

Non-Orthogonal Multiple Access in Multi-Cell Networks: Theory, Performance, and Practical Challenges

Wonjae Shin, Mojtaba Vaezi, Byungju Lee, David J. Love, Jungwoo Lee, and H. Vincent Poor

Abstract

Non-orthogonal multiple access (NOMA) is a potential enabler for the development of 5G and beyond wireless networks. By allowing multiple users to share the same time and frequency, NOMA can scale up the number of served users, increase the spectral efficiency, and improve user-fairness compared to existing orthogonal multiple access (OMA) techniques. While single-cell NOMA has drawn significant attention recently, much less attention has been given to multi-cell NOMA. This article discusses the opportunities and challenges of NOMA in a multi-cell environment. As the density of base stations and devices increases, inter-cell interference becomes a major obstacle in multi-cell networks. As such, identifying techniques that combine interference management approaches with NOMA is of great significance. After discussing the theory behind NOMA, this paper provides an overview of the current literature and discusses key implementation and research challenges, with an emphasis on multi-cell NOMA.

I. WHAT DRIVES NOMA?

The next generation of wireless networks will require a paradigm shift in order to support massive numbers of devices with diverse data rate and latency requirements. Particularly, the increasing demand for Internet of Things (IoT) devices poses challenging requirements on 5G wireless systems. Two key features of 5G are expected to be a latency of 1ms, compared to 10 ms in the 3rd Generation Partnership Project (3GPP) Long-Term Evolution (LTE), and support for 10 Gbps throughput.

To fulfill these requirements, numerous potential technologies have been introduced over the last few years. Among them is non-orthogonal multiple access (NOMA) [1], a technique to serve multiple users via a single wireless resource. NOMA can be realized in the power, code, or other domains [2], [3]. Code domain NOMA uses user-specific spreading sequences for sharing the entire resource, whereas power domain NOMA exploits the channel gain differences between the users for multiplexing via power allocation. Power domain NOMA can improve wireless communication in the following benefits:

W. Shin and J. Lee are with the Department of Electrical and Computer Engineering, Seoul National University, Seoul, Korea (e-mail: {wonjae.shin, junglee}@snu.ac.kr). M. Vaezi and H. V. Poor are with the Department of Electrical Engineering, Princeton University, Princeton, NJ, USA (e-mail: {mvaezi, poor}@princeton.edu). B. Lee (corresponding author) and D. J. Love are with School of Electrical and Computer Engineering, Purdue University, West Lafayette, IN, USA (e-mail: {byungjulee, djlove}@purdue.edu). (*Wonjae Shin and Mojtaba Vaezi contributed equally to this work.*)

- *Massive Connectivity*: There appears to be a reasonable consensus that NOMA is essential for massive connectivity. This is because the number of served users in all orthogonal multiple access (OMA) techniques is inherently limited by the number of resource blocks. In contrast, NOMA theoretically can serve an *arbitrary* number of users in each resource block by superimposing all users' signals. In this sense, NOMA can be tailored to typical IoT applications where a large number of devices sporadically try to transmit small packets.
- *Low Latency*: Latency requirements for 5G application are rather diverse. Unfortunately, OMA cannot guarantee such broad delay requirements because no matter how many bits a device wants to transmit the device must wait until an unoccupied resource block is available. On the contrary, NOMA supports flexible scheduling since it can accommodate a *variable* number of devices depending on the application that is being used and the perceived quality of service (QoS) of the device.
- *High Spectral Efficiency*: NOMA also surpasses OMA in terms of spectral efficiency and user-fairness. As will be seen in Section II, NOMA is the theoretically *optimal* way of using spectrum for both uplink and downlink communications. This is because every NOMA user can enjoy the whole bandwidth, whereas OMA users are limited to a smaller fraction of spectrum which is inversely proportional to the number of users. In addition, NOMA can also be combined with other emerging technologies, such as massive multiple-input multiple-output (MIMO) and mmWave technologies, to further support higher throughput.

In view of the above benefits, NOMA has drawn much attention from both academia and industry. However, much of the work in this context is limited to single-cell analysis, where there is no co-channel interference caused by an adjacent base station (BS). To verify the benefits of NOMA in a more realistic setting, it is necessary to consider a multi-cell network. Specifically, as wireless networks get denser and denser, inter-cell interference (ICI) becomes a major obstacle to achieving the benefits of NOMA. In this regard, we consider NOMA in a multi-cell environment for this article. We first discuss the theory behind NOMA and an overview of the literature of NOMA. We then explain the main implementation issues and research challenges, with particular interest on multi-cell NOMA. Finally, the system-level performance evaluation of multi-cell NOMA solutions will be provided before concluding the article.

II. THEORY BEHIND NOMA

Analysis of cellular communication can generally be classified as either *downlink* or *uplink*. In the downlink channel, the BS simultaneously transmits signals to multiple users, whereas in the uplink channel multiple users transmit data to the same BS.

From an information-theoretic perspective, the downlink and uplink are modeled by the *broadcast channel* (BC) and *multiple access channel* (MAC), respectively. The basic premise behind single-cell NOMA in the power domain is to reap the benefits promised by the theory of multi-user channels [4]. As such, we review what information theory promises for these channels, both in the single-cell and multi-cell settings. In particular, we seek to answer the following two questions in this section: 1) what are the highest achievable throughputs for these multi-user channels? and 2) how can a system achieve such rates?

