect to a Q-learning model, and provides larger stability in case of non-stationary traffic. Moreover, they have indicated that the proposed action exploration method reduces the convergence time.

To conclude this section, let us highlight that, although DRL has allowed great progresses in the context of system optimization in a stochastic environment, many challenges are still open, such as enabling distributed optimization in the context of multiple competitive or collaborative agents, or designing fast and low complex methods to update the learned policy after notable change in the system, which have not been observed during the training phase. More challenges faced by RL with respect to the energy efficiency are described in the next section.

Table XIII summarizes the main contributions presented in this section on ML and data driven energy efficient optimization.

X. OPEN RESEARCH DIRECTIONS

In this section, we identify lines of research, which according to the authors' understanding, still require further efforts to aid increasing the energy efficiency of 5G and future RANs.

A. *Multi-cell Energy Efficiency Theoretical Modelling*

As discussed throughout the survey, serving the end-users' QoS demands with the minimum power consumption is key to energy savings, and while the bounds and trade-offs to drive such optimisation in a single-cell case may be well understood, the fundamental understanding of energy efficiency in multi-cell RANs is still limited, due to complexity issues.

Large-scale multi-cell RANs are intricate to model (see Section V-B), and as a result, for tractability reasons, most of the current theoretical understanding on energy efficiency for wide-area networks have been derived based on, for example, the performance of the average UE in uniform networks with simplistic channel, operational and BS power consumption models [134] [135]. These models do not generally capture, however, all relevant features, such

Table XIII
SUMMARY OF MAIN ML AND DATA DRIVEN ENERGY EFFICIENT OPTIMIZATION LITERATURE.

Paper	Focus	Main contribution/idea
[178]	Cellular traffic forecasting	Encoder-decoder model based on gated dilated causal convolution
[181]	Traffic congestion prediction	Tree-based model with convolutional layers to capture the spatial information
[182]	Cellular traffic forecasting	Hybrid model consisting of LSTM unit and auto-encoder to capture temporal and spatial information, respectively
[185]	Cellular traffic forecasting	Densely connected CNNs used to capture the spatio-temporal characteristics of the traffic
[186]	Cellular traffic forecasting	Combination of ConvLSTM and CNN to model both temporal and long-distance spatial dependencies
[178]	Cellular traffic forecasting	Decomposition of the total data traffic volume into in-tower and inter-tower traffic to explicitly account for the spatial dependency
[188]	Cellular traffic forecasting	Cross-domain data used to enhance the accuracy of traffic forecasting
[189]	Cellular traffic forecasting	Model re-usability enhanced by considering a model based on auto-encoders
[141]	BS sleeping policy optimization	Q-learning method adopted to obtain the optimal sleeping time duration for the BSs
[192]	BS sleeping policy optimization	Fuzzy Q-learning method adopted to deal with the complexity of the DTX optimization problem
[196]	BS sleeping policy optimization	Action-wise experience replay proposed as a method for breaking the correlation between actions
[198]	BS sleeping policy optimization	Actor-critic DRL method adopted to control the small cell (de)activation in dense HetNets

as BS and UE distributions, NLoS and LoS transmissions, directional channels, and antenna correlations. Novel theoretical analyses embracing such complexities are thus required to characterise still unknown energy efficiency trade-offs, which may exist, and allow further technology breakthroughs in real deployments.

In this line, the work in [200]–[202] have represented a step forward, accounting for non-uniform BS and UE distributions in RANs, while being able to estimate local performance, i.e., not only the performance of the average UE, but also its distribution. On the same note, the work in [203] has characterised the cell load distribution for a given traffic density, and studied how this affects the network performance. These frameworks, however, are still in their infancy, and have been mostly applied, up to now, to the analysis of simpler single-antenna small cell networks. Further research is needed for their application to more complex RANs and advanced features, such as mMIMO.

With regard to channel models, to give another example, it is widely accepted that most mMIMO performance bounds used today work well when the useful signal coefficients behave almost deterministically, i.e., they have a non-zero mean and a small variance. However, the research in [204] has recently proven that under highly directional channels, for instance, the

channel hardening effect does not so clearly appear, and that new bounds are thus required, in terms of both capacity and energy efficiency, for a more accurate performance evaluation in these more realistic RANs setups.

It is also important to mention that there is a gap in the literature with respect to sophisticated energy efficiency performance analyses, via detailed numerical and/or system-level simulation tools, able to capture the complexity of large-scale multi-cell RANs. Gaining understanding of the interplay among complex features such as mMIMO, CA, and coordinated transmissions, together with their power consumption, which is hard to derive through a pure theoretical analysis, can provide new road maps to fundamental RANs deployment and operation. The intuition gained via these tools can also lead to new theoretical research avenues.

B. Energy Efficiency-driven Network Planning Tools

Deploying and operating a large-scale multi-cell RANs is expensive, and thus requires careful network dimensioning and planning to ensure an optimum radio resource utilisation, e.g., spectrum and bandwidth, number of BSs, their location, architecture and transmit power, number of antennas, and transceivers per BS, to cite a few [205], [206]. Importantly, RANs planning tools must ensure that the deployed system has a sufficient amount of radio resources, and can use it in an effective manner, to achieve the required level of network performance at the appropriate cost. Unfortunately, however, such tools are mostly network capacity-driven as of today, and not yet designed to derive optimal energy efficient deployments.

Once the RANs is planned and deployed, such implementation also imposes hard constraints on the future network performance and its energy consumption. It is thus of imperative importance to equip MNOs with sophisticated RANs planning tools for wide-area network design with energy efficiency at heart. For example, should an MNO deploy less BSs and CCs with larger mMIMO arrays, or in contrast, use more BS and CCs with smaller mMIMO arrays? The applicability of optimisation algorithms in the network planning phase to find such type of practical answers is crucial, and should rely on accurate topological descriptions of the deployment scenario, knowledge of current site deployments and performance, as well as UE and required traffic distributions and accurate BS power consumption models, among others [60].

It is also important to highlight that optimisation performance strictly depends on a reasonable trade-off between modelling accuracy and complexity, and thus MNOs should chose their network models carefully on a per problem basis. This reinforces the need for flexible and efficient

numerical radio propagation and system-level simulation tools as well as tailored optimization theories and algorithms, which have energy efficiency as key driver [207].

More developments in this area are needed at a professional level to make sure future RANs are optimally dimensioned, and energy waste is avoided.

C. Multi-carrier and Heterogeneous Network Analysis

In most scenarios, new 3GPP NR RANs deployments will coexist with existing ones, e.g., 3GPP LTE. In some cases, these deployments may be orthogonal in frequency with, e.g., 3GPP NR in the 3.5 GHz band and 3GPP LTE in the 2 GHz one. In some other cases, due to the scarcity of spectrum, 3GPP NR deployments will have to take place in the same spectrum already used by 3GPP LTE [29], [56]. In the latter case, the fundamental tool to enable such 3GPP NR/LTE spectrum coexistence is the dynamic time scheduling of both 3GPP NR and LTE, for which the 3GPP NR specification provides tools [208].

Given that 3GPP NR and LTE sites have very different characteristics (coverage, bandwidth, antennas, etc.), leading to distinct performance and energy consumption, it would also be desirable to inter-work and (de)activate these two technologies in a coordinated manner, while satisfying end-users' QoS demands with the minimum energy consumption. In some cases, 3GPP LTE may operate at lower carrier frequencies than 3GPP NR, and thus be able to provide a better blanket coverage at a smaller energy consumption. This may be the most energy efficient at low load periods. In contrast, at medium loads, 3GPP NR may be sufficient to provide the required capacity to the active UEs, and 3GPP LTE can be deactivated. If operated at different frequencies, at high loads, both technologies may be aggregated via dual connectivity [209].

In scenarios with the mentioned frequency imbalance, and because the BSs can avail of larger transmit power than the UEs, it may also make sense to simultaneously activate both technologies, and use downlink/uplink split [210], i.e., 3GPP NR for downlink transmissions and 3GPP LTE for the uplink [56]. However, this should only be done where and when necessary, under energy efficient conditions, avoiding potential energy waste due to having both technologies activated.

The optimisation of the inter-working between 3GPP NR and 3GPP LTE is also of high importance when 3GPP NR appears in the form of small cells [211] or millimetre wave [212] access points. These types of cells have a much smaller coverage radius, and can be (de)activated to provide boosted capacity where and when needed. Since some implementations of this type of cell, e.g., millimetre wave, may be power hungry, the usage of coordinated 5G sleep modes across

a large number of this type of cells is critical. This can be facilitated through the separation of the data plane and control plane [213], where the latter is continuously provided by underlying macrocells to ensure robust connectivity and mobility support, while capacity cells (i.e., small cells) allow for enhanced capacity and high rate data transmissions, locally and on demand.

In general, there is a lack of studies covering the understanding and optimisation of these technology inter-working practical use cases from a energy efficiency point of view.

