



OTFS-NOMA: An Efficient Approach for Exploiting Heterogenous User Mobility Profiles

DOI:

[10.1109/TCOMM.2019.2932934](https://doi.org/10.1109/TCOMM.2019.2932934)

Document Version

Accepted author manuscript

[Link to publication record in Manchester Research Explorer](#)

Citation for published version (APA):

Ding, Z., Schober, R., Fan, P., & Poor, H. V. (2019). OTFS-NOMA: An Efficient Approach for Exploiting Heterogenous User Mobility Profiles. *I E E Transactions on Communications*.
<https://doi.org/10.1109/TCOMM.2019.2932934>

Published in:

I E E Transactions on Communications

Citing this paper

Please note that where the full-text provided on Manchester Research Explorer is the Author Accepted Manuscript or Proof version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version.

General rights

Copyright and moral rights for the publications made accessible in the Research Explorer are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Takedown policy

If you believe that this document breaches copyright please refer to the University of Manchester's Takedown Procedures [<http://man.ac.uk/04Y6Bo>] or contact openresearch@manchester.ac.uk providing relevant details, so we can investigate your claim.



AUTHOR QUERIES

AUTHOR PLEASE ANSWER ALL QUERIES

PLEASE NOTE: We cannot accept new source files as corrections for your article. If possible, please annotate the PDF proof we have sent you with your corrections and upload it via the Author Gateway. Alternatively, you may send us your corrections in list format. You may also upload revised graphics via the Author Gateway.

AQ:1 = Please supply index terms/keywords for your paper. To download the IEEE Taxonomy, go to http://www.ieee.org/documents/taxonomy_v101.pdf.

AQ:2 = Author: Please confirm or add details for any funding or financial support for the research of this article.

AQ:3 = Please provide the postal code for The University of Manchester, Friedrich-Alexander-University Erlangen-Nurnberg, and Southwest Jiaotong University.

AQ:4 = Note that if you require corrections/changes to tables or figures, you must supply the revised files, as these items are not edited for you.

AQ:5 = Please provide the department name for Ref. [37].

AQ:6 = Current affiliation in the biography of Zhiguo Ding does not match the First Footnote. Please check and correct where needed.

OTFS-NOMA: An Efficient Approach for Exploiting Heterogenous User Mobility Profiles

Zhiguo Ding¹, Senior Member, IEEE, Robert Schober, Fellow, IEEE,
Pingzhi Fan, Fellow, IEEE, and H. Vincent Poor², Fellow, IEEE

Abstract—This paper considers a challenging communication scenario, in which users have heterogenous mobility profiles, e.g., some users are moving at high speeds and some users are static. A new non-orthogonal multiple-access (NOMA) transmission protocol that incorporates orthogonal time frequency space (OTFS) modulation is proposed. Thereby, users with different mobility profiles are grouped together for the implementation of NOMA. The proposed OTFS-NOMA protocol is shown to be applicable to both uplink and downlink transmission, where sophisticated transmit and receive strategies are developed to remove inter-symbol interference and harvest both multi-path and multi-user diversity. Analytical results demonstrate that both the high-mobility and the low-mobility users benefit from the application of OTFS-NOMA. In particular, the use of NOMA allows the spreading of the high-mobility users' signals over a large amount of time-frequency resources, which enhances the OTFS resolution and improves the detection reliability. In addition, OTFS-NOMA ensures that low-mobility users have access to bandwidth resources which in conventional OTFS-orthogonal multiple access (OTFS-OMA) would be solely occupied by the high-mobility users. Thus, OTFS-NOMA improves the spectral efficiency and reduces latency.

Index Terms—XXXXX.

I. INTRODUCTION

NON-ORTHOGONAL multiple access (NOMA) has been recognized as a paradigm shift for the design of multiple access techniques for the next generation of wireless networks [1]–[4]. Many existing works on NOMA have

focused on scenarios with low-mobility users, where users with different channel conditions or quality of service (QoS) requirements are grouped together for the implementation of NOMA. For example, in power-domain NOMA, a base station serves two users simultaneously [5], [6]. In particular, the base station first orders the users according to their channel conditions, where the ‘weak user’ which has a poorer connection to the base station is generally allocated more transmission power and the other user, referred to as the ‘strong user’, is allocated less power. As such, the two users can be served in the same time-frequency resource, which improves the spectral efficiency compared to orthogonal multiple access (OMA). In the case that users have similar channel conditions, grouping users with different QoS requirements can facilitate the implementation of NOMA and effectively exploit the potential of NOMA [7]–[9]. Various existing studies have shown that the NOMA principle can be applied to different communication networks, such as millimeter-wave networks [10], [11], massive multiple-input multiple-output (MIMO) systems [12], [13], hybrid multiple access systems [14], [15], visible light communication networks [16], [17], and mobile edge computing [18]. We also note that various standardization efforts have been made to facilitate the implementation of NOMA in practical systems. For example, a study for the application of NOMA for downlink transmission, termed multi-user superposition transmission (MUST), was carried out for the 3rd Generation Partnership Project (3GPP) Release 14, where 15 different forms of MUST were proposed and compared [19]. After this study was completed, MUST was formally included in 3GPP Release 15 which is also referred to as Evolved Universal Terrestrial Radio Access (E-UTRA) [20]. A study for the application of NOMA for uplink transmission has been recently carried out for 3GPP Release 16, where more than 20 different forms of NOMA have been proposed by various companies [21].

This paper considers the application of NOMA to a challenging communication scenario, where users have heterogeneous mobility profiles. Different from the existing works in [22], [23], the use of orthogonal time frequency space (OTFS) modulation is considered in this paper because of its superior performance in scenarios with doubly-dispersive channels [24]–[26]. Recall that the key idea of OTFS is to use the delay-Doppler plane, where the users' signals are

Manuscript received April 4, 2019; revised June 5, 2019; accepted July 14, 2019. The work of Z. Ding was supported by the UK Engineering and Physical Sciences Research Council under grant number EP/P009719/2 and by H2020-MSCA-RISE-2015 under grant number 690750. The work of P. Fan was supported by the National Natural Science Foundation of China under grant number 61731017, and the 111 Project (No.111-2-14). The work of H. V. Poor was supported by the U.S. National Science Foundation under Grants CCF-093970 and CCF-1513915. The associate editor coordinating the review of this article and approving it for publication was V. Raghavan. (Corresponding author: Zhiguo Ding.)

Z. Ding is with the Department of Electrical Engineering, Princeton University, Princeton, NJ 08544 USA, and also with the School of Electrical and Electronic Engineering, The University of Manchester, Manchester, U.K. (e-mail: zhiguo.ding@manchester.ac.uk).

R. Schober is with the Institute for Digital Communications, Friedrich-Alexander-University Erlangen-Nurnberg (FAU), Erlangen, Germany (e-mail: robert.schober@fau.de).

P. Fan is with the Institute of Mobile Communications, Southwest Jiaotong University, Chengdu, China (e-mail: pingzhifan@foxmail.com).

H. V. Poor is with the Department of Electrical Engineering, Princeton University, Princeton, NJ 08544 USA (e-mail: poor@princeton.edu).

Digital Object Identifier 10.1109/TCOMM.2019.2932934

orthogonally placed. Compared to conventional modulation schemes, such as orthogonal frequency-division multiplexing (OFDM), OTFS offers the benefit that the time-invariant channel gains in the delay-Doppler plane can be utilized, which simplifies channel estimation and signal detection in high-mobility scenarios. The impact of pulse-shaping waveforms on the performance of OTFS was studied in [27], and the design of interference cancellation and iterative detection for OTFS was investigated in [28]. The diversity gain achieved by OTFS was studied in [29], and the application of OTFS to multiple access was proposed in [30]. In [31] and [32], the concept of OTFS was combined with MIMO, which revealed that the use of spatial degrees of freedom can further enhance the performance of OTFS.

This paper considers the application of OTFS to NOMA communication networks, where the coexistence of NOMA and OTFS is investigated. In particular, this paper makes the following contributions:

- 1) A spectrally efficient OTFS-NOMA transmission protocol is proposed by grouping users with different mobility profiles for the implementation of NOMA. On the one hand, users with high mobility are served in the delay-Doppler plane, and their signals are modulated by OTFS. On the other hand, users with low mobility are served in the time-frequency plane, and their signals are modulated in a manner similar to conventional OFDM.
- 2) The proposed new OTFS-NOMA protocol is applied to both uplink and downlink transmission, where different rate and power allocation policies are used to suppress multiple access interference. In addition, sophisticated equalization techniques, such as the frequency-domain zero-forcing linear equalizer (FD-LE) and the decision feedback equalizer (FD-DFE), are employed to remove the inter-symbol interference in the delay-Doppler plane. The impact of the developed equalization techniques on OTFS-NOMA is analyzed by using the outage probability as the performance criterion. Strategies to harvest multi-path diversity and multi-user diversity are also introduced, which can further improve the outage performance of OTFS-NOMA transmission.
- 3) The developed analytical results demonstrate that both the high-mobility and the low-mobility users benefit from the proposed OTFS-NOMA scheme. The use of NOMA allows the high-mobility users' signals to be spread over a large amount of time-frequency resources without degrading the spectral efficiency. As a result, the OTFS resolution, which determines whether the users' channels can be accurately located in the delay-Doppler plane, is enhanced significantly, and therefore, the reliability of detecting the high-mobility users' signals is improved. We note that, in OTFS-OMA, enhancing the OTFS resolution implies that a large amount of time and frequency resources are solely occupied by the high-mobility users, which reduces the overall spectral efficiency since the high-mobility users' channel conditions are typically weaker than those of the low-mobility users. In contrast, the use of OTFS-NOMA ensures that the low-mobility users can access the

bandwidth resources which would be solely occupied by the high-mobility users in the OMA mode. Hence, OTFS-NOMA improves spectral efficiency and reduces latency, as with OTFS-OMA the low-mobility users may have to wait for a long time before the scarce bandwidth resources occupied by the high-mobility users become available. In addition, we note that for the low-mobility users, using OFDM yields the same reception reliability as using OTFS, as pointed out in [33]. Therefore, the proposed OTFS-NOMA scheme, which serves the low-mobility users in the time-frequency plane and modulates the low-mobility users' signals in a manner similar to OFDM, offers the same reception reliability as OTFS-OMA, which serves the low-mobility users in the delay-Doppler plane and modulates the low-mobility users' signals by OTFS. However, OTFS-NOMA has the benefit of reduced system complexity because the use of the complicated OTFS transforms is avoided.

II. FOUNDATIONS OF OTFS-NOMA

A. Time-Frequency Plane and Delay-Doppler Plane

The key idea of OTFS-NOMA is to efficiently use both the time-frequency plane and the delay-Doppler plane. A discrete time-frequency plane is obtained by sampling at intervals of T s and Δf Hz as follows:

$$\Lambda_{\text{TF}} = \{(nT, m\Delta f), n=0, \dots, N-1, m=0, \dots, M-1\}, \quad (1)$$

and the corresponding discrete delay-Doppler plane is given by

$$\Lambda_{\text{DD}} = \left\{ \left(\frac{k}{NT}, \frac{l}{M\Delta f} \right), k=0, \dots, N-1, l=0, \dots, M-1 \right\}, \quad (2)$$

where N and M denote the total number of time intervals and the total number of frequency subchannels, respectively. The choices for T and Δf are determined by the channel characteristics, as will be explained in the following subsection.

B. Channel Model

This paper considers a multi-user communication network in which one base station communicates with $(K+1)$ users, denoted by U_i , $0 \leq i \leq K$. Denote U_i 's channel response in the delay-Doppler plane by $h_i(\tau, \nu)$, where τ denotes the delay and ν denotes the Doppler shift. OTFS uses the sparsity feature of a wireless channel in the delay-Doppler plane, i.e., there are a small number of propagation paths between a transmitter and a receiver [24], [25], [28], which means that $h_i(\tau, \nu)$ can be expressed as follows:

$$h_i(\tau, \nu) = \sum_{p=0}^{P_i} h_{i,p} \delta(\tau - \tau_{i,p}) \delta(\nu - \nu_{i,p}), \quad (3)$$

where $(P_i + 1)$ denotes the number of propagation paths, and $h_{i,p}$, $\tau_{i,p}$, and $\nu_{i,p}$ denote the complex Gaussian channel gain,¹ the delay, and the Doppler shift associated with the

¹The Gaussian assumption has been commonly used in the OTFS literature [26]–[29] since each channel gain (or each tap of the delay-Doppler impulse response) represents a cluster of reflectors with specific delay and Doppler characteristics.

178 p -th propagation path, respectively. We assume that the $h_{i,p}$,
 179 $0 \leq p \leq P_i$, are independent and identically distributed (i.i.d.)
 180 random variables,² i.e., $h_{i,p} \sim CN\left(0, \frac{1}{P_i+1}\right)$, which means
 181 $\sum_{p=0}^{P_i} \mathcal{E}\{|h_{i,p}|^2\} = 1$, where $\mathcal{E}\{\cdot\}$ denotes the expectation
 182 operation. The discrete delay and Doppler tap indices for the
 183 p -th path of $h_i(\tau, \nu)$, denoted by $l_{\tau_{i,p}}$ and $k_{\nu_{i,p}}$, respectively,
 184 are given by [28]

$$185 \quad \tau_{i,p} = \frac{l_{\tau_{i,p}} + \tilde{l}_{\tau_{i,p}}}{M\Delta f}, \quad \nu_{i,p} = \frac{k_{\nu_{i,p}} + \tilde{k}_{\nu_{i,p}}}{NT}, \quad (4)$$

186 where $\tilde{l}_{\tau_{i,p}}$ and $\tilde{k}_{\nu_{i,p}}$ denote the fractional delay and the
 187 fractional Doppler shift, respectively.

188 The construction of Λ_{TF} and Λ_{DD} needs to ensure that T
 189 is not smaller than the maximal delay spread, and Δf is not
 190 smaller than the largest Doppler shift, i.e., $T \geq \max\{\tau_{i,p}, 0 \leq$
 191 $p \leq P_i, 0 \leq i \leq K\}$ and $\Delta f \geq \max\{\nu_{i,p}, 0 \leq p \leq P_i, 0 \leq$
 192 $i \leq K\}$. In addition, the choices of N and M affect the
 193 OTFS resolution, which determines whether $h_i(\tau, \nu)$ can be
 194 accurately located in the discrete delay-Doppler plane. In par-
 195 ticular, M and N need to be sufficiently large to approximately
 196 achieve ideal OTFS resolution, which ensures that $\tilde{l}_{\tau_{i,p}} =$
 197 $\tilde{k}_{\nu_{i,p}} = 0$, such that the interference caused by fractional delay
 198 and Doppler shift is effectively suppressed [24].

199 C. General Principle of OTFS-NOMA

200 To facilitate the illustration of the general principle of
 201 OTFS-NOMA, we first briefly describe OTFS-OMA, the
 202 benchmark scheme used in this paper. In OTFS-OMA, there
 203 is no spectrum sharing between the high-mobility users and
 204 the low-mobility users, i.e., if OTFS is used to serve the high-
 205 mobility users, the NT time intervals and the $M\Delta f$ frequency
 206 subchannels are occupied by the high-mobility users and the
 207 low-mobility users cannot be served in these resource blocks.
 208 The general principle of the proposed OTFS-NOMA scheme is
 209 to exploit both the delay-Doppler plane and the time-frequency
 210 plane, where users with heterogenous mobility profiles are
 211 grouped together and served simultaneously. On the one hand,
 212 for the users with high mobility, their signals are placed in
 213 the delay-Doppler plane, which means that the time-invariant
 214 channel gains in the delay-Doppler plane can be exploited. It is
 215 worth pointing out that in order to ensure that the channels
 216 can be located in the delay-Doppler plane, both N and M
 217 need to be large, which is a disadvantage of OTFS-OMA,
 218 since a significant number of frequency channels (e.g., $M\Delta f$)
 219 are occupied for a long time (e.g., NT) by the high-mobility
 220 users whose channel conditions can be quite weak. The use of
 221 OTFS-NOMA facilitates spectrum sharing and hence ensures
 222 that the high-mobility users' signals can be spread over a
 223 large amount of time-frequency resources without degrading
 224 the spectral efficiency.

225 On the other hand, for the users with low mobility, their
 226 signals are placed in the time-frequency plane. The inter-
 227 ference between the users with different mobility profiles

²In order to simplify the performance analysis, we assume that the users' channels are i.i.d. In practice, it is likely that the high-mobility users' channel conditions are worse than the low-mobility users' channel conditions. This channel difference is beneficial for the implementation of NOMA, and hence can further increase the performance gain of OTFS-NOMA over OTFS-OMA.

228 is managed by using the principle of NOMA. As a result,
 229 compared to OTFS-OMA, OTFS-NOMA improves the overall
 230 spectral efficiency since it encourages spectrum sharing among
 231 users with different mobility profiles and avoids that the
 232 bandwidth resources are solely occupied by the high-mobility
 233 users which might have weak channel conditions. In addition,
 234 the complexity of detecting the low-mobility users' signals is
 235 reduced, compared to OTFS-OMA which serves all users in
 236 the delay-Doppler plane.

237 In this paper, we assume that, among $(K + 1)$ users,
 238 U_0 is a user with high mobility, and the remaining K users,
 239 U_i for $1 \leq i \leq K$, are low-mobility users, which are
 240 referred to as 'NOMA' users.³ For OTFS-OMA, we assume
 241 that U_0 solely occupies all NM resource blocks in Λ_{DD} .
 242 In OTFS-NOMA, U_i , for $1 \leq i \leq K$, are opportunistic
 243 NOMA users and their signals are placed in Λ_{TF} . The design
 244 of downlink OTFS-NOMA transmission will be discussed
 245 in detail in Sections III, IV, and V. The application of
 246 OTFS-NOMA for uplink transmission will be considered in
 247 Section VI only briefly, due to space limitations.

248 III. DOWNLINK OTFS-NOMA - SYSTEM MODEL

249 In this section, the OTFS-NOMA downlink transmission
 250 protocol is described. In particular, assume that the base
 251 station sends NM symbols to U_0 , denoted by $x_0[k, l]$, $k \in$
 252 $\{0, \dots, N-1\}$, $l \in \{0, \dots, M-1\}$. By using the inverse
 253 symplectic finite Fourier transform (ISFFT), the high-mobility
 254 user's symbols placed in the delay-Doppler plane are converted
 255 to NM symbols in the time-frequency plane as follows [24]:

$$256 \quad X_0[n, m] = \frac{1}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} x_0[k, l] e^{j2\pi\left(\frac{kn}{N} - \frac{ml}{M}\right)}, \quad (5)$$

257 where $n \in \{0, \dots, N-1\}$ and $m \in \{0, \dots, M-1\}$. We note
 258 that the NM time-frequency signals can be viewed as N
 259 OFDM symbols containing M signals each. We assume that a
 260 rectangular window is applied to the transmitted and received
 261 signals.

262 The NOMA users' signals are placed directly in the time-
 263 frequency plane, and are superimposed with the high-mobility
 264 user's signals, $X_0[n, m]$. With NM orthogonal resource
 265 blocks available in the time-frequency plane, there are different
 266 ways for the K users to share the resource blocks. For
 267 illustration purposes, we assume that M users are selected
 268 from the K opportunistic NOMA users,⁴ where each NOMA

³We note that the principle of OTFS-NOMA can be extended to the case where multiple high-mobility users are served in the delay-Doppler plane. In this case, the NM signals in the delay-Doppler plane belong to different high-mobility users and OTFS is used as a type of multiple access technique [24], [30]. For downlink transmission, this change has no impact on the proposed detection schemes and the analytical results developed in this paper. For uplink transmission, the results developed in this paper are applicable to the case with multiple high-mobility users if the adaptive-rate transmission scheme proposed in Section VI is employed.

