

Securing Downlink Massive MIMO-NOMA Networks with Artificial Noise

Ming Zeng, Nam-Phong Nguyen, Octavia A. Dobre, and H. Vincent Poor

Abstract

In this paper, we focus on securing the confidential information of massive multiple-input multiple-output (MIMO) non-orthogonal multiple access (NOMA) networks by exploiting artificial noise (AN). An uplink training scheme is first proposed with minimum mean squared error estimation at the base station. Based on the estimated channel state information, the base station precodes the confidential information and injects the AN. Following this, the ergodic secrecy rate is derived for downlink transmission. An asymptotic secrecy performance analysis is also carried out for a large number of transmit antennas and high transmit power at the base station, respectively, to highlight the effects of key parameters on the secrecy performance of the considered system. Based on the derived ergodic secrecy rate, we propose the joint power allocation of the uplink training phase and downlink transmission phase to maximize the sum secrecy rates of the system. Besides, from the perspective of security, another optimization algorithm is proposed to maximize the energy efficiency. The results show that the combination of massive MIMO technique and AN greatly benefits NOMA networks in term of the secrecy performance. In addition, the effects of the uplink training phase and clustering process on the secrecy performance are revealed. Besides, the proposed optimization algorithms are compared with other baseline algorithms through simulations, and their superiority is validated. Finally, it is shown that the proposed system outperforms the conventional massive MIMO orthogonal multiple access in terms of the secrecy performance.

Index Terms

This work was supported in part by the Natural Sciences and Engineering Research Council of Canada (NSERC) through its Discovery program, and in part by the U.S. National Science Foundation under Grants CCF-093970 and CCF-1513915.

M. Zeng, N.-P. Nguyen, and O. A. Dobre are with Memorial University, Canada. (e-mail: {mzeng, nnguyen, odobre}.mun.ca). N.-P. Nguyen is also with Hanoi University of Science and Technology, Vietnam (email: phong.nguyennam@hust.edu.vn).

H. V. Poor is with the Electrical Engineering Department, Princeton University, Princeton, NJ, USA (e-mail: poor@princeton.edu).

All authors contributed equally to the article.

Non-orthogonal multiple access (NOMA), massive multiple-input multiple-output (MIMO), physical layer security, artificial noise (AN).

I. INTRODUCTION

The development of Internet-of-Things demands massive connectivity over the limited radio spectrum. This requires the next generation wireless networks deploy new multiple access technologies with better spectral efficiency [1]. Recently, non-orthogonal multiple access (NOMA) has been introduced as a solution for this challenge [2], [3]. Power-domain NOMA allows multiple users to share the same time-frequency resource simultaneously by using superposition coding and advanced interference cancellation techniques, such as successive interference cancellation (SIC) [4]–[6]. As a result, NOMA can enhance the capacity of a network in both spatial and temporal dimensions [7]–[10]. However, from the security viewpoint, sharing the same time-frequency resource among users imposes secrecy challenges.

Traditionally, the security issues have been handled at the higher layers using encryption approaches. However, the development of computing technologies and the tremendous growth in the number of wireless devices have surfaced the vulnerability of the conventional encryption methods [11]. As a result, physical layer security (PLS) has been introduced as an additional protecting layer to the conventional encryption methods for securing confidential information [12]. The principle of PLS is to take advantage of the randomness of the wireless channels to restrain the illegitimate side from overhearing the legitimate users [13]. The community has shown a great interest in applying PLS to NOMA networks. In [14], the authors investigated the secrecy outage probability (SOP) of NOMA relay networks with two types of relay, i.e., amplify-and-forward and decode-and-forward. The paper revealed that in the high signal-to-noise ratio regime, the SOP of the considered NOMA relay network converges to a constant value. In [15], the secrecy performance of a stochastic NOMA network was considered, by modelling its users' locations using stochastic geometry. The results showed that the secrecy diversity order of the considered system is determined by that of the user pair with a poorer channel. In [16], the authors derived a closed-form solution for maximizing the secrecy sum rate of the NOMA while taking the users' quality of service requirements into consideration. In [17], the authors investigated a NOMA system in the presence of an external eavesdropper. The SOP of the considered system was derived and used to optimize the decoding order, transmission

rates, and allocated power. These studies have laid the initial foundation for exploiting PLS in NOMA networks.

Recently, massive multiple-input multiple-output (MIMO) has become one of the key technologies for 5G network [18]–[20]. By deploying hundreds of antennas at the base station (BS) to serve tens of users, massive MIMO exploits the high spatial resolution and large array gain to greatly enhance the throughput, spectral efficiency, and energy efficiency (EE) [21]–[23]. Massive MIMO networks are suggested to operate in time division duplex to address pilot contamination by exploiting channel reciprocity [18]. In massive MIMO networks, the BS can obtain the knowledge of the channel state information (CSI) via uplink training sequences of the users and employ this knowledge to precode the transmit data. The combination of massive MIMO and NOMA seems to be naturally matched since it can offer a great performance enhancement for a large number of users [24]. However, there are some challenges of this combination. Since the number of orthogonal sequences for the uplink training phase is limited, the massive number of users has to be grouped in clusters. In a cluster, users share the same training sequence. As a consequence, the quality of the uplink training phase can be compromised. Therefore, the spatial resolution is decreased, which can lead to leakage of the confidential information. There have been several studies of PLS for massive MIMO-NOMA networks. In [25], the authors have investigated the secrecy performance of a NOMA massive MIMO network in the presence of an active eavesdropper. The inter-user interference was utilized to enhance the secrecy performance of the network. Artificial noise (AN) has proven its effectiveness to secure the legitimate side from malicious attempts [26], [27]. Recently, in [28], the authors have proposed a joint alignment of multi-user constellations and AN to secure the massive MIMO-NOMA networks. By using a water filling power allocation between the constellation and AN, the error rate of the legitimate user is eliminated with a large number of antennas at the receiver, while the error rate of the eavesdropper approaches a floor when the number of eavesdropper's antennas is large. So far, it is the only work that deploys AN in NOMA networks. Therefore, the role of AN in massive MIMO-NOMA networks is far from being well-understood.

In this paper, we propose an AN-based PLS method for the massive MIMO-NOMA networks in the presence of a passive eavesdropper. In order to secure the downlink transmission, the BS uses its knowledge of CSI to precode the confidential information and inject the AN, which is different from [25]. Besides, because of the high complexity of the uplink training phase in the massive MIMO-NOMA networks, the AN approaches in [26], [27] are not suitable. Therefore,

in this paper, the AN is injected in the null-space of the effective channels of the clusters in the downlink transmission phase. To emphasize the role of the uplink training process on the secrecy performance of the considered system, the CSI knowledge at the BS is the result of an estimation process that is more practical than the assumption of perfect CSI in other existing work on PLS for massive MIMO-NOMA networks. To the best of our knowledge, this is the first work using AN to secure massive MIMO-NOMA networks when taking imperfect channel estimation into account. The contributions of this paper can be summarized as follows:

- We demonstrate a framework to analyze the secrecy performance of an AN-aided massive MIMO-NOMA network while taking the imperfect channel estimation into consideration. In particular, the ergodic secrecy rates for users are derived. The asymptotic expressions of the legitimate and illegitimate rates for a large number of antennas and high transmit power at the BS are also obtained. Note that the AN-aided massive MIMO-OMA network is a special case of the proposed system. The analysis expressions can be applied directly with the number of users in each cluster being equal to one.
- The results reveal that by using a sufficiently large number of antennas at the BS, the AN only affects the eavesdropper. In addition, when the transmit power at the BS is sufficiently high, the secrecy performance of a user depends on the AN, the intra-cluster interference, and the channel estimation error of its cluster.
- In order to further exploit the interference and AN, we study the maximization of the sum ergodic secrecy rate (SE) and the maximization of the EE in terms of the ergodic secrecy rates. In this work, the EE is defined as the sum ergodic secrecy rate over the total transmit power, which includes both the uplink and downlink powers. For the SE maximization problem, we first decompose it into two sub-problems, i.e., uplink and downlink power allocation (PA), based on alternating optimization. Then, we address each sub-problem using difference of convex (DC) programming. The EE maximization problem is of fractional form, and can be transformed into a series of SE maximization problems, which can be solved accordingly. Numerical results show that the proposed algorithms can significantly enhance the performance of the considered system, compared with other baseline algorithms.

The rest of this paper is organized as follows. The system and channel models are described in Section II. The analytical expressions for the ergodic secrecy rates of the considered system are developed in Section III. In Section IV, the optimization problems are proposed, and the

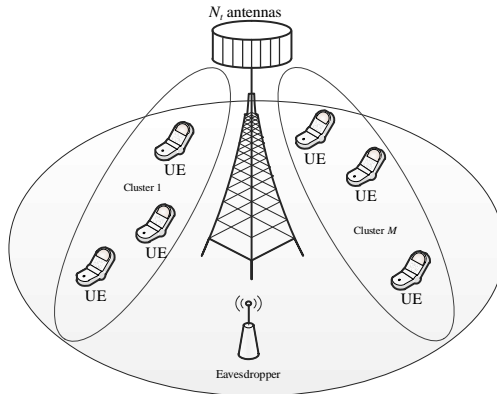


Fig. 1: System model.

solutions are discussed in Section V. The numerical results and discussions are presented in Section VI. Finally, we conclude the paper in Section VII.

Notations: Superscript $(\cdot)^H$ stands for the conjugate transpose. The expectation operation and Frobenius norm are denoted by $\mathbb{E}\{\cdot\}$ and $\|\cdot\|$, respectively. \mathbf{I}_{N_t} denotes the N_t -dimensional identity matrix. $\mathcal{CN}(\mu, \sigma^2)$ indicates complex normal distribution with μ mean and σ^2 variance.

II. SYSTEM AND CHANNEL MODELS

As shown in Fig. 1, we consider the downlink transmission in a massive MIMO-NOMA system, which includes one N_t -antenna BS, multiple single-antenna end users (UEs) that are grouped into M clusters with K_m users, $m = \{1, \dots, M\}$, in the m -th cluster, and one passive single-antenna eavesdropper. Before performing the downlink transmission, the BS needs the network's CSI to precode the information and inject the AN. Besides, the users also require knowledge of the precoding to decode the confidential information. Therefore, the BS and users exchange CSI and precoding knowledge in the training phases.

A. Training Phases

1) *Uplink training:* During one coherence interval duration of T samples, the users simultaneously send training sequences to the BS. Users in the same cluster employ the same training sequence. In order to prevent the training sequence of each cluster from interfering with each other, all clusters are assigned mutually orthogonal training sequences of length τ samples,

where $T \geq \tau \geq M$. The j -th cluster training sequence is denoted by a $\tau \times 1$ vector Φ_j , where $\Phi_j^H \Phi_i = 0, \forall i \neq j$, $\Phi_j^H \Phi_j = 1$. The received training signal at the BS is

$$\mathbf{Y} = \sum_{m=1}^M \sum_{k=1}^{K_m} \sqrt{P_{m,k} \beta_{m,k} \tau} \mathbf{h}_{m,k} \Phi_m^H + \mathbf{N}, \quad (1)$$

where $P_{m,k}$ is the transmit power of the k -th UE of the m -th cluster, $\beta_{m,k}$ is the large-scale fading, $\mathbf{h}_{m,k}$ is the small-scale fading, $\mathbf{h}_{m,k} \sim \mathcal{CN}(0, \mathbf{I}_{N_t})$, and the elements of $\mathbf{N} \sim \mathcal{CN}(0, 1)$ represent the additive white Gaussian noise (AWGN). Since Φ_m is known at the BS, the BS pre-processes the received signal as follows:

$$\begin{aligned} \underbrace{\mathbf{Y} \Phi_m^H}_{\tilde{\mathbf{y}}_m} &= \sum_{k=1}^{K_m} \sqrt{P_{m,k} \beta_{m,k} \tau} \mathbf{h}_{m,k} + \underbrace{\mathbf{N} \Phi_m^H}_{\tilde{\mathbf{n}}_m} \\ &= \sqrt{\sum_{k=1}^{K_m} P_{m,k} \beta_{m,k} \tau} \mathbf{h}_m + \tilde{\mathbf{n}}_m, \end{aligned} \quad (2)$$

where $\mathbf{h}_m = \frac{\sum_{k=1}^{K_m} \sqrt{P_{m,k} \beta_{m,k} \tau} \mathbf{h}_{m,k}}{\sqrt{\sum_{k=1}^{K_m} P_{m,k} \beta_{m,k} \tau}}$ is the effective channel for the m -th cluster.