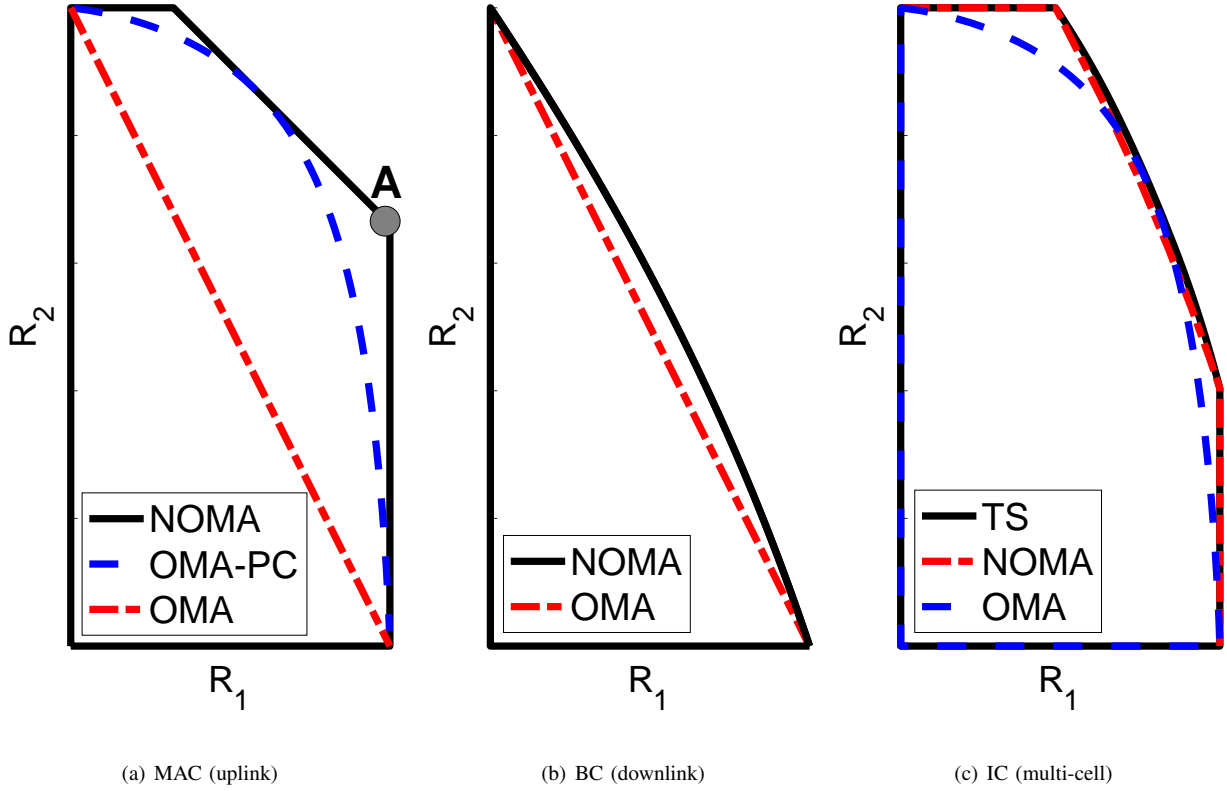


Fig. 1. Best achievable regions by OMA and NOMA in the multiple access channel (MAC), broadcast channel (BC), and interference channel (IC).

A. Single-Cell NOMA

The capacity regions of the two-user MAC and BC are achieved via NOMA, where both users' signals are transmitted at the same time and in the same frequency band [5]. The curves labeled by NOMA in Fig. 1(a) and Fig. 1(b) represent the MAC and BC capacity regions, respectively. Except for a few points, OMA is strictly suboptimal as can be seen from the figures. To gain more insight, we describe how the above regions are obtained. For OMA we consider a time division multiple access (TDMA) technique where α fraction of time ($0 \leq \alpha \leq 1$) is dedicated to user 1 and $\bar{\alpha} \triangleq 1 - \alpha$ fraction of time is dedicated to user 2. In this paper, $\mathcal{C}(x) \triangleq \frac{1}{2} \log_2(1 + x)$ and $\gamma_i = |h_i|^2 P$ is the received signal-to-noise ratio (SNR) for user i , where h_i is the channel gain, P is the transmitter power, and the noise power is normalized to unity.

1) *Uplink (MAC)*: Using OMA, each user sees a single-user channel in its dedicated fraction of time and, thus, $R_1 = \alpha \mathcal{C}(\gamma_1)$ and $R_2 = \bar{\alpha} \mathcal{C}(\gamma_2)$ are achievable. If power control is applied, these rates can be boosted to $R_1 = \alpha \mathcal{C}(\frac{\gamma_1}{\alpha})$ and $R_2 = \bar{\alpha} \mathcal{C}(\frac{\gamma_2}{\bar{\alpha}})$. In the case of NOMA, both users concurrently transmit, and their signals interfere with each other at the BS. The BS can use *successive interference cancellation* (SIC) to achieve any point in the NOMA region, which is the capacity region of this channel [4]. Particularly, to achieve point A the BS first decodes user 2's signal treating the other signal as noise. This results in $R_2 = \mathcal{C}(\frac{\gamma_2}{\gamma_1 + 1})$. The BS then removes user 2's

signal and decodes user 1's signal free of interference; i.e., $R_1 = \mathcal{C}(\gamma_1)$. From Fig. 1(a) it is seen that the gap between the NOMA and OMA regions becomes larger if power control is not used in OMA.