⁴The same M users can be scheduled as long as the users' channels do not change in the delay-Doppler plane. Otherwise, a new set of M users may be selected from the K opportunistic users. We also note that the number of the opportunistic users is assumed to be larger than the number of the frequency subchannels ($K \geq M$), which can be justified by a spectrum crunch scenario, i.e., there are not sufficient bandwidth resources available to support a large number of mobile devices.

user is to occupy one frequency subchannel and receive N information bearing symbols, denoted by $x_i(n)$, for $1 \leq i \leq M$ and $0 \leq n \leq N - 1$. The criterion employed for user scheduling and its impact on the performance of OTFS-NOMA will be discussed in Section V. Denote the time-frequency signals to be sent to U_i by $X_i[n, m]$, $1 \leq i \leq M$. The following mapping scheme is used in this paper⁵:

$$X_i[n, m] = \begin{cases} x_i(n) & \text{if } m = i - 1 \\ 0 & \text{otherwise,} \end{cases} \quad (6)$$

for $1 \leq i \leq M$ and $0 \leq n \leq N - 1$.

The base station superimposes U_0 's time-frequency signals with the NOMA users' signals as follows:

$$X[n, m] = \frac{\gamma_0}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} x_0[k, l] e^{j2\pi\left(\frac{kn}{N} - \frac{ml}{M}\right)} + \sum_{i=1}^M \gamma_i X_i[n, m], \quad (7)$$

where γ_i denotes the NOMA power allocation coefficient of user i , and $\sum_{i=0}^M \gamma_i^2 = 1$.

The transmitted signal at the base station is obtained by applying the Heisenberg transform to $X[n, m]$. By assuming perfect orthogonality between the transmit and receive pulses, the received signal at U_i in the time-frequency plane can be modelled as follows [24], [25], [28]:

$$Y_i[n, m] = H_i[n, m]X[n, m] + W_i[n, m], \quad (8)$$

where $W_i(n, m)$ is the white Gaussian noise in the time-frequency plane, and $H_i(n, m) = \iint h_i(\tau, \nu) e^{j2\pi\nu n T} e^{-j2\pi(\nu + m\Delta f)\tau} d\tau d\nu$.

IV. DOWNLINK OTFS-NOMA - DETECTING THE HIGH-MOBILITY USER'S SIGNALS

For the proposed downlink OTFS-NOMA scheme, U_0 directly detects its signals in the delay-Doppler plane by treating the NOMA users' signals as noise. In particular, in order to detect U_0 's signals, the symplectic finite Fourier transform (SFFT) is applied to $Y_0[n, m]$ to obtain the delay-Doppler estimates as follows:

$$\begin{aligned} y_0[k, l] &= \frac{1}{NM} \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} Y_0[n, m] e^{-j2\pi\left(\frac{nk}{N} - \frac{ml}{M}\right)} \\ &= \frac{1}{NM} \sum_{q=0}^M \gamma_q \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} x_q[n, m] h_{w,0}\left(\frac{k-n}{NT}, \frac{l-m}{M\Delta f}\right) \\ &\quad + z_0[k, l], \end{aligned} \quad (9)$$

where q denotes the user index, $z_0[k, l]$ is complex Gaussian noise, $x_q[k, l]$, $1 \leq q \leq M$, denotes the delay-Doppler representation of $X_q[n, m]$ and can be obtained by applying the SFFT to $X_q[n, m]$, the channel $h_{w,0}(\nu', \tau')$ is given by

$$h_{w,0}(\nu', \tau') = \iint h_i(\tau, \nu) w(\nu' - \nu, \tau' - \tau) e^{-j2\pi\nu\tau} d\tau d\nu, \quad (10)$$

⁵We note that mapping schemes different from (6) can also be used. For example, if N users are scheduled and each user is to occupy one time slot and receives an OFDM-like symbol containing M signals, we can set $X_i[n, m] = x_i(m)$, for $n = i - 1$.

and $w(\nu, \tau) = \sum_{c=0}^{N-1} \sum_{d=0}^{M-1} e^{-j2\pi(\nu c T - \tau d \Delta f)}$. To simplify the analysis, the power of the complex-Gaussian distributed noise is assumed to be normalized, i.e., $z_i[k, l] \sim CN(0, 1)$, where $CN(a, b)$ denotes a complex Gaussian distributed random variable with mean a and variance b .

By applying the channel model in (3), the relationship between the transmitted signals and the observations in the delay-Doppler plane can be expressed as follows [24], [25], [28]:

$$y_0[k, l] = \sum_{q=0}^M \gamma_q \sum_{p=0}^{P_0} h_{0,p} x_q[(k - k_{\nu_0,p})_N, (l - l_{\tau_0,p})_M] + z_0[k, l], \quad (11)$$

where $(\cdot)_N$ denotes the modulo N operator. As in [29]–[31], we assume that N and M are sufficiently large to ensure that both $\tilde{k}_{\nu_0,p}$ and $\tilde{l}_{\tau_0,p}$ are zero, i.e., there is no interference caused by fractional delay or fractional Doppler shift. We note that for OTFS-OMA, increasing N and M can significantly reduce spectral efficiency, whereas the use of large N and M becomes possible for OTFS-NOMA because of the spectrum sharing of users with different mobility profiles.

Define $\mathbf{y}_{0,k} = [y_0[k, 0] \cdots y_0[k, M-1]]^T$ and $\mathbf{y}_0 = [\mathbf{y}_{0,0}^T \cdots \mathbf{y}_{0,N-1}^T]^T$. Similarly, the signal vector \mathbf{x}_i and the noise vector \mathbf{z}_0 are constructed from $x_i[k, l]$ and $z_0[k, l]$, respectively. Based on (11), the system model can be expressed in matrix form as follows:

$$\mathbf{y}_0 = \gamma_0 \mathbf{H}_0 \mathbf{x}_0 + \underbrace{\sum_{q=1}^M \gamma_q \mathbf{H}_0 \mathbf{x}_q}_{\text{Interference and noise terms}} + \mathbf{z}_0, \quad (12)$$

where \mathbf{H}_0 is a block-circulant matrix and defined as follows:

$$\mathbf{H}_0 = \begin{bmatrix} \mathbf{A}_{0,0} & \mathbf{A}_{0,N-1} & \cdots & \mathbf{A}_{0,2} & \mathbf{A}_{0,1} \\ \mathbf{A}_{0,1} & \mathbf{A}_{0,0} & \ddots & \mathbf{A}_{0,3} & \mathbf{A}_{0,2} \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \mathbf{A}_{0,N-2} & \mathbf{A}_{0,N-3} & \ddots & \mathbf{A}_{0,0} & \mathbf{A}_{0,N-1} \\ \mathbf{A}_{0,N-1} & \mathbf{A}_{0,N-2} & \ddots & \mathbf{A}_{0,1} & \mathbf{A}_{0,0} \end{bmatrix}, \quad (13)$$

and each submatrix $\mathbf{A}_{0,n}$ is an $M \times M$ circulant matrix whose structure is determined by (11).

Example: Consider a special case with $N = 4$ and $M = 3$, and U_0 's channel is given by

$$h_0(\tau, \nu) = h_{0,0} \delta(\tau) \delta(\nu) + h_{0,1} \delta\left(\tau - \frac{1}{M\Delta f}\right) \delta\left(\nu - \frac{3}{NT}\right), \quad (14)$$

which means $k_0 = 0$, $k_1 = 3$, $l_0 = 0$, $l_1 = 1$. Therefore, the block-circulant matrix is given by

$$\mathbf{H}_0 = \begin{bmatrix} \mathbf{A}_{0,0} & \mathbf{A}_{0,3} & \mathbf{A}_{0,2} & \mathbf{A}_{0,1} \\ \mathbf{A}_{0,1} & \mathbf{A}_{0,0} & \mathbf{A}_{0,3} & \mathbf{A}_{0,2} \\ \mathbf{A}_{0,2} & \mathbf{A}_{0,1} & \mathbf{A}_{0,0} & \mathbf{A}_{0,3} \\ \mathbf{A}_{0,3} & \mathbf{A}_{0,2} & \mathbf{A}_{0,1} & \mathbf{A}_{0,0} \end{bmatrix}, \quad (15)$$

where $\mathbf{A}_{0,0} = h_{0,0}\mathbf{I}_3$, $\mathbf{A}_{0,1} = \mathbf{A}_{0,2} = \mathbf{0}_{3 \times 3}$ and $\mathbf{A}_{0,3} =$

$$\begin{bmatrix} 0 & 0 & h_{0,1} \\ h_{0,1} & 0 & 0 \\ 0 & h_{0,1} & 0 \end{bmatrix}.$$

Remark 1: It is well known that an $n \times n$ circulant matrix can be diagonalized by the $n \times n$ discrete Fourier transform (DFT) and inverse DFT matrices, denoted by \mathbf{F}_n and \mathbf{F}_n^{-1} , respectively, i.e., the columns of the DFT matrix are the eigenvectors of the circulant matrix. We note that directly applying the DFT factorization to \mathbf{H}_0 is not possible, since \mathbf{H}_0 is not a circulant matrix, but a block circulant matrix.

Because of the structure of \mathbf{H}_0 , inter-symbol interference still exists in the considered OTFS-NOMA system, and equalization is needed. We consider two equalization approaches, FD-LE and FD-DFE, which were both originally developed for single-carrier transmission with cyclic prefix [34], [35].

A. Design and Performance of FD-LE

The proposed FD-LE consists of two steps. Let \otimes denote the Kronecker product. The first step is to multiply the observation vector \mathbf{y}_0 by $\mathbf{F}_N \otimes \mathbf{F}_M^H$, which leads to the result in the following proposition.

Proposition 1: By applying the detection matrix $\mathbf{F}_N \otimes \mathbf{F}_M^H$ to observation vector \mathbf{y}_0 , the received signals for OTFS-NOMA downlink transmission can be written as follows:

$$\tilde{\mathbf{y}}_0 = \mathbf{D}_0(\mathbf{F}_N \otimes \mathbf{F}_M^H) \left(\gamma_0 \mathbf{x}_0 + \sum_{q=1}^M \gamma_q \mathbf{x}_q \right) + \tilde{\mathbf{z}}_0, \quad (16)$$

where $\tilde{\mathbf{y}}_0 = (\mathbf{F}_N \otimes \mathbf{F}_M^H) \mathbf{y}_0$, $\tilde{\mathbf{z}}_0 = (\mathbf{F}_N \otimes \mathbf{F}_M^H) \mathbf{z}_0$, \mathbf{D}_0 is a diagonal matrix whose $(kM + l + 1)$ -th main diagonal element is given by

$$D_0^{k,l} = \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} a_{0,n}^{m,1} e^{j2\pi \frac{lm}{M}} e^{-j2\pi \frac{kn}{N}}, \quad (17)$$

for $0 \leq k \leq N-1$, $0 \leq l \leq M-1$, and $a_{0,n}^{m,1}$ is the element located in the $(nM + m + 1)$ -th row and the first column of \mathbf{H}_0 .

Proof: Please refer to Appendix A. \square

With the simplified signal model shown in (16), the second step of FD-LE is to apply $(\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-1} \mathbf{D}_0^{-1}$ to $\tilde{\mathbf{y}}_0$. Thus, \mathbf{U}_0 's received signal is given by

$$\check{\mathbf{y}}_0 = \gamma_0 \mathbf{x}_0 + \underbrace{\sum_{q=1}^M \gamma_q \mathbf{x}_q + (\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-1} \mathbf{D}_0^{-1} \tilde{\mathbf{z}}_0}_{\text{Interference and noise terms}}, \quad (18)$$

where $\check{\mathbf{y}}_0 = (\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-1} \mathbf{D}_0^{-1} \tilde{\mathbf{y}}_0$. To simplify the analysis, we assume that the powers of all users' information-bearing signals are identical, which means that the transmit signal-to-noise ratio (SNR) can be defined as $\rho = \mathcal{E}\{|x_0[k, l]|^2\} = \mathcal{E}\{|x_i(n)|^2\}$, since the noise power is assumed to be normalized.⁶ The following lemma provides the signal-to-interference-plus-noise ratio (SINR) achieved by FD-LE.

⁶Following steps in the proof for Proposition 1 to show the statistical property of $\tilde{\mathbf{z}}_0$ in (66), we can also show that $W_i[n, m] \sim \mathcal{CN}(0, 1)$ if $z_i[k, l] \sim \mathcal{CN}(0, 1)$.

Lemma 1: Assume that $\gamma_i = \gamma_1$, for $1 \leq i \leq N$. By using FD-LE, the SINRs for detecting all $x_0[k, l]$, $0 \leq k \leq N-1$ and $0 \leq l \leq M-1$, are identical and given by

$$\text{SINR}_{0,kl}^{\text{LE}} = \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \frac{1}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} |D_0^{\tilde{k}, \tilde{l}}|^{-2}}. \quad (19)$$

Proof: Please refer to Appendix B. \square

Remark 2: The proof of Lemma 1 shows that $\sum_{i=0}^M \gamma_i^2 = 1$ can be simplified as $\gamma_0^2 + \gamma_i^2 = 1$ for $1 \leq i \leq M$, which is the motivation for assuming $\gamma_i = \gamma_1$. Following steps similar to those in the proofs for Proposition 1 and Lemma 1, one can show that directly applying \mathbf{H}_0^{-1} to the observation vector yields the same SINR. However, the proposed FD-LE scheme can be implemented more efficiently since $(\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-1} = \mathbf{F}_N^H \otimes \mathbf{F}_M$ and \mathbf{D}_0 is a diagonal matrix. Hence, the inversion of a full $NM \times NM$ matrix is avoided.

In this paper, the outage probability and the outage rate are used as performance criteria, since the outage probability can provide a tight bound on the probability of erroneous detection and is general in the sense that it does not depend on particular channel coding and modulation schemes used [36]. The outage probability achieved by FD-LE is given by $P(\log(1 + \text{SINR}_{0,kl}^{\text{LE}}) < R_0)$, where R_i , $0 \leq i \leq M$, denotes \mathbf{U}_i 's target data rate. It is difficult to analyze the outage probability for the following two reasons. First, the $D_0^{k,l}$, $k \in \{0, \dots, N-1\}$, $l \in \{0, \dots, M-1\}$, are not statistically independent, and second, the distribution of a sum of the inverse of exponentially distributed random variables is difficult to characterize. The following lemma provides an asymptotic result for the outage probability based on the SINR provided in Lemma 1.

Lemma 2: If $\gamma_0^2 > \gamma_1^2 \epsilon_0$, the diversity order achieved by FD-LE is one, where $\epsilon_0 = 2^{R_0} - 1$. Otherwise, the outage probability is always one.

Proof: Please refer to Appendix C. \square

Remark 3: Recall that the diversity order achieved by OTFS-OMA, where the high-mobility user, \mathbf{U}_0 , solely occupies the bandwidth resources, is also one. Therefore, the use of OTFS-NOMA ensures that the additional M low-mobility users are served without compromising \mathbf{U}_0 's diversity order, which improves the spectral efficiency compared to OTFS-OMA.

B. Design and Performance of FD-DFE

Different from FD-LE, which is a linear equalizer, FD-DFE is based on the idea of feeding back previously detected symbols. Since both \mathbf{x}_0 and \mathbf{x}_q , $q \geq 1$, experience the same fading channel, we first define $\mathbf{x} = \gamma_0 \mathbf{x}_0 + \sum_{q=1}^M \gamma_q \mathbf{x}_q$, which are the signals to be recovered by FD-DFE. Given the received signal vector shown in (12), the outputs of FD-DFE are given by

$$\hat{\mathbf{x}} = \mathbf{P}_0 \mathbf{y}_0 - \mathbf{G}_0 \check{\mathbf{x}}, \quad (20)$$

where $\check{\mathbf{x}}$ contains the decisions made on the symbols \mathbf{x} , \mathbf{P}_0 is the feed-forward part of the equalizer, and \mathbf{G}_0 is the feedback part of the equalizer. Similar to [34], [35], we use the following choices for \mathbf{P}_0 and \mathbf{G}_0 : $\mathbf{P}_0 = \mathbf{L}_0 (\mathbf{H}_0^H \mathbf{H}_0)^{-1} \mathbf{H}_0^H$, $\mathbf{G}_0 = \mathbf{L}_0 - \mathbf{I}_{NM}$, where \mathbf{L}_0 is a lower triangular matrix

with its main diagonal elements being ones in order to ensure causality of the feedback signals. With the above choices for \mathbf{P}_0 and \mathbf{G}_0 , U_0 's signals can be detected as follows:

$$\hat{\mathbf{x}} = \mathbf{L}_0(\mathbf{H}_0^H \mathbf{H}_0)^{-1} \mathbf{H}_0^H \mathbf{y}_0 - (\mathbf{L}_0 - \mathbf{I}_{NM}) \hat{\mathbf{x}}. \quad (21)$$

For FD-DFE, \mathbf{L}_0 is obtained from the Cholesky decomposition of \mathbf{H}_0 , i.e., $\mathbf{H}_0^H \mathbf{H}_0 = \mathbf{L}_0^H \mathbf{\Lambda}_0 \mathbf{L}_0$, where \mathbf{L}_0 is the desirable lower triangular matrix, and $\mathbf{\Lambda}_0$ is a diagonal matrix. Therefore, the estimates of \mathbf{x}_0 can be rewritten as follows:

$$\hat{\mathbf{x}} = \mathbf{x} + \mathbf{L}_0(\mathbf{H}_0^H \mathbf{H}_0)^{-1} \mathbf{H}_0^H \mathbf{z}_0 \quad (22)$$

$$= \gamma_0 \mathbf{x}_0 + \underbrace{\sum_{q=1}^M \gamma_q \mathbf{x}_q + \mathbf{L}_0(\mathbf{H}_0^H \mathbf{H}_0)^{-1} \mathbf{H}_0^H \mathbf{z}_0}_{\text{Interference and noise terms}}, \quad (23)$$

where perfect decision-making is assumed, i.e., $\check{\mathbf{x}} = \mathbf{x}$, and there is no error propagation [35], [37], [38]. We note that (23) yields an upper bound on the reception reliability of FD-DFE when error propagation cannot be completely avoided.

Following steps similar to those in the proof of Lemma 1, the covariance matrix for the interference-plus-noise term can be found as follows:

$$\mathbf{C}_{\text{cov}} = \rho \gamma_1^2 \mathbf{I}_{MN} + \mathbf{L}_0(\mathbf{H}_0^H \mathbf{H}_0)^{-1} \mathbf{L}_0^H = \rho \gamma_1^2 \mathbf{I}_{MN} + \mathbf{\Lambda}_0^{-1}, \quad (24)$$

where the last step follows from the fact that \mathbf{L}_0 is obtained from the Cholesky decomposition of \mathbf{H}_0 . Therefore, the SINR for detecting $x_0[k, l]$ can be expressed as follows:

$$\text{SINR}_{0,kl} = \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \lambda_{0,kl}^{-1}}, \quad (25)$$

where $\lambda_{0,kl}$ is the $(kM+l+1)$ -th element on the main diagonal of $\mathbf{\Lambda}_0$.