The BS uses the minimum mean squared error (MMSE) technique to estimate \mathbf{h}_m .¹ The estimate of \mathbf{h}_m is [25]

$$\hat{\mathbf{h}}_m = \frac{\sqrt{\sum_{k=1}^{K_m} P_{m,k} \beta_{m,k} \tau}}{1 + \sum_{k=1}^{K_m} P_{m,k} \beta_{m,k} \tau} \tilde{\mathbf{y}}_m. \quad (3)$$

The relation between $\mathbf{h}_{m,k}$ and $\hat{\mathbf{h}}_m$ is

$$\mathbf{h}_{m,k} = \sqrt{\rho_{m,k}} \hat{\mathbf{h}}_m + \sqrt{1 - \rho_{m,k}} \boldsymbol{\varepsilon}_{m,k}, \quad (4)$$

where $\boldsymbol{\varepsilon}_{m,k} \sim \mathcal{CN}(0, \mathbf{I}_{N_t})$ is the error vector, which is independent of $\hat{\mathbf{h}}_m$. Besides, $\rho_{m,k} = \frac{P_{m,k} \beta_{m,k} \tau}{1 + \sum_{i=1}^{K_m} P_{m,i} \beta_{m,i} \tau}$ [25].

Remark 1: For each cluster, the error of the estimation process depends on the uplink transmit power of each user, the number of users in a cluster, the large-scale fading, and the length of the training sequences. This error can be reduced by decreasing the number of users in a cluster. However, this leads to an increase in the number of clusters, and further yields more orthogonal training sequences, which are limited in certain cases, e.g., crowded stadium, busy city center, etc.

¹The use of MMSE has been widely adopted in massive MIMO system [19], [20].

After the estimation process, the BS uses the estimates of the cluster's effective channels to precode. In this paper, we assume that the BS employs the maximal ratio transmission (MRT) precoder, which is simple and nearly optimal in massive MIMO networks [20]. The precoder is defined as

$$\mathbf{w}_m = \frac{\hat{\mathbf{h}}_m}{\|\hat{\mathbf{h}}_m\|}. \quad (5)$$

2) *Downlink training*: The downlink training phase is similar to the uplink training phase, except that the BS uses the obtained precoder to beam the downlink pilots to the clusters. Since the downlink pilots are known at the users, these users can estimate accurately their effective channel gains, i.e., $|\sqrt{\beta_{m,k}}\mathbf{h}_{m,k}^H\mathbf{w}_m|^2$. We assume that the estimation process at users is perfect.² Without loss of generality, the users' effective channel gains of the m -th cluster are ordered as follows:

$$|\sqrt{\beta_{m,1}}\mathbf{h}_{m,1}^H\mathbf{w}_m|^2 \geq \dots \geq |\sqrt{\beta_{m,K_m}}\mathbf{h}_{m,K_m}^H\mathbf{w}_m|^2. \quad (6)$$

During this phase, the eavesdropper also obtains its effective channel gain, i.e., $|\sqrt{\beta_E}\mathbf{g}^H\mathbf{w}_m|^2$, where β_E is the large-scale fading and \mathbf{g} is the small-scale fading vector corresponding to the eavesdropper.

B. NOMA Downlink Transmission

In order to perform NOMA downlink transmission, the BS conducts superposition coding for each cluster. The superposition coding for the m -th cluster is as follows:

$$x_m = \sum_{k=1}^{K_m} \sqrt{Q_{m,k}}s_{m,k}, \quad (7)$$

where $Q_{m,k}$ is the transmit power allocated to UE $_{m,k}$, and $s_{m,k}$ is the corresponding transmitted signal, satisfying $\mathbb{E}\{|s_{m,k}|^2\} = 1$. For securing the confidential information, the BS injects AN into the transmitted signals. The BS combines all cluster signals as follows:

$$\mathbf{x} = \sum_{m=1}^M (\mathbf{w}_m x_m + \sqrt{Q_{m,0}}\mathbf{z}_m \lambda_m), \quad (8)$$

²This assumption is reasonable since it has been proven that at a sufficiently high transmit power, the error of the channel estimation process at the receiver is sufficiently small and can be neglected [27].

where \mathbf{w}_m and \mathbf{z}_m are the precoding vector and AN vector for the m -th cluster, respectively, $\hat{\mathbf{h}}_m^H \mathbf{z}_m = 0$, $\|\mathbf{z}_m\|^2 = 1$; $Q_{m,0}$ is the power allocated for the AN and λ_m is the AN signal of the m -th cluster, $\mathbb{E}\{|\lambda_m|^2\} = 1$.

The received signal at the UE $_{m,k}$ is

$$\begin{aligned}
y_{m,k} &= \underbrace{\sqrt{\beta_{m,k}} \mathbf{h}_{m,k}^H \sqrt{Q_{m,k}} \mathbf{w}_m s_{m,k}}_{\text{Desired signal}} \\
&+ \underbrace{\sqrt{\beta_{m,k}} \mathbf{h}_{m,k}^H \left(\sum_{i=1, i \neq k}^{K_m} \sqrt{Q_{m,i}} \mathbf{w}_m s_{m,i} + \sqrt{Q_{m,0}} \mathbf{z}_m \lambda_m \right)}_{\text{Intra-cluster interference and AN}} \\
&+ \underbrace{\sqrt{\beta_{m,k}} \mathbf{h}_{m,k}^H \sum_{j=1, j \neq m}^M \left(\sum_{i=1}^{K_j} \sqrt{Q_{j,i}} \mathbf{w}_j s_{j,i} + \sqrt{Q_{j,0}} \mathbf{z}_j \lambda_j \right)}_{\text{Inter-cluster interference and AN}} \\
&+ n_{m,k}, \tag{9}
\end{aligned}$$

where $n_{m,k} \sim \mathcal{CN}(0, 1)$ is the AWGN at UE $_{m,k}$.

The eavesdropper tries to intercept the confidential information of UE $_{m,k}$. The received signal at the eavesdropper is

$$\begin{aligned}
y_{m,k}^e &= \underbrace{\sqrt{\beta_E} \mathbf{g}^H \sqrt{Q_{m,k}} \mathbf{w}_m s_{m,k}}_{\text{Desired signal}} \\
&+ \underbrace{\sqrt{\beta_E} \mathbf{g}^H \left(\sum_{i=1, i \neq k}^{K_m} \sqrt{Q_{m,i}} \mathbf{w}_m s_{m,i} + \sqrt{Q_{m,0}} \mathbf{z}_m \lambda_m \right)}_{\text{Intra-cluster interference and AN}} \\
&+ \underbrace{\sqrt{\beta_E} \mathbf{g}^H \sum_{j=1, j \neq m}^M \left(\sum_{i=1}^{K_j} \sqrt{Q_{j,i}} \mathbf{w}_j s_{j,i} + \sqrt{Q_{j,0}} \mathbf{z}_j \lambda_j \right)}_{\text{Inter-cluster interference and AN}} \\
&+ n_e, \tag{10}
\end{aligned}$$

where $n_e \sim \mathcal{CN}(0, 1)$ is the AWGN at the eavesdropper.

III. SECRECY PERFORMANCE ANALYSIS

In this section, we derive the ergodic secrecy rate of UE $_{m,k}$ from its ergodic legitimate rate and its corresponding ergodic eavesdropping rate.

A. Ergodic Secrecy Rate

The ergodic secrecy rate of UE_{*m,k*} is

$$\begin{aligned} R_{m,k}^{sec} &= \mathbb{E} \{ [R_{m,k} - R_{m,k}^e]^+ \} \\ &\approx [\mathbb{E} \{ R_{m,k} \} - \mathbb{E} \{ R_{m,k}^e \}]^+, \end{aligned} \quad (11)$$

where $[x]^+ = \max(x, 0)$. This approximation is reasonable in massive MIMO systems owing to the channel hardening property [29]. The achievable rate of UE_{*m,k*} is

$$\bar{R}_{m,k} = \mathbb{E} \{ R_{m,k} \} \approx \left(1 - \frac{\tau}{T}\right) \log_2(1 + \bar{\gamma}_{m,k}), \quad (12)$$

where $\mathbb{E} \{ \cdot \}$ denotes the expectation operator and $\bar{\gamma}_{m,k} = \frac{\kappa_{m,k}}{\sum_{t=1}^3 \mathfrak{S}_{m,k,t+1}}$, with

$$\begin{aligned} &\kappa_{m,k} \\ &= \left| \mathbb{E} \left\{ \sqrt{Q_{m,k} \beta_{m,k}} \mathbf{h}_{m,k}^H \mathbf{w}_m \right\} \right|^2 \\ &= Q_{m,k} \beta_{m,k} \left| \mathbb{E} \left\{ (\sqrt{\rho_{m,k}} \hat{\mathbf{h}}_m^H \mathbf{w}_m + \sqrt{1 - \rho_{m,k}} \boldsymbol{\varepsilon}_{m,k}^H \mathbf{w}_m) \right\} \right|^2 \\ &\stackrel{(a)}{=} Q_{m,k} \beta_{m,k} \rho_{m,k} \left| \mathbb{E} \left\{ \left\| \hat{\mathbf{h}}_m \right\| \right\} \right|^2 \\ &\stackrel{(b)}{=} Q_{m,k} \beta_{m,k} \rho_{m,k} \frac{\Gamma^2(N_t + \frac{1}{2})}{\Gamma^2(N_t)} \\ &\stackrel{(c)}{\approx} Q_{m,k} \beta_{m,k} \rho_{m,k} N_t, \end{aligned} \quad (13)$$

where step (a) holds true because $\mathbb{E} \{ \boldsymbol{\varepsilon}_{m,k}^H \mathbf{w}_m \} = \mathbb{E} \{ \boldsymbol{\varepsilon}_{m,k}^H \} \mathbb{E} \{ \mathbf{w}_m \} = 0$, $\Gamma(\cdot)$ is the gamma function, step (b) is based on the fact that $\left\| \hat{\mathbf{h}}_m \right\|$ has a scaled Chi distribution with $2N_t$ degrees of freedom by a factor of $\frac{1}{\sqrt{2}}$ [27]. Therefore, $\mathbb{E} \left\{ \left\| \hat{\mathbf{h}}_m \right\| \right\} = \frac{\Gamma(N_t + \frac{1}{2})}{\Gamma(N_t)}$, and step (c) is obtained by using the approximation $\frac{\Gamma^2(N_t + \frac{1}{2})}{\Gamma^2(N_t)} \xrightarrow{N_t \rightarrow \infty} N_t$ [30].

Further, $\mathfrak{S}_{m,k,i}$ for $i = \{1, 2, 3\}$ in the expression of $\bar{\gamma}_{m,k}$ are given in (14), (15) and (16), respectively on the top of the next page. Note that step (a) in (15) is obtained because $\hat{\mathbf{h}}_m^H \mathbf{z}_m = 0$ and $\boldsymbol{\varepsilon}_{m,k}$ is independent of \mathbf{z}_m .

It can be seen that $\mathfrak{S}_{m,k,1}$ denotes the desired signal leakage due to the imperfect uplink channel estimation, while $\mathfrak{S}_{m,k,2}$ represents the intra-cluster interference after SIC and the AN leakage. In addition, $\mathfrak{S}_{m,k,3}$ expresses the inter-cluster interference and AN.