D. Data-Driven Optimisation

State-of-the-art data-driven approaches for network traffic prediction are mainly related to measurements biased by the observation point, i.e., the BSs, which usually does not report the effective traffic demand but rather the cell throughput or resource usage, which depend on the network deployment, interference, current RRM parameters, and running energy saving schemes. The use of this potentially biased measurements may be an issue when adopting prediction algorithms as enablers for improving the mobile network energy efficiency. In fact, the implementation of any energy saving scheme will impact the load distribution across the network, and make related forecasting unreliable. As a result, further research is encouraged with regard to the estimation and prediction of unbiased metrics, whose characterization is not affected by the algorithms implemented using these observations.

Moreover, there is currently a lack of understanding on how prediction errors, e.g. traffic forecasting errors, affect the gains provided by energy saving schemes. In this line, most of the current literature focuses on measuring the performance of the prediction in the metric space, (function of the difference between the predicted value and ground-truth, e.g., physical resource block (PRB) usage, throughput). However, it is generally not straightforward to derive how improvements in the prediction of such metrics help to minimise the RANs power consumption. It is thus recommended that prediction accuracy is investigated also in terms of energy savings when developing data-driven optimisation schemes [214].

As discussed in Section IX, there is also a lack of end-to-end ML frameworks that jointly use supervised learning and RL to characterise and optimise the 5G system. Specifically, supervised learning models can provide multi-step traffic predictions, achieving a comprehensive forecast of future status of the mobile network environment. Using this information can help RL algorithms to converge faster to optimal operational policies, and enhance the performance of the exploration

phase, e.g., optimally deciding the moment to (de)activate BS functionalities without affecting network performance.

Importantly, it should be stressed that recent progresses in the areas of computational processing and data storage, as well as the increased availability of big data, have made the use of AI more practical than ever in many challenging fields. However, the acquisition of large data sets in RANs is currently challenging, and their processing energy demanding, which limits the opportunities for implement data-driven optimization in 5G RANs. To address this issue, the use of joint data-driven and model-based approaches is being widely explored [215], [216], and is becoming the foundation of new optimization mechanisms for large-scale multi-cell RANs. Moreover, these techniques can be used for bootstrapping ML models, thus reducing their need for data, computational complexity, power consumption, and latency.

E. Green AI

Most of the recent breakthroughs led by novel ML solutions have been possible thanks to the ever-increasing computational capacity of dedicated hardware platforms. The work in [217] has highlighted that, in the last decade, while ML models evolved from AlexNet—an image recognition DNN presented in 2012 [218]—to AlphaZero—a RL algorithm proposed in 2018 [219]—, the associated computational cost trend increased by 300,000x. In this line, the work in [220] has also analyzed the energy consumption issues arising from the need of exponentially larger computational resources to continue marginally increasing the model accuracy, and has estimated that the carbon footprint of the current brute force trend is environmentally unfriendly.

The training of DNNs on mobile devices in a distributed, computationally and energy efficient manner is also an ongoing research topic, which can brake down the aforementioned complexity. Collaborative learning schemes, such as federated learning [221], should be considered to mitigate the energy inefficiencies resulting from traditional, centralised ML approaches. Moreover, notable efforts are being performed towards hardware design and software accelerators, which make possible to also move part of the ML process to the UE itself. The more computing is performed on the mobile device, the less data needs to be sent to the cloud and/or network, enabling a reduction of the energy consumption due to the minimization of data exchange (see [222], [223] and references therein). In addition, tailored *early stopping* can also be used to terminate the training process when a near-optimal solution has been found thus reducing the

required number of iterations —and the associated energy consumption— needed to train the ML model [224].

However, to pave the way for the successful integration of ML techniques to drive the modelling and optimization of 5G RANs and beyond, it is key to consider, not only accuracy key performance indicators, but also computational and energy consumption aspects during the design, training and exploitation phases of such data-driven optimisation algorithms. To achieve this goal, while designing an ML model, in addition to accuracy and optimisation-related metrics, computational efficiency and energy consumption must also be considered. Using these metrics, ML architectures that converge faster and/or need to be updated less frequently can be designed to optimise RANs, prioritising energy consumption over model accuracy. One possible approach towards this goal is to investigate yet undiscovered ML paradigms. For instance, as discussed in [215], [216], model-based RL solutions, which exploit a-priori expert knowledge to characterise physically/mathematically the system evolution, can be integrated with model-free RL architectures, which interact with the environment to identify the optimal policy. Similarly, as explained in Section IX-B, RL can leverage information about the future status of the system obtained by supervised learning forecasting to speed up the training speed and converge towards improved operating conditions.

XI. CONCLUSIONS

In this paper, we have provided an overview of the state-of-the-art on the fundamental understanding and practical considerations of the energy efficiency challenge in 5G networks. We have surveyed in detail the available BS power consumption models and metrics for the optimization of the energy efficiency. We have also reviewed the impact on energy efficiency of four 3GPP NR key enabling technologies, i.e., mMIMO, the lean carrier design, ASMs, and ML, presenting the findings in the literature with respect to their available bounds, trade-offs and/or practical achievable energy savings. Importantly, we have highlighted the need for adapting the network resources to meet the end-users' QoS demands, while minimizing network power consumption, and have surveyed the related research differentiating among different algorithm time scales and classes (i.e. micro-sleeps, carrier and channel shutdown). We have also stressed the role that spatio-temporal predictions and online optimisation via ML will play in the previous network power consumption minimization task, and discussed state-of-the-art ML related approaches in

such energy efficiency field. To conclude, we have also provided discussion around the lines of research that need further work to make 5G networks greener.

As a final note, given the enabling effect of that telecommunications systems and the impact that they can have in meeting the requirements for a sustainable development, we encourage the research community to continue making progress toward a sustainable communication system.

LIST OF ACRONYMS

3GPP	third generation partnership project
4G	fourth generation
5G	fifth generation
5GC	5G core network
AC	alternating current
AEE	area energy efficiency
AI	artificial intelligence
ANN	artificial neural network
APC	area power consumption
ARIMA	autoregressive integrated moving average
ASE	area spectral efficiency
ASM	advanced sleep mode
BBU	baseband unit
BCH	broadcast channel
BS	base station
BWP	bandwidth part
CA	carrier aggregation
CC	component carrier
CO_{2e}	carbon dioxide equivalent
CNN	convolutional neural network
CoMP	coordinated multi-point
ConvLSTM	convolutional LSTM
CPU	central processing unit
CRAN	centralized radio access network
C-RAN	cloud radio access network

CRS	cell-specific reference symbol
CSI	channel state information
CSI-RS	channel state information-reference signals
CU	central unit
D2D	device to device
DC	direct current
DCI	downlink control information
DL	downlink
DNN	deep neural network
DQL	deep Q-learning
DRAN	distributed radio access network
DRL	deep reinforcement learning
DRX	discontinuous reception
DTX	discontinuous transmission
DU	distributed unit
eMBB	enhanced mobile broadband
ETSI	European telecommunications standards institute
FDD	frequency division duplexing
FE	front-end
FeMBMS	further evolved multimedia broadcast multicast service
FFT	fast Fourier transform
FIS	fuzzy inference system
GHG	greenhouse gas
gNB	next generation NodeB
GNN	graph neural network
GPP	general purpose processor
HetNet	heterogeneous network
HPPP	homogeneous Poison point process
ICT	information and communication technology
IMT	international mobile telecommunications
ITU	international telecommunication union
KPI	key performance indicator

KNN	<i>k</i> -nearest neighbours
LAA	licensed Assisted Access
LoS	line-of-sight
LSTM	long short term memory cell
LTE	long term evolution
LTE-A	long term evolution advanced
MAC	medium access control
MAE	Mean Absolute Error
MBSFN	multicast-broadcast single-frequency network
MCS	modulation and coding scheme
MDAF	management data analytics function
MIMO	multiple-input multiple-output
ML	machine learning
mMIMO	massive multiple-input multiple-output
mMTC	massive machine type communication
MNO	mobile network operator
MRC	maximal ratio combining
MRT	maximum ratio transmission
NEE	network energy efficiency
NF	network function
NFV	network functions virtualization
NG	next generation
NG-RAN	next generation radio access network
NLoS	non-line-of-sight
NR	new radio
NRMSE	normalised root mean square error
NWDAF	network data analytics function
OAM	operation, administration and maintenance
OFDM	orthogonal frequency division multiplexing
OSI	open systems interconnection
PA	power amplifier
PCell	primary cell

PDCP	packet data convergence protocol
PHY	physical layer
POI	point of interest
PRB	physical resource block
PSS	primary synchronisation channel
QoS	quality of service
RAN	radio access network
RAT	radio access technology
RE	resource efficiency
RF	radio frequency
RL	reinforcement learning
RLC	radio link control
RLF	radio link failure
RNN	recurrent neural network
RRC	radio resource control
RRH	remote radio head
RRM	radio resource management
SARIMA	seasonal autoregressive integrated moving average
SCell	secondary cell
SDO	standard development organization
SI	system information
SIB1	systeminformationblocktype1
SINR	signal to interference plus noise ratio
SISO	single input single output
SNR	signal to noise ratio
SRS	sounding reference signals
SS	synchronization signal
SSB	synchronization Signal/PBCH block
SIB1	system information block 1
SSS	secondary synchronisation channel
SVR	support vector regression
TDD	time division duplexing