Remark 4: We note that there is a fundamental difference between the two equalization schemes. One can observe from (19) that the SINRs achieved by FD-LE for different $x_0[k, l]$ are identical. However, for FD-DFE, different symbols experience different effective fading gains, $\lambda_{0,kl}$. Therefore, FD-DFE can realize unequal error protection for data streams with different priorities. This comes at the price of a higher computational complexity.

We further note that the use of FD-DFE also ensures that multi-path diversity can be harvested, as shown in the following. The outage performance analysis for FD-DFE requires knowledge of the distribution of the effective channel gains, $\lambda_{0,kl}$. Because of the implicit relationship between $\mathbf{\Lambda}_0$ and \mathbf{H}_0 , a general expression for the outage probability achieved by FD-DFE is difficult to obtain. However, analytical results can be developed for special cases to show that the use of FD-DFE can realize the maximal multi-path diversity.

In particular, the SINR for $x_0[N-1, M-1]$ is a function of $\lambda_{0,(N-1)(M-1)}$ which is the last element on the main diagonal of $\mathbf{\Lambda}_0$. Recall that $\mathbf{\Lambda}_0$ is obtained via Cholesky decomposition, i.e., $\mathbf{H}_0^H \mathbf{H}_0 = \mathbf{L}_0^H \mathbf{\Lambda}_0 \mathbf{L}_0$. Because \mathbf{L}_0 is a lower triangular matrix, $\lambda_{0,(N-1)(M-1)}$ is equal to the element of $\mathbf{H}_0^H \mathbf{H}_0$ located in the NM -th column and the NM -th row,

which means

$$\lambda_{0,(N-1)(M-1)} = \sum_{p=0}^{P_0} |h_{0,p}|^2. \quad (26)$$

Since the channel gains are i.i.d. and follow $h_{0,p} \sim \mathcal{CN}(0, \frac{1}{P_0+1})$, the probability density function (pdf) of $\sqrt{P_0+1} \lambda_{0,(N-1)(M-1)}$ is given by

$$f(x) = \frac{1}{P_0!} e^{-x} x^{P_0}. \quad (27)$$

By using the above pdf, the outage probability and the diversity order can be obtained by some algebraic manipulations, as shown in the following corollary.

Corollary 1: Assume $\gamma_0^2 > \gamma_1^2 \epsilon_0$. The use of FD-DFE realizes the following outage probability for detection of $x_0[N-1, M-1]$:

$$P_{N-1, M-1}^0 = \frac{1}{P_0!} g\left(P_0+1, \frac{\epsilon_0(P_0+1)}{\rho(\gamma_0^2 - \gamma_1^2 \epsilon_0)}\right), \quad (28)$$

where $g(\cdot)$ denotes the incomplete Gamma function. The full multi-path diversity order, P_0+1 , is achievable for $x_0[N-1, M-1]$.

Remark 5: The results in Corollary 1 can be extended to OTFS-OMA with FD-DFE straightforwardly. We also note that diversity gains larger than one are not achievable with FD-LE as shown in Lemma 2, which is one of the disadvantages of FD-LE compared to FD-DFE.

Remark 6: We note that not all NM data streams can benefit from the full diversity gain. The simulation results provided in Section VII (Fig. 2) show that the diversity orders achievable for $x_0[k, l]$, $k < N-1$ and $l < M-1$, are smaller than that for $x_0[N-1, M-1]$, and the diversity order for $x_0[0, 0]$ is one, i.e., the same value as for FD-LE. We further note that the diversity result in Corollary 1 is obtained by assuming that there is no error propagation, i.e., it is assumed that when detecting the i -th element of \mathbf{x} in (21), the first $(i-1)$ elements of \mathbf{x} have already been correctly detected. Because of this assumption, the diversity gain developed in Corollary 1 is an upper bound on the diversity gain achieved by FD-DFE. If the assumption does not hold, the diversity orders for $x_0[k, l]$ will be capped by the worst case, i.e., the diversity gain for $x_0[0, 0]$ which is one.

Remark 7: FD-DFE entails a higher implementation complexity than FD-LE, as explained in the following. The complexity of FD-LE is mainly caused by computing the inversion of $\mathbf{H}_0^H \mathbf{H}_0$. However, for FD-DFE, \mathbf{L}_0 needs to be computed, in addition to $(\mathbf{H}_0^H \mathbf{H}_0)^{-1}$, as shown in (21). Recall that \mathbf{L}_0 is obtained from the Cholesky decomposition of the $NM \times NM$ matrix \mathbf{H}_0 , which entails a computational complexity of $\mathcal{O}(N^3 M^3)$. Therefore, the computational complexity of FD-DFE is higher than that of FD-LE, but FD-DFE offers a performance gain in terms of reception reliability compared to FD-LE, as shown in Section VII.

V. DOWNLINK OTFS-NOMA - DETECTING THE NOMA USERS' SIGNALS

Successive interference cancellation (SIC) will be carried out by the NOMA users, where each NOMA user first decodes

540 the high mobility user's signal in the delay-Doppler plane
541 and then decodes its own signal in the time-frequency plane.
542 The two stages of SIC are discussed in the following two
543 subsections, respectively.

544 A. Stage I of SIC

545 Following steps similar to the ones in the previous section,
546 each NOMA user also observes the mixture of the $(M + 1)$
547 users' signals in the delay-Doppler plane as follows:

$$548 \quad \mathbf{y}_i = \gamma_0 \mathbf{H}_i \mathbf{x}_0 + \underbrace{\sum_{q=1}^M \gamma_q \mathbf{H}_i \mathbf{x}_q}_{\text{Interference and noise terms}} + \mathbf{z}_i, \quad (29)$$

549 where \mathbf{H}_i and \mathbf{z}_i are defined similar to \mathbf{H}_0 and \mathbf{z}_0 , respectively.

550 We assume that the low-mobility NOMA users do not
551 experience Doppler shift, and therefore, their channels can be
552 simplified as follows:

$$553 \quad h_i(\tau) = \sum_{p=0}^{P_i} h_{i,p} \delta(\tau - \tau_{i,p}), \quad (30)$$

554 for $1 \leq i \leq K$, which means that each NOMA user's
555 channel matrix, \mathbf{H}_i , $1 \leq i \leq N$, is a block-diagonal matrix,
556 i.e., $\mathbf{A}_{i,0}$ is a non-zero circulant matrix and $\mathbf{A}_{i,n} = \mathbf{0}_{M \times M}$,
557 for $1 \leq n \leq N - 1$. Therefore, each NOMA user can
558 divide its observation vector into N equal-length sub-vectors,
559 i.e., $\mathbf{y}_i = [\mathbf{y}_{i,0}^T \cdots \mathbf{y}_{i,N-1}^T]^T$, which yields the following
560 simplified system model:

$$561 \quad \mathbf{y}_{i,n} = \gamma_0 \mathbf{A}_{i,0} \mathbf{x}_{0,n} + \sum_{q=1}^M \gamma_q \mathbf{A}_{i,0} \mathbf{x}_{q,n} + \mathbf{z}_{i,n}, \quad (31)$$

562 where, similar to $\mathbf{y}_{i,n}$, $\mathbf{x}_{i,n}$ and $\mathbf{z}_{i,n}$ are obtained from \mathbf{x}_i
563 and \mathbf{z}_i , respectively. Therefore, unlike the high-mobility user,
564 the NOMA users can perform their signal detection based on
565 reduced-size observation vectors, which reduces the computa-
566 tional complexity.

567 Since $\mathbf{A}_{i,0}$ is a circulant matrix, the two equalization
568 approaches used in the previous section are still applicable.
569 First, we consider the use of FD-LE. Following the same steps
570 as in the proof for Proposition 1, in the first step of FD-LE, the
571 DFT matrix is applied to the reduced-size observation vector,
572 which yields the following:

$$573 \quad \tilde{\mathbf{y}}_{i,n} = \tilde{\mathbf{D}}_i \mathbf{F}_M^H \left(\gamma_0 \mathbf{x}_{0,n} + \sum_{q=1}^M \gamma_q \mathbf{x}_{q,n} \right) + \tilde{\mathbf{z}}_{i,n}, \quad (32)$$

574 where $\tilde{\mathbf{y}}_{i,n} = \mathbf{F}_M^H \mathbf{y}_{i,n}$ and $\tilde{\mathbf{z}}_{i,n} = \mathbf{F}_M^H \mathbf{z}_{i,n}$. Compared to \mathbf{D}_i
575 in Proposition 1 which is an $NM \times NM$ matrix, $\tilde{\mathbf{D}}_i$ is an
576 $M \times M$ diagonal matrix, and its $(l + 1)$ -th main diagonal
577 element is given by $\tilde{D}_i^l = \sum_{m=0}^{M-1} a_{i,0}^{m,1} e^{j2\pi \frac{lm}{M}}$, for $0 \leq l \leq$
578 $M - 1$, where $a_{i,0}^{m,1}$ is the element located in the $(m + 1)$ -th
579 row and the first column of $\mathbf{A}_{i,0}$. Unlike conventional OFDM,
580 which uses \mathbf{F}_M at the receiver, \mathbf{F}_M^H is used here. Because
581 $\mathbf{F}_M^H \mathbf{A}_{i,0} \mathbf{F}_M = [\mathbf{F}_M \mathbf{A}_{i,0}^* \mathbf{F}_M^H]^*$, the sign of the exponent of
582 the exponential component of \tilde{D}_i^l is different from that in the
583 conventional case.

584 In the second step of FD-LE, $\mathbf{F}_M \tilde{\mathbf{D}}_i^{-1}$ is applied to $\tilde{\mathbf{y}}_{i,n}$.
585 Following steps similar to the ones in the proof for Lemma 1,
586 the SINR for detecting $x_0[k, l]$ can be obtained as follows:

$$587 \quad \text{SINR}_{0,kl}^{i,\text{LE}} = \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \frac{1}{M} \sum_{l=0}^{M-1} |\tilde{D}_i^l|^{-2}}. \quad (33)$$

588 We note that $\text{SINR}_{0,k_1 l}^{i,\text{LE}} = \text{SINR}_{0,k_2 l}^{i,\text{LE}}$, for $k_1 \neq k_2$, due to the
589 time invariant nature of the channels.

590 If FD-DFE is used, the corresponding SINR for detecting
591 $x_0[k, l]$ is given by

$$592 \quad \text{SINR}_{0,kl}^{i,\text{DFE}} = \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \tilde{\lambda}_{0,l}^{-1}}, \quad (34)$$

593 where $\tilde{\lambda}_{0,l}$ is obtained from the Cholesky decomposition
594 of $\mathbf{A}_{i,0}$. The details for the derivation of (34) are omitted
595 here due to space limitations.

596 B. Stage II of SIC

597 Assume that \mathbf{U}_0 's NM signals can be decoded and removed
598 successfully, which means that, in the time-frequency plane,
599 the NOMA users observe the following:

$$600 \quad Y_i[n, m] = \sum_{q=1}^M \gamma_q H_i[n, m] X_q[n, m] + W_i[n, m] \\ 601 \quad = \gamma_1 H_i[n, m] x_{m+1}(n) + W_i[n, m], \quad (35)$$

602 where the last step follows from the mapping scheme used
603 in (6) and it is assumed that all NOMA users employ the
604 same power allocation coefficient. We note that \mathbf{U}_i is only
605 interested in $Y_i[n, i - 1]$, $0 \leq n \leq N - 1$. Therefore, \mathbf{U}_i 's n -th
606 information bearing signal, $x_i(n)$, can be detected by applying
607 a one-tap equalizer as follows:

$$608 \quad \hat{x}_i(n) = \frac{Y_i[n, i - 1]}{\gamma_1 H_i[n, i - 1]}, \quad (36)$$

609 which means that the SNR for detecting $x_i(n)$ is given by

$$610 \quad \text{SNR}_{i,n} = \rho \gamma_1^2 |\tilde{D}_i^{i-1}|^2, \quad (37)$$

611 since $W_i[n, i - 1]$ is white Gaussian noise and $H_i[n, i - 1] =$
612 \tilde{D}_i^{i-1} . We note that $\text{SNR}_{i,n_1} = \text{SNR}_{i,n_2}$, for $n_1 \neq n_2$, which
613 is due to the time-invariant nature of the channel.

614 Without loss of generality, assume that the same target data
615 rate R_i is used for $x_i(n)$, $0 \leq n \leq N - 1$. Therefore, the outage
616 probability for $x_i(n)$ is given by

$$617 \quad \text{P}_{i,n}^{\text{LE}} \\ 618 \quad = 1 - \text{P} \left(\text{SNR}_{i,n} > \epsilon_i, \text{SINR}_{0,kl}^{i,\text{LE}} > \epsilon_0, \forall l \right) \\ 619 \quad = 1 - \text{P} \left(\rho \gamma_1^2 |\tilde{D}_i^{i-1}|^2 > \epsilon_i, \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \frac{1}{M} \sum_{l=0}^{M-1} |\tilde{D}_i^l|^{-2}} > \epsilon_0 \right), \quad (38)$$

620 if FD-LE is used in the first stage of SIC. If FD-DFE is used
621 in the first stage of SIC, the outage probability for $x_i(n)$ is
622 given by

$$623 \quad \text{P}_{i,n}^{\text{DFE}} = 1 - \text{P} \left(\text{SNR}_{i,n} > \epsilon_i, \text{SINR}_{0,kl}^{i,\text{DFE}} > \epsilon_0, \forall l \right) \\ 624 \quad = 1 - \text{P} \left(\rho \gamma_1^2 |\tilde{D}_i^{i-1}|^2 > \epsilon_i, \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \tilde{\lambda}_{0,l}^{-1}} > \epsilon_0, \forall l \right), \quad (39)$$

where $\epsilon_i = 2^{R_i} - 1$. Again because of the correlation between the random variables $|\tilde{D}_i^l|^{-2}$ and $\tilde{\lambda}_{0,l}$, the exact expressions for the outage probabilities are difficult to obtain. Alternatively, the achievable diversity order is analyzed in the following subsections.

1) *Random User Scheduling*: If the M users are randomly selected from the K available users, which means that each $|\tilde{D}_i^l|^2$ is complex Gaussian distributed. For the FD-LE case, the outage probability, $P_{i,n}^{\text{LE}}$, can be upper bounded as follows:

$$P_{i,n}^{\text{LE}} \leq 1 - \text{P} \left(\rho\gamma_1^2 |\tilde{D}_i^{\min}|^2 > \epsilon_i, \frac{\rho\gamma_0^2}{\rho\gamma_1^2 + |\tilde{D}_i^{\min}|^{-2}} > \epsilon_0 \right), \quad (40)$$

where $|\tilde{D}_i^{\min}|^2 = \min\{|\tilde{D}_i^m|^2, 0 \leq m \leq M-1\}$. The upper bound on the outage probability in (40) can be rewritten as follows:

$$P_{i,n}^{\text{LE}} \leq 1 - \text{P} \left(|\tilde{D}_i^{\min}|^2 > \bar{\epsilon} \right), \quad (41)$$

where $\bar{\epsilon} = \max \left\{ \frac{\epsilon_0}{\rho(\gamma_0^2 - \gamma_1^2 \epsilon_0)}, \frac{\epsilon_i}{\rho\gamma_1^2} \right\}$. As a result, an upper bound on the outage probability can be obtained as follows:

$$P_{i,n}^{\text{LE}} \leq \text{P} \left(|\tilde{D}_i^{\min}|^2 < \bar{\epsilon} \right) \leq MP \left(|\tilde{D}_i^0|^2 < \bar{\epsilon} \right) \doteq \frac{1}{\rho}, \quad (42)$$

where $\text{P}^o \doteq \rho^{-d}$ denotes exponential equality, i.e., $d = -\lim_{\rho \rightarrow \infty} \frac{\log \text{P}^o}{\log \rho}$ [36]. Therefore, the following corollary can be obtained.

Corollary 2: For random user scheduling and FD-LE, a diversity order of 1 is achievable at the NOMA users.

Our simulation results in Section VII show that a diversity order of 1 is also achievable for FD-DFE, although we do not have a formal proof for this conclusion, yet.

2) *Realizing Multi-User Diversity*: The diversity order of OTFS-NOMA can be improved by carrying out opportunistic user scheduling, which yields multi-user diversity gains. For illustration purpose, we propose a greedy user scheduling policy, where a single NOMA user is scheduled to transmit in all resource blocks of the time-frequency plane. From the analysis of the random scheduling case we deduce that $|\tilde{D}_i^{\min}|^2$ is critical to the outage performance. Therefore, the scheduled NOMA user, denoted by U_{i^*} , is selected based on the following criterion:

$$i^* = \arg \max_{i \in \{1, \dots, K\}} \left\{ |\tilde{D}_i^{\min}|^2 \right\}. \quad (43)$$

By using the assumption that the users' channel gains are independent and following steps similar to the ones in the proof for Lemma 2, the following corollary can be obtained in a straightforward manner.

Corollary 3: For FD-LE, the user scheduling strategy shown in (43) realizes the maximal multi-user diversity gain, K .

Remark 8: The reason why a multi-user diversity gain of K can be realized by the proposed scheduling strategy is explained in the following. Recall that the SINR for FD-LE to detect $x_0[k, l]$ is $\text{SINR}_{0,kl}^{i,\text{LE}} = \frac{\rho\gamma_0^2}{\rho\gamma_1^2 + \frac{1}{M} \sum_{i=0}^{M-1} |\tilde{D}_i^l|^{-2}}$. If this SINR is too small, the first stage of SIC will fail and an

outage event will occur. To improve the SINR, it is important to ensure that for a scheduled user, its weakest channel gain, $|\tilde{D}_i^{\min}|^2 = \min\{|\tilde{D}_i^m|^2, 0 \leq m \leq M-1\}$, is not too small. The used scheduling strategy shown in (43) is essentially a max-min strategy and ensures that the user with the strongest $|\tilde{D}_i^{\min}|^2$ is selected from the K candidates, which effectively exploits multi-user diversity.

We note that the user scheduling strategy shown in (43) is also useful for improving the performance of FD-DFE, as shown in Section VII.

VI. UPLINK OTFS-NOMA TRANSMISSION

The design of uplink OTFS-NOMA is similar to that of downlink OTFS-NOMA, and due to space limitations, we mainly focus on the difference between the two cases in this section. Again, we assume that U_0 is grouped with M NOMA users, selected from the K available users. U_0 's NM signals are placed in the delay-Doppler plane, and are denoted by $x_0[k, l]$, where $0 \leq k \leq N-1$ and $0 \leq l \leq M-1$. The corresponding time-frequency signals, $X_0[n, m]$, are obtained by applying ISFFT to $x_0[k, l]$. On the other hand, the NOMA users' signals, $x_i(n)$, are mapped to time-frequency signals, $X_i[n, m]$, according to (6).

Following steps similar to the ones for the downlink case, the base station's observations in the time-frequency plane are given by

$$\begin{aligned} Y[n, m] &= \sum_{q=0}^M H_q[n, m] X_q[n, m] + W[n, m] \\ &= \frac{H_0(n, m)}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} x_0[k, l] e^{j2\pi(\frac{kn}{N} - \frac{ml}{M})} \\ &\quad + \sum_{q=1}^M H_q[n, m] X_q[n, m] + W[n, m], \end{aligned} \quad (44)$$

where $W[n, m]$ is the Gaussian noise at the base station in the time-frequency plane. We assume that all users employ the same transmit pulse as well as the same transmit power. The base station applies SIC to first detect the NOMA users' signals in the time-frequency plane, and then tries to detect the high-mobility user's signals in the delay-Doppler plane, as shown in the following two subsections.