Remark 2: Note that perfect SIC is assumed to obtain $\mathfrak{S}_{m,k,2}$. That is, the k -th user first decodes and subtracts the interfering signals from the K_m -th to the $(k+1)$ -th user in sequence, and then

$$\begin{aligned}
\mathfrak{S}_{m,k,1} &= Q_{m,k}\beta_{m,k} \left(\mathbb{E} \{ |\mathbf{h}_{m,k}^H \mathbf{w}_m|^2 \} - (\mathbb{E} \{ \mathbf{h}_{m,k}^H \mathbf{w}_m \})^2 \right) \\
&= Q_{m,k}\beta_{m,k} \left(\mathbb{E} \left\{ \left| \sqrt{\rho_{m,k}} \hat{\mathbf{h}}_m^H \mathbf{w}_m + \sqrt{1 - \rho_{m,k}} \boldsymbol{\varepsilon}_{m,k}^H \mathbf{w}_m \right|^2 \right\} - (\mathbb{E} \{ \mathbf{h}_{m,k}^H \mathbf{w}_m \})^2 \right) \\
&= Q_{m,k}\beta_{m,k} \left(\rho_{m,k} \mathbb{E} \left\{ \left| \hat{\mathbf{h}}_m^H \mathbf{w}_m \right|^2 \right\} + (1 - \rho_{m,k}) \mathbb{E} \left\{ \left| \boldsymbol{\varepsilon}_{m,k}^H \mathbf{w}_m \right|^2 \right\} - (\mathbb{E} \{ \mathbf{h}_{m,k}^H \mathbf{w}_m \})^2 \right) \\
&= Q_{m,k}\beta_{m,k} \left(\rho_{m,k} N_t + 1 - \rho_{m,k} - \rho_{m,k} \frac{\Gamma^2(N_t + \frac{1}{2})}{\Gamma^2(N_t)} \right) \\
&= Q_{m,k}\beta_{m,k}(1 - \rho_{m,k}), \tag{14}
\end{aligned}$$

$$\begin{aligned}
\mathfrak{S}_{m,k,2} &= \mathbb{E} \left\{ \beta_{m,k} \left(\sum_{i=1}^{k-1} Q_{m,i} |\mathbf{h}_{m,k}^H \mathbf{w}_m|^2 + Q_{m,0} |\mathbf{h}_{m,k}^H \mathbf{z}_m|^2 \right) \right\} \\
&= \beta_{m,k} \left(\sum_{i=1}^{k-1} Q_{m,i} \mathbb{E} \{ |\mathbf{h}_{m,k}^H \mathbf{w}_m|^2 \} + Q_{m,0} \mathbb{E} \{ |\mathbf{h}_{m,k}^H \mathbf{z}_m|^2 \} \right) \\
&\stackrel{(a)}{=} \beta_{m,k} \left[\sum_{i=1}^{k-1} Q_{m,i} (\rho_{m,k} N_t + 1 - \rho_{m,k}) + Q_{m,0} (1 - \rho_{m,k}) \right], \tag{15}
\end{aligned}$$

$$\begin{aligned}
\mathfrak{S}_{m,k,3} &= \mathbb{E} \left\{ \beta_{m,k} \sum_{j=1, j \neq m}^M \left(\sum_{i=1}^{K_j} Q_{j,i} |\mathbf{h}_{m,k}^H \mathbf{w}_j|^2 + Q_{j,0} |\mathbf{h}_{m,k}^H \mathbf{z}_j|^2 \right) \right\} \\
&= \beta_{m,k} \sum_{j=1, j \neq m}^M \sum_{i=0}^{K_j} Q_{j,i}. \tag{16}
\end{aligned}$$

demodulates its desired signal $s_{m,k}$. In other words, the residual intra-cluster interference is only from the users with stronger channel gains, i.e., the first user to the $(k - 1)$ -th user. In practice, owing to channel estimation error, hardware limitation, low signal quality, and so on, the decoding error of the weak interfering signal may occur. Consequently, there exists residual interference from the weak users after SIC, namely imperfect SIC. This residual interference is similar to the intra-cluster interference. As shown in [31]–[33], the residual interference can be modeled as a linear function of the power of the interfering signal, and the coefficient of imperfect SIC can be obtained through long-term measurements. As a result, the ergodic secrecy rate in the presence of imperfect SIC can be directly derived by adding the term of residual interference in

$$\bar{R}_{m,k}^e = \left(1 - \frac{\tau}{T}\right) \log_2 \left(1 + \frac{Q_{m,k}\beta_{\mathbb{E}}}{\beta_{\mathbb{E}} \sum_{i=0, i \neq k}^{K_m} Q_{m,i} + \beta_{\mathbb{E}} \sum_{j=1, j \neq m}^M \sum_{i=0}^{K_j} Q_{j,i} + 1}\right). \quad (18)$$

$\mathfrak{S}_{m,k,2}$.

The ergodic eavesdropping rate corresponding to $\text{UE}_{m,k}$ is

$$\bar{R}_{m,k}^e = \mathbb{E} \{R_{m,k}^e\} \approx \left(1 - \frac{\tau}{T}\right) \log_2(1 + \bar{\gamma}_{m,k}^e), \quad (17)$$

where $\bar{\gamma}_{m,k}^e = \frac{\kappa_{m,k}^e}{\sum_{t=1}^2 \mathfrak{S}_{m,k,t}^e + 1}$, with

$$\begin{aligned} \kappa_{m,k}^e &= Q_{m,k}\beta_{\mathbb{E}}\mathbb{E} \{|\mathbf{g}^H \mathbf{w}_m|^2\} = Q_{m,k}\beta_{\mathbb{E}}, \\ \mathfrak{S}_{m,k,1}^e &= \sum_{i=1, i \neq k}^{K_m} Q_{m,i}\beta_{\mathbb{E}}\mathbb{E} \{|\mathbf{g}^H \mathbf{w}_m|^2\} \\ &\quad + Q_{m,0}\beta_{\mathbb{E}}\mathbb{E} \{|\mathbf{g}^H \mathbf{z}_m|^2\} \\ &= \beta_{\mathbb{E}} \sum_{i=0, i \neq k}^{K_m} Q_{m,i}, \\ \mathfrak{S}_{m,k,2}^e &= \beta_{\mathbb{E}} \sum_{j=1, j \neq m}^M \left(\sum_{i=1}^{K_j} Q_{j,i}\mathbb{E} \{|\mathbf{g}^H \mathbf{w}_j|^2\} \right. \\ &\quad \left. + Q_{j,0}\mathbb{E} \{|\mathbf{g}^H \mathbf{z}_m|^2\} \right) \\ &= \beta_{\mathbb{E}} \sum_{j=1, j \neq m}^M \sum_{i=0}^{K_j} Q_{j,i}. \end{aligned}$$

Therefore, $\bar{R}_{m,k}^e$ can be simplified as (18) on the top of the next page.³

By comparing the intra-cluster interference terms in $\bar{R}_{m,k}$ and $\bar{R}_{m,k}^e$, i.e., $\mathfrak{S}_{m,k,2}$ and $\mathfrak{S}_{m,k,1}^e$, we can observe that the intra-cluster interference has less impact on the legitimate users owing to SIC. This helps to achieve a higher secrecy rate.

³It is possible to extend this work to the case of multiple eavesdroppers or multi-antenna eavesdropper since (18) can be applied to each eavesdropper or each antenna of a multi-antenna eavesdropper. The secrecy performance in these cases is determined by the strongest eavesdropper or the strongest eavesdropping antenna.

B. Asymptotic Secrecy Performance

In this subsection, increasing the number of antennas and the transmit power at the BS are respectively studied to reveal insights into the considered system.

1) *Large Number of Antennas at the BS*: We first investigate the impact of a large number of antennas at the BS on the secrecy performance. From (18), we can observe that the eavesdropping rate is independent of the number of antennas at the BS. When this number is large, the legitimate rate is expressed as

$$\bar{R}_{m,k} \stackrel{N_t \rightarrow \infty}{=} \left(1 - \frac{\tau}{T}\right) \log_2 \left(1 + \frac{Q_{m,k}}{\sum_{i=1}^{k-1} Q_{m,i}}\right). \quad (19)$$

Remark 3: When the number of antennas at the BS is sufficiently large, the secrecy rate converges to a constant value. At the legitimate side, the effect of imperfect CSI, fading, inter-cluster interference, and AN leakage is negligible because of channel hardening. The legitimate rate depends only on the intra-cluster transmit powers. Meanwhile, the eavesdropping rate suffers from noise, interferences, and fading. Obviously, by using AN, the secrecy performance can be guaranteed in this scenario.

2) *High Transmit Power at the BS*: In order to reveal the impact of the transmit power at the BS, the transmit power for each user is set proportional to the maximum transmit power of the BS, i.e., $Q_{m,k} = \sigma_{m,k} Q_{\max}$, where Q_{\max} is the maximum transmit power at the BS and $\sum_{m=1}^M \sum_{k=1}^{K_m} \sigma_{m,k} = 1$. When Q_{\max} is large, the legitimate rate and the eavesdropping rate are respectively approximated as (20) and (21) on the top of the next page.

Remark 4: When the transmit power at the BS is high, we can observe that:

- The secrecy rate converges to a constant value. This value is independent of fading and the maximum transmit power.
- The legitimate rate and the eavesdropping rate suffer from the same amount of inter-cluster interference and inter-cluster AN. In other words, the secrecy rate is independent of the inter-cluster interference and inter-cluster AN.
- The eavesdropper is affected by the AN more heavily than the legitimate user. This effect depends on the uplink training process. Recalling Remark 1, we can conclude that the secrecy performance depends on the number of available orthogonal pilots.

$$\bar{R}_{m,k} \stackrel{Q_{\max} \rightarrow \infty}{=} \left(1 - \frac{\tau}{T}\right) \times \log_2 \left(1 + \frac{\sigma_{m,k} \rho_{m,k} N_t}{\sigma_{m,k} (1 - \rho_{m,k}) + \left[\sum_{i=1}^{k-1} \sigma_{m,i} (\rho_{m,k} N_t + 1 - \rho_{m,k}) + \sigma_{m,0} (1 - \rho_{m,k}) \right] + \sum_{j=1, j \neq m}^M \sum_{i=0}^{K_j} \sigma_{i,j}} \right), \quad (20)$$

$$\bar{R}_{m,k}^e \stackrel{Q_{\max} \rightarrow \infty}{=} \left(1 - \frac{\tau}{T}\right) \log_2 \left(1 + \frac{\sigma_{m,k}}{\sum_{i=0, i \neq k}^{K_m} \sigma_{m,i} + \sum_{j=1, j \neq m}^M \sum_{i=0}^{K_j} \sigma_{i,j}} \right). \quad (21)$$

IV. OPTIMIZATION PROBLEMS

In this section, we consider the optimization of the uplink and downlink PA to fully exploit the potential of the proposed secure massive MIMO-NOMA network. Two system level criteria are respectively considered, i.e., the SE maximization and the EE maximization.

A. SE Maximization

First, we aim to maximize the SE for the considered system, which is formulated as

$$\max_{\mathbf{P}, \mathbf{Q}} \sum_{m=1}^M \sum_{k=1}^{K_m} R_{m,k}^{sec} \quad (22a)$$

$$\text{s.t. } 0 \leq P_{m,k} \leq P_{m,k}^{\max}, m \in \{1, \dots, M\}, \quad (22b)$$

$$k \in \{1, \dots, K_m\},$$

$$Q_{m,k} \geq 0, m \in \{1, \dots, M\}, k \in \{0, \dots, K_m\}, \quad (22c)$$

$$\sum_{m=1}^M \sum_{k=0}^{K_m} Q_{m,k} \leq Q_{\max}, \quad (22d)$$

where $\mathbf{P} \in \mathcal{R}^{M \times K_m}$ and $\mathbf{Q} \in \mathcal{R}^{M \times (K_m+1)}$ denote the matrix for the uplink and downlink power, respectively. Equations (22b) and (22d) represent the maximum transmit power constraint for each user in uplink and the total power constraint in downlink, respectively. Note that there exists a one-to-one mapping between $P_{m,k}$ and $\rho_{m,k}$.