2) *Downlink (BC)*: In the downlink, OMA can only achieve $R_1 = \alpha\mathcal{C}(\gamma_1)$ and $R_2 = \bar{\alpha}\mathcal{C}(\gamma_2)$. However, making use of a NOMA scheme can strictly increase this rate region as shown in Fig. 1(b). In particular, the capacity region of this channel is known and can be achieved using *superposition coding* at the BS. For decoding, the user with the stronger channel uses SIC to decode its signal free of interference, i.e., $R_1 = \mathcal{C}(\beta\gamma_1)$, while the other user is capable of decoding at a rate of $R_2 = \mathcal{C}(\frac{\bar{\beta}\gamma_2}{\beta\gamma_2+1})$, where β is the fraction of the BS power allocated to user 1's data and $\bar{\beta} = 1 - \beta$. By varying β from 0 to 1, any rate pair (R_1, R_2) on the boundary of capacity region of the BC (NOMA region) can be achieved.

The fact that the capacity region of downlink NOMA is known enables us to find the optimum power allocation corresponding to any point (R_1, R_2) on the boundary of the capacity region. In fact, all we need to know to achieve such a rate pair is to find what fraction of the BS power should be allocated to each user. Corresponding to each (R_1, R_2) there is a $0 \leq \beta \leq 1$ such that βP and $\bar{\beta}P$ are the optimal powers for user 1 and user 2, respectively, where P is the BS power. Conversely, every β generates a point on the boundary of the capacity region.

The above argument implies that NOMA can improve *user-fairness* smoothly and in an optimal way by flexible power allocation. Suppose that a user has a poor channel condition. To boost this user's rate and improve user-fairness the BS can simply increase the fraction of power allocated to this user. We can look at this problem from yet another perspective. To increase the rate of such a user, we can maximize the weighted sum-rate $\mu R_1 + R_2$ where a high weight (μ) is given to such a user. This is because, to maximize $\mu R_1 + R_2$ for any $\mu \geq 0$ there exists an optimal power allocation strategy, determined by β . Seeing that $\mu > 1$ ($\mu < 1$) corresponds to the case where user 1 has higher (lower) weight than user 2, to improve the user-fairness we can assign an appropriate weight to the important user and find the corresponding β .

3) *K-User Uplink/Downlink*: In the above, we described coding strategies for the two-user uplink/downlink channels. Interestingly, very similar coding schemes are still capacity-achieving for the K -user MAC and BC, i.e., superposition coding with SIC gives the largest region for the K -user BC. Similarly, to achieve the capacity region of the K -user MAC, the users transmit their signals concurrently and the BS applies SIC, as described in [4, Section 6.1.4]. These schemes are based on NOMA as they allow multiple users to transmit at the same time and frequency. Additionally, OMA is strictly suboptimal [4].

B. Multi-Cell NOMA

In a multi-cell setting, these problems are more involved and simple channel models are insufficient. Unfortunately, capacity-achieving schemes are unknown. However, the achievable rate regions for the interference channel indicate the superiority of NOMA to OMA, as shown in Fig. 1(c).

1) *Interference Channel (IC)*: The capacity region of the two-user IC is not known in general; however, it is known that OMA is strictly sub-optimal. The Han-Kobayashi (HK) scheme [5] is the best known achievable scheme for the IC. In its basic form, the HK scheme employs rate-splitting and superposition coding at each transmitter. Since it uses superposition coding, the basic HK implies a NOMA. In general, the HK scheme applies time-sharing

to improve the basic HK region and can be seen as a combination of NOMA and OMA [6]. The HK scheme that combines NOMA and OMA gives the largest rate region [6], as shown in Fig. 1(c). In this figure, OMA refers to TDMA whereas NOMA refers to the basic HK scheme in which time-sharing is not applied. The third curve, labeled TS, is based on the HK scheme with time-sharing (TS) in which two time slots are used: in one time slot both users are active while in the other time slot only one of them is transmitting. As can be seen from this figure, both NOMA and OMA are suboptimal when compared with the case where NOMA and OMA are combined with.

2) *Interfering MAC and BC*: Consider a mutually interfering two-cell network in the uplink, where each cell includes one MAC. Assume that only one of the transmitters of each MAC (typically the closest one to the cell-edge) is interfering with the BS of the other MAC. In this network, the interfering transmitters can employ HK coding, similar to that used in the IC, while the non-interfering transmitters in each MAC employ single-user coding. This NOMA-based transmission results in an inner bound which is within a one-bit gap of the capacity region [7]. Likewise, one can use interfering BC to model a mutually interfering two-cell downlink network.

Despite years of intensive research, finding optimal uplink and downlink transmit/receive strategies for multi-cell networks remains rather elusive. In fact, as discussed earlier, even for a much simpler case of the two-user IC, the optimal coding strategy is still unknown. Nonetheless, fundamental results from information theory as a whole suggest that NOMA-based techniques result in a superior rate region when compared with OMA.

It should be highlighted that, despite the above insight from information theory, OMA techniques have been used in the cellular networks from 1G to 4G, mainly to avoid interference due to its simplicity.¹ In addition, the lack of understanding of optimal strategies for multi-cell networks has motivated pragmatic approaches in which interference is simply treated as noise.

III. SINGLE-CELL NOMA: A REVIEW

As explained in Section II, the basic theory of NOMA has been around for several decades. However, a new wave of research on NOMA has been motivated by the advance of processors which make it possible to implement SIC at the user equipment. Saito et al. [1] first observed the potential of NOMA for 5G systems. They showed that NOMA can improve system throughput and user-fairness over orthogonal frequency division multiple access (OFDMA). Since then, NOMA has attracted considerable attention from both industry and academia. To make this concept more practical, several issues like user-pairing, power allocation, and SIC implementation issues have been studied in [8]. NOMA has also been considered in the 3GPP LTE-A systems under the name of multi-user superposition transmission [9].

The performance of NOMA can be further boosted in multi-antenna networks. MIMO-NOMA solutions exploit multiplexing and diversity gains to improve outage probability and throughput, by converting the MIMO channel into multiple parallel channels [10].