A. Stage I of SIC

The base station will first try to detect the NOMA users' signals in the time-frequency plane by treating the signals from U_0 as noise, which is the first stage of SIC.

By using (6), $x_i(n)$ can be estimated as follows:

$$\begin{aligned} \hat{x}_i(n) &= \frac{Y[n, i-1]}{H_i[n, i-1]} \\ &= x_i[n] + \frac{H_0[n, i-1] X_0[n, i-1] + W[n, i-1]}{H_i[n, i-1]}. \end{aligned} \quad (45)$$

Define an $NM \times 1$ vector, $\bar{\mathbf{x}}_0$, whose $(nM + m + 1)$ -th element is $X_0[n, m]$. Recall that $X_0[n, m]$ is obtained from the ISFFT of $x_0[k, l]$, i.e.,

$$\bar{\mathbf{x}}_0 = (\mathbf{F}_N^H \otimes \mathbf{F}_M) \mathbf{x}_0, \quad (46)$$

721 which means $X_0[n, m]$ follows the same distribution as
 722 $x_0[k, l]$. By applying steps similar to those in the proof for
 723 Lemma 1, the SINR for detecting $x_i(n)$ is given by

$$724 \quad \text{SINR}_{i,n} = \frac{\rho |H_i[n, i-1]|^2}{\rho |H_0[n, i-1]|^2 + 1}. \quad (47)$$

725 Unlike downlink OTFS-NOMA, there are two possible
 726 strategies for uplink OTFS-NOMA to combat multiple access
 727 interference, as shown in the following two subsections.

728 1) *Adaptive-Rate Transmission*: One strategy to combat
 729 multiple access interference is to impose the following con-
 730 straint on $x_i(n)$:

$$731 \quad R_{i,n} \leq \log \left(1 + \frac{\rho |H_i[n, i-1]|^2}{\rho |H_0[n, i-1]|^2 + 1} \right), \quad (48)$$

732 which means that the first stage of SIC is guaranteed to be
 733 successful. Therefore, the M low-mobility users are served
 734 without affecting U_0 's outage probability, i.e., the use of
 735 NOMA is transparent to U_0 .

736 Because U_i 's data rate is adaptive, outage events when
 737 decoding $x_i(n)$ do not happen, which means that an appropri-
 738 ate criterion for the performance evaluation is the ergodic rate.
 739 Recall that $H_i[n, i-1] = \tilde{D}_i^{i-1}$ and $H_0[n, i-1] = D_0^{n, i-1}$.
 740 Therefore, U_i 's ergodic rate is given by

$$741 \quad \mathcal{E}\{R_{i,n}\} \leq \mathcal{E} \left\{ \log \left(1 + \frac{\rho |\tilde{D}_i^{i-1}|^2}{\rho |D_0^{n, i-1}|^2 + 1} \right) \right\}. \quad (49)$$

742 We note that the ergodic rate of uplink OTFS-NOMA can
 743 be further improved by modifying the user scheduling strategy
 744 proposed in (43), as shown in the following. Particularly,
 745 denote the NOMA user which is scheduled to transmit in the
 746 m -th frequency subchannel by $U_{i_m^*}$, and this user is selected
 747 by using the following criterion:

$$748 \quad i_m^* = \arg \max_{i \in \{1, \dots, K\}} \left\{ |\tilde{D}_i^m|^2 \right\}. \quad (50)$$

749 We note that a single user might be scheduled on multiple
 750 frequency channels, which reduces user fairness.

751 Because the integration of the logarithm function appearing
 752 in (49) leads to non-insightful special functions, we will use
 753 simulations to evaluate the ergodic rate of OTFS-NOMA in
 754 Section VII.

755 2) *Fixed-Rate Transmission*: If the NOMA users do not
 756 have the capabilities to adapt their transmission rates, they
 757 have to use fixed data rates R_i for transmission, which means
 758 that outage events can happen and the achieved outage perfor-
 759 mance is analyzed in the following. For illustration purposes,
 760 we focus on the case when the user scheduling strategy shown
 761 in (50) is used.

762 The outage probability for detecting $x_{i_m^*}(n)$ is given by

$$763 \quad P_{i_m^*, n} = \text{P} \left(\log \left(1 + \frac{\rho |\tilde{D}_{i_m^*}^{i_m^*-1}|^2}{\rho |D_0^{n, i_m^*-1}|^2 + 1} \right) < R_{i_m^*} \right). \quad (51)$$

764 Following steps similar to the ones in the proof for Lemma 2,
 765 we can show that $|\tilde{D}_{i_m^*}^{i_m^*-1}|^2$ and $|D_0^{n, i_m^*-1}|^2$ are independent,

and the use of the user scheduling scheme in (50) simplifies
 the outage probability as follows:

$$766 \quad P_{i_m^*, n} = \text{P} \left(\log \left(1 + \frac{\rho |\tilde{D}_{i_m^*}^{i_m^*-1}|^2}{\rho |D_0^{n, i_m^*-1}|^2 + 1} \right) < R_{i_m^*} \right) \quad 767$$

$$768 \quad = \int_0^\infty \left(1 - e^{-\frac{\epsilon_{i_m^*}(1+\rho y)}{\rho}} \right)^K e^{-y} dy, \quad (52) \quad 769$$

770 where we use the fact that the cumulative distribution function
 771 of $|\tilde{D}_{i_m^*}^{i_m^*-1}|^2$ is $(1 - e^{-x})^K$ because of the adopted user
 772 scheduling strategy.

773 The outage probability can be further simplified as follows:

$$774 \quad P_{i_m^*, n} = \sum_{k=0}^K \binom{K}{k} (-1)^k \int_0^\infty e^{-\frac{k\epsilon_{i_m^*}(1+\rho y)}{\rho}} e^{-y} dy \quad 775$$

$$776 \quad = \sum_{k=0}^K \binom{K}{k} (-1)^k e^{-\frac{k\epsilon_{i_m^*}}{\rho}} \frac{1}{k\epsilon_{i_m^*} + 1}. \quad (53) \quad 777$$

778 At high SNR, the outage probability can be approximated
 779 as follows:

$$780 \quad P_{i_m^*, n} \approx \sum_{k=0}^K \binom{K}{k} (-1)^k \frac{1}{k\epsilon_{i_m^*} + 1}, \quad (54) \quad 781$$

782 which is no longer a function of ρ , i.e., the outage probability
 783 has an error floor at high SNR. This is due to the fact that
 784 $U_{i_m^*}$ is subject to strong interference from U_0 .

785 However, we can show that the error floor experienced by
 786 $U_{i_m^*}$ can be reduced by increasing K , i.e., inviting more oppor-
 787 tunistic users for NOMA transmission. In particular, assuming
 788 $K\epsilon_{i_m^*} \rightarrow 0$, the outage probability can be approximated as
 789 follows:

$$790 \quad P_{i_m^*, n} \approx \sum_{k=0}^K \binom{K}{k} (-1)^k (1 + k\epsilon_{i_m^*})^{-1} \quad 791$$

$$792 \quad \approx \sum_{k=0}^K \binom{K}{k} (-1)^k \sum_{l=0}^\infty (-1)^l k^l \epsilon_{i_m^*}^l, \quad (55) \quad 793$$

794 where we use the fact that $(1+x)^{-1} = \sum_{l=0}^\infty (-1)^l x^l$, $|x| < 1$.
 795 Therefore, the error floor at high SNR can be approximated
 796 as follows:

$$797 \quad P_{i_m^*, n} \approx \sum_{l=0}^\infty (-1)^l \epsilon_{i_m^*}^l \sum_{k=0}^K \binom{K}{k} (-1)^k k^l \quad 798$$

$$799 \quad \approx (-1)^K \epsilon_{i_m^*}^K (-1)^K K! = K! \epsilon_{i_m^*}^K, \quad (56) \quad 800$$

801 where we use the identities $\sum_{k=0}^K \binom{K}{k} (-1)^k k^l = 0$, for $l < K$
 802 and $\sum_{k=0}^K \binom{K}{k} (-1)^k k^K = (-1)^K K!$.

803 The conclusion that increasing K reduces the error floor
 804 can be confirmed by defining $f(k) = k! \epsilon_{i_m^*}^k$ and using the
 805 following fact:

$$806 \quad f(k) - f(k+1) = k! \epsilon_{i_m^*}^k (1 - (k+1)\epsilon_{i_m^*}) > 0, \quad (57)$$

807 where it is assumed that $k\epsilon_{i_m^*} \rightarrow 0$.

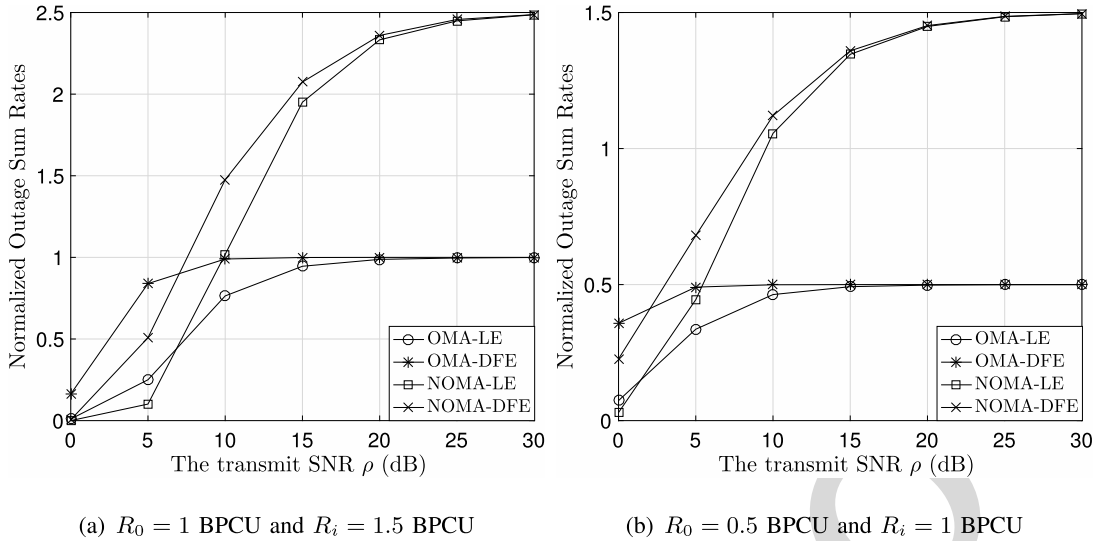


Fig. 1. Impact of OTFS-NOMA on the downlink sum rates. $M = N = K = 16$. $P_0 = P_i = 3$. BPCU denotes bit per channel use. $\gamma_0^2 = \frac{3}{4}$ and $\gamma_i^2 = \frac{1}{4}$ for $i > 0$. Random user scheduling is used.

B. Stage II of SIC

If adaptive transmission is used, the NOMA users' signals can be detected successfully during the first stage of SIC. Therefore, they can be removed from the observations at the base station, i.e., $\bar{Y}[n, m] = Y[n, m] - \sum_{q=1}^N H_q(n, m)X_q[n, m]$, and SFFT is applied to obtain the delay-Doppler observations as follows:

$$y_0[k, l] = \frac{1}{NM} \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} \bar{Y}[n, m] e^{-j2\pi(\frac{nk}{N} - \frac{ml}{M})}$$

$$= \sum_{p=0}^{P_0} h_{0,p} x_0[(k - k_{\mu_{0,p}})_N, (l - l_{\tau_{0,p}})_M] + z[k, l], \quad (58)$$

where $z[k, l]$ denote additive noise. U_0 's signals can be detected by applying either of the two considered equalization approaches, and the same performance as for OTFS-OMA can be realized. The analytical development is similar to the downlink case, and hence is omitted due to space limitations.

However, if fixed-rate transmission is used, the uplink outage events for decoding $x_0[k, l]$ are different from the downlink ones, as shown in the following. Particularly, the use of FD-LE yields the following SINR expression for decoding $x_0[k, l]$:

$$\text{SINR}_{0,kl}^{\text{LE}} = \frac{\rho}{\frac{1}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} |D_0^{k,l}|^{-2}}. \quad (59)$$

If FD-DFE is used, the SNR for detection of $x_0[k, l]$ is given by

$$\text{SINR}_{0,kl}^{\text{DFE}} = \rho \lambda_{0,kl}. \quad (60)$$

Therefore, the outage probability for detecting $x_0[k, l]$ is given by

$$P_{kl} = 1 - \text{P}(\text{SINR}_{0,kl}^{\text{DFE/LE}} > \epsilon_0, \text{SNR}_{i,n} > \epsilon_i \forall i, n)$$

$$\geq 1 - \text{P}(\text{SNR}_{i,n} > \epsilon_i \forall i, n) \geq \text{P}(\text{SNR}_{1,0} < \epsilon_i).$$

Since $\text{P}(\text{SNR}_{1,0} < \epsilon_i)$ has an error floor as shown in the previous subsection, the uplink outage probability for detection

TABLE I
DELAY-DOPPLER PROFILE FOR U_0 'S CHANNEL

Propagation path index (p)	0	1	2	3
Delay ($\tau_{0,p}$) μs	8.33	25	41.67	58.33
Delay tap index ($l_{\tau_{0,p}}$)	2	6	10	14
Doppler ($\nu_{0,p}$) Hz	0	0	468.8	468.8
Doppler tap index ($k_{\nu_{0,p}}$)	0	0	1	1

of U_0 's signals does not go to zero even if $\rho \rightarrow \infty$, which is different from the downlink case. Therefore, if fixed-rate transmission is used, adding the M low-mobility users into the bandwidth, which would be solely occupied by U_0 in OTFS-OMA, improves connectivity but degrades U_0 's performance.

VII. NUMERICAL STUDIES

In this section, the performance of OTFS-NOMA is evaluated via computer simulations. Similar to [26]–[28], we first define the delay-Doppler profile for U_0 's channel as shown in Table I, where $P_0 = 3$ and the subchannel spacing is $\Delta f = 7.5$ kHz. Therefore, the maximal speed corresponding to the largest Doppler shift $\nu_{0,3} = 468.8$ Hz is 126.6 km/h if the carrier frequency is $f_c = 4$ GHz. On the other hand, the NOMA users' channels are assumed to be time invariant with $P_i = 3$ propagation paths, i.e., $\tau_{i,p} = 0$ for $p \geq 4, i \geq 1$. For all the users' channels, we assume that $\sum_{p=0}^{P_i} \mathcal{E}\{|h_{i,p}|^2\} = 1$ and $|h_{i,p}|^2 \sim CN(0, \frac{1}{P_i+1})$. For the fixed rate transmission scheme, a simple choice for power allocation ($\gamma_0^2 = \frac{3}{4}$ and $\gamma_i^2 = \frac{1}{4}$ for $i > 0$) is considered. The performance of OTFS-NOMA could be further improved by optimizing γ_i according to the users' channel conditions and QoS requirements.

In Fig. 1, downlink OTFS-NOMA transmission is evaluated by using the normalized outage sum rate as the performance criterion which is defined as $\frac{1}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} (1 - P_{0,kl}) R_0$

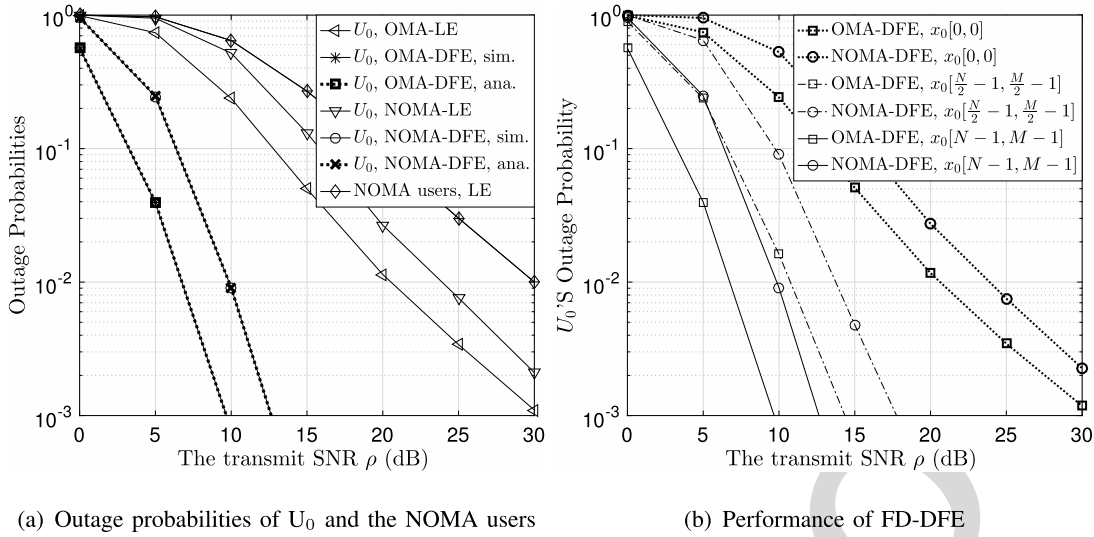


Fig. 2. The outage performance of downlink OTFS-OMA and OTFS-NOMA. $M = N = K = 16$. $P_0 = P_i = 3$. $\gamma_0^2 = \frac{3}{4}$ and $\gamma_i^2 = \frac{1}{4}$ for $i > 0$. $R_0 = 0.5$ BPCU and $R_i = 1$ BPCU. In Fig. 2(a), for FD-DFE, the performance of $x_0[N-1, M-1]$ is shown. Random user scheduling is used.

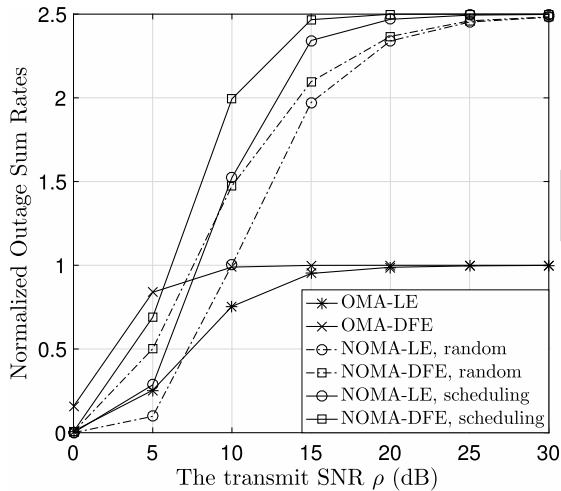


Fig. 3. Impact of user scheduling on the downlink outage sum rates. $P_0 = P_i = 3$. $R_0 = 1$ BPCU and $R_i = 1.5$ BPCU. $M = N = K = 16$, $\gamma_0^2 = \frac{3}{4}$ and $\gamma_i^2 = \frac{1}{4}$ for $i > 0$.