TTI	transmission time interval
UE	user equipment
UL	uplink
UMTS	universal mobile telecommunication system
UN	United Nations
URLLC	ultra-reliable low-latency communication
VNF	virtual network function
VRAN	virtualized radio access network
WSEE	weighted sum of the energy efficiencies
WPEE	weighted product of the energy efficiencies
WMEE	weighted minimum of the energy efficiencies
ZF	zero forcing

REFERENCES

- [1] J. Rifkin, *The Third Industrial Revolution: How Lateral Power Is Transforming Energy, the Economy, and the World*, 1st ed. St. Martin's Press, Sep. 2011.
- [2] S. R. Weart, *The Discovery of Global Warming: Revised and Expanded Edition*, 2nd ed. Harvard University Press, Oct. 2008.
- [3] Worldometer, "Current world population," Available at <https://www.worldometers.info/world-population/> (2020/10/27).
- [4] NASA, "Global climate change: Vital signs of the planet," Available at <https://climate.nasa.gov/vital-signs/carbon-dioxide/> (2020/10/27).
- [5] United Nations Environment Programme, "UN Emissions Gap Report 2019," Tech. Rep., Nov. 2019. [Online]. Available: <https://wedocs.unep.org/bitstream/handle/20.500.11822/30797/EGR2019.pdf>
- [6] T. F. Stocker *et al.*, "Climate Change 2013. The Physical Science Basis. Working Group I Contribution to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change - Abstract for decision-makers," Tech. Rep., Oct. 2013. [Online]. Available: <https://www.osti.gov/etdeweb/biblio/22221318>
- [7] NASA Goddard Institute for Space Studies, GISTEMP Team, "GISS Surface Temperature Analysis (GISTEMP), version 4," Available at <data.giss.nasa.gov/gistemp/> (2020/10/27).
- [8] N. Lenssen, G. Schmidt, J. Hansen, M. Menne, A. Persin, R. Ruedy, and D. Zyss, "Improvements in the GISTEMP uncertainty model," *J. Geophys. Res. Atmos.*, vol. 124, no. 12, 2019.
- [9] Secretariat of the United Nations Conference on Trade and Development, "UN Trade and Development report 2019," Tech. Rep., Aug. 2019. [Online]. Available: https://unctad.org/en/PublicationsLibrary/tdr2019_en.pdf
- [10] A. Cazenave, "How fast are the ice sheets melting?" *Science*, vol. 314, no. 5803, pp. 1250–1252, Nov. 2006.
- [11] R. Seager, M. Ting, I. Held, Y. Kushnir, J. Lu, G. Vecchi, H. Huang *et al.*, "Model projections of an imminent transition to a more arid climate in Southwestern North America," *Science*, vol. 316, no. 5828, pp. 1181–1184, Nov. 2007.
- [12] G. B. Bonan, "Forests and climate Change: Forcings, feedbacks, and the climate benefits of forests," *Science*, vol. 320, no. 5882, pp. 1444–1449, Jun. 2008.

- [13] A. Clement, R. Burgman, and J. Norris, “Observational and model evidence for positive low-level cloud feedback,” *Science*, vol. 325, no. 5939, pp. 460–464, Jul. 2009.
- [14] GSMA, “2019 Mobile Industry SDG Impact Report,” Tech. Rep., Sep. 2019. [Online]. Available: <https://www.gsma.com/betterfuture/wp-content/uploads/2019/10/2019-09-24-a60d6541465e86561f37f0f77ebee0f7-1.pdf>
- [15] —, “The Enablement Effect: The impact of mobile communications technologies on carbon emission reductions,” Tech. Rep., Feb. 2020. [Online]. Available: https://www.gsma.com/betterfuture/wp-content/uploads/2019/12/GSMA_Enablement_Effect.pdf
- [16] Huawei Technologies Co., Ltd., “Green 5G: Building a sustainable world,” Tech. Rep., Aug. 2020. [Online]. Available: <https://www.huawei.com/en/public-policy/green-5g-building-a-sustainable-world>
- [17] I.-R. M.2083, “Framework and overall objectives of the future development of IMT for 2020 and beyond,” Tech. Rep., Sep. 2015. [Online]. Available: <https://www.itu.int/rec/R-REC-M.2083>
- [18] ITU-T, “Summary of SMART 2020 Report,” Tech. Rep., Aug. 2008. [Online]. Available: <https://www.itu.int/md/T05-FG.ICT-C-0004/en>
- [19] —, “Sustainable ICT in corporate organizations,” Tech. Rep., 2012. [Online]. Available: https://www.itu.int/dms_pub/itu-t/opb/tut/T-TUT-ICT-2012-10-PDF-E.pdf
- [20] C40, “Carbon emissions are already falling in 30 cities,” Available at <https://www.bloomberg.com/news/articles/2019-10-09/c40-the-cities-where> (2019/10/09).
- [21] SMARTer2030, “ICT Solutions for 21st Century Challenges,” Tech. Rep., 2015. [Online]. Available: http://smarter2030.gesi.org/downloads/Full_report.pdf
- [22] GSMA, “Energy Efficiency,” Tech. Rep., May 2019. [Online]. Available: <https://www.gsma.com/futurenetworks/wiki/energy-efficiency-2/>
- [23] A. Abrol and R. K. Jha, “Power Optimization in 5G Networks: A Step Towards Green Communication,” *IEEE Access*, vol. 4, pp. 1355–1374, 2016.
- [24] M. Usama and M. Erol-Kantarci, “A Survey on Recent Trends and Open Issues in Energy Efficiency of 5G,” *Sensors (Basel)*, vol. 19, no. 14, Jul. 2019.
- [25] J. Lorincz, A. Capone, and J. Wu, “Greener, Energy-Efficient and Sustainable Networks: State-Of-The-Art and New Trends,” *Sensors (Basel)*, vol. 19, no. 22, Nov. 2019.
- [26] X. Cao, L. Liu, Y. Cheng, and X. Shen, “Towards energy-efficient wireless networking in the big data era: A survey,” *IEEE Communications Surveys Tutorials*, vol. 20, no. 1, pp. 303–332, 2018.
- [27] L. Pierucci, “The quality of experience perspective toward 5G technology,” *IEEE Wireless Comm.*, vol. 22, no. 4, pp. 10–16, 2015.
- [28] 3GPP, “3gpp specification series: 38 series.”
- [29] E. Dahlman, S. Parkvall, and J. Skold, *5G NR: The Next Generation Wireless Access Technology*. Academic Press, Aug. 2018.
- [30] J. Wu, S. Guo, H. Huang, W. Liu, and Y. Xiang, “Information and communications technologies for sustainable development goals: State-of-the-art, needs and perspectives,” *IEEE Communications Surveys Tutorials*, vol. 20, no. 3, pp. 2389–2406, 2018.
- [31] R. Mahapatra, Y. Nijssure, G. Kaddoum, N. Ul Hassan, and C. Yuen, “Energy efficiency tradeoff mechanism towards wireless green communication: A survey,” *IEEE Communications Surveys Tutorials*, vol. 18, no. 1, pp. 686–705, 2016.
- [32] Y. Chen, S. Zhang, S. Xu, and G. Y. Li, “Fundamental trade-offs on green wireless networks,” *IEEE Communications Magazine*, vol. 49, no. 6, pp. 30–37, 2011.