¹For 3G, wideband code-division multiple access (WCDMA) was adopted, wherein *orthogonal* channelization codes are used within a cell, yet *quasi-orthogonal* scrambling codes are used to reduce the inter-cell interference.

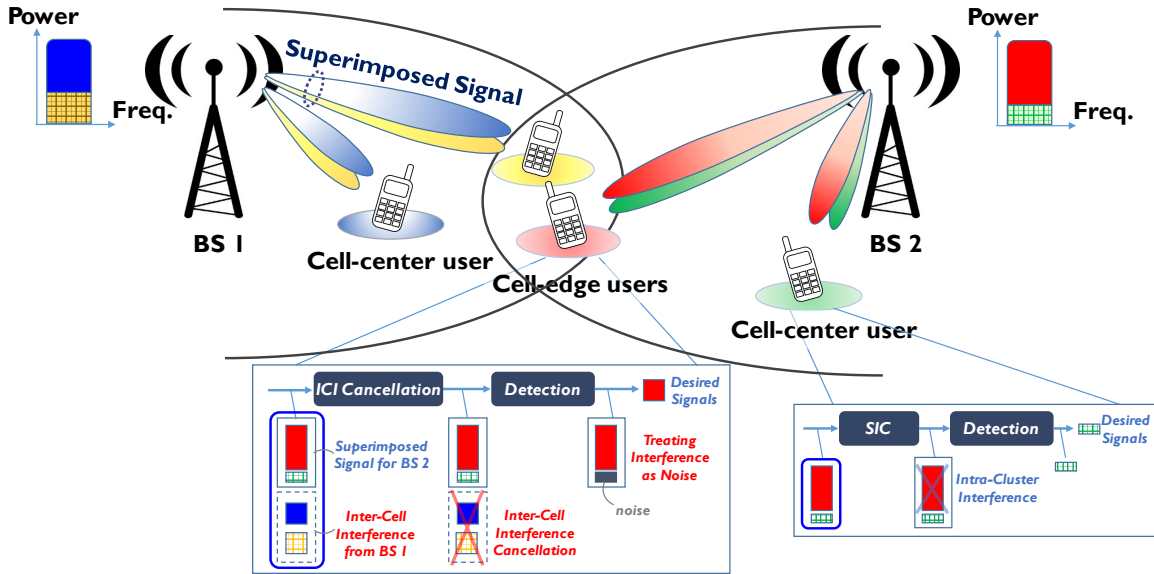


Fig. 2. An illustration of multi-cell NOMA networks.

IV. MULTI-CELL NOMA SOLUTIONS

In this section, we discuss recent research that combines interference management approaches with NOMA, called multi-cell NOMA. As illustrated in Fig. 2, ICI is the main issue in multi-cell NOMA networks, as it reduces a cell-edge user's performance. This is in contrast with single-cell NOMA, which aims at improving the user-fairness. Multi-cell techniques are used to harness the effect of ICI.

Multi-cell techniques can be categorized into coordinated scheduling/beamforming (CS/CB) and joint processing (JP) [11]. This classification is based on whether the data messages desired at the users should be shared among multiple BSs or not. These techniques can be combined with NOMA. For NOMA-CS/CB, data for a user is only available at and transmitted from a single BS. In contrast, NOMA-JP relies on data sharing among more than one BS.

A. NOMA with Joint Processing

In NOMA-JP, the users' data symbols are available at more than one BS. Based on the number of active BSs that serve a user, we can further divide NOMA-JP into two classes: NOMA-joint transmission (JT) and NOMA-dynamic cell selection (DCS).

1) *NOMA-JT*: This approach requires multiple BSs to simultaneously serve a user using a shared wireless resource instead of acting as interference to each other. This significantly improves the quality of the received signal at cell-edge users at the cost of slightly diminished rates for cell-center users. This cooperative setting is similar to a single-cell NOMA as the ICI for cell-edge users can be completely canceled using network MIMO techniques [12]. Such an approach usually relies on global channel state information (CSI) at all transmitters, which results in excessive backhaul overhead. To overcome the CSI sharing overhead for NOMA-JT, a coordinated superposition coding (CSC)

scheme for a two-cell downlink network was introduced in [13]. In this scheme, each cell-center user is served by its corresponding BS while the cell-edge user is served by both BSs, as shown in Fig. 3. Specifically, two BSs transmit Alamouti coded signals to a cell-edge user to achieve a higher transmission rate, while each BS also transmits signals to the cell-center user. It has been shown that the coordination between two cells allows NOMA to provide a common cell-edge user with a reasonable transmission rate without sacrificing cell-center users' rates. Let P and P_c be the powers of the cell-center and cell-edge users' messages per cell, respectively. Assume that $\gamma_{i,m} = |h_{i,m}|^2$ for $i \in \{1, 2, c\}$ and $m \in \{1, 2\}$, where $h_{j,m}$ and $h_{c,m}$ denote the channel coefficients to the cell-center user in cell j and the common cell-edge user from BS m , $\forall j, m \in \{1, 2\}$, respectively. The sum-rate of NOMA-JT given by $R_1 + R_2 + R_c$ (sum of rates for cell-center users R_1 and R_2 , and a common cell-edge user R_c) where $R_1 = \mathcal{C}(\frac{\gamma_{1,1}P}{\gamma_{1,2}P+1})$, $R_2 = \mathcal{C}(\frac{\gamma_{2,2}P}{\gamma_{2,1}P+1})$, and $R_c = \min\{\mathcal{C}(\frac{(\gamma_{1,1}+\gamma_{1,2})P_c}{(\gamma_{1,1}+\gamma_{1,2})P+1})$, $\mathcal{C}(\frac{(\gamma_{2,1}+\gamma_{2,2})P_c}{(\gamma_{2,1}+\gamma_{2,2})P+1})$, $\mathcal{C}(\frac{(\gamma_{c,1}+\gamma_{c,2})P_c}{(\gamma_{c,1}+\gamma_{c,2})P+1})\}$. Note that the last term comes from the condition that cell-edge user's message has to be decoded by that user and also cell-center users in both cells in order to operate SIC.