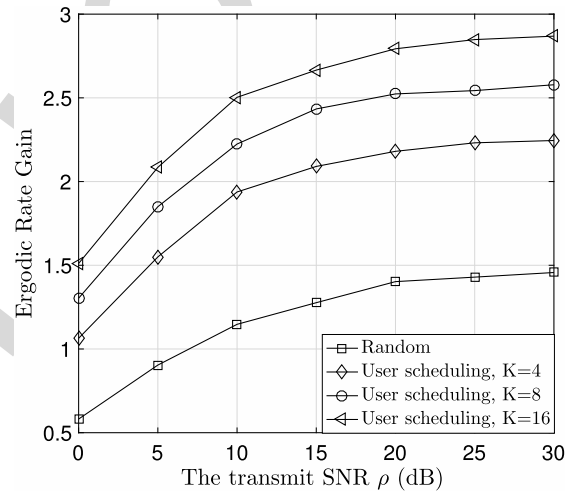


Fig. 4. The ergodic rate gain of OTFS-NOMA over OTFS-OMA. The NOMA users adapt their data rates according to (48). $P_0 = P_i = 3$. $M = N = 16$.

and $\frac{1}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} (1 - P_{0,kl}) R_0 + \frac{1}{NM} \sum_{i=1}^M \sum_{n=0}^{N-1} (1 - P_{i,n}) R_i$ for OTFS-OMA and OTFS-NOMA, respectively. Fig. 1 shows that the use of OTFS-NOMA can significantly improve the sum rate at high SNR for both considered choices of R_0 and R_i . The reason for this performance gain is the fact that the maximal sum rate achieved by OTFS-OMA is capped by R_0 , whereas OTFS-NOMA can provide sum rates up to $R_0 + R_i$. Comparing Fig. 1(a) to Fig. 1(b), one can observe that the performance loss of OTFS-NOMA at low SNR can be mitigated by reducing the target data rates, since reducing the target rates improves the probability of successful SIC. Furthermore, both figures show that FD-DFE outperforms FD-LE in the entire considered range of SNRs; however, we note that the performance gain of FD-DFE over FD-LE is achieved at the expense of increased computational complexity.

In Fig. 2, the outage probabilities achieved by downlink OTFS-OMA and OTFS-NOMA are shown. As can be seen

from Fig. 2(a), the diversity order achieved with FD-LE for detection of $x_0[k, l]$ is one, as expected from Lemma 2. As discussed in Section IV-B, one advantage of FD-DFE over FD-LE is that FD-DFE facilitates multi-path fading diversity gains, whereas FD-LE is limited to a diversity gain of one. This conclusion is confirmed by Fig. 2(a), where the analytical results developed in Corollary 1 are also verified. Fig. 2(b) shows the outage probabilities achieved by FD-DFE for different $x_0[k, l]$. As shown in the figure, the lowest outage probability is obtained for $x_0[N-1, M-1]$, whereas the outage probability of $x_0[0, 0]$ is the largest, which is due to the fact that, in FD-DFE, different signals $x_0[k, l]$ are affected by different effective channel gains, $\lambda_{0,kl}$. Another important observation from the figures is that the FD-LE outage probability is the same as the FD-DFE outage probability for detection of $x_0[0, 0]$, which fits the intuition that for FD-DFE the reliability of the first decision ($x_0[0, 0]$) is the same as that of FD-LE. For the same reason, FD-LE and FD-DFE yield similar performance for detection of the NOMA users'

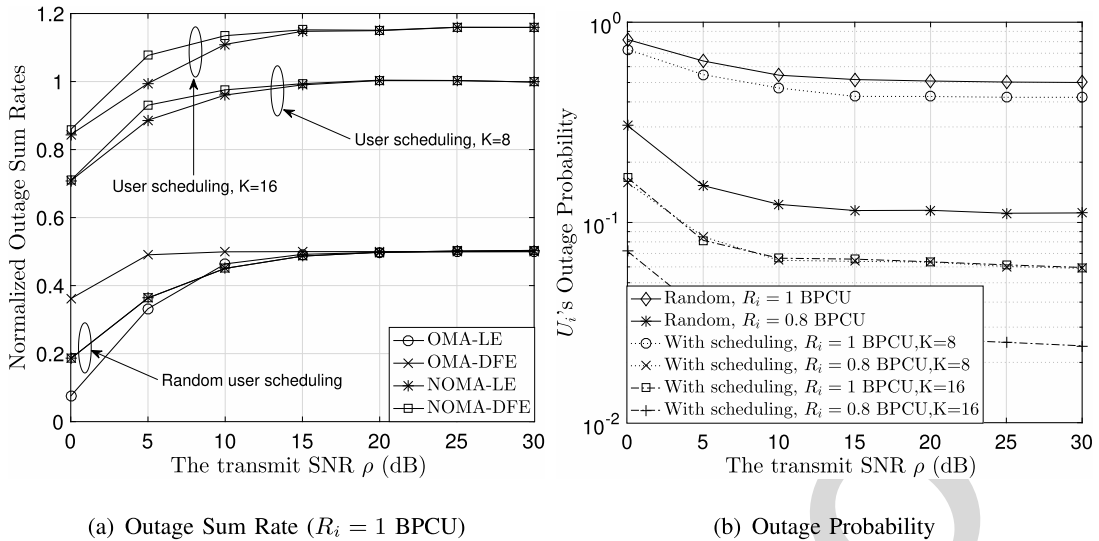


Fig. 5. The performance of uplink OTFS-NOMA. Fixed-rate transmission is used by the NOMA users. $M = N = 16$. $P_0 = P_i = 3$. $R_0 = 0.5$ BPCU. $\gamma_0^2 = \frac{3}{4}$ and $\gamma_i^2 = \frac{1}{4}$ for $i > 0$.

893 signals, since the FD-DFE outage performance is dominated
 894 by the reliability for detection of $x_0[0, 0]$, and hence is the
 895 same as that of FD-LE.

896 In addition to multi-path diversity, another degree of free-
 897 dom available in the considered OTFS-NOMA downlink
 898 scenario is multi-user diversity, which can be harvested by
 899 applying user scheduling as discussed in Section V-B. Fig. 3
 900 demonstrates the benefits of exploiting multi-user diversity.
 901 With random user scheduling, at low SNR, the performance
 902 of OTFS-NOMA is worse than that of OTFS-OMA, which
 903 is also consistent with Fig. 1. By increasing the number
 904 of users participating in OTFS-NOMA, the performance of
 905 OTFS-NOMA can be improved, particularly at low and moder-
 906 ate SNR. For example, for FD-LE, the performance of
 907 OTFS-NOMA approaches that of OTFS-OMA at low SNR
 908 by exploiting multi-user diversity, and for FD-DFE, an extra
 909 gain of 0.5 BPCU can be achieved at moderate SNR.

910 In Figs. 4 and 5, the performance of uplink OTFS-NOMA is
 911 evaluated. As discussed in Section VI, the NOMA users have
 912 two choices for their transmission rates, namely adaptive and
 913 fixed rate transmission. The use of adaptive rate transmission
 914 can ensure that the implementation of NOMA is transparent
 915 to U_0 , which means that U_0 's QoS requirements are strictly
 916 guaranteed. Since U_0 achieves the same performance for
 917 OTFS-NOMA and OTFS-OMA when adaptive rate transmis-
 918 sion is used, we only focus on the NOMA users' performance,
 919 where the ergodic rate in (49) is used as the criterion.
 920 We note that this ergodic rate is the net performance gain of
 921 OTFS-NOMA over OTFS-OMA, which is the reason why the
 922 vertical axis in Fig. 4 is labeled 'Ergodic Rate Gain'. When
 923 the M users are randomly selected from the K NOMA users,
 924 the ergodic rate gain is moderate, e.g., 1.5 bit per channel
 925 use (BPCU) at $\rho = 30$ dB. By applying the scheduling strategy
 926 proposed in (50), the ergodic rate gain can be significantly
 927 improved, e.g., nearly by a factor of two compared to the
 928 random case with $K = 16$ and $\rho = 30$ dB.

929 Fig. 5 focuses on the case with fixed rate transmission, and
 930 similar to Fig. 1, the normalized outage sum rate is used as
 931 performance criterion in Fig. 5(a). One can observe that with
 932 random user scheduling, the sum rate of OTFS-NOMA is similar
 933 to that of OTFS-OMA. This is due to the fact that no inter-
 934 ference mitigation strategy, such as power or rate allocation,
 935 is used for NOMA uplink transmission, which means that U_0
 936 and the NOMA users cause strong interference to each other
 937 and SIC failure may happen frequently. By applying the user
 938 scheduling strategy proposed in (50), the channel conditions of
 939 the scheduled users become quite different, which facilitates
 940 the implementation of SIC. This benefit of user scheduling
 941 can be clearly observed in Fig. 5(a), where NOMA achieves
 942 a significant gain over OMA although advanced power or rate
 943 allocation strategies are not used. Fig. 5(a) also shows that the
 944 difference between the performance of FD-LE and FD-DFE is
 945 insignificant for the uplink case. This is due to the fact that the
 946 outage events during the first stage of SIC dominate the outage
 947 performance, and they are not affected by whether FD-LE or
 948 FD-DFE is employed. Another important observation from
 949 Fig. 5(a) is that the maximal sum rate $R_0 + R_i$ cannot be
 950 realized, even at high SNR. The reason for this behaviour is
 951 the existence of the error floor for the NOMA users' outage
 952 probabilities, as shown in Fig. 5(b). The analytical results
 953 provided in Section V-B show that increasing K can reduce
 954 the error floor, which is confirmed by Fig. 5(b).

955 VIII. CONCLUSION

956 In this paper, we have proposed OTFS-NOMA uplink and
 957 downlink transmission schemes, where users with different
 958 mobility profiles are grouped together for the implemen-
 959 tation of NOMA. The analytical results developed in the
 960 paper demonstrate that both the high-mobility and the low-
 961 mobility users benefit from the application of OTFS-NOMA.
 962 In particular, the use of NOMA enables the spreading of
 963 the signals of a high-mobility user over a large amount

of time-frequency resources, which enhances the OTFS resolution and improves the detection reliability. In addition, OTFS-NOMA ensures that the low-mobility users have access to the bandwidth resources which would be solely occupied by the high-mobility users in OTFS-OMA. Hence, OTFS-NOMA improves the spectral efficiency and reduces latency. An interesting topic for future works is studying the impact of non-zero fractional delays and fractional Doppler shifts on the performance of the developed OTFS-NOMA protocol. Furthermore, in this paper, the users' channel gains (the taps of the delay-Doppler impulse response) have been assumed to be Gaussian distributed, and an important direction for future research is to investigate the impact of other types of channel distributions on the performance of OTFS-NOMA. Moreover, the combination of emerging spectrally efficient 5G solutions, such as 5G New Radio Bandwidth Part (5G-NR-BWP) [39], [40] and software-controlled metasurfaces [41], with OTFS-NOMA is also a promising topic for future research.

APPENDIX A PROOF FOR PROPOSITION 1

Intuitively, the use of $\mathbf{F}_N \otimes \mathbf{F}_M^H$ is analogous to the application of the ISFFT which transforms signals from the delay-Doppler plane to the time-frequency plane, where inter-symbol interference is removed, i.e., the user's channel matrix is diagonalized. The following proof confirms this intuition and reveals how the diagonalized channel matrix is related to the original block circulant matrix. We first apply $\mathbf{F}_N \otimes \mathbf{I}_M$ to \mathbf{y}_0 , which yields the following:

$$\begin{aligned} & (\mathbf{F}_N \otimes \mathbf{I}_M)\mathbf{y}_0 \\ &= (\mathbf{F}_N \otimes \mathbf{I}_M)\mathbf{H}_0 \left(\gamma_0 \mathbf{x}_0 + \sum_{q=1}^M \gamma_q \mathbf{x}_q \right) + (\mathbf{F}_N \otimes \mathbf{I}_M)\mathbf{z}_0 \\ &= \text{diag} \left\{ \sum_{n=0}^{N-1} \mathbf{A}_{0,n} e^{-j \frac{2\pi l n}{N}}, 0 \leq l \leq N-1 \right\} (\mathbf{F}_N \otimes \mathbf{I}_M) \\ & \quad \times \left(\gamma_0 \mathbf{x}_0 + \sum_{q=1}^M \gamma_q \mathbf{x}_q \right) + (\mathbf{F}_N \otimes \mathbf{I}_M)\mathbf{z}_0, \end{aligned} \quad (61)$$

where $\text{diag}\{\mathbf{B}_1, \dots, \mathbf{B}_N\}$ denotes a block-diagonal matrix with \mathbf{B}_n , $1 \leq n \leq N$, on its main diagonal. Note that $\sum_{n=0}^{N-1} \mathbf{A}_{0,n} e^{-j \frac{2\pi l n}{N}}$, $0 \leq l \leq N-1$, is a sum of N $M \times M$ circulant matrices, each of which can be further diagonalized by \mathbf{F}_M . Therefore, we can apply $\mathbf{I}_N \otimes \mathbf{F}_M^H$ to $(\mathbf{F}_N \otimes \mathbf{I}_M)\mathbf{y}_0$, which yields the following:

$$\begin{aligned} & (\mathbf{I}_N \otimes \mathbf{F}_M^H)(\mathbf{F}_N \otimes \mathbf{I}_M)\mathbf{y}_0 \\ &= \text{diag} \left\{ \sum_{n=0}^{N-1} \mathbf{A}_{0,n} e^{-j \frac{2\pi l n}{N}}, 0 \leq l \leq N-1 \right\} \\ & \quad \times (\mathbf{F}_N \otimes \mathbf{I}_M)(\mathbf{I}_N \otimes \mathbf{F}_M^H) \left(\gamma_0 \mathbf{x}_0 + \sum_{q=1}^M \gamma_q \mathbf{x}_q \right) \\ & \quad + (\mathbf{I}_N \otimes \mathbf{F}_M^H)(\mathbf{F}_N \otimes \mathbf{I}_M)\mathbf{z}_0, \end{aligned} \quad (62)$$

where $\mathbf{A}_{0,n}$ is a diagonal matrix, $\mathbf{A}_{0,n} = \text{diag} \left\{ \sum_{m=0}^{M-1} a_{0,n}^{m,1} e^{j \frac{2\pi t m}{M}}, 0 \leq t \leq M-1 \right\}$, and $a_{0,n}^{m,1}$ is the element located in the m -th row and first column of $\mathbf{A}_{0,n}$.

By applying a property of the Kronecker product, $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{A}\mathbf{C}) \otimes (\mathbf{B}\mathbf{D})$, the received signals can be simplified as follows:

$$\begin{aligned} & (\mathbf{F}_N \otimes \mathbf{F}_M^H)\mathbf{y}_0 \\ &= \text{diag} \left\{ \sum_{n=0}^{N-1} \mathbf{A}_{0,n} e^{-j \frac{2\pi l n}{N}}, 0 \leq l \leq N-1 \right\} (\mathbf{F}_N \otimes \mathbf{F}_M^H) \\ & \quad \times \left(\gamma_0 \mathbf{x}_0 + \sum_{q=1}^M \gamma_q \mathbf{x}_q \right) + (\mathbf{F}_N \otimes \mathbf{F}_M^H)\mathbf{z}_0, \end{aligned} \quad (63)$$

where the $(kM + l + 1)$ -th element on the main diagonal of \mathbf{D}_0 is $D_0^{k,l}$ as defined in the proposition. The proof for the proposition is complete.

APPENDIX B PROOF FOR LEMMA 1

In order to facilitate the SINR analysis, the system model in (18) is further simplified. Define $\tilde{X}[n, m] = \sum_{i=1}^M X_i[n, m]$. With the mapping scheme used in (6), the NOMA users' signals are interleaved and orthogonally placed in the time-frequency plane, i.e., $\tilde{X}[n, m]$ is simply U_{m+1} 's n -th signal, $x_{m+1}(n)$. Denote the outcome of the SFFT of $\tilde{X}[n, m]$ by $\tilde{x}[k, l]$, which yields the following transform:

$$\tilde{x}[k, l] = \frac{1}{\sqrt{NM}} \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} \tilde{X}[n, m] e^{-j2\pi \left(\frac{nk}{N} - \frac{ml}{M} \right)}. \quad (64)$$

Denote the $NM \times 1$ vector collecting the $\tilde{x}[k, l]$ by $\tilde{\mathbf{x}}$ and the $NM \times 1$ vector collecting the $\tilde{X}[n, m]$ by $\tilde{\mathbf{x}}$, which means that (64) can be rewritten as follows:

$$\tilde{\mathbf{x}} = (\mathbf{F}_N \otimes \mathbf{F}_M^H)\tilde{\mathbf{x}}. \quad (65)$$

Therefore, the model for the received signals in (18) can be re-written as follows:

$$\begin{aligned} \check{\mathbf{y}}_0 &= \gamma_0 \mathbf{x}_0 + \gamma_1 \tilde{\mathbf{x}} + (\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-1} \mathbf{D}_0^{-1} \tilde{\mathbf{z}}_i \\ &= \gamma_0 \mathbf{x}_0 + \gamma_1 \underbrace{(\mathbf{F}_N \otimes \mathbf{F}_M^H)\tilde{\mathbf{x}} + (\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-1} \mathbf{D}_0^{-1} \tilde{\mathbf{z}}_0}_{\text{Interference and noise terms}}, \end{aligned} \quad (66)$$

where we have used the assumption that $\gamma_i = \gamma_1$, for $1 \leq i \leq N$. Note that the power of the information-bearing signals is simply $\gamma_0^2 \rho$, and therefore, the key step to obtain the SINR is to find the covariance matrix of the interference-plus-noise term.

We first show that $\tilde{\mathbf{z}}_0 \triangleq (\mathbf{F}_N \otimes \mathbf{F}_M^H)\mathbf{z}_0$ is still a complex Gaussian vector, i.e., $\tilde{\mathbf{z}}_i \sim CN(0, \mathbf{I}_{NM})$. Recall that \mathbf{z}_0 contains NM i.i.d. complex Gaussian random variables. Furthermore, $\mathbf{F}_N \otimes \mathbf{F}_M^H$ is a unitary matrix as shown in the following:

$$\begin{aligned} & (\mathbf{F}_N \otimes \mathbf{F}_M^H)(\mathbf{F}_N \otimes \mathbf{F}_M^H)^H \stackrel{(a)}{=} (\mathbf{F}_N \otimes \mathbf{F}_M^H)(\mathbf{F}_N^H \otimes \mathbf{F}_M) \\ & \stackrel{(b)}{=} (\mathbf{F}_N \mathbf{F}_N^H) \otimes (\mathbf{F}_M^H \mathbf{F}_M) \\ & = \mathbf{I}_{NM}, \end{aligned} \quad (67)$$

where step (a) follows from the fact that $(\mathbf{A} \otimes \mathbf{B})^H = \mathbf{A}^H \otimes \mathbf{B}^H$ and step (b) follows from the fact that

1051 $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD})$. Therefore, $(\mathbf{F}_N \otimes \mathbf{F}_M^H) \mathbf{z}_0 \sim$
 1052 $CN(0, \mathbf{I}_{NM})$ given the fact that $\mathbf{z}_0 \sim CN(0, \mathbf{I}_{NM})$ and a
 1053 unitary transformation of a Gaussian vector is still a Gaussian
 1054 vector.