- [33] G. Y. Li, Z. Xu, C. Xiong, C. Yang, S. Zhang, Y. Chen, and S. Xu, "Energy-efficient wireless communications: tutorial, survey, and open issues," *IEEE Wireless Communications*, vol. 18, no. 6, pp. 28–35, 2011.
- [34] Z. Hasan, H. Boostanimehr, and V. K. Bhargava, "Green cellular networks: A survey, some research issues and challenges," *IEEE Communications Surveys Tutorials*, vol. 13, no. 4, pp. 524–540, 2011.
- [35] D. Feng, C. Jiang, G. Lim, L. J. Cimini, G. Feng, and G. Y. Li, "A survey of energy-efficient wireless communications," *IEEE Communications Surveys and Tutorials*, vol. 15, no. 1, pp. 167–178, 2013.
- [36] A. De Domenico, E. Calvanese Strinati, and A. Capone, "Green cellular networks: A survey, some research issues and challenges," *Computer Communications*, vol. 37, pp. 5–24, 2014.
- [37] L. Budzisz, F. Ganji, G. Rizzo, M. Ajmone Marsan, M. Meo, Y. Zhang, G. Koutitas, L. Tassiulas, S. Lambert, B. Lannoo, M. Pickavet, A. Conte, I. Haratcherev, and A. Wolisz, "Dynamic resource provisioning for energy efficiency in wireless access networks: A survey and an outlook," *IEEE Communications Surveys Tutorials*, vol. 16, no. 4, pp. 2259–2285, 2014.
- [38] C.-L. I, C. Rowell, S. Han, Z. Xu, G. Li, and Z. Pan, "Toward green and soft: a 5g perspective," *IEEE Communications Magazine*, vol. 52, no. 2, pp. 66–73, 2014.
- [39] Q. Wu, G. Y. Li, W. Chen, D. W. K. Ng, and R. Schober, "An overview of sustainable green 5g networks," *IEEE Wireless Communications*, vol. 24, no. 4, pp. 72–80, 2017.
- [40] S. Zhang, Q. Wu, S. Xu, and G. Y. Li, "Fundamental green tradeoffs: Progresses, challenges, and impacts on 5g networks," *IEEE Communications Surveys and Tutorials*, vol. 19, no. 1, pp. 33–56, 2017.
- [41] A. Mughees, M. Tahir, M. A. Sheikh, and A. Ahad, "Towards energy efficient 5g networks using machine learning: Taxonomy, research challenges, and future research directions," *IEEE Access*, vol. 8, pp. 187 498–187 522, 2020.
- [42] Y. R. Li, M. Chen, J. Xu, L. Tian, and K. Huang, "Power saving techniques for 5g and beyond," *IEEE Access*, vol. 8, pp. 108 675–108 690, 2020.
- [43] Q. Nadeem, A. Kammoun, and M. Alouini, "Elevation Beamforming With Full Dimension MIMO Architectures in 5G Systems: A Tutorial," *IEEE Communications Surveys Tutorials*, vol. 21, no. 4, pp. 3238–3273, 2019.
- [44] F. Rusek *et al.*, "Scaling up MIMO: Opportunities and challenges with very large arrays," *IEEE Signal Process. Mag.*, vol. 30, no. 1, pp. 40–60, Jan. 2013.
- [45] X. Lin, J. Li, R. Baldemair, J. T. Cheng, S. Parkvall, D. C. Larsson, H. Koorapaty, M. Frenne, S. Falahati, A. Grovlen, and K. Werner, "5G New Radio: Unveiling the Essentials of the Next Generation Wireless Access Technology," *IEEE Communications Standards Magazine*, vol. 3, no. 3, pp. 30–37, 2019.
- [46] M. Fuentes, J. L. Carcel, C. Dietrich, L. Yu, E. Garro, V. Pauli, F. I. Lazarakis, O. Grøndalen, O. Bulakci, J. Yu, W. Mohr, and D. Gomez-Barquero, "5g new radio evaluation against imt-2020 key performance indicators," *IEEE Access*, vol. 8, pp. 110 880–110 896, 2020.
- [47] M. Paolini, Senza Fili, "AI and machine learning: Why now? Network optimization in the age of 5G," 2019. Accessed on 18/10/2020. [Online]. Available: <https://www.intel.com/content/dam/www/public/us/en/documents/reports/ai-and-5g-report.pdf>
- [48] D. C. Nguyen *et al.*, "Wireless AI: Enabling an AI-Governed Data Life Cycle," *IEEE Commun. Surveys Tutorials*, 2020.
- [49] E. Pateromichelakis, F. Moggio, C. Mannweiler, P. Arnold, M. Shariat, M. Einhaus, Q. Wei, O. Bulakci, and A. De Domenico, "End-to-End Data Analytics Framework for 5G Architecture," *IEEE Access*, vol. 7, pp. 40 295–40 312, 2019.
- [50] R. Ferrús, O. Sallent, and J. Perez-Romero, "Data Analytics Architectural Framework for Smarter Radio Resource Management in 5G Radio Access Networks," *IEEE Communications Magazine*, vol. 58, no. 5, pp. 98–104, 2020.
- [51] T. L. Marzetta, "Noncooperative Cellular Wireless with Unlimited Numbers of Base Station Antennas," *IEEE Transactions on Wireless Communications*, vol. 9, no. 11, pp. 3590–3600, 2010.

- [52] T. L. Marzetta, E. G. Larsson, H. Yang, and H. Q. Ngo, *Fundamentals of Massive MIMO*. Cambridge University Press, 2016.
- [53] E. Björnson, J. Hoydis, and L. Sanguinetti, “Massive mimo networks: Spectral, energy, and hardware efficiency,” *Foundations and Trends® in Signal Processing*, vol. 11, no. 3-4, pp. 154–655, 2017.
- [54] T. L. Marzetta, “How Much Training is Required for Multiuser MIMO?” in *2006 Fortieth Asilomar Conference on Signals, Systems and Computers*, 2006, pp. 359–363.
- [55] J. Jose, A. Ashikhmin, T. L. Marzetta, and S. Vishwanath, “Pilot Contamination and Precoding in Multi-Cell TDD Systems,” *IEEE Transactions on Wireless Communications*, vol. 10, no. 8, pp. 2640–2651, 2011.
- [56] H. Holma, D. A. Toskala, and T. Nakamura, *5G Technology : 3GPP New Radio*. John Wiley & Sons Ltd., Feb. 2020.
- [57] E. Dahlman, S. Parkvall, and J. Skold, *4G: LTE/LTE-Advanced for Mobile Broadband*. Academic Press, Oct. 2013.
- [58] 3GPP TSG RAN, “TR 37.910, Study on new radio access technology: Radio access architecture and interfaces,” *VI4.0.0*, Mar. 2017.
- [59] O. Blume, D. Zeller, and U. Barth, “Approaches to energy efficient wireless access networks,” in *2010 4th International Symposium on Communications, Control and Signal Processing (ISCCSP)*, 2010, pp. 1–5.
- [60] Ajay R. Mishra, *Fundamentals of Network Planning and Optimisation 2G/3G/4G: Evolution to 5G*, 2nd ed. John Wiley & Sons Ltd., Jul. 2018.
- [61] P. Frenger and R. Tano, “More Capacity and Less Power: How 5G NR Can Reduce Network Energy Consumption,” in *2019 IEEE 89th Vehicular Technology Conference (VTC2019-Spring)*, 2019, pp. 1–5.
- [62] Z. Niu, “TANGO: traffic-aware network planning and green operation,” *IEEE Wireless Communications*, vol. 18, no. 5, pp. 25–29, 2011.
- [63] E. Oh, K. Son, and B. Krishnamachari, “Dynamic base station switching-on/off strategies for green cellular networks,” *IEEE Transactions on Wireless Communications*, vol. 12, no. 5, pp. 2126–2136, 2013.
- [64] 3GPP TSG RAN, “TR 37.816, Study on RAN-centric data collection and utilization for LTE and NR,” *VI6.0.0*, Sep. 2019.
- [65] N. C. Luong, D. T. Hoang, S. Gong, D. Niyato, P. Wang, Y. Liang, and D. I. Kim, “Applications of deep reinforcement learning in communications and networking: A survey,” *IEEE Communications Surveys Tutorials*, vol. 21, no. 4, pp. 3133–3174, 2019.
- [66] S. Hu, Y. Ouyang, Y. Yao, M. H. Fallah, and W. Lu, “A study of lte network performance based on data analytics and statistical modeling,” in *2014 23rd Wireless and Optical Communication Conference (WOCC)*, 2014, pp. 1–6.
- [67] A. Tang, R. Tam, A. Cadrin-Chênevert, W. Guest, J. Chong, J. Barfett, L. Chepelev, R. Cairns, J. R. Mitchell, M. D. Cicero *et al.*, “Canadian association of radiologists white paper on artificial intelligence in radiology,” *Canadian Association of Radiologists Journal*, vol. 69, no. 2, pp. 120–135, 2018.
- [68] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. The MIT Press, 2018. [Online]. Available: <http://incompleteideas.net/book/the-book-2nd.html>
- [69] N. C. Thompson, K. Greenewald, K. Lee, and G. F. Manso, “The computational limits of deep learning,” 2020.
- [70] 3GPP TSG RAN, “TR 38.801, Study on self evaluation towards IMT-2020 submission,” *VI6.1.0*, Sept. 2010.
- [71] 3GPP TSG SA, “TR 32.972, Telecommunication management; Study on system and functional aspects of energy efficiency in 5G networks,” *VI6.1.0*, Sept. 2019.
- [72] G. Auer, V. Giannini, C. Desset, I. Godor, P. Skillermark, M. Olsson, M. A. Imran, D. Sabella, M. J. Gonzalez, O. Blume, and A. Fehske, “How much energy is needed to run a wireless network?” *IEEE Wireless Communications*, vol. 18, no. 5, pp. 40–49, 2011.