B. EE Maximization

We also consider maximization of EE, defined as the sum ergodic secrecy rate over the total transmit power, which includes both the uplink and downlink power [34]–[36]. Moreover, for uplink and downlink power, both fixed circuit power and dynamic transmit power are considered [37], [38]. We denote the overall circuit power of the system as P_f . Then, the EE is given as

$$\eta_{EE} = \frac{\sum_{m=1}^M \sum_{k=1}^{K_m} R_{m,k}^{sec}}{\sum_{m=1}^M \sum_{k=1}^{K_m} P_{m,k} + \sum_{m=1}^M \sum_{k=0}^{K_m} Q_{m,k} + P_f}. \quad (23)$$

Accordingly, the EE optimization problem can be expressed as

$$\max_{\mathbf{P}, \mathbf{Q}} \eta_{EE}, \text{ s.t. (22b) – (22d)}. \quad (24)$$

V. PROPOSED SOLUTIONS

A. SE Maximization

Problem (22) is clearly non-convex, owing to the non-convex objective function. Moreover, it can be seen that the uplink power \mathbf{P} and downlink power \mathbf{Q} are coupled in the objective function. This coupling makes (22) difficult to handle. To address it, we propose to decompose the original problem into the following two sub-problems:

1) *Uplink Power Allocation for Channel Estimation:* For this sub-problem, we assume that the downlink power is appropriately allocated to the users and the AN, i.e., \mathbf{Q} is known and given. Then, the original problem can be simplified as

$$\max_{\mathbf{P}} \sum_{m=1}^M \sum_{k=1}^{K_m} R_{m,k}^{sec}, \text{ s.t. (22b)}. \quad (25)$$

2) *Downlink Power Allocation for Data Transmission:* Likewise, here we assume that the uplink power is appropriately allocated to the users, i.e., \mathbf{P} is known and given. Then, the original problem is re-expressed as

$$\max_{\mathbf{Q}} \sum_{m=1}^M \sum_{k=1}^{K_m} R_{m,k}^{sec}, \text{ s.t. (22c), (22d)}. \quad (26)$$

$$\begin{aligned}
f &= (1 - \frac{\tau}{T}) \sum_{m=1}^M \sum_{k=1}^{K_m} \underbrace{\log_2 \left((a_1 + a_2)\beta_{m,k}\tau P_{m,k} + a_3\tau \sum_{i=1}^{K_m} \beta_{m,i}P_{m,i} + a_3 \right)}_{f_1(\mathbf{P})} \\
&\quad - (1 - \frac{\tau}{T}) \sum_{m=1}^M \sum_{k=1}^{K_m} \underbrace{\log_2 \left(a_2\beta_{m,k}\tau P_{m,k} + a_3\tau \sum_{i=1}^{K_m} \beta_{m,i}P_{m,i} + a_3 \right)}_{f_2(\mathbf{P})}. \tag{28}
\end{aligned}$$

For sub-problem (1), since \mathbf{Q} is given, it can be seen that $\bar{R}_{m,k}^e$ is a constant. Then, we only need to consider $\bar{R}_{m,k}$. After some mathematical manipulations, $\bar{R}_{m,k}$ can be expressed as

$$\begin{aligned}
\bar{R}_{m,k} &= (1 - \frac{\tau}{T}) \log_2 \left(1 + \frac{\kappa_{m,k}}{\sum_{t=1}^3 \mathfrak{S}_{m,k,t} + 1} \right) \\
&= (1 - \frac{\tau}{T}) \\
&\quad \times \log_2 \left(1 + \frac{a_1\beta_{m,k}\tau P_{m,k}}{a_2\beta_{m,k}\tau P_{m,k} + a_3\tau \sum_{i=1}^{K_m} \beta_{m,i}P_{m,i} + a_3} \right), \tag{27}
\end{aligned}$$

where $a_1 = Q_{m,k}\beta_{m,k}N_t$, $a_2 = \beta_{m,k}[(N_t - 1) \sum_{i=1}^{k-1} Q_{m,i} - Q_{m,0} - Q_{m,k}]$, $a_3 = \beta_{m,k} \sum_{i=0}^k Q_{m,i} + \beta_{m,k} \sum_{j \neq m} \sum_{i=0}^{K_j} Q_{j,i} + 1$.

On this basis, we further transform $f = \sum_{m=1}^M \sum_{k=1}^{K_m} \bar{R}_{m,k}$ as (28) on the top of the next page.

Note that $(1 - \frac{\tau}{T})$ is a constant, which does not affect the solution and can be removed. Then, (25) can be re-expressed as

$$\max_{\mathbf{P}} \sum_{m=1}^M \sum_{k=1}^{K_m} f_1(\mathbf{P}) - f_2(\mathbf{P}), \text{ s.t. (22b)}, \tag{29}$$

where both functions $f_1(\mathbf{P})$ and $f_2(\mathbf{P})$ are concave. Thus, the objective $\sum_{m=1}^M \sum_{k=1}^{K_m} f_1(\mathbf{P}) - f_2(\mathbf{P})$ is a DC function. The gradient of f_2 at $P_{j,i}, \forall j \in \{1, \dots, M\}, i \in \{1, \dots, K_j\}$ is given

Algorithm 1: Proposed Power Allocation Algorithm for Sum Rate Maximization

```

1 Initialize  $\varepsilon \leftarrow 10^{-3}$ ; Initialize feasible downlink power  $\mathbf{Q}^{(0)}$ ;
2 repeat{Outer iteration}
3   Uplink power allocation:
4   repeat{Inner iteration}
5      $\mathbf{P}^{(l)} \leftarrow \max_{\mathbf{P}} \sum_{m=1}^M \sum_{k=1}^{K_m} [f_1(\mathbf{P}) - f_2(\mathbf{P}^{(l-1)}) - (P_{m,k} - P_{m,k}^{(l-1)}) \times$ 
6        $\sum_{j=1}^M \sum_{i=1}^{K_j} \nabla f_2(P_{j,i}^{(l-1)})]$  s.t. (22b)
7   until  $\mathbf{P}^{(l)}$  converges;
8    $R_{\text{sum}}^u \leftarrow \sum_{m=1}^M \sum_{k=1}^{K_m} R_{m,k}^{\text{sec}}$ ;
9   Downlink power allocation:
10  repeat{Inner iteration}
11     $\mathbf{Q}^{(l)} \leftarrow \max_{\mathbf{Q}} \sum_{m=1}^M \sum_{k=0}^{K_m} [g_1(\mathbf{Q}) + g_3(\mathbf{Q}) - g_2(\mathbf{Q}^{(l-1)}) - g_4(\mathbf{Q}^{(l-1)}) -$ 
12       $(Q_{m,k} - Q_{m,k}^{(l-1)}) \times \sum_{j=1}^M \sum_{i=0}^{K_j} \nabla g_2(Q_{j,i}^{(l-1)}) + \nabla g_4(Q_{j,i}^{(l-1)})]$ 
13      s.t. (22c), (22d)
14    until  $\mathbf{Q}^{(l)}$  converges;
15     $R_{\text{sum}}^d \leftarrow \sum_{m=1}^M \sum_{k=1}^{K_m} R_{m,k}^{\text{sec}}$ ;
16    Compute  $\varepsilon^* \leftarrow R_{\text{sum}}^d - R_{\text{sum}}^u$ .
17  until  $\varepsilon^* \leq \varepsilon$ ;
18  if  $R_{m,k}^{\text{sec}} < 0, \forall m \in \{1, \dots, M\}, k \in \{1, \dots, K_m\}$ 
19     $R_{m,k}^{\text{sec}} \leftarrow 0$ ;
20  end
21   $R_{\text{sum}} \leftarrow \sum_{m=1}^M \sum_{k=1}^{K_m} R_{m,k}^{\text{sec}}$ .

```

by

$$\nabla f_2(P_{j,i}) = \begin{cases} \frac{(a_2+a_3)\beta_{m,k}\tau/\ln 2}{a_2\beta_{m,k}\tau P_{m,k}+a_3\tau \sum_{i=1}^{K_m} \beta_{m,i}P_{m,i}+a_3}, & j = m, i = k, \\ \frac{a_3\beta_{m,i}\tau/\ln 2}{a_2\beta_{m,k}\tau P_{m,k}+a_3\tau \sum_{i=1}^{K_m} \beta_{m,i}P_{m,i}+a_3}, & j = m, i \neq k, \\ 0, & j \neq m. \end{cases}$$

The following procedure generates a sequence $\{\mathbf{P}^{(l)}\}$ of improved feasible solutions [39], [40]. Initialized from a feasible $\{\mathbf{P}^{(0)}\}$, $\{\mathbf{P}^{(l)}\}$ is obtained as the optimal solution of the following

$$\begin{aligned}
\bar{R}_{m,k} &= (1 - \frac{\tau}{T}) \log_2 \left(1 + \frac{b_1 Q_{m,k}}{b_2 Q_{m,k} + b_3 \sum_{i=1}^{k-1} Q_{m,i} + b_2 Q_{m,0} + \beta_{m,k} \sum_{j \neq m} \sum_{i=0}^{K_j} Q_{j,i} + 1} \right) \\
&= (1 - \frac{\tau}{T}) \log_2 \underbrace{\left((b_1 + b_2) Q_{m,k} + b_3 \sum_{i=1}^{k-1} Q_{m,i} + b_2 Q_{m,0} + \beta_{m,k} \sum_{j \neq m} \sum_{i=0}^{K_j} Q_{j,i} + 1 \right)}_{g_1(\mathbf{Q})} \\
&\quad - (1 - \frac{\tau}{T}) \log_2 \underbrace{\left(b_2 Q_{m,k} + b_3 \sum_{i=1}^{k-1} Q_{m,i} + b_2 Q_{m,0} + \beta_{m,k} \sum_{j \neq m} \sum_{i=0}^{K_j} Q_{j,i} + 1 \right)}_{g_2(\mathbf{Q})} \tag{31}
\end{aligned}$$

convex problem at the l -th iteration:

$$\begin{aligned}
\max_{\mathbf{P}} \quad & \sum_{m=1}^M \sum_{k=1}^{K_m} \left[f_1(\mathbf{P}) - f_2(\mathbf{P}^{(l-1)}) - \right. \\
& \left. (P_{m,k} - P_{m,k}^{(l-1)}) \times \sum_{j=1}^M \sum_{i=1}^{K_j} \nabla f_2(P_{j,i}^{(l-1)}) \right] \\
\text{s.t.} \quad & (22b). \tag{30}
\end{aligned}$$

Note that (30) can be efficiently solved by available convex software packages [41]. Moreover, since there exists no inter-cluster interference, the sum rate maximization can be done in parallel for each cluster, i.e., the system sum rate maximization equals to the cluster sum rate maximization.

After solving the above problem, we can obtain the value for \mathbf{P} . Accordingly, we can obtain $\rho_{m,k}$. On this basis, for sub-problem (2), after some mathematical manipulations, $\bar{R}_{m,k}$ can be expressed as (31) on the top of the next page. Note that in (31), $b_1 = \rho_{m,k} \beta_{m,k} N_t$, $b_2 = \beta_{m,k} (1 - \rho_{m,k})$, and $b_3 = \beta_{m,k} (\rho_{m,k} N_t + 1 - \rho_{m,k})$.

The gradient of g_2 at $Q_{j,i}, \forall j \in \{1, \dots, M\}, i \in \{0, \dots, K_j\}$ is given by (32) on the top of the next page.

$$\nabla g_2(Q_{j,i}) = \begin{cases} \frac{b_2}{b_2 Q_{m,k} + b_3 \sum_{i=1}^{k-1} Q_{m,i} + b_2 Q_{m,0} + \beta_{m,k} \sum_{j \neq m} \sum_{i=0}^{K_j} Q_{j,i} + 1}, & j = m, i = k \text{ or } 0, \\ \frac{b_3}{b_2 Q_{m,k} + b_3 \sum_{i=1}^{k-1} Q_{m,i} + b_2 Q_{m,0} + \beta_{m,k} \sum_{j \neq m} \sum_{i=0}^{K_j} Q_{j,i} + 1}, & j = m, i = 1, \dots, k-1, \\ \beta_{m,k}, & j \neq m, \\ 0, & \text{otherwise.} \end{cases} \quad (32)$$

Next, let us consider $-\bar{R}_{m,k}^e$, which can be re-written as

$$\begin{aligned} -\bar{R}_{m,k}^e &= (1 - \frac{\tau}{T}) \log_2 \left(\underbrace{\beta_E \sum_{i \neq k} Q_{m,i} + \beta_E \sum_{j \neq m} \sum_{i=0}^{K_j} Q_{j,i} + 1}_{g_3(\mathbf{Q})} \right) \\ &\quad - (1 - \frac{\tau}{T}) \log_2 \left(\underbrace{\beta_E \sum_{i=0}^{K_m} Q_{m,i} + \beta_E \sum_{j \neq m} \sum_{i=0}^{K_j} Q_{j,i} + 1}_{g_4(\mathbf{Q})} \right). \end{aligned} \quad (33)$$

The gradient of g_4 at $Q_{j,i}, \forall j \in \{1, \dots, M\}, i \in \{0, \dots, K_j\}$ is given by

$$\nabla g_4(Q_{j,i}) = \frac{\beta_E / \ln 2}{\beta_E \sum_{i=0}^{K_m} Q_{m,i} + \beta_E \sum_{j \neq m} \sum_{i=0}^{K_j} Q_{j,i} + 1}. \quad (34)$$

The following procedure generates a sequence $\{\mathbf{Q}^{(l)}\}$ of improved feasible solutions [39], [40]. Initialized from a feasible $\{\mathbf{Q}^{(0)}\}$, $\{\mathbf{Q}^{(l)}\}$ is obtained as the optimal solution of the following convex problem at the l -th iteration:

$$\begin{aligned} &\max_{\mathbf{Q}} \sum_{m=1}^M \sum_{k=0}^{K_m} [g_1(\mathbf{Q}) + g_3(\mathbf{Q}) - g_2(\mathbf{Q}^{(l-1)}) - g_4(\mathbf{Q}^{(l-1)}) - \\ &\quad (Q_{m,k} - Q_{m,k}^{(l-1)}) \times \sum_{j=1}^M \sum_{i=0}^{K_j} \nabla g_2(Q_{j,i}^{(l-1)}) + \nabla g_4(Q_{j,i}^{(l-1)})] \\ &\text{s.t. (22c), (22d).} \end{aligned} \quad (35)$$

Note that (35) can also be efficiently solved by available convex software packages [41].