2) *NOMA-DCS*: In this case, the user's data is shared among multiple BSs, but it is transmitted only from one selected BS. Note that the transmitting BS can be dynamically changed over time by using order statistics. Suppose $|h_{c,2}|^2 > |h_{c,1}|^2$; then, BS 2 becomes the sole serving BS for a cell-edge user until the order statistics is not changed. That is, only BS 2 employs NOMA strategy to support a pair of cell-edge and cell-center users at the same time while BS 1 serves only its corresponding cell-center user (see Fig. 3). Since BS 1 employs OMA instead of NOMA, rate expressions for NOMA-JT, except for R_2 , should be modified for NOMA-DCS as $R_1 = \mathcal{C}(\frac{\gamma_1 P}{\gamma_1(P+P_c)+1})$ and $R_c = \min\{\mathcal{C}(\frac{\gamma_2 P_c}{(\gamma_1+\gamma_2)P+1})$, $\mathcal{C}(\frac{\gamma_2^c P_c}{(\gamma_1^c+\gamma_2^c)P+1})\}$ since user 1 does not use SIC for NOMA transmission.

B. NOMA with Coordinated Scheduling/Beamforming

The designs of CS/CB for NOMA differ from those of JP in that the users' data are not shared among the BSs. However, the cooperating BSs still need to exchange global CSI and cooperative scheduling information via a standardized interface named X2. This may result in a non-negligible overhead especially for high mobility cell-edge users. This subsection briefly discusses how to apply CS or CB to NOMA to tackle ICI problem. The illustration of NOMA-CS/CB is shown in Fig. 3.

1) *NOMA-CB*: In this case, data for a user is only available at one serving BS, and the beamforming decision is made with coordination that relies on global CSI. The authors in [14] proposed two novel *interference alignment* (IA)-based CB methods in which two BSs jointly optimize their beamforming vectors in order to improve the data rates of cell-edge users by removing ICI. Both algorithms aim to choose the transmit/receive beamforming vectors to satisfy the zero ICI as well as zero inter-cluster interference. These algorithms are termed interfering channel alignment (ICA)-based CB and IA-based CB. The former requires global CSI at the BS. However, the latter only requires the knowledge of cell-edge users' serving channel at the BS but with a slightly large number of antennas to compensate for the lack of interfering channels' knowledge. In particular, when the number of users is sufficiently large, it turns out that the number of extra antennas required for the latter scheme becomes negligible. Moreover, the transmit and receive beamforming vectors for uplink multi-cell NOMA also can be directly obtained by uplink-downlink duality.