- [73] D. López-Pérez, M. Ding, H. Claussen, and A. H. Jafari, "Towards 1 Gbps/UE in Cellular Systems: Understanding Ultra-Dense Small Cell Deployments," *IEEE Communications Surveys & Tutorials*, vol. 17, no. 4, pp. 2078–2101, 2015.
- [74] B. Chen, J. Chen, Y. Gao, and J. Zhang, "Coexistence of lte-laa and wi-fi on 5 ghz with corresponding deployment scenarios: A survey," *IEEE Communications Surveys Tutorials*, vol. 19, no. 1, pp. 7–32, 2017.
- [75] Hua Wang, C. Rosa, and K. I. Pedersen, "Dedicated carrier deployment in heterogeneous networks with inter-site carrier aggregation," in *2013 IEEE Wireless Communications and Networking Conference (WCNC)*, 2013, pp. 756–760.
- [76] C. Rosa, K. Pedersen, H. Wang, P. Michaelsen, S. Barbera, E. Malkamäki, T. Henttonen, and B. Sébire, "Dual connectivity for lte small cell evolution: functionality and performance aspects," *IEEE Communications Magazine*, vol. 54, no. 6, pp. 137–143, 2016.
- [77] K. I. Pedersen, F. Frederiksen, C. Rosa, H. Nguyen, L. G. U. Garcia, and Y. Wang, "Carrier aggregation for LTE-advanced: functionality and performance aspects," *IEEE Communications Magazine*, vol. 49, no. 6, pp. 89–95, 2011.
- [78] G. Yu, Q. Chen, R. Yin, H. Zhang, and G. Y. Li, "Joint Downlink and Uplink Resource Allocation for Energy-Efficient Carrier Aggregation," *IEEE Transactions on Wireless Communications*, vol. 14, no. 6, pp. 3207–3218, 2015.
- [79] S. Tombaz, P. Frenger, F. Athley, E. Semaan, C. Tidestav, and A. Furuskar, "Energy Performance of 5G-NX Wireless Access Utilizing Massive Beamforming and an Ultra-Lean System Design," in *2015 IEEE Global Communications Conference (GLOBECOM)*, 2015, pp. 1–7.
- [80] A. De Domenico, R. Gerzaguét, N. Cassiau, A. Clemente, R. D'Errico, C. Dehos, J. L. Gonzalez, D. Ktenas, L. Manat, V. Savin, and A. Siligaris, "Making 5g millimeter-wave communications a reality [industry perspectives]," *IEEE Wireless Communications*, vol. 24, no. 4, pp. 4–9, 2017.
- [81] E. Björnson, L. Sanguinetti, J. Hoydis, and M. Debbah, "Optimal design of energy-efficient multi-user mimo systems: Is massive mimo the answer?" *IEEE Transactions on Wireless Communications*, vol. 14, no. 6, pp. 3059–3075, 2015.
- [82] B. Debaillie, C. Desset, and F. Louagie, "A flexible and future-proof power model for cellular base stations," in *IEEE 81st Vehicular Technology Conference (VTC Spring)*, 2015, pp. 1–7.
- [83] M. M. A. Hossain, C. Cavdar, E. Björnson, and R. Jäntti, "Energy Saving Game for Massive MIMO: Coping With Daily Load Variation," *IEEE Transactions on Vehicular Technology*, vol. 67, no. 3, pp. 2301–2313, 2018.
- [84] E. Björnson, L. Sanguinetti, J. Hoydis, and M. Debbah, "Designing multi-user mimo for energy efficiency: When is massive mimo the answer?" in *2014 IEEE Wireless Communications and Networking Conference (WCNC)*, 2014, pp. 242–247.
- [85] M. Fiorani, S. Tombaz, J. Martensson, B. Skubic, L. Wosinska, and P. Monti, "Modeling energy performance of C-RAN with optical transport in 5G network scenarios," *IEEE/OSA Journal of Optical Communications and Networking*, vol. 8, no. 11, pp. B21–B34, 2016.
- [86] C. Rowell and Rohde & Schwarz, "Pillars of 5G: Spectral & Energy Efficiency," *Microwave Journal*, pp. 182–184, 2020.
- [87] A. Checko, A. P. Avramova, M. S. Berger, and H. L. Christiansen, "Evaluating c-ran fronthaul functional splits in terms of network level energy and cost savings," *Journal of Communications and Networks*, vol. 18, no. 2, pp. 162–172, 2016.
- [88] M. Shehata, A. Elbanna, F. Musumeci, and M. Tornatore, "Multiplexing gain and processing savings of 5g radio-access-network functional splits," *IEEE Transactions on Green Communications and Networking*, vol. 2, no. 4, pp. 982–991, 2018.
- [89] N. Nikaein, E. Schiller, R. Favraud, R. Knopp, I. Alyafawi, and T. Braun, *Towards a Cloud-Native Radio Access Network*. Cham: Springer International Publishing, 2017, pp. 171–202.
- [90] T. Sigwele, A. S. Alam, P. Pillai, and Y. F. Hu, "Energy-efficient cloud radio access networks by cloud based workload consolidation for 5g," *Journal of Network and Computer Applications*, vol. 78, pp. 1 – 8, 2017.

- [91] A. Younis, T. X. Tran, and D. Pompili, "Bandwidth and Energy-Aware Resource Allocation for Cloud Radio Access Networks," *IEEE Transactions on Wireless Communications*, vol. 17, no. 10, pp. 6487–6500, 2018.
- [92] D. Sabella, A. De Domenico, E. Katranaras, M. A. Imran, M. di Girolamo, U. Salim, M. Lalam, K. Samdanis, and A. Maeder, "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.
- [93] T. Werthmann, H. Grob-Lipski, and M. Proebster, "Multiplexing gains achieved in pools of baseband computation units in 4g cellular networks," in *2013 IEEE 24th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*, 2013, pp. 3328–3333.
- [94] J. K. Chaudhary, A. Kumar, J. Bartelt, and G. Fettweis, "C-ran employing xran functional split: Complexity analysis for 5g nr remote radio unit," in *2019 European Conference on Networks and Communications (EuCNC)*, 2019, pp. 580–585.
- [95] S. Khatibi, K. Shah, and M. Roshdi, "Modelling of computational resources for 5G RAN," in *Proc. EuCNC*, Jun. 2018, pp. 1–5.
- [96] L. Wang and S. Zhou, "On the fronthaul statistical multiplexing gain," *IEEE Communications Letters*, vol. 21, no. 5, pp. 1099–1102, 2017.
- [97] M. Fiorani, S. Tombaz, J. Mårtensson, B. Skubic, L. Wosinska, and P. Monti, "Energy performance of C-RAN with 5G-NX radio networks and optical transport," in *IEEE International Conference on Communications (ICC)*, 2016, pp. 1–6.
- [98] A. Zappone and E. Jorswieck, "Energy efficiency in wireless networks via fractional programming theory," *Foundations and Trends® in Communications and Information Theory*, vol. 11, no. 3-4, pp. 185–396, 2015.
- [99] ETSI TC EE, "ES 203 228, Environmental Engineering (EE); Assessment of mobile network energy efficiency," *V1.3.1*, Oct. 2020.
- [100] 3GPP TSG RAN, "TR 38.913, Study on Scenarios and Requirements for Next Generation Access Technologies," *V16.0.0*, Jul. 2020.
- [101] 3GPP TSG SA, "TR 21.866, Study on Energy Efficiency Aspects of 3GPP Standards," *V15.0.0*, Jun. 2017.
- [102] Z. Yan, M. Peng, and C. Wang, "Economical Energy Efficiency: An Advanced Performance Metric for 5G Systems," *IEEE Wireless Communications*, vol. 24, no. 1, pp. 32–37, 2017.
- [103] 3GPP TSG SA, "TS 28.554, Management and orchestration; 5G end to end Key Performance Indicators (KPI)," *V17.3.0*, June 2021.
- [104] H. Peters, "Axiomatic bargaining game theory," *Mathematical Programming and Operations Research*, W. Leinfellner and G. Eberlein, Eds. Kluwer Academic Publishers, vol. 9, 1992.
- [105] A. Ghosh, A. Maeder, M. Baker, and D. Chandramouli, "5g evolution: A view on 5g cellular technology beyond 3gpp release 15," *IEEE Access*.
- [106] H. Q. Ngo, E. G. Larsson, and T. L. Marzetta, "Energy and spectral efficiency of very large multiuser mimo systems," *IEEE Transactions on Communications*, vol. 61, no. 4, pp. 1436–1449, 2013.
- [107] E. Björnson, J. Hoydis, M. Kountouris, and M. Debbah, "Massive mimo systems with non-ideal hardware: Energy efficiency, estimation, and capacity limits," *IEEE Transactions on Information Theory*, vol. 60, no. 11, pp. 7112–7139, 2014.
- [108] J. Hoydis, S. ten Brink, and M. Debbah, "Massive mimo in the ul/dl of cellular networks: How many antennas do we need?" *IEEE Journal on Selected Areas in Communications*, vol. 31, no. 2, pp. 160–171, 2013.
- [109] J. Tang, D. K. C. So, E. Alsusa, and K. A. Hamdi, "Resource efficiency: A new paradigm on energy efficiency and spectral efficiency tradeoff," *IEEE Transactions on Wireless Communications*, vol. 13, no. 8, pp. 4656–4669, 2014.