Now we have solved the two sub-problems. We repeat them after each other until convergence. Then, for those users with negative rates, we set their rates to zero following the $[\cdot]^+$ operation. The specific procedure is summarized in Algorithm 1.

$$\max_{\mathbf{P}, \mathbf{Q}} \sum_{m=1}^M \sum_{k=1}^{K_m} R_{m,k}^{sec} - \lambda^{(l-1)} \left(\sum_{m=1}^M \sum_{k=1}^{K_m} P_{m,k} + \sum_{m=1}^M \sum_{k=0}^{K_m} Q_{m,k} + P_f \right), \text{ s.t. (22b) - (22d)}. \quad (36)$$

B. EE Maximization

It is clear that (24) belongs to a fractional problem, which can be transformed into a series of parametric subtractive-form subproblems as (36) on the top of the next page based on Dinkelbach algorithm [42].

Note in (36), $\lambda^{(l-1)}$ is a non-negative parameter. Starting from $\lambda^{(0)} = 0$, $\lambda^{(l)}$ can be updated by $\lambda^{(l)} = \frac{\sum_{m=1}^M \sum_{k=1}^{K_m} R_{m,k}^{sec(l)}}{\sum_{m=1}^M \sum_{k=1}^{K_m} P_{m,k}^{(l)} + \sum_{m=1}^M \sum_{k=0}^{K_m} Q_{m,k}^{(l)} + P_f}$, where $R_{m,k}^{sec(l)}$, $P_{m,k}^{(l)}$ and $Q_{m,k}^{(l)}$ are the updated rates and power after solving (36). As shown in [42], $\lambda^{(l)}$ keeps growing as l increases. When $\lambda^{(l)} - \lambda^{(l-1)}$ is smaller than a certain threshold, e.g., 10^{-3} , the iterations terminate, and the obtained $\lambda^{(l)}$ is the maximum EE of (24).

Then, the problem lies in how to solve (36) for a given λ . It is clear that (36) is similar to the sum rate maximization problem (22), except for the extra linear part in the objective function. Adding a linear part does not affect the way of solving the problem, and thus, we can apply the proposed sum rate maximization here directly. The specific procedure is summarized in Algorithm 2.

C. Complexity and Convergence

The proposed SE maximization algorithm includes inner and outer iterations. For the inner iteration, i.e., the DC programming, its convergence has been shown in [39], [40]. For the outer iteration, on one hand, the SE increases or remains unchanged for both the uplink and downlink PA; on the other hand, there exists an upper bound for the SE. Therefore, the outer iteration terminates within a limited number of iterations, i.e., the proposed SE maximization algorithm always converges.

The proposed EE maximization algorithm also includes inner and outer iterations. For the inner iteration, i.e., the SE maximization, its convergence has been shown above. For the outer iteration, i.e., the fractional programming, it always converges to the stationary and optimal solution [42]. Therefore, the proposed EE maximization algorithm always converges.

Now, we discuss the computational complexity of the proposed algorithms. First, we look at the proposed SE maximization algorithm. Denote the number of iterations for solving the

Algorithm 2: Energy-Efficient Power Allocation Algorithm

```

1 Initialize  $\varepsilon \leftarrow 10^{-6}$ ,  $\lambda \leftarrow 0$ , Initialize feasible power  $\mathbf{P}^{(0)}$ .
2 repeat{Outer iteration}
3   repeat{Inner iteration}
4      $\mathbf{P}^{(l)}, \mathbf{Q}^{(l)} \leftarrow$ 
5        $\max \sum_{m=1}^M \sum_{k=1}^{K_m} R_{m,k}^{sec} - \lambda(\sum_{m=1}^M \sum_{k=1}^{K_m} P_{m,k} + \sum_{m=1}^M \sum_{k=0}^{K_m} Q_{m,k} + P_f)$ 
6       s.t. (22b), (22c), (22d)
7     until  $\mathbf{P}^{(l)}, \mathbf{Q}^{(l)}$  converge;
8     Compute  $\varepsilon^* \leftarrow \sum_{m=1}^M \sum_{k=1}^{K_m} R_{m,k}^{sec(l)} - \lambda(\sum_{m=1}^M \sum_{k=1}^{K_m} P_{m,k}^{(l)} + \sum_{m=1}^M \sum_{k=0}^{K_m} Q_{m,k}^{(l)} + P_f)$ .
9     Update  $\lambda \leftarrow \frac{\sum_{m=1}^M \sum_{k=1}^{K_m} R_{m,k}^{sec(l)}}{\sum_{m=1}^M \sum_{k=1}^{K_m} P_{m,k}^{(l)} + \sum_{m=1}^M \sum_{k=0}^{K_m} Q_{m,k}^{(l)} + P_f}$ .
10  until  $\varepsilon^* \leq \varepsilon$ ;
11  if  $R_{m,k}^{sec(l)} < 0, \forall m \in \{1, \dots, M\}, k \in \{1, \dots, K_m\}$ 
12     $R_{m,k}^{sec(l)} \leftarrow 0$ ;
13  end
14   $\eta_{EE} \leftarrow \frac{\sum_{m=1}^M \sum_{k=1}^{K_m} R_{m,k}^{sec(l)}}{\sum_{m=1}^M \sum_{k=1}^{K_m} P_{m,k}^{(l)} + \sum_{m=1}^M \sum_{k=0}^{K_m} Q_{m,k}^{(l)} + P_f}$ .

```

uplink and downlink PA as I_1 and I_2 , respectively. The corresponding number of dual variables for solving (30) and (35) is denoted as D_1 and D_2 , respectively. Then, if the number of outer iteration is I_3 , the overall computational complexity of the proposed SE maximization algorithm is $O(I_3(I_1 D_1^2 + I_2 D_2^2))$. Next, we consider the proposed EE maximization problem. Denote its outer iteration as I_4 , then it can be easily shown that the overall computational complexity of the proposed EE maximization algorithm is $O(I_4 I_3 (I_1 D_1^2 + I_2 D_2^2))$.

VI. NUMERICAL RESULTS

In this section, we firstly investigate the behavior of the system without PA to highlight the effects of key parameters on the secrecy performance in subsection VI-A. The effectiveness of our proposed PA algorithms is then evaluated in subsection VI-B.

A. Fixed PA

Without loss of generality, we consider the following scenario. The total transmit power is allocated 80% for information transmission and 20% for AN. The effect of varying the AN

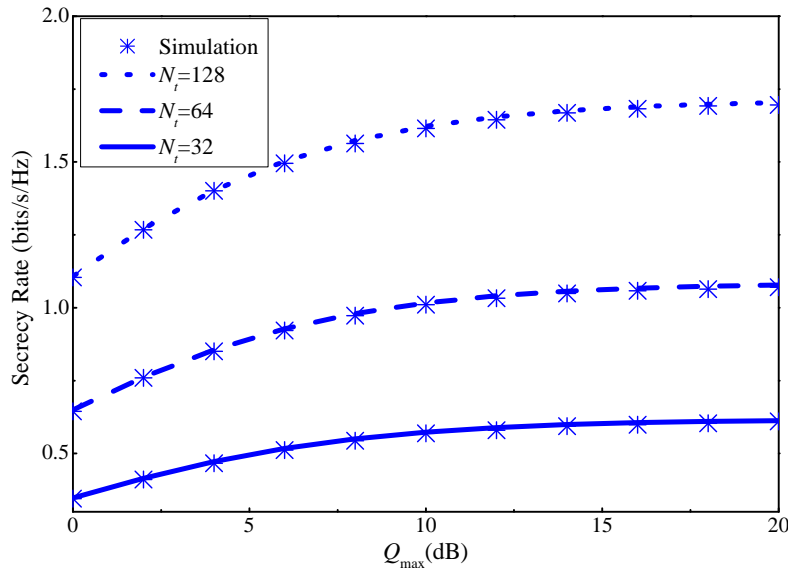


Fig. 2: Secrecy rate at the 2^{nd} user in the 5^{th} cluster versus the total transmit power at the BS, for different numbers of transmit antennas.

power allocation will be shown later in Fig. 3. The power is equally assigned to each user, and the AN power for each cluster is the same. $T = 300$ units and $\tau = M$ units. $\beta_{m,k}$ for each user is a random value between 0 and 100 and satisfies the condition $\beta_{m,1} \geq \dots \geq \beta_{m,K_m}$, while that for the illegitimate user is fixed to $\beta_E = 10$. Unless explicitly mentioned, this setup is kept throughout the section.

Without loss of generality, the ergodic secrecy rate of the 2^{nd} user in the 5^{th} cluster is selected to show in Fig. 2. The number of cluster is $M = 10$ and the number of users in a cluster is $K = 2$.⁴ It can be seen that the approximation in (11) and the simulation results match very well. Throughout the numerical results section, this approximation will be used. Besides, when the total transmit power at the BS increases, the secrecy rate at a user converges to a constant value. This is because of the interference and AN within the cluster and from other clusters. In addition, we can also observe that an increase in the number of antennas at the BS can lift the secrecy performance. The reason is that by increasing the number of antennas, the spatial transmitting beams become sharper, which leads to a decrease in inter-cluster interference and AN leakage, and an increase in the desired signal. The next figure will reveal how to take advantage of this

⁴The subscript m in K_m is dropped since the same number of users is considered in clusters.

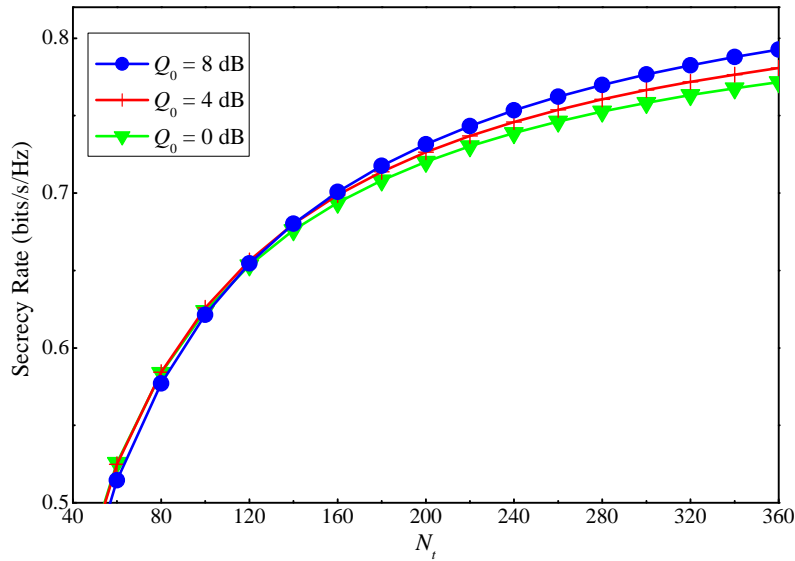


Fig. 3: Secrecy rate at the 2^{nd} user in the 5^{th} cluster versus the number of antennas at the BS, for different AN powers.

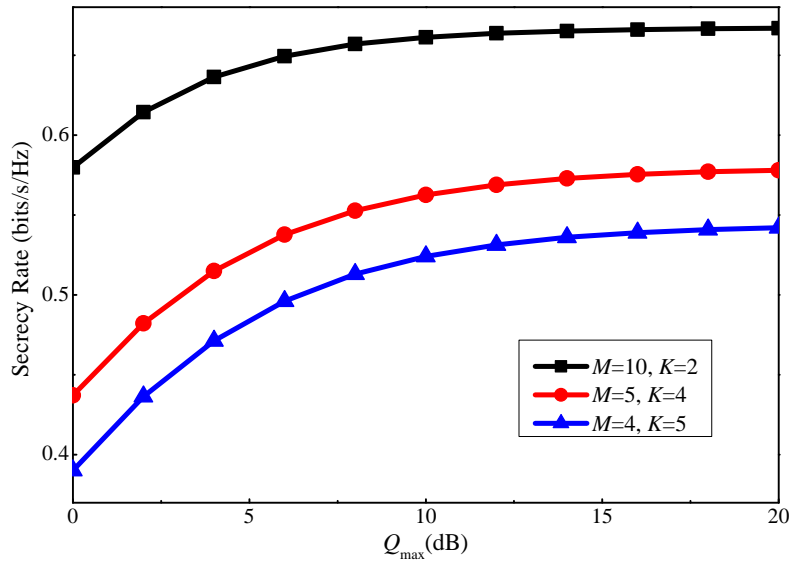


Fig. 4: Secrecy rate at the 2^{nd} user of the 2^{nd} cluster versus the total transmit power at the BS, for different clustering scenarios.

property to enhance secrecy performance.