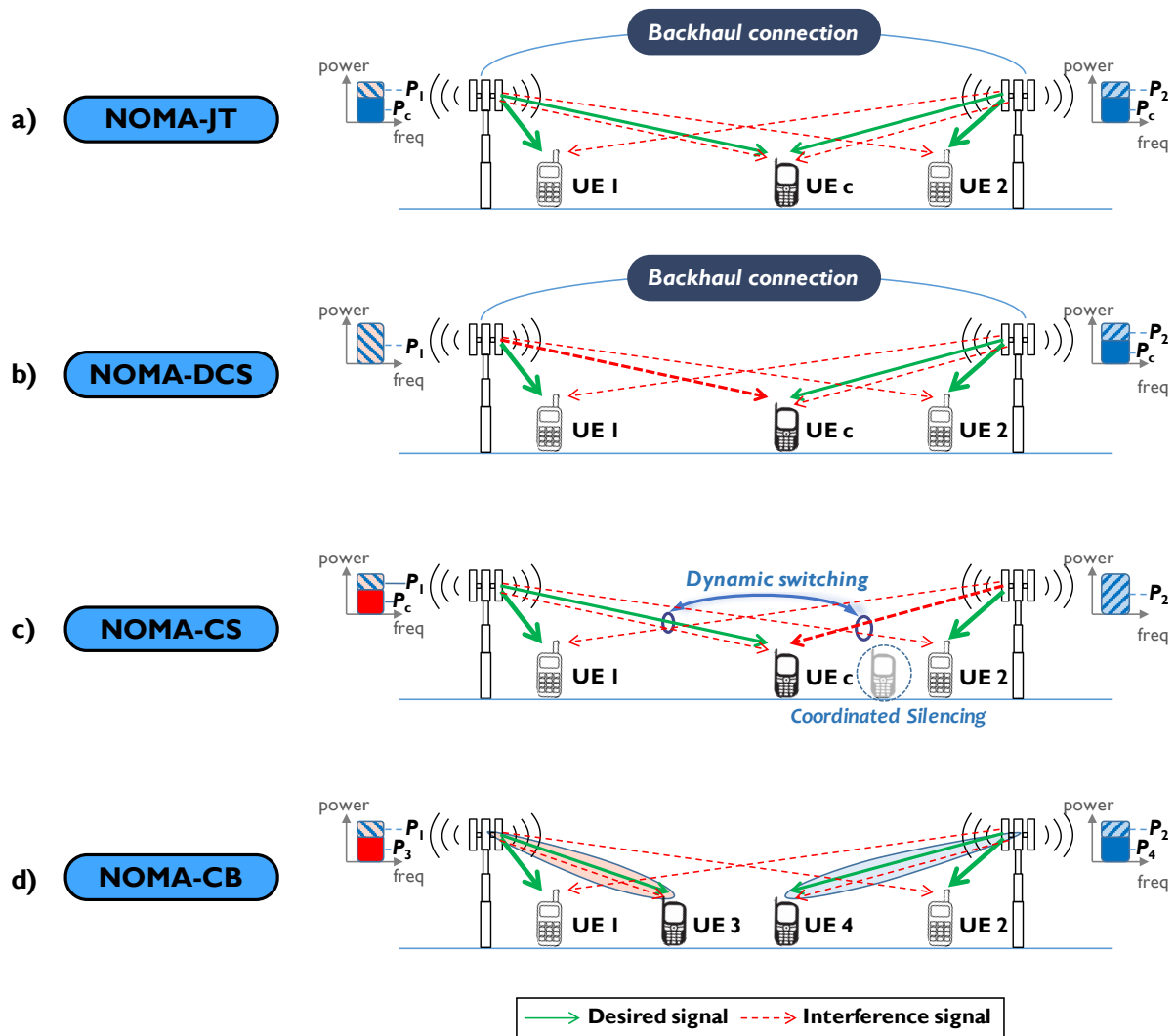


Fig. 3. Multi-cell NOMA solutions: a) NOMA-JT, b) NOMA-DCS, c) NOMA-CS, d) NOMA-CB

2) *NOMA-CS*: The key idea of *NOMA-CS* is to allow geographically separated BSs to coordinate scheduling to serve NOMA users with less ICI so as to ensure the proper QoS of cell-edge users. To guarantee the required data rate of the cell-edge users in a cell, the adjacent BS may decide not to transmit a superimposed message to a set of NOMA users, but just a dedicated message to a single cell-center user as in Fig. 3. However, such a *NOMA-CS* scheme is formulated as a combinatorial optimization, which is NP-hard. Therefore, a simple scheduling algorithm is indispensable in order to determine a set of NOMA users scheduled in each BS within a certain scheduling interval.

A summary of different multi-cell NOMA techniques is provided in Table 1. In addition, we compare the number of supported users by different NOMA schemes according to the number of clusters in each cell and the number of BS/user antennas. For comparison, we consider two-cell scenarios and assume that each BS and user has K antennas. Each cell consists of K clusters each having two users. It should be highlighted that single-cell NOMA

can support $2K$ users [10]. In contrast, single-cell OMA can serve only K users since the number of served users is limited by the number of antennas at the BS [4].

V. PRACTICAL CHALLENGES FOR MULTI-CELL NOMA

A. SIC Implementation Issues

As seen in Section II, SIC is at the heart of NOMA, and NOMA achieves the capacity region of the downlink and uplink channels (in a single-cell network) and the best rate region in the multi-cell setting. SIC, however, suffers from several practical issues, such as:

1) *Hardware Complexity*: SIC implies that each user has to decode information intended for all other users before its own in the SIC decoding order. This causes the complexity of decoding to scale with the number of users in the cell. To reduce the complexity, we can divide users in to multiple clusters and apply encoding/decoding within each cluster. Then, the complexity would be reasonable enough to be handled thanks to the advance of processor technologies during past decades. In fact, 3GPP LTE-A recently includes a new category of relatively complex user terminals, named network assisted interference cancellation and suppression (NAICS).

2) *Error propagation*: Error propagation means that if an error occurs in decoding a certain user's signal, all other users after this user in the SIC decoding order will be affected and their signals are likely to be decoded incorrectly. The side effect can be compensated by using stronger codes provided that the number of users is not very large.²

B. Imperfect CSI

Without perfect CSI at the user side it is not possible to completely remove the effect of the other users' signals from the received signal, which results in error propagation. Moreover, without perfect CSI about the interfering links at the BS, a joint precoder that guarantees no ICI is not known yet. In this regard, new beamforming designs which are robust to CSI errors must be developed for multi-cell NOMA.

²By making use of implementable near-capacity achieving AWGN channel codes (such as LDPC codes), we can get closer to the capacity region in practice.

TABLE I
A COMPARISON OF DIFFERENT MULTI-CELL NOMA SOLUTIONS.

	NOMA-CS	NOMA-CB	NOMA-DCS	NOMA-JT
# of transmission points	1	1	1 (Dynamic)	≥ 2
Shared information	CSI, Scheduling	CSI, BF	CSI, Data	(CSI), Data, BF
Backhaul type	Non-ideal	Non-ideal	Ideal	Ideal
Total # of supported users	$\ll 4K$	$4(K - 1)$	$3K$	$3K$ (or $4K$)
References		[14]	[13]	[12], [13]

C. Multi-User Power Allocation and Clustering

Power allocation and clustering are important factors that determine the performance gain of NOMA. To explain the effect of these factors, consider the simple case of single-cell two-user NOMA and assume that βP and $(1-\beta)P$ are the powers allocated to user 1 and user 2. As described in Section II, by varying β different points on the NOMA curve can be achieved. Therefore, β determines the rates for the users. This implies that with power allocation we can manage the system throughput and user-fairness. If there are more than two users in one cell, from the theory we know that all users' signals should be superimposed together; i.e., having one cluster maximizes the system throughput. However, in practice, having only one cluster can result in a serious performance degradation due to SIC error when there are many users in each cell. A suboptimal, but more practical, solution is to have multiple clusters per cell. However, it is still very hard to find the optimal clustering for a given number of clusters and the optimal solution is unknown. In multi-cell networks, ICI comes in which makes clustering, and power allocation even harder. It should be noted that in the multi-cell case, even when no SIC error assumed, the optimal clustering and power allocation solutions are not known. Therefore, clustering and power allocation algorithms with reasonable complexity and good performance are inevitable to implement NOMA in practical cellular systems.