- [110] Y. Hao, Q. Ni, H. Li, and S. Hou, "Energy and spectral efficiency tradeoff with user association and power coordination in massive mimo enabled hetnets," *IEEE Communications Letters*, vol. 20, no. 10, pp. 2091–2094, 2016.
- [111] Y. Huang, S. He, J. Wang, and J. Zhu, "Spectral and energy efficiency tradeoff for massive mimo," *IEEE Transactions on Vehicular Technology*, vol. 67, no. 8, pp. 6991–7002, 2018.
- [112] L. You, J. Xiong, A. Zappone, W. Wang, and X. Gao, "Spectral efficiency and energy efficiency tradeoff in massive mimo downlink transmission with statistical csit," *IEEE Transactions on Signal Processing*, vol. 68, pp. 2645–2659, 2020.
- [113] Z. Liu, W. Du, and D. Sun, "Energy and spectral efficiency tradeoff for massive mimo systems with transmit antenna selection," *IEEE Transactions on Vehicular Technology*, vol. 66, no. 5, pp. 4453–4457, 2017.
- [114] J. Choi, D. J. Love, and P. Bidigare, "Downlink training techniques for fdd massive mimo systems: Open-loop and closed-loop training with memory," *IEEE Journal of Selected Topics in Signal Processing*, vol. 8, no. 5, pp. 802–814, 2014.
- [115] X. Gao, B. Jiang, X. Li, A. B. Gershman, and M. R. McKay, "Statistical eigenmode transmission over jointly correlated mimo channels," *IEEE Transactions on Information Theory*, vol. 55, no. 8, pp. 3735–3750, 2009.
- [116] E. Björnson and E. G. Larsson, "How energy-efficient can a wireless communication system become?" in *2018 52nd Asilomar Conference on Signals, Systems, and Computers*. IEEE, 2018, pp. 1252–1256.
- [117] E. V. Belmega and S. Lasaulce, "Energy-efficient precoding for multiple-antenna terminals," *IEEE Transactions on Signal Processing*, vol. 59, no. 1, pp. 329–340, 2011.
- [118] V. S. Varma, S. Lasaulce, M. Debbah, and S. E. Elayoubi, "an energy-efficient framework for the analysis of mimo slow fading channels," *IEEE Transactions on Signal Processing*.
- [119] L. Liu, G. Miao, and J. Zhang, "Energy-efficient scheduling for downlink multi-user mimo," in *2012 IEEE International Conference on Communications (ICC)*, 2012, pp. 4394–4394.
- [120] J. Mao, J. Gao, Y. Liu, and G. Xie, "Simplified semi-orthogonal user selection for mu-mimo systems with zfbf," *IEEE Wireless Communications Letters*, vol. 1, no. 1, pp. 42–45, 2012.
- [121] X. Zhang, S. Zhou, Z. Niu, and X. Lin, "An energy-efficient user scheduling scheme for multiuser mimo systems with rf chain sleeping," in *2013 IEEE Wireless Communications and Networking Conference (WCNC)*, 2013, pp. 169–174.
- [122] H. Li, H. Zhang, D. Li, Y. Liu, and A. Nallanathan, "Joint antenna selection and user scheduling in downlink multi-user mimo systems," in *2018 IEEE 4th International Conference on Computer and Communications (ICCC)*, 2018, pp. 1072–1076.
- [123] H. Lee and A. Paulraj, "Mimo systems based on modulation diversity," *IEEE Transactions on Communications*, vol. 58, no. 12, pp. 3405–3409, 2010.
- [124] A. He, S. Srikanteswara, K. K. Bae, T. R. Newman, J. H. Reed, W. H. Tranter, M. Sajadieh, and M. Verhelst, "Power consumption minimization for mimo systems — a cognitive radio approach," *IEEE Journal on Selected Areas in Communications*, vol. 29, no. 2, pp. 469–479, 2011.
- [125] M. Di Renzo, H. Haas, A. Ghayeb, S. Sugiura, and L. Hanzo, "Spatial modulation for generalized mimo: Challenges, opportunities, and implementation," *Proceedings of the IEEE*, vol. 102, no. 1, pp. 56–103, 2014.
- [126] M. Haenggi, *Stochastic geometry for wireless networks*. Cambridge University Press, 2012.
- [127] J. Andrews, F. Baccelli, and R. Ganti, "A tractable approach to coverage and rate in cellular networks," *IEEE Transactions on Communications*, vol. 59, no. 11, pp. 3122–3134, Nov. 2011.
- [128] H. Dhillon, R. Ganti, F. Baccelli, and J. Andrews, "Modeling and analysis of K-tier downlink heterogeneous cellular networks," *IEEE Journal on Selected Areas in Communications*, vol. 30, no. 3, pp. 550–560, Apr. 2012.
- [129] T. Bai and R. W. Heath, "Coverage and Rate Analysis for Millimeter-Wave Cellular Networks," *IEEE Transactions on Wireless Communications*, vol. 14, no. 2, pp. 1100–1114, Feb. 2015.

- [130] M. Ding, P. Wang, D. López-Pérez, G. Mao, and Z. Lin, “Performance impact of LoS and NLoS transmissions in dense cellular networks,” *IEEE Transactions on Wireless Communications*, vol. 15, no. 3, pp. 2365–2380, Mar. 2016.
- [131] T. Bai and R. W. Heath, “Analyzing uplink sinr and rate in massive mimo systems using stochastic geometry,” *IEEE Transactions on Communications*, vol. 64, no. 11, pp. 4592–4606, 2016.
- [132] Q. Zhang, H. H. Yang, T. Q. S. Quek, and J. Lee, “Heterogeneous cellular networks with los and nlos transmissions—the role of massive mimo and small cells,” *IEEE Transactions on Wireless Communications*, vol. 16, no. 12, pp. 7996–8010, 2017.
- [133] P. Parida and H. S. Dhillon, “Stochastic geometry-based uplink analysis of massive mimo systems with fractional pilot reuse,” *IEEE Transactions on Wireless Communications*, vol. 18, no. 3, pp. 1651–1668, 2019.
- [134] E. Björnson, L. Sanguinetti, and M. Kountouris, “Deploying dense networks for maximal energy efficiency: Small cells meet massive mimo,” *IEEE Journal on Selected Areas in Communications*, vol. 34, no. 4, pp. 832–847, 2016.
- [135] Y. Xin, D. Wang, J. Li, H. Zhu, J. Wang, and X. You, “Area Spectral Efficiency and Area Energy Efficiency of Massive MIMO Cellular Systems,” *IEEE Transactions on Vehicular Technology*, vol. 65, no. 5, pp. 3243–3254, 2016.
- [136] E. Björnson, L. Sanguinetti, and M. Kountouris, “Designing wireless broadband access for energy efficiency: Are small cells the only answer?” in *2015 IEEE International Conference on Communication Workshop (ICCW)*, 2015, pp. 136–141.
- [137] A. Zappone, L. Sanguinetti, G. Bacci, E. Jorswieck, and M. Debbah, “Energy-efficient power control: A look at 5g wireless technologies,” *IEEE Transactions on Signal Processing*, vol. 64, no. 7, pp. 1668–1683, 2016.
- [138] A. Zappone, E. Björnson, L. Sanguinetti, and E. Jorswieck, “Globally optimal energy-efficient power control and receiver design in wireless networks,” *IEEE Transactions on Signal Processing*, vol. 65, no. 11, pp. 2844–2859, 2017.
- [139] P. Frenger, P. Moberg, J. Malmudin, Y. Jading, and I. Godor, “Reducing energy consumption in lte with cell dtx,” in *2011 IEEE 73rd Vehicular Technology Conference (VTC Spring)*, 2011, pp. 1–5.
- [140] W. Yang, K. Lin, and H. Wei, “5g on-demand si acquisition framework and performance evaluation,” *IEEE Access*, vol. 7, pp. 163 245–163 261, 2019.
- [141] F. E. Salem, T. Chahed, Z. Altman, and A. Gati, “Traffic-aware Advanced Sleep Modes management in 5G networks,” in *2019 IEEE Wireless Communications and Networking Conference (WCNC)*, 2019, pp. 1–6.
- [142] F. E. Salem, “Management of advanced sleep modes for energy-efficient 5G networks,” Theses, Institut Polytechnique de Paris, Dec. 2019. [Online]. Available: <https://tel.archives-ouvertes.fr/tel-02500618>
- [143] R. Tano, M. Tran, and P. Frenger, “KPI Impact on 5G NR Deep Sleep State Adaption,” in *IEEE 90th Vehicular Technology Conference (VTC2019-Fall)*, 2019, pp. 1–5.
- [144] S. Tombaz, P. Frenger, M. Olsson, and A. Nilsson, “Energy performance of 5g-nx radio access at country level,” in *2016 IEEE 12th International Conference on Wireless and Mobile Computing, Networking and Communications (WiMob)*, 2016, pp. 1–6.
- [145] Orange, “RP-190326, NR enhancements for Advanced Sleep Modes,” *3GPP RAN #83*, Mar. 2019.
- [146] 3GPP TSG RAN Rel-18 workshop, “RWS-210447, Network energy saving and green operation for NR,” Jun. 2021.
- [147] 3GPP TSG RAN, “TR 36.927, Potential solutions for energy saving for E-UTRAN,” *V16.0.0*, Jul. 2020.
- [148] T. Kim, Y. Kim, Q. Lin, F. Sun, J. Fu, Y. Kim, A. Papasakellariou, H. Ji, and J. Lee, “Evolution of Power Saving Technologies for 5G New Radio,” *IEEE Access*, vol. 8, pp. 198 912–198 924, 2020.
- [149] M. Feng, S. Mao, and T. Jiang, “Base Station ON-OFF Switching in 5G Wireless Networks: Approaches and Challenges,” *IEEE Wireless Communications*, vol. 24, no. 4, pp. 46–54, 2017.
- [150] —, “BOOST: Base Station on-off Switching Strategy for Green Massive MIMO HetNets,” *IEEE Transactions on Wireless Communications*, vol. 16, no. 11, pp. 7319–7332, 2017.
- [151] 3GPP TSG RAN, “TS 38.300, NR and NG-RAN Overall Description,” *V16.6.0*, Sep. 2019.