1055 Therefore, the covariance matrix of the interference-plus-
 1056 noise term is given by

$$1057 \mathbf{C}_{\text{cov}} = \gamma_1^2 \mathcal{E} \left\{ (\mathbf{F}_N \otimes \mathbf{F}_M^H) \check{\mathbf{x}} \check{\mathbf{x}}^H (\mathbf{F}_N \otimes \mathbf{F}_M^H)^H \right\} \\
 1058 + \mathcal{E} \left\{ (\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-1} \mathbf{D}_0^{-1} \check{\mathbf{z}}_0 \check{\mathbf{z}}_0^H \mathbf{D}_0^{-H} (\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-H} \right\}. \quad (68)$$

1061 Recall that the $(nM + m + 1)$ -th element of $\check{\mathbf{x}}$ is $\tilde{X}[n, m]$
 1062 which is equal to $x_{m+1}(n)$. Therefore, the covariance matrix
 1063 can be further simplified as follows:

$$1064 \mathbf{C}_{\text{cov}} = \gamma_1^2 \rho (\mathbf{F}_N \otimes \mathbf{F}_M^H) (\mathbf{F}_N \otimes \mathbf{F}_M^H)^H \\
 1065 + (\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-1} \mathbf{D}_0^{-1} \mathbf{D}_0^{-H} (\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-H} \\
 1066 = \gamma_1^2 \rho \mathbf{I}_{MN} + (\mathbf{F}_N^H \otimes \mathbf{F}_M) \mathbf{D}_0^{-1} \mathbf{D}_0^{-H} (\mathbf{F}_N \otimes \mathbf{F}_M^H), \quad (69)$$

1068 where the noise power is assumed to be normalized.

1069 Following the same steps as in the proof of Proposi-
 1070 tion 1, we learn that, by construction, $(\mathbf{F}_N^H \otimes \mathbf{F}_M) \mathbf{D}_0^{-1} \mathbf{D}_0^{-H}$
 1071 $(\mathbf{F}_N \otimes \mathbf{F}_M^H)$ is also a block-circulant matrix, which means
 1072 that the elements on the main diagonal of $(\mathbf{F}_N^H \otimes \mathbf{F}_M)$
 1073 $\mathbf{D}_0^{-1} \mathbf{D}_0^{-H} (\mathbf{F}_N \otimes \mathbf{F}_M^H)$ are identical. Without loss of gener-
 1074 ality, denote the diagonal elements of $(\mathbf{F}_N^H \otimes \mathbf{F}_M) \mathbf{D}_0^{-1} \mathbf{D}_0^{-H}$
 1075 $(\mathbf{F}_N \otimes \mathbf{F}_M^H)$ by ϕ . Therefore, ϕ can be found by using the
 1076 trace of the matrix as follows:

$$1077 \phi = \frac{1}{NM} \text{Tr} \left\{ (\mathbf{F}_N^H \otimes \mathbf{F}_M) \mathbf{D}_0^{-1} \mathbf{D}_0^{-H} (\mathbf{F}_N \otimes \mathbf{F}_M^H) \right\} \\
 1078 = \frac{1}{NM} \text{Tr} \left\{ (\mathbf{F}_N \otimes \mathbf{F}_M^H) (\mathbf{F}_N^H \otimes \mathbf{F}_M) \mathbf{D}_0^{-1} \mathbf{D}_0^{-H} \right\} \\
 1079 = \frac{1}{NM} \text{Tr} \left\{ \mathbf{D}_0^{-1} \mathbf{D}_0^{-H} \right\} = \frac{1}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} |D_0^{k,l}|^{-2}. \quad (70)$$

1080 Therefore, the SINR for detection of $x_0[k, l]$ is given by

$$1081 \text{SINR}_{0,kl}^{LE} = \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \phi}, \quad (71)$$

1082 and the proof is complete.

1083 APPENDIX C 1084 PROOF FOR LEMMA 2

1085 The lemma is proved by first developing upper and lower
 1086 bounds on the outage probability, and then showing that both
 1087 bounds have the same diversity order.

1088 An upper bound on $\text{SINR}_{0,kl}$ is given by

$$1089 \text{SINR}_{0,kl} = \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \frac{1}{NM} \sum_{\tilde{k}=0}^{N-1} \sum_{\tilde{l}=0}^{M-1} |D_0^{\tilde{k}, \tilde{l}}|^{-2}} \\
 1090 \leq \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \frac{1}{NM} |D_0^{0,0}|^{-2}}. \quad (72)$$

Therefore, the outage probability, denoted by $P_{0,kl}$, can be
 lower bounded as follows:

$$1093 P_{0,kl} \geq \mathbb{P} \left(\frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \frac{1}{NM} |D_0^{0,0}|^{-2}} < \epsilon_0 \right) \\
 1094 = \mathbb{P} \left(|D_0^{0,0}|^2 < \frac{\epsilon_0}{NM \rho (\gamma_0^2 - \gamma_1^2 \epsilon_0)} \right), \quad (73)$$

where we assume that $\gamma_0^2 > \gamma_1^2 \epsilon_0$. Otherwise, the outage
 probability is always one.

To evaluate the lower bound on the outage probability,
 the distribution of $D_0^{u,v}$ is required. Recall from (16) that $D_0^{u,v}$
 is the $((v-1)M + u)$ -th main diagonal element of \mathbf{D}_0 and
 can be expressed as follows:

$$1101 D_0^{u,v} = \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} a_{0,n}^{m,1} e^{j2\pi \frac{um}{M}} e^{-j2\pi \frac{vn}{N}}, \quad (74)$$

which is the ISFFT of $a_{0,n}^{m,1}$. Therefore, we have the following
 property:

$$1102 \tilde{\mathbf{D}}_0 = \sqrt{NM} \mathbf{F}_M^H \mathbf{A}_0 \mathbf{F}_N, \quad (75)$$

where the element in the u -th row and the v -th column of $\tilde{\mathbf{D}}_0$
 is $D_0^{u,v}$ and the element in the m -th row and the n -th column
 of \mathbf{A}_0 is $a_{0,n}^{m,1}$.

The matrix-based expression shown in (75) can be vector-
 ized as follows:

$$1103 \text{Diag}(\mathbf{D}_0) = \text{vec}(\tilde{\mathbf{D}}_0) = \sqrt{NM} \text{vec}(\mathbf{F}_M^H \mathbf{A}_0 \mathbf{F}_N) \\
 1104 = \sqrt{NM} (\mathbf{F}_N \otimes \mathbf{F}_M^H) \text{vec}(\mathbf{A}_0), \quad (76)$$

where $\text{Diag}(\mathbf{A})$ denotes a vector collecting all elements on
 the main diagonal of \mathbf{A} and we use the facts that $(\mathbf{C}^T \otimes$
 $\mathbf{A}) \text{vec}(\mathbf{B}) = \text{vec}(\mathbf{D})$ if $\mathbf{ABC} = \mathbf{D}$, and $\mathbf{F}_N^T = \mathbf{F}_N$.

We note that $\text{vec}(\mathbf{A}_0)$ contains only $(P_0 + 1)$ non-zero
 elements, where the remaining elements are zero. Therefore,
 each element on the main diagonal of \mathbf{D}_0 is a superposition
 of $(P_0 + 1)$ i.i.d. random variables, $h_{i,p} \sim CN(0, \frac{1}{P_0+1})$.
 We further note that the coefficients for the superposition are
 complex exponential constants, i.e., the magnitude of each
 coefficient is one. Therefore, each element on the main di-
 agonal of \mathbf{D}_0 is still complex Gaussian distributed, i.e., $D_0^{u,v} \sim$
 $CN(0, 1)$, which means that the lower bound on the outage
 probability shown in (73) can be expressed as follows:

$$1125 P_{0,kl} \geq 1 - e^{-\frac{\epsilon_0}{NM \rho (\gamma_0^2 - \gamma_1^2 \epsilon_0)}} \doteq \frac{1}{\rho}. \quad (77)$$

On the other hand, an upper bound on the outage probability
 is given by

$$1128 P_{0,kl} \leq \mathbb{P} \left(\frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \frac{1}{NM} \sum_{\tilde{k}=0}^{N-1} \sum_{\tilde{l}=0}^{M-1} |D_0^{\min}|^{-2}} < \epsilon_0 \right), \quad (78)$$

where $|D_0^{\min}| = \min\{|D_0^{k,l}|, \forall l \in \{0, \dots, M-1\}, k \in$
 $\{0, \dots, N-1\}\}$.

Therefore, the outage probability can be upper bounded as
 follows:

$$1134 P_{0,kl} \leq \mathbb{P} \left(|D_0^{\min}|^2 < \frac{\epsilon_0}{\rho (\gamma_0^2 - \gamma_1^2 \epsilon_0)} \right). \quad (79)$$

It is important to point out that the $|D_0^{k,l}|^2$, $l \in \{0, \dots, M-1\}$, $k \in \{0, \dots, N-1\}$, are identically but not independently distributed. This correlation property is shown as follows. The covariance matrix of the effective channel gains, i.e., the elements on the main diagonal of \mathbf{D}_0 , is given by

$$\begin{aligned} & \mathcal{E} \{ \text{Diag}(\mathbf{D}_0) \text{Diag}(\mathbf{D}_0)^H \} \\ &= NM \mathcal{E} \{ (\mathbf{F}_N \otimes \mathbf{F}_M^H) \text{vec}(\mathbf{A}_0) \text{vec}(\mathbf{A}_0)^H (\mathbf{F}_N \otimes \mathbf{F}_M^H)^H \} \\ &= NM (\mathbf{F}_N \otimes \mathbf{F}_M^H) \mathcal{E} \{ \text{vec}(\mathbf{A}_0) \text{vec}(\mathbf{A}_0)^H \} (\mathbf{F}_N \otimes \mathbf{F}_M^H)^H. \end{aligned} \quad (80)$$

Because the channel gains, $h_{0,p}$, are i.i.d., $\mathcal{E} \{ \text{vec}(\mathbf{A}_0) \text{vec}(\mathbf{A}_0)^H \}$ is a diagonal matrix, where only (P_0+1) of its main diagonal elements are non-zero. Following the same steps as in the proof for Proposition 1, one can show that the product of $(\mathbf{F}_N \otimes \mathbf{F}_M^H)$, a diagonal matrix, and $(\mathbf{F}_N \otimes \mathbf{F}_M^H)^H$ yields a block circulant matrix, which means that $\mathcal{E} \{ \text{Diag}(\mathbf{D}_0) \text{Diag}(\mathbf{D}_0)^H \}$ is a block-circulant matrix, not a diagonal matrix. Therefore, the $|D_0^{k,l}|^2$, $l \in \{0, \dots, M-1\}$, $k \in \{0, \dots, N-1\}$, are correlated, and not independent.

Although the $|D_0^{k,l}|^2$ are not independent, an upper bound on $P_{0,kl}$ can be still found as follows:

$$\begin{aligned} P_{0,kl} &\leq \text{P} \left(|D_0^{\min}|^2 < \frac{\epsilon_0}{\rho(\gamma_0^2 - \gamma_1^2 \epsilon_0)} \right) \\ &\leq \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} \text{P} \left(|D_0^{k,l}|^2 < \frac{\epsilon_0}{\rho(\gamma_0^2 - \gamma_1^2 \epsilon_0)} \right) \\ &\leq MNP \left(|D_0^{0,0}|^2 < \frac{\epsilon_0}{\rho(\gamma_0^2 - \gamma_1^2 \epsilon_0)} \right) \\ &= MN \left(1 - e^{-\frac{\epsilon_0}{\rho(\gamma_0^2 - \gamma_1^2 \epsilon_0)}} \right) \doteq \frac{1}{\rho}. \end{aligned} \quad (81)$$

Since both the upper and lower bounds on the outage probability have the same diversity order, the proof of the lemma is complete.

REFERENCES

[1] Z. Ding *et al.*, "Application of non-orthogonal multiple access in LTE and 5G networks," *IEEE Commun. Mag.*, vol. 55, no. 2, pp. 185–191, Feb. 2017.

[2] W. Shin, M. Vaezi, B. Lee, D. J. Love, J. Lee, and H. V. Poor, "Non-orthogonal multiple access in multi-cell networks: Theory, performance, and practical challenges," *IEEE Commun. Mag.*, vol. 55, no. 10, pp. 176–183, Oct. 2017.

[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," *IEEE Commun. Surveys Tuts.*, vol. 19, no. 2, pp. 721–742, 2nd Quart., 2017.

[4] H. Sari, F. Vanhaverbeke, and M. Moeneclaey, "Multiple access using two sets of orthogonal signal waveforms," *IEEE Commun. Lett.*, vol. 4, no. 1, pp. 4–6, Jan. 2000.

[5] Y. Saito, A. Benjebbour, Y. Kishiyama, and T. Nakamura, "System-level performance evaluation of downlink non-orthogonal multiple access (NOMA)," in *Proc. IEEE 24th Annu. Int. Symp. Pers., Indoor, Mobile Radio Commun. (PIMRC)*, London, U.K., Sep. 2013, pp. 611–615.

[6] Z. Ding, Z. Yang, P. Fan, and H. V. Poor, "On the performance of non-orthogonal multiple access in 5G systems with randomly deployed users," *IEEE Signal Process. Lett.*, vol. 21, no. 12, pp. 1501–1505, Dec. 2014.

[7] Z. Ding, P. Fan, and H. V. Poor, "Impact of user pairing on 5G nonorthogonal multiple-access downlink transmissions," *IEEE Trans. Veh. Technol.*, vol. 65, no. 8, pp. 6010–6023, Aug. 2016.

[8] A. Brighente and S. Tomasin, "Power allocation for non-orthogonal millimeter wave systems with mixed traffic," *IEEE Trans. Wireless Commun.*, vol. 18, no. 1, pp. 432–443, Jan. 2019.

[9] Y. Wu, L. P. Qian, K. Ni, C. Zhang, and X. Shen, "Delay-minimization nonorthogonal multiple access enabled multi-user mobile edge computation offloading," *IEEE J. Sel. Topics Signal Process.*, vol. 13, no. 3, pp. 392–407, Jun. 2019.

[10] Z. Ding, P. Fan, and H. V. Poor, "Random beamforming in millimeter-wave NOMA networks," *IEEE Access*, vol. 5, pp. 7667–7681, 2017.

[11] Y. Zhou, V. W. S. Wong, and R. Schober, "Coverage and rate analysis of millimeter wave NOMA networks with beam misalignment," *IEEE Trans. Wireless Commun.*, vol. 17, no. 12, pp. 8211–8227, Dec. 2018.

[12] Z. Ding and H. V. Poor, "Design of massive-MIMO-NOMA with limited feedback," *IEEE Signal Process. Lett.*, vol. 23, no. 5, pp. 629–633, May 2016.

[13] R. Chopra, C. R. Murthy, H. A. Suraweera, and E. G. Larsson, "Analysis of nonorthogonal training in massive MIMO under channel aging with SIC receivers," *IEEE Signal Process. Lett.*, vol. 26, no. 2, pp. 282–286, Feb. 2019.

[14] A. Maatouk, E. Çalıřkan, M. Koca, M. Assaad, G. Gui, and H. Sari, "Frequency-domain NOMA with two sets of orthogonal signal waveforms," *IEEE Commun. Lett.*, vol. 22, no. 5, pp. 906–909, May 2018.

[15] Y. Sun, D. W. K. Ng, Z. Ding, and R. Schober, "Optimal joint power and subcarrier allocation for full-duplex multicarrier non-orthogonal multiple access systems," *IEEE Trans. Commun.*, vol. 65, no. 3, pp. 1077–1091, Mar. 2017.

[16] L. Yin, W. O. Popoola, X. Wu, and H. Haas, "Performance evaluation of non-orthogonal multiple access in visible light communication," *IEEE Trans. Commun.*, vol. 64, no. 12, pp. 5162–5175, Dec. 2016.

[17] H. Marshoud, V. M. Kapinas, G. K. Karagiannidis, and S. Muhaidat, "Non-orthogonal multiple access for visible light communications," *IEEE Photon. Technol. Lett.*, vol. 28, no. 1, pp. 51–54, Jan. 2016.

[18] Z. Ding, P. Fan, and H. V. Poor, "Impact of non-orthogonal multiple access on the offloading of mobile edge computing," *IEEE Trans. Commun.*, vol. 67, no. 1, pp. 375–390, Jan. 2019.

[19] *Study on Downlink Multiuser Superposition Transmission for LTE*, document TR 36.859, 3GPP, Mar. 2015.

[20] *Technical Specification Group Radio Access Network; Evolved Universal Terrestrial Radio Access (E-UTRA); Physical Channels and Modulation (Release 15)*, document TS 36.211, 3GPP, Jan. 2019.

[21] *Study on Non-Orthogonal Multiple Access (NOMA) for NR (Release 16)*, document TR 38.812, 3GPP, Dec. 2018.

[22] B. Di, L. Song, Y. Li, and Z. Han, "V2X meets NOMA: Non-orthogonal multiple access for 5G-enabled vehicular networks," *IEEE Wireless Commun.*, vol. 24, no. 6, pp. 14–21, Dec. 2017.

[23] Y. Chen, L. Wang, Y. Ai, B. Jiao, and L. Hanzo, "Performance analysis of NOMA-SM in vehicle-to-vehicle massive MIMO channels," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 12, pp. 2653–2666, Dec. 2017.

[24] R. Hadani and A. Monk, "OTFS: A new generation of modulation addressing the challenges of 5G," 2018, *arXiv:1802.02623*. [Online]. Available: <https://arxiv.org/abs/1802.02623>

[25] R. Hadani *et al.*, "Orthogonal time frequency space modulation," in *Proc. IEEE Wireless Commun. Netw. Conf. (WCNC)*, San Francisco, CA, USA, Mar. 2017, pp. 1–6.

[26] K. R. Murali and A. Chockalingam, "On OTFS modulation for high-Doppler fading channels," in *Proc. Inf. Theory Appl. Workshop (ITA)*, San Diego, CA, USA, Feb. 2018, pp. 1–10.

[27] P. Raviteja, Y. Hong, E. Viterbo, and E. Biglieri, "Practical pulse-shaping waveforms for reduced-cyclic-prefix OTFS," *IEEE Trans. Veh. Tech.*, vol. 68, no. 1, pp. 957–961, Jan. 2019.

[28] P. Raviteja, K. T. Phan, Y. Hong, and E. Viterbo, "Interference cancellation and iterative detection for orthogonal time frequency space modulation," *IEEE Trans. Wireless Commun.*, vol. 17, no. 10, pp. 6501–6515, Oct. 2018.

[29] G. D. Surabhi, R. M. Augustine, and A. Chockalingam, "On the diversity of uncodet OTFS modulation in doubly-dispersive channels," 2018, *arXiv:1808.07747*. [Online]. Available: <https://arxiv.org/abs/1808.07747>

[30] V. Khammammetti and S. K. Mohammed, "OTFS-based multiple-access in high Doppler and delay spread wireless channels," *IEEE Wireless Commun. Lett.*, vol. 8, no. 2, pp. 528–531, Apr. 2019.

[31] A. RezaazadehReyhani, A. Farhang, M. Ji, R. R. Chen, and B. Farhang-Boroujeny, "Analysis of discrete-time MIMO OFDM-based orthogonal time frequency space modulation," in *Proc. IEEE Int. Conf. Communicat. (ICC)*, Kansas City, MO, USA, May 2018, pp. 1–6.