In Fig. 3, we demonstrate the advantage of combining AN and massive MIMO technique in NOMA networks. In this setup, the transmit power assigned to each user is 10 dB, and the

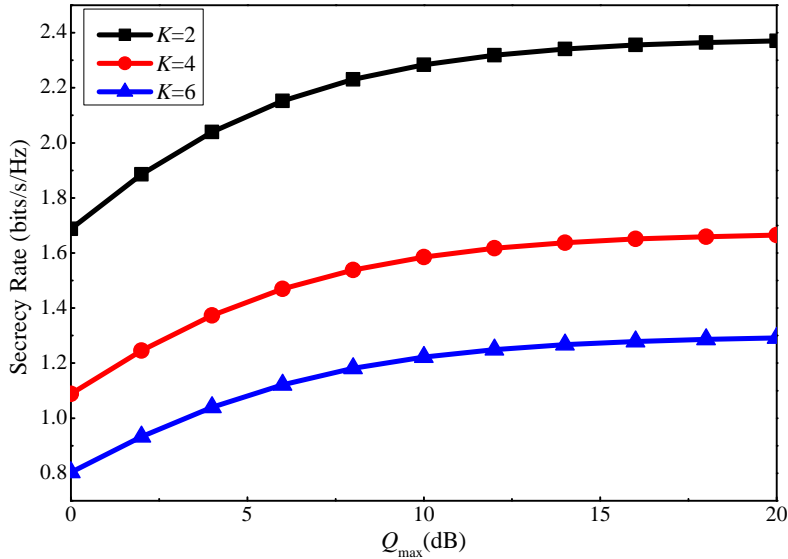


Fig. 5: The secrecy rate at the the 2^{nd} cluster versus the total transmit power at the BS, for a fixed number of clusters and different numbers of users.

AN power is varied as $\{0, 4, 8\}$ dB. The number of clusters is $M = 10$ with $K = 2$ users in a cluster. We can observe that when the number of transmit antennas is sufficiently large, the more the AN allocated power is, the better the secrecy performance at the 2^{nd} user is. The main reason is that for the legitimate side, the channel hardening property of massive MIMO technique helps reducing the AN leakage and the inter-cluster interference at each cluster. Meanwhile, the secrecy performance of the eavesdropper decreases when the AN power increases.

Figures 4 and 5 depict the effect of clustering on the secrecy performance. In Fig. 4, the total number of users is 20, which are clustered into three scenarios: $\{M = 10, K = 2\}$, $\{M = 5, K = 4\}$, and $\{M = 4, K = 5\}$. The total transmit power for each scenario is the same. The results show that the smaller the number of users in a cluster is, the better the secrecy performance at a user is. Meanwhile, in Fig. 5, the scenario of limited number of orthogonal sequences is shown. In this scenario, we assume that the number of available orthogonal sequences is 10, therefore, the number of clusters is $M = 10$. The number of users in a cluster is varied as $K = \{2, 4, 6\}$ to highlight its effect on the secrecy performance of a cluster. It is observed that although the transmit power for each user is identical and the AN power for each cluster is the same, the cluster with more users has smaller total secrecy rate than the ones with a smaller number of users. The reason is that when the number of users

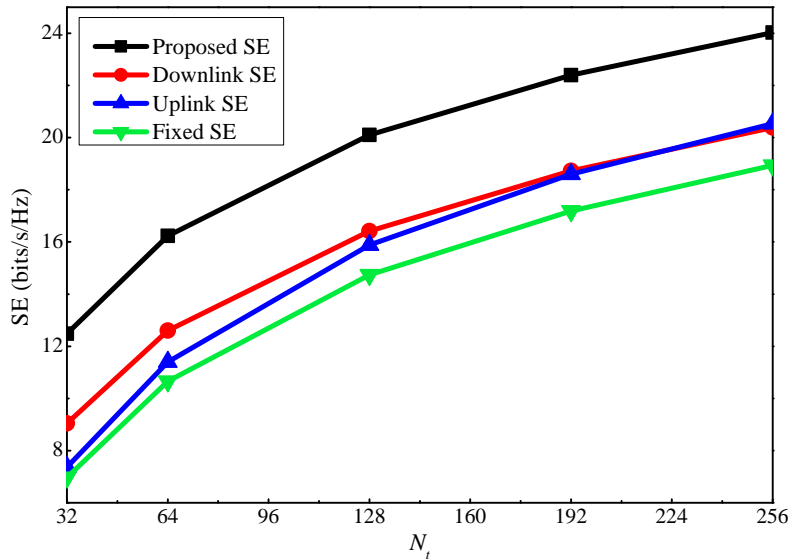


Fig. 6: SE comparison between the proposed algorithm and baseline algorithms.

in a cluster is small, the error of the uplink training process at this cluster is also small. As a consequence, the beam of the BS for this cluster is more precise, followed by a decrease in intra-cluster interference and AN leakage. This also reduces the imposed interference from this cluster to other clusters. In other words, for a better secrecy performance of each user and cluster, it is crucial to keep the number of users in a cluster small (minimum is two users for NOMA networks).

B. Optimized PA

In the following, we investigate the effectiveness of the proposed SE and EE maximization algorithms. We consider a scenario with four clusters, each with three users, i.e., $M = 4$ and $K = 3$. The simulation parameters are as follows: $Q_{\max} = 20$ dB, $P_{\max} = 0$ dB. $T = 300$ units. The large scale channel gain $\beta_{m,k}$ for each user is a random value between 0 and 100, while that for the illegitimate user is fixed to $\beta_E = 10$.

First, we investigate the effectiveness of the proposed SE maximization algorithm, referred to as Proposed SE. We compare it with three baseline algorithms, as follows: Downlink SE, which allocates the maximum uplink power to each user, and on this basis, performs PA for the downlink transmission as the Proposed SE. In contrast, the Uplink SE first allocates 80% of the total downlink power to the users equally, and 20% of the total downlink power to the AN

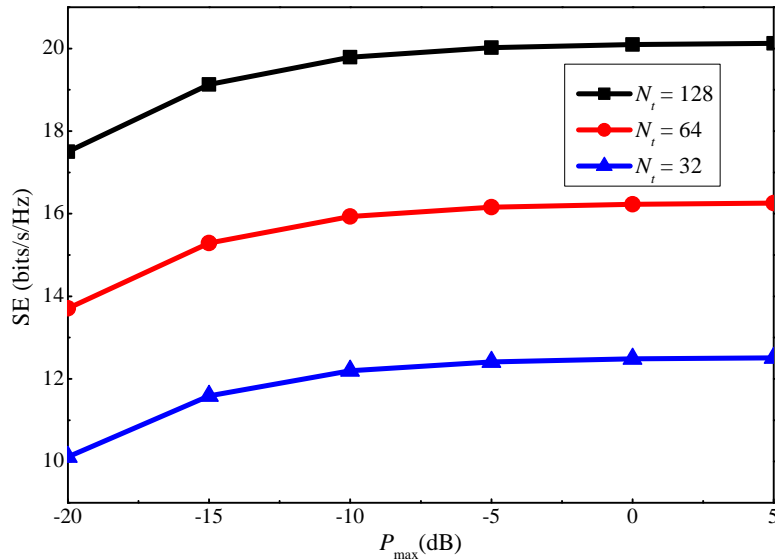


Fig. 7: SE versus the maximum uplink power, for different numbers of BS antennas.

equally. Then, uplink power is optimized as the Proposed SE. Fixed SE allocates the maximum uplink power for each user, equal downlink power allocation among the users, and the AN as in the above subsection. As shown in Fig. 6, the SE provided by all four algorithms grows with the number of transmit antennas. Moreover, among them, it can be seen that Proposed SE achieves the best performance, followed by Downlink SE, Uplink SE, and Fixed PA. This fully reveals the necessity of performing power optimization for the considered system. Furthermore, both uplink and downlink PA are required to achieve the best performance. Nonetheless, by comparing Downlink SE and Uplink SE, we can conclude that an appropriate allocation of the downlink power may play a larger role in the current setting.

To further show the effect of the uplink and downlink power on the achieved SE, Figs. 7 and 8 plot the SE versus the maximum uplink and downlink power, respectively. $N_t = \{32, 64, 128\}$ is respectively considered in each case. It is clear that the SE increases with both the maximum uplink and downlink powers. The former is because increasing the maximum uplink power leads to a more precise channel estimation result, which improves the beamforming sharpness and thus, the SE. The latter is because more power is available for data transmission. However, after a certain point, the increase becomes minor for both power values. This can be explained by the logarithmic relation between the power and user rate. Moreover, for the downlink power, increasing it also leads to a larger illegitimate rate and intra-cluster interference. Besides, by

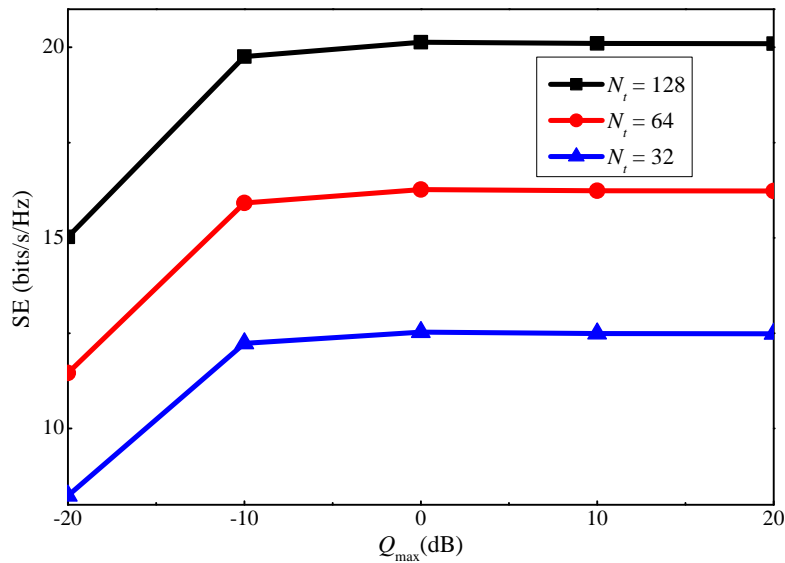


Fig. 8: SE versus the maximum downlink power, for different numbers of BS antennas.

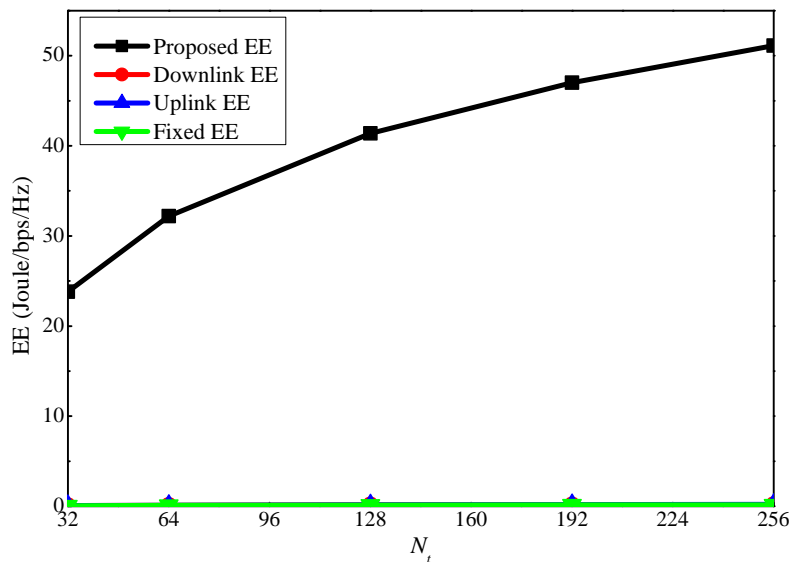


Fig. 9: EE comparison between the proposed algorithm and other baseline algorithms.

comparing the three antenna scenarios, we can conclude that increasing the number of antennas can significantly increase the SE.