- [152] L. Chiaraviglio, F. Cuomo, M. Listanti, E. Manzia, and M. Santucci, "Fatigue-Aware Management of Cellular Networks Infrastructure with Sleep Modes," *IEEE Transactions on Mobile Computing*, vol. 16, no. 11, pp. 3028–3041, 2017.
- [153] H. Celebi, Y. Yapici, I. Guvenç, and H. Schulzrinne, "Load-Based On/Off Scheduling for Energy-Efficient Delay-Tolerant 5G Networks," *IEEE Transactions on Green Communications and Networking*, vol. 3, no. 4, pp. 955–970, 2019.
- [154] N. Saxena, A. Roy, and H. Kim, "Traffic-Aware Cloud RAN: A Key for Green 5G Networks," *IEEE Journal on Selected Areas in Communications*, vol. 34, no. 4, pp. 1010–1021, 2016.
- [155] Z. Niu, Y. Wu, J. Gong, and Z. Yang, "Cell zooming for cost-efficient green cellular networks," *IEEE Communications Magazine*, vol. 48, no. 11, pp. 74–79, 2010.
- [156] A. Conte, A. Feki, L. Chiaraviglio, D. Ciullo, M. Meo, and M. A. Marsan, "Cell wilting and blossoming for energy efficiency," *IEEE Wireless Communications*, vol. 18, no. 5, pp. 50–57, 2011.
- [157] H. Jiang, S. Yi, L. Wu, H. Leung, Y. Wang, X. Zhou, Y. Chen, and L. Yang, "Data-Driven Cell Zooming for Large-Scale Mobile Networks," *IEEE Transactions on Network and Service Management*, vol. 15, no. 1, pp. 156–168, 2018.
- [158] 3GPP TSG SA, "TS 28.310, Management and orchestration; Energy efficiency of 5G," *V16.2.0*, Sept. 2020.
- [159] 3GPP TSG RAN, "R3-206885, Support of inter-system inter-RAT energy saving," Nov. 2020.
- [160] K. N. R. S. V. Prasad, E. Hossain, and V. K. Bhargava, "Energy efficiency in massive mimo-based 5g networks: Opportunities and challenges," *IEEE Wireless Communications*, vol. 24, no. 3, pp. 86–94, 2017.
- [161] S. Asaad, A. M. Rabiei, and R. R. Müller, "Massive MIMO With Antenna Selection: Fundamental Limits and Applications," *IEEE Transactions on Wireless Communications*, vol. 17, no. 12, pp. 8502–8516, 2018.
- [162] X. Gao, O. Edfors, J. Liu, and F. Tufvesson, "Antenna selection in measured massive MIMO channels using convex optimization," in *IEEE Globecom Workshops (GC Wkshps)*, 2013, pp. 129–134.
- [163] J. Xu and L. Qiu, "Energy efficiency optimization for mimo broadcast channels," *IEEE Transactions on Wireless Communications*, vol. 12, no. 2, pp. 690–701, 2013.
- [164] N. P. Le, F. Safaei, and L. C. Tran, "Antenna selection strategies for mimo-ofdm wireless systems: An energy efficiency perspective," *IEEE Transactions on Vehicular Technology*, vol. 65, no. 4, pp. 2048–2062, 2016.
- [165] B. Makki, A. Ide, T. Svensson, T. Eriksson, and M. Alouini, "A genetic algorithm-based antenna selection approach for large-but-finite mimo networks," *IEEE Transactions on Vehicular Technology*, vol. 66, no. 7, pp. 6591–6595, 2017.
- [166] B.-J. Lee, S.-L. Ju, N.-I. Kim, and K.-S. Kim, "Enhanced Transmit-Antenna Selection Schemes for Multiuser Massive MIMO Systems," *Wireless Communications and Mobile Computing*, Jun. 2017.
- [167] M. Arash, E. Yazdian, Mohammad, S. Fazel, G. Brante, and M. A. Imran, "Employing antenna selection to improve energy-efficiency in massive MIMO systems," *Wireless Communications and Mobile Computing*, Aug. 2017.
- [168] K. Senel, E. Björnson, and E. G. Larsson, "Joint Transmit and Circuit Power Minimization in Massive MIMO With Downlink SINR Constraints: When to Turn on Massive MIMO?" *IEEE Transactions on Wireless Communications*, vol. 18, no. 3, pp. 1834–1846, 2019.
- [169] W. Pramudito, E. Alsusa, A. Al-Dweik, K. A. Hamdi, D. K. C. So, and T. L. Marzetta, "Load-Aware Energy Efficient Adaptive Large Scale Antenna System," *IEEE Access*, vol. 8, pp. 82 592–82 606, 2020.
- [170] Y. Li, "Deep Reinforcement Learning," *CoRR*, vol. abs/1810.06339, Oct. 2018. [Online]. Available: <http://arxiv.org/abs/1810.06339>
- [171] T. O'Shea and J. Hoydis, "An introduction to deep learning for the physical layer," *IEEE Transactions on Cognitive Communications and Networking*, vol. 3, no. 4, pp. 563–575, 2017.
- [172] D. L. Jagerman, B. Melamed, and W. Willinger, "Stochastic modeling of traffic processes," *Frontiers in queueing*, pp. 271–370, 1997.
- [173] R. J. Hyndman and G. Athanasopoulos, *Forecasting: principles and practice*. OTexts, 2018.