D. Operation with FFR

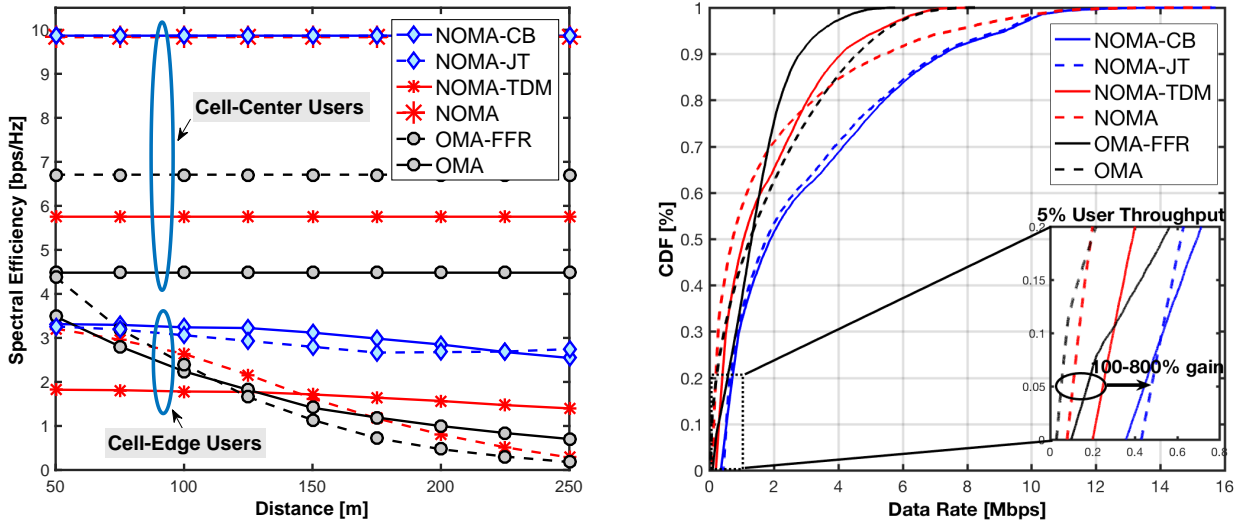
The basic idea of fractional frequency reuse (FFR) is to split a cell's bandwidth into multiple subbands and *orthogonally* allocate subbands for the cell-edge regions of the adjacent cells. This concept is in contrast with NOMA wherein orthogonalization is avoided due to its suboptimality. Despite being theoretically suboptimal, FFR is important as it offers a simple approach for ICI management without requiring CSI. Thus, it is important to investigate methods that can bring NOMA and FFR-based networks together. A simple idea to make use of both FFR and NOMA is to apply NOMA in the cell-center band and cell-edge band separately, which would pair cell-center users together (in the cell-center band) and cell-edge users together (in the cell-edge band). However, such users are not expected to have very different channel conditions, and any NOMA gain may not be noteworthy. Another idea is to pair a user from the cell-center region with a user from the cell-edge region in the cell-edge band to avoid ICI. Such a pairing will reduce cell-edge users' rates as their specific bands can be shared by the cell-center users too, which sacrifices the cell-edge users' rates. This, in turn, deteriorates user-fairness.

E. Security

The fact that in a NOMA-based transmission the user with better channel condition is able to decode the other user's signal brings new security concerns. Upper-layer security approaches (e.g., cryptographic) are still relevant since only the legitimate user has a key to decode its message. Nonetheless, physical layer security schemes are of interest but cannot be easily applied to the new environment.

VI. PERFORMANCE OF MULTI-CELL NOMA

In order to observe the potential gain of NOMA, numerical analysis is performed under a realistic multi-cell environment. In our simulation, we consider a two-cell downlink cellular network. As a performance metric, we



(a) Spectral efficiency as a function of cell-edge user's location

(b) Individual data rate CDF for random user deployments

Fig. 4. Performance comparison of different transmission schemes in multi-cell downlink networks

use the cumulative distribution function (CDF) of the user throughput, and the individual user throughputs for the cell-center and cell-edge users. Detailed simulation parameters are provided in Table II. In our simulation, following schemes are considered: OMA, OMA-FFR, NOMA, NOMA-TDM, NOMA-JT, and NOMA-CB. In OMA-FFR, FFR is used in controlling ICI on top of OMA transmission. Due to the effect of FFR, the cell-edge users experience no ICI while the cell-center users receive interference from other cell. In OMA, single-cell operation [10] is applied by treating all ICI as noise. Compared to OMA-FFR, ICI is a significant issue especially for cell-edge users, resulting in a severe SNR loss. Since it is inherently difficult to apply FFR in NOMA-based schemes as discussed in Section V, NOMA-TDM and NOMA schemes are considered. NOMA-TDM refers to a NOMA scheme that allows users in different cells to share one resource block via some form of orthogonalization, but NOMA simply acts as a single-cell operation [10] by treating the ICI as noise. In NOMA-CB and NOMA-JT, two BSs jointly optimize their beamforming vectors in order to mitigate ICI [13], [14].