Next, we investigate the proposed EE algorithm. Here $P_f = -5$ dB. We first compare the proposed EE maximization algorithm with the other three baseline algorithms when $Q_{\max} = 20$ dB and $P_{\max} = 0$ dB. According to Fig. 9, the EE for the other algorithms is quite small

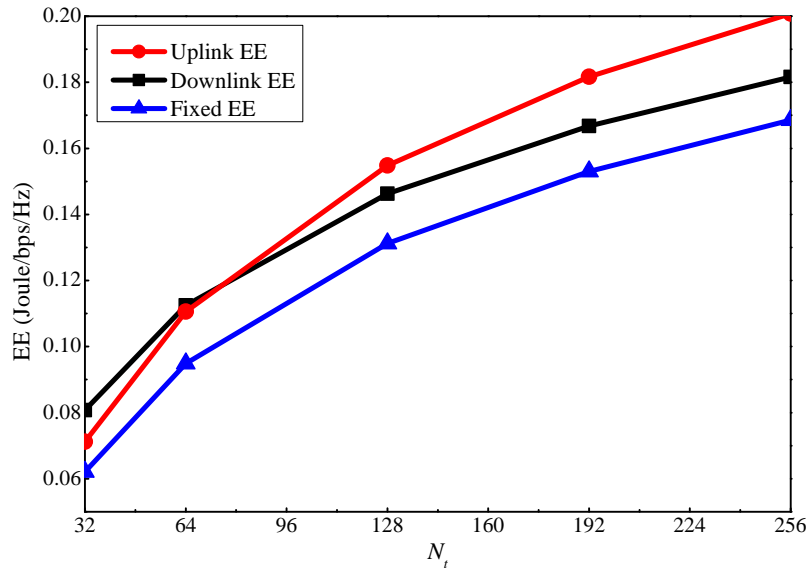


Fig. 10: EE for the three baseline algorithms.

compared with the proposed algorithm. This is because when $Q_{\max} = 20$ dB and $P_{\max} = 0$ dB, the power level is quite high, and thus, a large part of the available power is not used to maximize the EE. However, for the three baseline algorithms, at least one of the uplink and downlink power is fully consumed according to the setting. This leads to low EE. Figure 10 only shows these three algorithms, and it can be seen that all of them increase with the antenna number as the proposed algorithm.

Similar to the SE, we also show how the EE varies with the maximum uplink and downlink power in Figs. 11 and 12, respectively. For both cases, the EE first grows with the maximum power constraint, and after a certain threshold, i.e., $P_{\max} = -20$ dB and $Q_{\max} = -10$ dB, it remains unchanged even if the maximum power constraint continues to grow. This is because the slow increases in the SE cannot compensate for the power increment when the power is high, and thus, no more power will be consumed by the users to maximize the EE. By comparing the EE figures with the sum rate ones, i.e., Fig. 7 versus Fig. 11, and Fig. 8 versus Fig. 12, we can observe that the EE reaches the turning point at a smaller power value than the sum rate. This is because after the sum rate increment over the power declines to a certain value, no more extra power is used to maximize the EE.

The baseline massive MIMO-OMA can be considered as a special case of the proposed massive MIMO-NOMA scheme with just one user in each cluster. Accordingly, the legitimate achievable

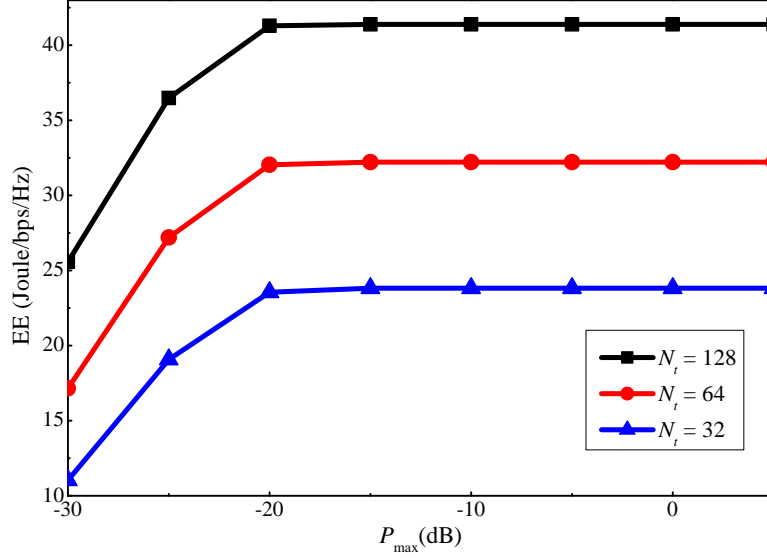


Fig. 11: EE versus the maximum uplink power, for different numbers of BS antennas.

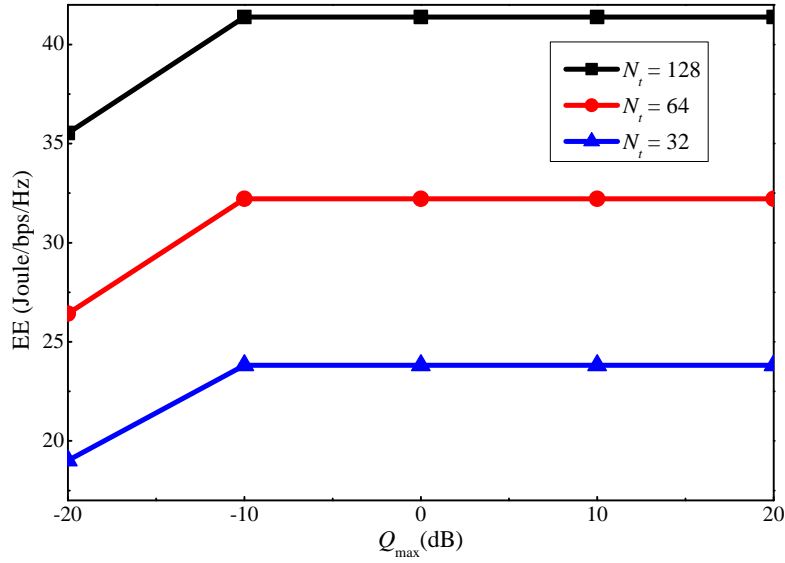


Fig. 12: EE versus the maximum downlink power, for different numbers of BS antennas.

rate of the m -th user is:

$$R_m^{OMA} = \left(1 - \frac{\tau}{T}\right) \log_2 \left(1 + \frac{\kappa_m}{\sum_{i=1}^3 I_{m,i} + 1}\right), \quad (37)$$

where $\kappa_m = Q_m \beta_m \rho_m N_t$, $I_{m,1} = Q_m \beta_m (1 - \rho_m)$, $I_{m,2} = \sum_{i \neq m}^M Q_i \beta_m$, $I_{m,3} = Q_{m,0} \beta_m (1 - \rho_m) +$

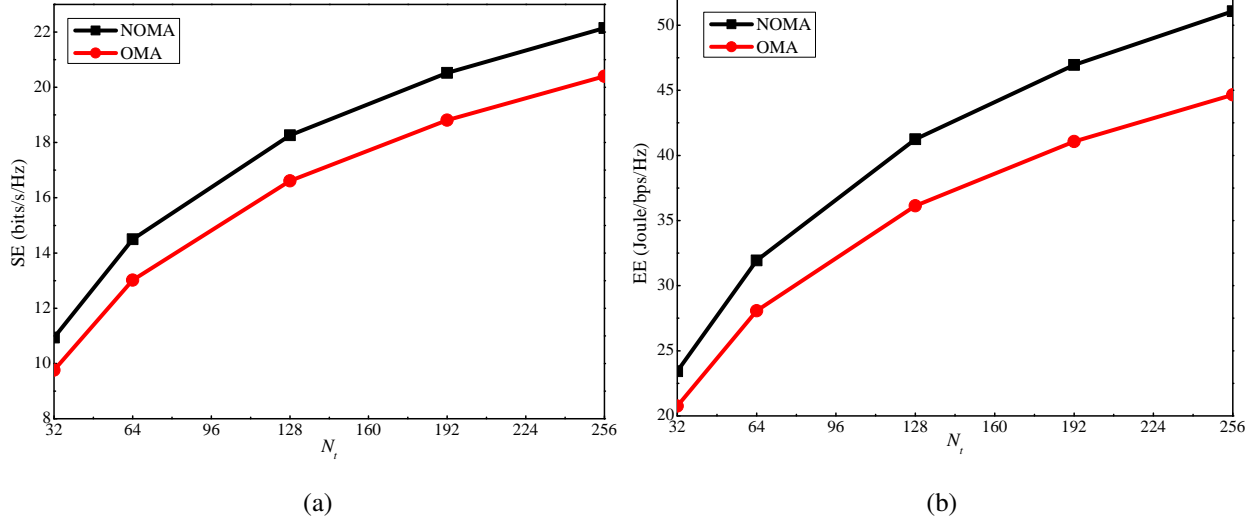


Fig. 13: Performance comparison for NOMA and OMA when the number of antenna varies:
(a) SE; (b) EE.

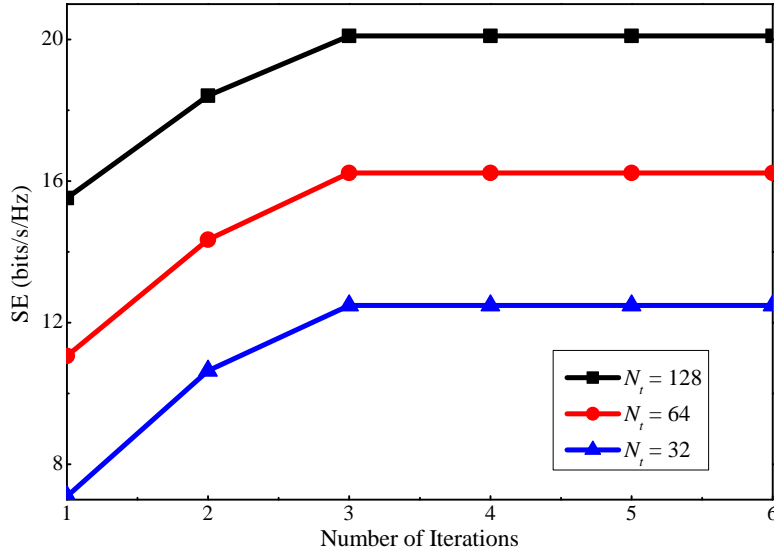


Fig. 14: Convergence of the proposed SE algorithm.

$\sum_{i \neq m}^M Q_{i,0} \beta_m$, Q_m is the downlink power for the m -th user, $Q_{m,0}$ is the AN power for the m -th user, β_m is the large scale fading of the m -th user, $\rho_m = \frac{P_m \beta_m \tau}{P_m \beta_m \tau + 1}$, τ is the length of training sequences that is the same as the NOMA case, and P_m is the uplink transmit power of the m -th

user. The achievable eavesdropping rate corresponding to the m -th user is:

$$R_{E,m}^{OMA} = \left(1 - \frac{\tau}{T}\right) \log_2 \left(1 + \frac{Q_m \beta_E}{\sum_{i \neq m}^M Q_i \beta_E + \sum_{i=1}^M Q_{i,0} \beta_E + 1}\right). \quad (38)$$

The achievable secrecy rate of the m -th user is

$$R_{S,m}^{OMA} = [R_m^{OMA} - R_{E,m}^{OMA}]^+. \quad (39)$$

In simulations, to compare the proposed massive MIMO-NOMA with the baseline massive MIMO-OMA, we consider a scenario with four clusters and two users in each cluster. TDMA is used for the baseline massive MIMO-OMA, and thus, each user in one cluster is only served half the time. Fig. 13 shows the corresponding SE and EE comparison between the considered schemes. It is clear that the proposed scheme outperforms the baseline massive MIMO-OMA when the number of antennas at the BS increases, which shows its superiority.

Finally, Figs. 14 and 15 show how many iterations are required for the proposed SE and EE maximization algorithms to converge, respectively. Note that here an iteration means solving either the uplink or the downlink DC programming problem, which requires to solve an average of five convex problems according to the simulation. Results for three different antenna numbers are presented when $Q_{\max} = 20$ dB and $P_{\max} = 0$ dB. It can be seen that a small number of iterations are required for the proposed SE and EE maximization algorithms to converge.