- [174] A. Azari, P. Papapetrou, S. Denic, and G. Peters, "Cellular traffic prediction and classification: A comparative evaluation of lstm and arima," in *Discovery Science*, P. Kralj Novak, T. Šmuc, and S. Džeroski, Eds. Springer International Publishing, 2019, pp. 129–144.
- [175] F. Martínez, M. P. Frías, M. D. Pérez, and A. J. Rivera, "A methodology for applying k-nearest neighbor to time series forecasting," *Artificial Intelligence Review*, vol. 52, no. 3, 2019.
- [176] C.-H. Wu, J.-M. Ho, and D.-T. Lee, "Travel-time prediction with support vector regression," *IEEE transactions on intelligent transportation systems*, vol. 5, no. 4, pp. 276–281, 2004.
- [177] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [178] X. Zhang and J. You, "A gated dilated causal convolution based encoder-decoder for network traffic forecasting," *IEEE Access*, vol. 8, pp. 6087–6097, 2020.
- [179] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, u. Kaiser, and I. Polosukhin, "Attention is all you need," ser. NIPS'17. Red Hook, NY, USA: Curran Associates Inc., 2017, p. 6000–6010.
- [180] X. Wang, Z. Zhou, F. Xiao, K. Xing, Z. Yang, Y. Liu, and C. Peng, "Spatio-temporal analysis and prediction of cellular traffic in metropolis," *IEEE Transactions on Mobile Computing*, vol. 18, no. 9, pp. 2190–2202, 2019.
- [181] M. S. Hossain and G. Muhammad, "A deep-tree-model-based radio resource distribution for 5g networks," *IEEE Wireless Communications*, vol. 27, no. 1, pp. 62–67, 2020.
- [182] J. Wang, J. Tang, Z. Xu, Y. Wang, G. Xue, X. Zhang, and D. Yang, "Spatiotemporal modeling and prediction in cellular networks: A big data enabled deep learning approach," in *IEEE INFOCOM 2017 - IEEE Conference on Computer Communications*, 2017, pp. 1–9.
- [183] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio, *Deep learning*. MIT press Cambridge, 2016, vol. 1, no. 2.
- [184] X. Shi, Z. Chen, H. Wang, D.-Y. Yeung, W.-K. Wong, and W.-c. Woo, "Convolutional lstm network: A machine learning approach for precipitation nowcasting," *Advances in neural information processing systems*, vol. 28, pp. 802–810, 2015.
- [185] C. Zhang, H. Zhang, D. Yuan, and M. Zhang, "Citywide cellular traffic prediction based on densely connected convolutional neural networks," *IEEE Communications Letters*, vol. 22, no. 8, pp. 1656–1659, 2018.
- [186] C. Zhang and P. Patras, "Long-term mobile traffic forecasting using deep spatio-temporal neural networks," in *Proceedings of the Eighteenth ACM International Symposium on Mobile Ad Hoc Networking and Computing*. New York, NY, USA: Association for Computing Machinery, 2018, p. 231–240. [Online]. Available: <https://doi.org/10.1145/3209582.3209606>
- [187] F. Xu, Y. Li, H. Wang, P. Zhang, and D. Jin, "Understanding mobile traffic patterns of large scale cellular towers in urban environment," *IEEE/ACM Transactions on Networking*, vol. 25, no. 2, pp. 1147–1161, 2017.
- [188] C. Zhang, H. Zhang, J. Qiao, D. Yuan, and M. Zhang, "Deep transfer learning for intelligent cellular traffic prediction based on cross-domain big data," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 6, pp. 1389–1401, 2019.
- [189] X. Wang, T. Yang, Y. Cui, Y. Jin, and H. Wang, "Bsenet: A data-driven spatio-temporal representation learning for base station embedding," *IEEE Access*, vol. 8, pp. 51 674–51 683, 2020.
- [190] P. V. Klaine, M. A. Imran, O. Onireti, and R. D. Souza, "A survey of machine learning techniques applied to self-organizing cellular networks," *IEEE Communications Surveys Tutorials*, vol. 19, no. 4, pp. 2392–2431, 2017.
- [191] T. Jaakkola, M. I. Jordan, and S. P. Singh, "On the convergence of stochastic iterative dynamic programming algorithms," *Neural Computation*, vol. 6, no. 6, pp. 1185–1201, 1994.
- [192] A. De Domenico and D. Kténas, "Reinforcement learning for interference-aware cell dtx in heterogeneous networks," in *2018 IEEE Wireless Communications and Networking Conference (WCNC)*, 2018, pp. 1–6.
- [193] L. Zadeh, "The concept of a linguistic variable and its application to approximate reasoning—i," *Information Sciences*, vol. 8, no. 3, pp. 199 – 249, 1975.

- [194] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [195] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, "Deep Reinforcement Learning: A Brief Survey," *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 26–38, 2017.
- [196] J. Liu, B. Krishnamachari, S. Zhou, and Z. Niu, "DeepNap: Data-Driven Base Station Sleeping Operations Through Deep Reinforcement Learning," *IEEE Internet of Things Journal*, vol. 5, no. 6, pp. 4273–4282, 2018.
- [197] T. Schaul, J. Quan, I. Antonoglou, and D. Silver, "Prioritized experience replay," in *International Conference on Learning Representations*, Puerto Rico, 2016.
- [198] J. Ye and Y. J. A. Zhang, "Drag: Deep reinforcement learning based base station activation in heterogeneous networks," *IEEE Transactions on Mobile Computing*, vol. 19, no. 9, pp. 2076–2087, 2020.
- [199] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," 2019.
- [200] M. Di Renzo, S. Wang, and X. Xi, "Inhomogeneous double thinning—modeling and analysis of cellular networks by using inhomogeneous poisson point processes," *IEEE Transactions on Wireless Communications*, vol. 17, no. 8, pp. 5162–5182, 2018.
- [201] Q. Cui, X. Yu, Y. Wang, and M. Haenggi, "The sir meta distribution in poisson cellular networks with base station cooperation," *IEEE Transactions on Communications*, vol. 66, no. 3, pp. 1234–1249, 2018.
- [202] S. Wang and M. Di Renzo, "On the mean interference-to-signal ratio in spatially correlated cellular networks," *IEEE Wireless Communications Letters*, vol. 9, no. 3, pp. 358–362, 2020.
- [203] G. Ghatak, A. De Domenico, and M. Coupechoux, "Coverage analysis and load balancing in hetnets with millimeter wave multi-rat small cells," *IEEE Transactions on Wireless Communications*, vol. 17, no. 5, pp. 3154–3169, 2018.
- [204] G. Caire, "On the ergodic rate lower bounds with applications to massive mimo," *IEEE Transactions on Wireless Communications*, vol. 17, no. 5, pp. 3258–3268, 2018.
- [205] M. J. Nawrocki, M. Dohler, and A. H. Aghvami, *Understanding UMTS Radio Network Modelling, Planning and Automated Optimisation: Theory and Practice*, 1st ed. John Wiley & Sons Ltd., Apr. 2006.
- [206] J. Laiho, A. Wacker, and T. Novosad, *Radio Network Planning and Optimisation for UMTS*, 2nd ed. John Wiley & Sons Ltd., Feb. 2006.
- [207] J. T. J. Penttinen, *5G Network Planning and Optimization*, 2019, pp. 255–269.
- [208] L. Wan, Z. Guo, and X. Chen, "Enabling efficient 5g nr and 4g lte coexistence," *IEEE Wireless Communications*, vol. 26, no. 1, pp. 6–8, 2019.
- [209] O. N. C. Yilmaz, O. Teyeb, and A. Orsino, "Overview of lte-nr dual connectivity," *IEEE Communications Magazine*, vol. 57, no. 6, pp. 138–144, 2019.
- [210] D. López-Pérez, J. Ling, B. H. Kim, V. Subramanian, S. Kanugovi, and M. Ding, "Boosted wifi through lte small cells: The solution for an all-wireless enterprise," in *2016 IEEE 27th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*, 2016, pp. 1–6.
- [211] D. López-Pérez, M. Ding, H. Claussen, and A. H. Jafari, "Towards 1 gbps/ue in cellular systems: Understanding ultra-dense small cell deployments," *IEEE Communications Surveys Tutorials*, vol. 17, no. 4, pp. 2078–2101, 2015.
- [212] M. Giordani, M. Polese, A. Roy, D. Castor, and M. Zorzi, "A tutorial on beam management for 3gpp nr at mmwave frequencies," *IEEE Communications Surveys Tutorials*, vol. 21, no. 1, pp. 173–196, 2019.
- [213] A. Mohamed, O. Onireti, M. A. Imran, A. Imran, and R. Tafazolli, "Control-data separation architecture for cellular radio access networks: A survey and outlook," *IEEE Communications Surveys Tutorials*, vol. 18, no. 1, pp. 446–465, 2016.

- [214] G. Vallero, D. Renga, M. Meo, and M. A. Marsan, “Greener RAN Operation Through Machine Learning,” *IEEE Transactions on Network and Service Management*, vol. 16, no. 3, pp. 896–908, 2019.
- [215] N. Shlezinger, J. Whang, Y. C. Eldar, and A. G. Dimakis, “Model-based deep learning,” 2020.
- [216] A. Zappone, M. Di Renzo, and M. Debbah, “Wireless networks design in the era of deep learning: Model-based, ai-based, or both?” *IEEE Transactions on Communications*, vol. 67, no. 10, pp. 7331–7376, 2019.
- [217] R. Schwartz, J. Dodge, N. A. Smith, and O. Etzioni, “Green ai,” *Commun. ACM*, vol. 63, no. 12, p. 54–63, Nov. 2020. [Online]. Available: <https://doi.org/10.1145/3381831>
- [218] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems*, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, Eds., vol. 25. Curran Associates, Inc., 2012, pp. 1097–1105. [Online]. Available: <https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf>
- [219] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel *et al.*, “A general reinforcement learning algorithm that masters chess, shogi, and go through self-play,” *Science*, vol. 362, no. 6419, pp. 1140–1144, 2018.
- [220] E. Strubell, A. Ganesh, and A. McCallum, “Energy and policy considerations for deep learning in NLP,” in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy: Association for Computational Linguistics, Jul. 2019, pp. 3645–3650. [Online]. Available: <https://www.aclweb.org/anthology/P19-1355>
- [221] P. Kairouz and H. B. McMahan, “Advances and open problems in federated learning,” *Foundations and Trends® in Machine Learning*, vol. 14, no. 1, pp. –, 2021. [Online]. Available: <http://dx.doi.org/10.1561/22000000083>
- [222] Y. Deng, “Deep Learning on Mobile Devices: a Review,” in *Mobile Multimedia/Image Processing, Security, and Applications*, vol. 10993. International Society for Optics and Photonics, 2019, p. 109930A.
- [223] J. Lee, N. Chirkov, E. Ignasheva, Y. Pisarchyk, M. Shieh, F. Riccardi, R. Sarokin, A. Kulik, and M. Grundmann, “On-device Neural Net Inference with Mobile GPUs,” *arXiv preprint arXiv:1907.01989*, 2019.
- [224] L. Prechelt, *Early Stopping — But When?* Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 53–67.