In Fig. 4(a), we plot the performance of the cell-center and cell-edge users. Generally, the performance of OMA and NOMA decreases significantly with the location of cell-edge users since ICI mitigation is not considered. On the other hand, NOMA-TDM and OMA-FFR divide resources to support multi-cell environment, thus the rates of cell-edge users are improved compared to the single-cell operation schemes, such as OMA and NOMA. NOMA-CB can fully exploit all the resources to support all the users, and shows a twice increased performance compared to NOMA-TDM. NOMA-JT shows a performance almost similar to that of NOMA-CB, but its gain increases as the cell-edge user gets closer to the boarder of the cell. This is because the cell-edge user can take advantage of the link from the neighboring BS to improve its SNR via data sharing. Note that the cell-edge user performance of OMA-FFR and OMA is even better than that of NOMA-CB when the location of cell-edge user is relatively close

TABLE II
SIMULATION PARAMETERS

Cell layout	2 Cells
Cell radius	0.25 Km
Path loss exponent	4
Channel model	Rayleigh fading model
Channel estimation	Ideal
Number of transmitter antennas	4
Number of receiver antennas	4
Number of clusters per cell	4
Number of users per cluster	2
Users' locations	Randomly generated and uniformly distributed within the cell
User pairing	cell-center user from the disc with radius 0.125 Km cell-edge user from the ring
Transmission power	10 W
Noise power spectral density	10^{-10} W/Hz
Maximum number of multiplexed UEs	1 (OMA), 2 (NOMA)

to the BS, due to the inherently remaining inter-user interference of cell-edge user from NOMA transmission [15]. This phenomenon accounts for a motivation behind user-pairing for NOMA to be implemented in practice.

In Fig. 4(b), we plot the CDF of the user throughput. It can be seen that NOMA-CB and NOMA-JT achieve the best performance for any throughput (e.g., 5%-tile CDF point and average for user throughput). This is because ICI is effectively controlled by exploiting the multi-cell NOMA transmissions. As a matter of fact, OMA and NOMA are advantageous to the cell-center users' throughput due to full usage of a resource block while these are limited by severe ICI especially for the cell-edge users. On the contrary, NOMA-TDM and OMA-FFR are deployed to overcome ICI by further splitting a resource block into two parts (i.e., two cells) with sacrificing cell-center users' throughput.

VII. CONCLUSION

In this article, we described the theory behind NOMA in single-cell (both uplink and downlink) and multi-cell networks. This is followed by an up-to-date literature review of interference management techniques that apply NOMA in multi-cell networks. Numerical results show the significance of interference cancellation in NOMA. We have also highlighted major practical issues and challenges that arise in the implementation of multi-cell NOMA.

REFERENCES

- [1] Y. Saito, Y. Kishiyama, A. Benjebbour, T. Nakamura, A. Li, and K. Higuchi, "Non-orthogonal multiple access (NOMA) for cellular future radio access," in *Proc. IEEE 77th Vehicular Technology Conference (VTC Spring)*, pp. 1–5, 2013.
- [2] L. Dai, B. Wang, Y. Yuan, S. Han, C.-L. I, and Z. Wang, "Non-orthogonal multiple access for 5G: Solutions, challenges, opportunities, and future research trends," *IEEE Communications Magazine*, vol. 53, no. 9, pp. 74–81, 2015.
- [3] S. M. R. Islam, N. Avazov, O. A. Dobre, and K. S. Kwak, "Power-domain non-orthogonal multiple access (NOMA) in 5G systems: Potentials and challenges," *arXiv preprint arXiv:1608.05783*, 2016.

- [4] D. Tse and P. Viswanath, *Fundamentals of Wireless Communication*. Cambridge university press, 2005.
- [5] A. El Gamal and Y. H. Kim, *Network Information Theory*. Cambridge University Press, 2011.
- [6] M. Vaezi and H. V. Poor, "Simplified Han-Kobayashi region for one-sided and mixed Gaussian interference channels," in *Proc. IEEE International Conference on Communications (ICC)*, pp. 1–6, 2016.
- [7] Y. Pang and M. Varanasi, "Approximate capacity region of the MAC-IC-MAC," *arXiv preprint arXiv:1604.02234*, 2016.
- [8] Z. Ding, P. Fan, and H. V. Poor, "Impact of user pairing on 5G non-orthogonal multiple access downlink transmissions," *IEEE Transactions on Vehicular Technology*, vol. 19, no. 8, pp. 1462–1465, 2015.
- [9] 3GPP TD RP-150496, "Study on Downlink Multiuser Superposition Transmission," Mar. 2015.
- [10] Z. Ding, F. Adachi, and H. V. Poor, "The application of MIMO to non-orthogonal multiple access," *IEEE Transactions on Wireless Communications*, vol. 15, no. 1, pp. 537–552, 2016.
- [11] D. Lee, H. Seo, B. Clerckx, E. Hardouin, D. Mazzaresse, S. Nagata, and K. Sayana, "Coordinated multipoint transmission and reception in LTE-advanced: deployment scenarios and operational challenges," *IEEE Communications Magazine*, vol. 50, no. 2, pp. 148–155, 2012.
- [12] S. Han, C.-L. I, Z. Xu, and Q. Sun, "Energy efficiency and spectrum efficiency co-design: From NOMA to network NOMA," *IEEE Multimedia Communications Technical Committee E-Letter*, vol. 9, no. 5, pp. 21–24, 2014.
- [13] J. Choi, "Non-orthogonal multiple access in downlink coordinated two-point systems," *IEEE Communications Letters*, vol. 18, no. 2, pp. 313–316, 2014.
- [14] W. Shin, M. Vaezi, B. Lee, D. J. Love, J. Lee, and H. V. Poor, "Coordinated beamforming for multi-cell MIMO-NOMA," *IEEE Communications Letters*, vol. 21, no. 1, pp. 84–87, 2017.
- [15] H. Tabassum, M. S. Ali, E. Hossain, M. Hossain, and D. I. Kim, "Non-orthogonal multiple access (NOMA) in cellular uplink and downlink: Challenges and enabling techniques," *arXiv preprint arXiv:1608.05783*, 2016.