VII. CONCLUSION

In this paper, an AN-aided scheme has been proposed to ensure secrecy in massive MIMO-NOMA networks. The ergodic secrecy rate and its asymptotic value have been derived to spotlight the roles of key parameters on the secrecy performance of the considered system. The results have revealed that with a sufficiently large number of transmit antennas at the BS, only the illegitimate side is affected by the AN. In addition, when the transmit power at the BS is high, the secrecy performance of a user is independent of the inter-cluster interference and AN and is determined by the uplink training process, which depends on the number of users in a cluster, the uplink transmit power, and the large-scale fading. Besides, the results also suggest to keep the number of users in a cluster small for a better secrecy performance at each user and cluster. Furthermore, numerical results validate that our proposed optimization algorithms can obtain significant improvements over the baseline algorithms, i.e., Uplink PA, Downlink PA and

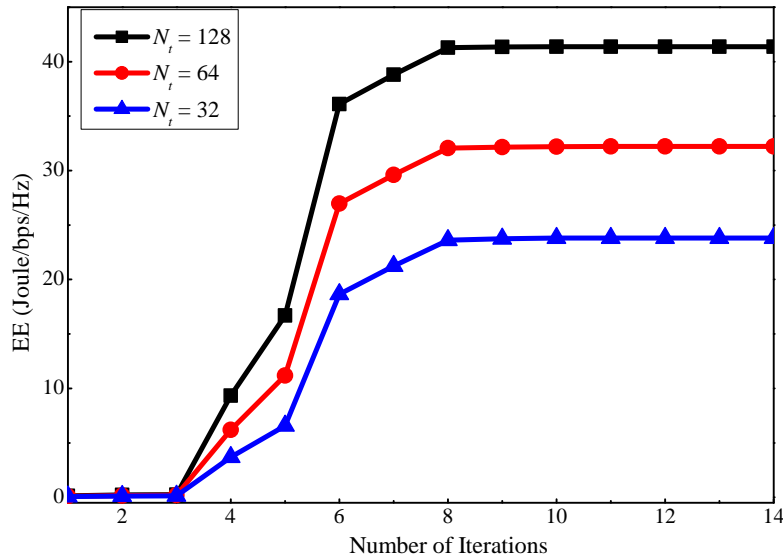


Fig. 15: Convergence of the proposed EE algorithm.

Fixed PA, in terms of the sum ergodic secrecy rate and energy efficiency. This fully reveals the necessity of performing power optimization for the considered system, and the effectiveness of the proposed algorithms. Finally, from the perspective of sum ergodic secrecy rate and its energy efficiency, our proposed system surpasses the conventional massive MIMO-OMA system.

REFERENCES

- [1] V. W. S. Wong et al., *Key Technologies for 5G Wireless Systems*. Cambridge, UK: Cambridge University Press, 2017.
- [2] S. M. R. Islam, M. Zeng, and O. A. Dobre, “NOMA in 5G systems: Exciting possibilities for enhancing spectral efficiency,” *IEEE 5G Tech. Focus*, vol. 1, no. 2, May 2017. [Online]. Available: 307 <http://5g.ieee.org/tech-focus>.
- [3] S. M. R. Islam, M. Zeng, O. A. Dobre, and K. Kwak, “Resource allocation for downlink noma systems: Key techniques and open issues,” *IEEE Wireless Commun. Mag.*, vol. 25, no. 2, pp. 40–47, April 2018.
- [4] S. M. R. Islam et al., “Power-domain non-orthogonal multiple access (NOMA) in 5G systems: Potentials and challenges,” *IEEE Commun. Surv. Tuts.*, vol. 19, no. 2, pp. 721–742, Second quarter 2017.
- [5] M. Zeng, G. I. Tsiropoulos, O. A. Dobre, and M. H. Ahmed, “Power allocation for cognitive radio networks employing non-orthogonal multiple access,” in *Proc IEEE Globecom*, Washington DC, USA, Dec. 2016.
- [6] Z. Wei, D. W. K. Ng, J. Yuan, and H. Wang, “Optimal resource allocation for power-efficient mc-noma with imperfect channel state information,” *IEEE Trans. Commun.*, vol. 65, no. 9, pp. 3944–3961, Sep. 2017.
- [7] M. Zeng, A. Yadav, O. A. Dobre, G. I. Tsiropoulos, and H. V. Poor, “Capacity comparison between MIMO-NOMA and MIMO-OMA with multiple users in a cluster,” *IEEE J. Sel. Areas Commun.*, vol. 35, no. 10, pp. 2413–2424, Oct. 2017.
- [8] —, “On the sum rate of MIMO-NOMA and MIMO-OMA systems,” *IEEE Wireless Commun. Lett.*, vol. 6, no. 4, pp. 534–537, Aug. 2017.

- [9] Z. Ding, F. Adachi, and H. V. Poor, "The application of MIMO to non-orthogonal multiple access," *IEEE Trans. Wireless Commun.*, vol. 15, no. 1, pp. 537–552, Jan. 2016.
- [10] M. Zeng, A. Yadav, O. A. Dobre, and H. V. Poor, "A fair individual rate comparison between MIMO-NOMA and MIMO-OMA," in *Proc IEEE Globecom Wkshps*, Singapore, Dec 2017, pp. 1–5.
- [11] N. Yang, L. Wang, G. Geraci, M. Elkashlan, J. Yuan, and M. D. Renzo, "Safeguarding 5G wireless communication networks using physical layer security," *IEEE Commun. Mag.*, vol. 53, no. 4, pp. 20–27, Apr. 2015.
- [12] A. D. Wyner, "The wire-tap channel," *Bell Syst. Tech. J.*, vol. 54, no. 8, pp. 1355–1387, Jan. 1975.
- [13] Y.-S. Shiu, S. Chang, H.-C. Wu, S. Huang, and H.-H. Chen, "Physical layer security in wireless networks: A tutorial," *IEEE Wireless Commun.*, vol. 18, no. 2, pp. 66–74, Apr. 2011.
- [14] J. Chen, L. Yang, and M.-S. Alouini, "Physical layer security for cooperative NOMA systems," *IEEE Trans. Veh. Technol.*, vol. 67, no. 5, pp. 4645–4649, May 2018.
- [15] Y. Liu, Z. Qin, M. Elkashlan, Y. Gao, and L. Hanzo, "Enhancing the physical layer security of non-orthogonal multiple access in large-scale networks," *IEEE Trans. Wireless Commun.*, vol. 16, no. 3, pp. 1656 – 1672, Mar. 2017.
- [16] Y. Zhang, H.-M. Wang, Q. Yang, and Z. Ding, "Secrecy sum rate maximization in non-orthogonal multiple access," *IEEE Commun. Lett.*, vol. 20, no. 5, pp. 930–933, May 2016.
- [17] B. He, A. Liu, N. Yang, and V. K. N. Lau, "On the design of secure non-orthogonal multiple access systems," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 10, pp. 2196–2206, Oct. 2017.
- [18] E. Larsson, O. Edfors, F. Tufvesson, and T. Marzetta, "Massive MIMO for next generation wireless systems," *IEEE Commun. Mag.*, vol. 52, no. 2, pp. 186–195, Feb. 2014.
- [19] T. L. Marzetta, "Noncooperative cellular wireless with unlimited numbers of base station antennas," *IEEE Trans. Commun.*, vol. 9, no. 11, pp. 3590–3600, Nov. 2010.
- [20] H. Q. Ngo, E. G. Larsson, and T. L. Marzetta, "Energy and spectral efficiency of very large multiuser MIMO systems," *IEEE Trans. Commun.*, vol. 61, no. 4, pp. 1436–1449, Apr. 2013.
- [21] W. Hao, M. Zeng, Z. Chu, S. Yang, and G. Sun, "Energy-efficient resource allocation for mmWave massive MIMO HetNets with wireless backhaul," *IEEE Access*, vol. 6, pp. 2457–2471, Feb. 2018.
- [22] W. Hao, M. Zeng, Z. Chu, and S. Yang, "Energy-efficient power allocation in millimeter wave massive MIMO with non-orthogonal multiple access," *IEEE Wireless Commun. Lett.*, vol. 6, no. 6, pp. 782–785, Dec. 2017.
- [23] M. Zeng, W. Hao, O. A. Dobre, and H. V. Poor, "Energy-efficient power allocation in uplink mmWave massive MIMO with NOMA," *IEEE Trans. Veh. Technol.*, pp. 1–1, 2019.
- [24] J. Ma, C. Liang, C. Xu, and L. Ping, "On orthogonal and superimposed pilot schemes in massive MIMO NOMA systems," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 12, pp. 2696 – 2707, Dec. 2017.
- [25] X. Chen, Z. Zhang, C. Zhong, D. W. K. Ng, and R. Jia, "Exploiting inter-user interference for secure massive non-orthogonal multiple access," *IEEE J. Sel. Areas Commun.*, vol. 36, no. 4, pp. 788–801, Apr. 2018.
- [26] J. Zhu, R. Schober, and V. K. Bhargava, "Linear precoding of data and artificial noise in secure massive MIMO systems," *IEEE Trans. Wireless Commun.*, vol. 15, no. 3, pp. 2245–2261, Mar. 2016.
- [27] N.-P. Nguyen, H. Q. Ngo, T. Q. Duong, H. D. Tuan, and K. Tourki, "Secure massive MIMO with the artificial noise-aided downlink training," *IEEE J. Sel. Areas Commun.*, vol. 36, no. 4, pp. 802 – 816, Apr. 2018.
- [28] Y.-Y. Zhang, J.-K. Zhang, and H.-Y. Yu, "Physically securing energy-based massive MIMO MAC via joint alignment of multi-user constellations and artificial noise," *IEEE J. Sel. Areas Commun.*, vol. 36, no. 4, pp. 829 – 844, Apr. 2018.
- [29] N.-P. Nguyen, H. Q. Ngo, T. Q. Duong, H. D. Tuan, and D. B. da Costa, "Full-duplex cyber-weapon with massive arrays," *IEEE Trans. Commun.*, vol. 65, no. 12, pp. 5544 – 5558, Aug. 2017.

- [30] I. S. Gradshteyn and I. M. Ryzhik, *Table of Integrals, Series, and Products*, 7th ed. San Diego, CA: Academic press, 2007.
- [31] H. Tabassum, E. Hossain, and J. Hossain, "Modeling and analysis of uplink non-orthogonal multiple access in large-scale cellular networks using poisson cluster processes," *IEEE Trans. Commun.*, vol. 65, no. 8, pp. 3555–3570, Aug 2017.
- [32] H. Sun, B. Xie, R. Q. Hu, and G. Wu, "Non-orthogonal multiple access with sic error propagation in downlink wireless mimo networks," in *Proc IEEE VTC*, Sep. 2016, pp. 1–5.
- [33] X. Chen, Z. Zhang, C. Zhong, R. Jia, and D. W. K. Ng, "Fully non-orthogonal communication for massive access," *IEEE Trans. Commun.*, vol. 66, no. 4, pp. 1717–1731, Apr. 2018.
- [34] M. Zeng, A. Yadav, O. A. Dobre, and H. V. Poor, "Energy-efficient power allocation for MIMO-NOMA with multiple users in a cluster," *IEEE Access*, vol. 6, pp. 5170–5181, 2018.
- [35] —, "Energy-efficient power allocation for hybrid multiple access systems," in *Proc IEEE ICC Wkshps*, Kansas City, MO, USA, May 2018, pp. 1–5.
- [36] W. Hao, Z. Chu, F. Zhou, S. Yang, G. Sun, and K. Wong, "Green communication for NOMA-based CRAN," *IEEE Internet of Things J.*, pp. 1–1, 2018.
- [37] M. Zeng, A. Yadav, O. A. Dobre, and H. V. Poor, "Energy-efficient power allocation for uplink NOMA," in *Proc IEEE Globecom*, Abu Dhabi, UAE, Dec. 2018, pp. 1–6.
- [38] —, "Energy-efficient joint User-RB association and power allocation for uplink hybrid NOMA-OMA," *IEEE Internet of Things J.*, pp. 1–1, 2019.
- [39] H. H. Kha, H. D. Tuan, and H. H. Nguyen, "Fast global optimal power allocation in wireless networks by local d.c. programming," *IEEE Trans. Wireless Commun.*, vol. 11, no. 2, pp. 510–515, Feb. 2012.
- [40] N. Vucic, S. Shi, and M. Schubert, "DC programming approach for resource allocation in wireless networks," in *Proc. Int. Symp. Modeling Optimization Mobile, Ad Hoc Wireless Netw.*, May 2010, pp. 380–386.
- [41] M. Grant and S. Boyd, "CVX: Matlab software for disciplined convex programming, version 1.21." Available: <http://cvxr.com/cvx>, Dec. 2010.
- [42] W. Dinkelbach, "On nonlinear fractional programming," *Manag. Sci.*, vol. 13, no. 7, pp. 3492–498, Mar. 1967.