- 1264 [32] M. K. Ramachandran and A. Chockalingam, "MIMO-OTFS in high-
1265 Doppler fading channels: Signal detection and channel estimation,"
1266 in *Proc. IEEE GLOBECOM*, Kansas City, MO, USA, Dec. 2018,
1267 pp. 206–212.
- 1268 [33] P. Raviteja, E. Viterbo, and Y. Hong, "OTFS performance on static
1269 multipath channels," *IEEE Wireless Commun. Lett.*, vol. 8, no. 3,
1270 pp. 745–748, Jun. 2019.
- 1271 [34] D. Falconer, S. L. Ariyavisitakul, A. Benyamin-Seeyar, and B. Eidson,
1272 "Frequency domain equalization for single-carrier broadband wireless
1273 systems," *IEEE Commun. Mag.*, vol. 40, no. 4, pp. 58–66, Apr. 2002.
- 1274 [35] J. Louveaux, L. Vandendorpe, and T. Sartenar, "Cyclic prefixed single
1275 carrier and multicarrier transmission: Bit rate comparison," *IEEE*
1276 *Commun. Lett.*, vol. 7, no. 4, pp. 180–182, Apr. 2003.
- 1277 [36] L. Zheng and D. N. C. Tse, "Diversity and multiplexing: A fundamental
1278 tradeoff in multiple-antenna channels," *IEEE Trans. Inf. Theory*, vol. 49,
1279 no. 5, pp. 1073–1096, May 2003.
- 1280 [37] B. Devillers, "Cyclic prefixed block transmission for wireless commu-
1281 nications: Performance analysis and optimization," Ph.D. dissertation,
1282 Université catholique de Louvain, Ottignies-Louvain-la-Neuve, Belgium,
1283 2009.
- 1284 [38] B. Devillers, J. Louveaux, and L. Vandendorpe, "About the diversity in
1285 cyclic prefixed single-carrier systems," *Phys. Commun.*, vol. 1, no. 4,
1286 pp. 266–276, Dec. 2008.
- 1287 [39] J. Jeon, "NR wide bandwidth operations," *IEEE Commun. Mag.*, vol. 56,
1288 no. 3, pp. 42–46, Mar. 2018.
- 1289 [40] C. Sexton, N. Marchetti, and L. A. DaSilva, "Customization and
1290 trade-offs in 5G RAN slicing," *IEEE Commun. Mag.*, vol. 57, no. 4,
1291 pp. 116–122, Apr. 2019.
- 1292 [41] C. Liaskos, S. Nie, A. Tsioliaridou, A. Pitsillides, S. Ioannidis, and
1293 I. Akyildiz, "A new wireless communication paradigm through software-
1294 controlled metasurfaces," *IEEE Commun. Mag.*, vol. 56, no. 9,
1295 pp. 162–169, Sep. 2018.

the Vodafone Foundation for Research in Mobile Communications, the
2006 UBC Killam Research Prize, the 2007 Wilhelm Friedrich Bessel
Research Award of the Alexander von Humboldt Foundation, the 2008
Charles McDowell Award for Excellence in Research from UBC, the 2011
Alexander von Humboldt Professorship, the 2012 NSERC E.W.R. Steacie
Fellowship, and the 2017 Wireless Communications Recognition Award by
the IEEE Wireless Communications Technical Committee. He was listed as
a 2017 Highly Cited Researcher by the Web of Science. He is also the Chair
of the Steering Committee of IEEE TRANSACTIONS ON MOLECULAR, BIO-
LOGICAL AND MULTI-SCALE COMMUNICATION, a member of the Editorial
Board of PROCEEDINGS OF THE IEEE, a Member-at-Large of the Board
of Governors of ComSoc, and the ComSoc Director of journals. He is also a
Distinguished Lecturer of the IEEE Communications Society (ComSoc). From
2012 to 2015, he served as the Editor-in-Chief of IEEE TRANSACTIONS ON
COMMUNICATIONS.



Pingzhi Fan (M'93–SM'99–F'15) received the M.Sc. degree in computer science from Southwest Jiaotong University, China, in 1987, and the Ph.D. degree in electronic engineering from Hull University, U.K., in 1994.

He was the Chief Scientist of the National 973 Research Project (MoST) from 2012 to 2016. He is currently a Professor and the Director of the Institute of Mobile Communications, Southwest Jiaotong University. He has been a Visiting Professor with Leeds University, U.K., since 1997, and has been a Guest Professor with Shanghai Jiaotong University since 1999. He has over 280 research papers published in various international journals and eight books (including edited). He is the inventor of 22 granted patents. His research interests include vehicular communications, wireless networks for big data, and signal design and coding. He is also a fellow of IET, CIE, and CIC. He was a recipient of the U.K. ORS Award in 1992 and the Outstanding Young Scientist Award (NSFC) in 1998. He has served as the general chair or the TPC chair of a number of international conferences. He is also the Founding Chair of IEEE VTS BJ Chapter, IEEE ComSoc CD Chapter, and IEEE Chengdu Section. He is also the guest editor or editorial member of several international journals. He has also served as the Board Member of IEEE Region 10, IET (IEE) Council, and IET Asia-Pacific Region. He is also an IEEE VTS Distinguished Lecturer (2015–2019).

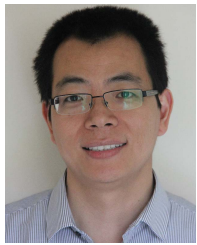


H. Vincent Poor (M'77–SM'82–F'87) received the Ph.D. degree in EECS from Princeton University in 1977.

From 1977 to 1990, he was on the faculty of the University of Illinois at Urbana–Champaign. Since 1990, he has been on the faculty at Princeton University, where he is currently the Michael Henry Strater University Professor of electrical engineering. From 2006 to 2016, he served as the Dean of the School of Engineering and Applied Science, Princeton University. He has also held visiting appointments at several other universities, including most recently at Berkeley and Cambridge. His research interests are in the areas of information theory and signal processing, and their applications in wireless networks, energy systems, and related fields. Among his publications in these areas is the recent book *Multiple Access Techniques for 5G Wireless Networks and Beyond*. (Springer, 2019). He is also a member of the National Academy of Engineering and the National Academy of Sciences. He is also a Foreign Member of the Chinese Academy of Sciences, the Royal Society, and other national and international academies. He was a recipient of the Marconi and Armstrong Awards of the IEEE Communications Society in 2007 and 2009, respectively. Recent recognition of his work includes the 2017 IEEE Alexander Graham Bell Medal, the 2019 ASEE Benjamin Garver Lamme Award, the D.Sc. (*honoris causa*) from Syracuse University awarded in 2017, and the D.Eng. (*honoris causa*) from the University of Waterloo awarded in 2019.

AQ:6

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318



Zhiguo Ding (S'03–M'05–SM'15) received the B.Eng. degree in electrical engineering from the Beijing University of Posts and Telecommunications in 2000 and the Ph.D. degree in electrical engineering from Imperial College London in 2005.

From July 2005 to April 2018, he was with Queen's University Belfast, Imperial College, Newcastle University, and Lancaster University. Since April 2018, he has been a Professor of communications with The University of Manchester. From October 2012 to September 2018, he was an

Academic Visitor with Princeton University. His research interests are 5G networks, game theory, cooperative and energy harvesting networks, and statistical signal processing. He was a recipient of the Best Paper Award at IET ICWMC-2009 and IEEE WCSP-2014, the EU Marie Curie Fellowship (2012–2014), the Top IEEE TVT Editor 2017, the 2018 IEEE Communication Society Heinrich Hertz Award, the 2018 IEEE Vehicular Technology Society Jack Neubauer Memorial Award, and the 2018 IEEE Signal Processing Society Best Signal Processing Letter Award. He was an Editor of IEEE WIRELESS COMMUNICATIONS LETTERS and IEEE COMMUNICATIONS LETTERS from 2013 to 2016. He has been serving as an Editor for IEEE TRANSACTIONS ON COMMUNICATIONS, IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY, and *Journal of Wireless Communications and Mobile Computing*.

1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334



Robert Schober (M'01–SM'08–F'10) received the Diploma (Univ.) and Ph.D. degrees in electrical engineering from the Friedrich-Alexander University of Erlangen-Nuremberg (FAU), Germany, in 1997 and 2000, respectively.

From 2002 to 2011, he was a Professor and the Canada Research Chair with The University of British Columbia (UBC), Vancouver, Canada. Since January 2012, he has been an Alexander von Humboldt Professor and the Chair for Digital Communication with FAU. His research interests fall

into the broad areas of communication theory, wireless communications, and statistical signal processing. He is also a fellow of the Canadian Academy of Engineering and the Engineering Institute of Canada. He was a recipient of several awards for his work, including the 2002 Heinz Maier-Leibnitz Award of the German Science Foundation (DFG), the 2004 Innovations Award of

1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399

AUTHOR QUERIES

AUTHOR PLEASE ANSWER ALL QUERIES

PLEASE NOTE: We cannot accept new source files as corrections for your article. If possible, please annotate the PDF proof we have sent you with your corrections and upload it via the Author Gateway. Alternatively, you may send us your corrections in list format. You may also upload revised graphics via the Author Gateway.

AQ:1 = Please supply index terms/keywords for your paper. To download the IEEE Taxonomy, go to http://www.ieee.org/documents/taxonomy_v101.pdf.

AQ:2 = Author: Please confirm or add details for any funding or financial support for the research of this article.

AQ:3 = Please provide the postal code for The University of Manchester, Friedrich-Alexander-University Erlangen-Nurnberg, and Southwest Jiaotong University.

AQ:4 = Note that if you require corrections/changes to tables or figures, you must supply the revised files, as these items are not edited for you.

AQ:5 = Please provide the department name for Ref. [37].

AQ:6 = Current affiliation in the biography of Zhiguo Ding does not match the First Footnote. Please check and correct where needed.

OTFS-NOMA: An Efficient Approach for Exploiting Heterogenous User Mobility Profiles

Zhiguo Ding¹, Senior Member, IEEE, Robert Schober, Fellow, IEEE,
Pingzhi Fan, Fellow, IEEE, and H. Vincent Poor², Fellow, IEEE

Abstract—This paper considers a challenging communication scenario, in which users have heterogenous mobility profiles, e.g., some users are moving at high speeds and some users are static. A new non-orthogonal multiple-access (NOMA) transmission protocol that incorporates orthogonal time frequency space (OTFS) modulation is proposed. Thereby, users with different mobility profiles are grouped together for the implementation of NOMA. The proposed OTFS-NOMA protocol is shown to be applicable to both uplink and downlink transmission, where sophisticated transmit and receive strategies are developed to remove inter-symbol interference and harvest both multi-path and multi-user diversity. Analytical results demonstrate that both the high-mobility and the low-mobility users benefit from the application of OTFS-NOMA. In particular, the use of NOMA allows the spreading of the high-mobility users' signals over a large amount of time-frequency resources, which enhances the OTFS resolution and improves the detection reliability. In addition, OTFS-NOMA ensures that low-mobility users have access to bandwidth resources which in conventional OTFS-orthogonal multiple access (OTFS-OMA) would be solely occupied by the high-mobility users. Thus, OTFS-NOMA improves the spectral efficiency and reduces latency.

Index Terms—XXXXX.

I. INTRODUCTION

NON-ORTHOGONAL multiple access (NOMA) has been recognized as a paradigm shift for the design of multiple access techniques for the next generation of wireless networks [1]–[4]. Many existing works on NOMA have

focused on scenarios with low-mobility users, where users with different channel conditions or quality of service (QoS) requirements are grouped together for the implementation of NOMA. For example, in power-domain NOMA, a base station serves two users simultaneously [5], [6]. In particular, the base station first orders the users according to their channel conditions, where the ‘weak user’ which has a poorer connection to the base station is generally allocated more transmission power and the other user, referred to as the ‘strong user’, is allocated less power. As such, the two users can be served in the same time-frequency resource, which improves the spectral efficiency compared to orthogonal multiple access (OMA). In the case that users have similar channel conditions, grouping users with different QoS requirements can facilitate the implementation of NOMA and effectively exploit the potential of NOMA [7]–[9]. Various existing studies have shown that the NOMA principle can be applied to different communication networks, such as millimeter-wave networks [10], [11], massive multiple-input multiple-output (MIMO) systems [12], [13], hybrid multiple access systems [14], [15], visible light communication networks [16], [17], and mobile edge computing [18]. We also note that various standardization efforts have been made to facilitate the implementation of NOMA in practical systems. For example, a study for the application of NOMA for downlink transmission, termed multi-user superposition transmission (MUST), was carried out for the 3rd Generation Partnership Project (3GPP) Release 14, where 15 different forms of MUST were proposed and compared [19]. After this study was completed, MUST was formally included in 3GPP Release 15 which is also referred to as Evolved Universal Terrestrial Radio Access (E-UTRA) [20]. A study for the application of NOMA for uplink transmission has been recently carried out for 3GPP Release 16, where more than 20 different forms of NOMA have been proposed by various companies [21].

This paper considers the application of NOMA to a challenging communication scenario, where users have heterogeneous mobility profiles. Different from the existing works in [22], [23], the use of orthogonal time frequency space (OTFS) modulation is considered in this paper because of its superior performance in scenarios with doubly-dispersive channels [24]–[26]. Recall that the key idea of OTFS is to use the delay-Doppler plane, where the users' signals are

Manuscript received April 4, 2019; revised June 5, 2019; accepted July 14, 2019. The work of Z. Ding was supported by the UK Engineering and Physical Sciences Research Council under grant number EP/P009719/2 and by H2020-MSCA-RISE-2015 under grant number 690750. The work of P. Fan was supported by the National Natural Science Foundation of China under grant number 61731017, and the 111 Project (No.111-2-14). The work of H. V. Poor was supported by the U.S. National Science Foundation under Grants CCF-093970 and CCF-1513915. The associate editor coordinating the review of this article and approving it for publication was V. Raghavan. (Corresponding author: Zhiguo Ding.)

Z. Ding is with the Department of Electrical Engineering, Princeton University, Princeton, NJ 08544 USA, and also with the School of Electrical and Electronic Engineering, The University of Manchester, Manchester, U.K. (e-mail: zhiguo.ding@manchester.ac.uk).

R. Schober is with the Institute for Digital Communications, Friedrich-Alexander-University Erlangen-Nurnberg (FAU), Erlangen, Germany (e-mail: robert.schober@fau.de).

P. Fan is with the Institute of Mobile Communications, Southwest Jiaotong University, Chengdu, China (e-mail: pingzhifan@foxmail.com).

H. V. Poor is with the Department of Electrical Engineering, Princeton University, Princeton, NJ 08544 USA (e-mail: poor@princeton.edu).

Digital Object Identifier 10.1109/TCOMM.2019.2932934

orthogonally placed. Compared to conventional modulation schemes, such as orthogonal frequency-division multiplexing (OFDM), OTFS offers the benefit that the time-invariant channel gains in the delay-Doppler plane can be utilized, which simplifies channel estimation and signal detection in high-mobility scenarios. The impact of pulse-shaping waveforms on the performance of OTFS was studied in [27], and the design of interference cancellation and iterative detection for OTFS was investigated in [28]. The diversity gain achieved by OTFS was studied in [29], and the application of OTFS to multiple access was proposed in [30]. In [31] and [32], the concept of OTFS was combined with MIMO, which revealed that the use of spatial degrees of freedom can further enhance the performance of OTFS.

This paper considers the application of OTFS to NOMA communication networks, where the coexistence of NOMA and OTFS is investigated. In particular, this paper makes the following contributions:

- 1) A spectrally efficient OTFS-NOMA transmission protocol is proposed by grouping users with different mobility profiles for the implementation of NOMA. On the one hand, users with high mobility are served in the delay-Doppler plane, and their signals are modulated by OTFS. On the other hand, users with low mobility are served in the time-frequency plane, and their signals are modulated in a manner similar to conventional OFDM.
- 2) The proposed new OTFS-NOMA protocol is applied to both uplink and downlink transmission, where different rate and power allocation policies are used to suppress multiple access interference. In addition, sophisticated equalization techniques, such as the frequency-domain zero-forcing linear equalizer (FD-LE) and the decision feedback equalizer (FD-DFE), are employed to remove the inter-symbol interference in the delay-Doppler plane. The impact of the developed equalization techniques on OTFS-NOMA is analyzed by using the outage probability as the performance criterion. Strategies to harvest multi-path diversity and multi-user diversity are also introduced, which can further improve the outage performance of OTFS-NOMA transmission.
- 3) The developed analytical results demonstrate that both the high-mobility and the low-mobility users benefit from the proposed OTFS-NOMA scheme. The use of NOMA allows the high-mobility users' signals to be spread over a large amount of time-frequency resources without degrading the spectral efficiency. As a result, the OTFS resolution, which determines whether the users' channels can be accurately located in the delay-Doppler plane, is enhanced significantly, and therefore, the reliability of detecting the high-mobility users' signals is improved. We note that, in OTFS-OMA, enhancing the OTFS resolution implies that a large amount of time and frequency resources are solely occupied by the high-mobility users, which reduces the overall spectral efficiency since the high-mobility users' channel conditions are typically weaker than those of the low-mobility users. In contrast, the use of OTFS-NOMA ensures that the low-mobility users can access the

bandwidth resources which would be solely occupied by the high-mobility users in the OMA mode. Hence, OTFS-NOMA improves spectral efficiency and reduces latency, as with OTFS-OMA the low-mobility users may have to wait for a long time before the scarce bandwidth resources occupied by the high-mobility users become available. In addition, we note that for the low-mobility users, using OFDM yields the same reception reliability as using OTFS, as pointed out in [33]. Therefore, the proposed OTFS-NOMA scheme, which serves the low-mobility users in the time-frequency plane and modulates the low-mobility users' signals in a manner similar to OFDM, offers the same reception reliability as OTFS-OMA, which serves the low-mobility users in the delay-Doppler plane and modulates the low-mobility users' signals by OTFS. However, OTFS-NOMA has the benefit of reduced system complexity because the use of the complicated OTFS transforms is avoided.

II. FOUNDATIONS OF OTFS-NOMA

A. Time-Frequency Plane and Delay-Doppler Plane

The key idea of OTFS-NOMA is to efficiently use both the time-frequency plane and the delay-Doppler plane. A discrete time-frequency plane is obtained by sampling at intervals of T s and Δf Hz as follows:

$$\Lambda_{\text{TF}} = \{(nT, m\Delta f), n=0, \dots, N-1, m=0, \dots, M-1\}, \quad (1)$$

and the corresponding discrete delay-Doppler plane is given by

$$\Lambda_{\text{DD}} = \left\{ \left(\frac{k}{NT}, \frac{l}{M\Delta f} \right), k=0, \dots, N-1, l=0, \dots, M-1 \right\}, \quad (2)$$

where N and M denote the total number of time intervals and the total number of frequency subchannels, respectively. The choices for T and Δf are determined by the channel characteristics, as will be explained in the following subsection.

B. Channel Model

This paper considers a multi-user communication network in which one base station communicates with $(K+1)$ users, denoted by U_i , $0 \leq i \leq K$. Denote U_i 's channel response in the delay-Doppler plane by $h_i(\tau, \nu)$, where τ denotes the delay and ν denotes the Doppler shift. OTFS uses the sparsity feature of a wireless channel in the delay-Doppler plane, i.e., there are a small number of propagation paths between a transmitter and a receiver [24], [25], [28], which means that $h_i(\tau, \nu)$ can be expressed as follows:

$$h_i(\tau, \nu) = \sum_{p=0}^{P_i} h_{i,p} \delta(\tau - \tau_{i,p}) \delta(\nu - \nu_{i,p}), \quad (3)$$

where $(P_i + 1)$ denotes the number of propagation paths, and $h_{i,p}$, $\tau_{i,p}$, and $\nu_{i,p}$ denote the complex Gaussian channel gain,¹ the delay, and the Doppler shift associated with the

¹The Gaussian assumption has been commonly used in the OTFS literature [26]–[29] since each channel gain (or each tap of the delay-Doppler impulse response) represents a cluster of reflectors with specific delay and Doppler characteristics.

178 p -th propagation path, respectively. We assume that the $h_{i,p}$,
 179 $0 \leq p \leq P_i$, are independent and identically distributed (i.i.d.)
 180 random variables,² i.e., $h_{i,p} \sim CN\left(0, \frac{1}{P_i+1}\right)$, which means
 181 $\sum_{p=0}^{P_i} \mathcal{E}\{|h_{i,p}|^2\} = 1$, where $\mathcal{E}\{\cdot\}$ denotes the expectation
 182 operation. The discrete delay and Doppler tap indices for the
 183 p -th path of $h_i(\tau, \nu)$, denoted by $l_{\tau_{i,p}}$ and $k_{\nu_{i,p}}$, respectively,
 184 are given by [28]

$$185 \quad \tau_{i,p} = \frac{l_{\tau_{i,p}} + \tilde{l}_{\tau_{i,p}}}{M\Delta f}, \quad \nu_{i,p} = \frac{k_{\nu_{i,p}} + \tilde{k}_{\nu_{i,p}}}{NT}, \quad (4)$$

186 where $\tilde{l}_{\tau_{i,p}}$ and $\tilde{k}_{\nu_{i,p}}$ denote the fractional delay and the
 187 fractional Doppler shift, respectively.

188 The construction of Λ_{TF} and Λ_{DD} needs to ensure that T
 189 is not smaller than the maximal delay spread, and Δf is not
 190 smaller than the largest Doppler shift, i.e., $T \geq \max\{\tau_{i,p}, 0 \leq$
 191 $p \leq P_i, 0 \leq i \leq K\}$ and $\Delta f \geq \max\{\nu_{i,p}, 0 \leq p \leq P_i, 0 \leq$
 192 $i \leq K\}$. In addition, the choices of N and M affect the
 193 OTFS resolution, which determines whether $h_i(\tau, \nu)$ can be
 194 accurately located in the discrete delay-Doppler plane. In par-
 195 ticular, M and N need to be sufficiently large to approximately
 196 achieve ideal OTFS resolution, which ensures that $\tilde{l}_{\tau_{i,p}} =$
 197 $\tilde{k}_{\nu_{i,p}} = 0$, such that the interference caused by fractional delay
 198 and Doppler shift is effectively suppressed [24].

199 C. General Principle of OTFS-NOMA

200 To facilitate the illustration of the general principle of
 201 OTFS-NOMA, we first briefly describe OTFS-OMA, the
 202 benchmark scheme used in this paper. In OTFS-OMA, there
 203 is no spectrum sharing between the high-mobility users and
 204 the low-mobility users, i.e., if OTFS is used to serve the high-
 205 mobility users, the NT time intervals and the $M\Delta f$ frequency
 206 subchannels are occupied by the high-mobility users and the
 207 low-mobility users cannot be served in these resource blocks.
 208 The general principle of the proposed OTFS-NOMA scheme is
 209 to exploit both the delay-Doppler plane and the time-frequency
 210 plane, where users with heterogenous mobility profiles are
 211 grouped together and served simultaneously. On the one hand,
 212 for the users with high mobility, their signals are placed in
 213 the delay-Doppler plane, which means that the time-invariant
 214 channel gains in the delay-Doppler plane can be exploited. It is
 215 worth pointing out that in order to ensure that the channels
 216 can be located in the delay-Doppler plane, both N and M
 217 need to be large, which is a disadvantage of OTFS-OMA,
 218 since a significant number of frequency channels (e.g., $M\Delta f$)
 219 are occupied for a long time (e.g., NT) by the high-mobility
 220 users whose channel conditions can be quite weak. The use of
 221 OTFS-NOMA facilitates spectrum sharing and hence ensures
 222 that the high-mobility users' signals can be spread over a
 223 large amount of time-frequency resources without degrading
 224 the spectral efficiency.

225 On the other hand, for the users with low mobility, their
 226 signals are placed in the time-frequency plane. The inter-
 227 ference between the users with different mobility profiles

²In order to simplify the performance analysis, we assume that the users' channels are i.i.d. In practice, it is likely that the high-mobility users' channel conditions are worse than the low-mobility users' channel conditions. This channel difference is beneficial for the implementation of NOMA, and hence can further increase the performance gain of OTFS-NOMA over OTFS-OMA.

228 is managed by using the principle of NOMA. As a result,
 229 compared to OTFS-OMA, OTFS-NOMA improves the overall
 230 spectral efficiency since it encourages spectrum sharing among
 231 users with different mobility profiles and avoids that the
 232 bandwidth resources are solely occupied by the high-mobility
 233 users which might have weak channel conditions. In addition,
 234 the complexity of detecting the low-mobility users' signals is
 235 reduced, compared to OTFS-OMA which serves all users in
 236 the delay-Doppler plane.

237 In this paper, we assume that, among $(K + 1)$ users,
 238 U_0 is a user with high mobility, and the remaining K users,
 239 U_i for $1 \leq i \leq K$, are low-mobility users, which are
 240 referred to as 'NOMA' users.³ For OTFS-OMA, we assume
 241 that U_0 solely occupies all NM resource blocks in Λ_{DD} .
 242 In OTFS-NOMA, U_i , for $1 \leq i \leq K$, are opportunistic
 243 NOMA users and their signals are placed in Λ_{TF} . The design
 244 of downlink OTFS-NOMA transmission will be discussed
 245 in detail in Sections III, IV, and V. The application of
 246 OTFS-NOMA for uplink transmission will be considered in
 247 Section VI only briefly, due to space limitations.

248 III. DOWNLINK OTFS-NOMA - SYSTEM MODEL

249 In this section, the OTFS-NOMA downlink transmission
 250 protocol is described. In particular, assume that the base
 251 station sends NM symbols to U_0 , denoted by $x_0[k, l]$, $k \in$
 252 $\{0, \dots, N-1\}$, $l \in \{0, \dots, M-1\}$. By using the inverse
 253 symplectic finite Fourier transform (ISFFT), the high-mobility
 254 user's symbols placed in the delay-Doppler plane are converted
 255 to NM symbols in the time-frequency plane as follows [24]:

$$256 \quad X_0[n, m] = \frac{1}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} x_0[k, l] e^{j2\pi\left(\frac{kn}{N} - \frac{ml}{M}\right)}, \quad (5)$$

257 where $n \in \{0, \dots, N-1\}$ and $m \in \{0, \dots, M-1\}$. We note
 258 that the NM time-frequency signals can be viewed as N
 259 OFDM symbols containing M signals each. We assume that a
 260 rectangular window is applied to the transmitted and received
 261 signals.

262 The NOMA users' signals are placed directly in the time-
 263 frequency plane, and are superimposed with the high-mobility
 264 user's signals, $X_0[n, m]$. With NM orthogonal resource
 265 blocks available in the time-frequency plane, there are different
 266 ways for the K users to share the resource blocks. For
 267 illustration purposes, we assume that M users are selected
 268 from the K opportunistic NOMA users,⁴ where each NOMA

³We note that the principle of OTFS-NOMA can be extended to the case where multiple high-mobility users are served in the delay-Doppler plane. In this case, the NM signals in the delay-Doppler plane belong to different high-mobility users and OTFS is used as a type of multiple access technique [24], [30]. For downlink transmission, this change has no impact on the proposed detection schemes and the analytical results developed in this paper. For uplink transmission, the results developed in this paper are applicable to the case with multiple high-mobility users if the adaptive-rate transmission scheme proposed in Section VI is employed.

⁴The same M users can be scheduled as long as the users' channels do not change in the delay-Doppler plane. Otherwise, a new set of M users may be selected from the K opportunistic users. We also note that the number of the opportunistic users is assumed to be larger than the number of the frequency subchannels ($K \geq M$), which can be justified by a spectrum crunch scenario, i.e., there are not sufficient bandwidth resources available to support a large number of mobile devices.

user is to occupy one frequency subchannel and receive N information bearing symbols, denoted by $x_i(n)$, for $1 \leq i \leq M$ and $0 \leq n \leq N - 1$. The criterion employed for user scheduling and its impact on the performance of OTFS-NOMA will be discussed in Section V. Denote the time-frequency signals to be sent to U_i by $X_i[n, m]$, $1 \leq i \leq M$. The following mapping scheme is used in this paper⁵:

$$X_i[n, m] = \begin{cases} x_i(n) & \text{if } m = i - 1 \\ 0 & \text{otherwise,} \end{cases} \quad (6)$$

for $1 \leq i \leq M$ and $0 \leq n \leq N - 1$.

The base station superimposes U_0 's time-frequency signals with the NOMA users' signals as follows:

$$X[n, m] = \frac{\gamma_0}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} x_0[k, l] e^{j2\pi(\frac{kn}{N} - \frac{ml}{M})} + \sum_{i=1}^M \gamma_i X_i[n, m], \quad (7)$$

where γ_i denotes the NOMA power allocation coefficient of user i , and $\sum_{i=0}^M \gamma_i^2 = 1$.

The transmitted signal at the base station is obtained by applying the Heisenberg transform to $X[n, m]$. By assuming perfect orthogonality between the transmit and receive pulses, the received signal at U_i in the time-frequency plane can be modelled as follows [24], [25], [28]:

$$Y_i[n, m] = H_i[n, m]X[n, m] + W_i[n, m], \quad (8)$$

where $W_i(n, m)$ is the white Gaussian noise in the time-frequency plane, and $H_i(n, m) = \iint h_i(\tau, \nu) e^{j2\pi\nu n T} e^{-j2\pi(\nu + m\Delta f)\tau} d\tau d\nu$.

IV. DOWNLINK OTFS-NOMA - DETECTING THE HIGH-MOBILITY USER'S SIGNALS

For the proposed downlink OTFS-NOMA scheme, U_0 directly detects its signals in the delay-Doppler plane by treating the NOMA users' signals as noise. In particular, in order to detect U_0 's signals, the symplectic finite Fourier transform (SFFT) is applied to $Y_0[n, m]$ to obtain the delay-Doppler estimates as follows:

$$\begin{aligned} y_0[k, l] &= \frac{1}{NM} \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} Y_0[n, m] e^{-j2\pi(\frac{nk}{N} - \frac{ml}{M})} \\ &= \frac{1}{NM} \sum_{q=0}^M \gamma_q \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} x_q[n, m] h_{w,0} \left(\frac{k-n}{NT}, \frac{l-m}{M\Delta f} \right) \\ &\quad + z_0[k, l], \end{aligned} \quad (9)$$

where q denotes the user index, $z_0[k, l]$ is complex Gaussian noise, $x_q[k, l]$, $1 \leq q \leq M$, denotes the delay-Doppler representation of $X_q[n, m]$ and can be obtained by applying the SFFT to $X_q[n, m]$, the channel $h_{w,0}(\nu', \tau')$ is given by

$$h_{w,0}(\nu', \tau') = \iint h_i(\tau, \nu) w(\nu' - \nu, \tau' - \tau) e^{-j2\pi\nu\tau} d\tau d\nu, \quad (10)$$

⁵We note that mapping schemes different from (6) can also be used. For example, if N users are scheduled and each user is to occupy one time slot and receives an OFDM-like symbol containing M signals, we can set $X_i[n, m] = x_i(m)$, for $n = i - 1$.

and $w(\nu, \tau) = \sum_{c=0}^{N-1} \sum_{d=0}^{M-1} e^{-j2\pi(\nu c T - \tau d \Delta f)}$. To simplify the analysis, the power of the complex-Gaussian distributed noise is assumed to be normalized, i.e., $z_i[k, l] \sim CN(0, 1)$, where $CN(a, b)$ denotes a complex Gaussian distributed random variable with mean a and variance b .

By applying the channel model in (3), the relationship between the transmitted signals and the observations in the delay-Doppler plane can be expressed as follows [24], [25], [28]:

$$y_0[k, l] = \sum_{q=0}^M \gamma_q \sum_{p=0}^{P_0} h_{0,p} x_q[(k - k_{\nu_{0,p}})_N, (l - l_{\tau_{0,p}})_M] + z_0[k, l], \quad (11)$$

where $(\cdot)_N$ denotes the modulo N operator. As in [29]–[31], we assume that N and M are sufficiently large to ensure that both $\tilde{k}_{\nu_{0,p}}$ and $\tilde{l}_{\tau_{0,p}}$ are zero, i.e., there is no interference caused by fractional delay or fractional Doppler shift. We note that for OTFS-OMA, increasing N and M can significantly reduce spectral efficiency, whereas the use of large N and M becomes possible for OTFS-NOMA because of the spectrum sharing of users with different mobility profiles.

Define $\mathbf{y}_{0,k} = [y_0[k, 0] \cdots y_0[k, M-1]]^T$ and $\mathbf{y}_0 = [\mathbf{y}_{0,0}^T \cdots \mathbf{y}_{0,N-1}^T]^T$. Similarly, the signal vector \mathbf{x}_i and the noise vector \mathbf{z}_0 are constructed from $x_i[k, l]$ and $z_0[k, l]$, respectively. Based on (11), the system model can be expressed in matrix form as follows:

$$\mathbf{y}_0 = \underbrace{\gamma_0 \mathbf{H}_0 \mathbf{x}_0 + \sum_{q=1}^M \gamma_q \mathbf{H}_0 \mathbf{x}_q}_{\text{Interference and noise terms}} + \mathbf{z}_0, \quad (12)$$

where \mathbf{H}_0 is a block-circulant matrix and defined as follows:

$$\mathbf{H}_0 = \begin{bmatrix} \mathbf{A}_{0,0} & \mathbf{A}_{0,N-1} & \cdots & \mathbf{A}_{0,2} & \mathbf{A}_{0,1} \\ \mathbf{A}_{0,1} & \mathbf{A}_{0,0} & \ddots & \mathbf{A}_{0,3} & \mathbf{A}_{0,2} \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \mathbf{A}_{0,N-2} & \mathbf{A}_{0,N-3} & \ddots & \mathbf{A}_{0,0} & \mathbf{A}_{0,N-1} \\ \mathbf{A}_{0,N-1} & \mathbf{A}_{0,N-2} & \ddots & \mathbf{A}_{0,1} & \mathbf{A}_{0,0} \end{bmatrix}, \quad (13)$$

and each submatrix $\mathbf{A}_{0,n}$ is an $M \times M$ circulant matrix whose structure is determined by (11).

Example: Consider a special case with $N = 4$ and $M = 3$, and U_0 's channel is given by

$$h_0(\tau, \nu) = h_{0,0} \delta(\tau) \delta(\nu) + h_{0,1} \delta \left(\tau - \frac{1}{M\Delta f} \right) \delta \left(\nu - \frac{3}{NT} \right), \quad (14)$$

which means $k_0 = 0$, $k_1 = 3$, $l_0 = 0$, $l_1 = 1$. Therefore, the block-circulant matrix is given by

$$\mathbf{H}_0 = \begin{bmatrix} \mathbf{A}_{0,0} & \mathbf{A}_{0,3} & \mathbf{A}_{0,2} & \mathbf{A}_{0,1} \\ \mathbf{A}_{0,1} & \mathbf{A}_{0,0} & \mathbf{A}_{0,3} & \mathbf{A}_{0,2} \\ \mathbf{A}_{0,2} & \mathbf{A}_{0,1} & \mathbf{A}_{0,0} & \mathbf{A}_{0,3} \\ \mathbf{A}_{0,3} & \mathbf{A}_{0,2} & \mathbf{A}_{0,1} & \mathbf{A}_{0,0} \end{bmatrix}, \quad (15)$$

where $\mathbf{A}_{0,0} = h_{0,0}\mathbf{I}_3$, $\mathbf{A}_{0,1} = \mathbf{A}_{0,2} = \mathbf{0}_{3 \times 3}$ and $\mathbf{A}_{0,3} =$

$$\begin{bmatrix} 0 & 0 & h_{0,1} \\ h_{0,1} & 0 & 0 \\ 0 & h_{0,1} & 0 \end{bmatrix}.$$

Remark 1: It is well known that an $n \times n$ circulant matrix can be diagonalized by the $n \times n$ discrete Fourier transform (DFT) and inverse DFT matrices, denoted by \mathbf{F}_n and \mathbf{F}_n^{-1} , respectively, i.e., the columns of the DFT matrix are the eigenvectors of the circulant matrix. We note that directly applying the DFT factorization to \mathbf{H}_0 is not possible, since \mathbf{H}_0 is not a circulant matrix, but a block circulant matrix.

Because of the structure of \mathbf{H}_0 , inter-symbol interference still exists in the considered OTFS-NOMA system, and equalization is needed. We consider two equalization approaches, FD-LE and FD-DFE, which were both originally developed for single-carrier transmission with cyclic prefix [34], [35].

A. Design and Performance of FD-LE

The proposed FD-LE consists of two steps. Let \otimes denote the Kronecker product. The first step is to multiply the observation vector \mathbf{y}_0 by $\mathbf{F}_N \otimes \mathbf{F}_M^H$, which leads to the result in the following proposition.

Proposition 1: By applying the detection matrix $\mathbf{F}_N \otimes \mathbf{F}_M^H$ to observation vector \mathbf{y}_0 , the received signals for OTFS-NOMA downlink transmission can be written as follows:

$$\tilde{\mathbf{y}}_0 = \mathbf{D}_0(\mathbf{F}_N \otimes \mathbf{F}_M^H) \left(\gamma_0 \mathbf{x}_0 + \sum_{q=1}^M \gamma_q \mathbf{x}_q \right) + \tilde{\mathbf{z}}_0, \quad (16)$$

where $\tilde{\mathbf{y}}_0 = (\mathbf{F}_N \otimes \mathbf{F}_M^H) \mathbf{y}_0$, $\tilde{\mathbf{z}}_0 = (\mathbf{F}_N \otimes \mathbf{F}_M^H) \mathbf{z}_0$, \mathbf{D}_0 is a diagonal matrix whose $(kM + l + 1)$ -th main diagonal element is given by

$$D_0^{k,l} = \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} a_{0,n}^{m,1} e^{j2\pi \frac{lm}{M}} e^{-j2\pi \frac{kn}{N}}, \quad (17)$$

for $0 \leq k \leq N-1$, $0 \leq l \leq M-1$, and $a_{0,n}^{m,1}$ is the element located in the $(nM + m + 1)$ -th row and the first column of \mathbf{H}_0 .

Proof: Please refer to Appendix A. \square

With the simplified signal model shown in (16), the second step of FD-LE is to apply $(\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-1} \mathbf{D}_0^{-1}$ to $\tilde{\mathbf{y}}_0$. Thus, \mathbf{U}_0 's received signal is given by

$$\check{\mathbf{y}}_0 = \gamma_0 \mathbf{x}_0 + \underbrace{\sum_{q=1}^M \gamma_q \mathbf{x}_q + (\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-1} \mathbf{D}_0^{-1} \tilde{\mathbf{z}}_0}_{\text{Interference and noise terms}}, \quad (18)$$

where $\check{\mathbf{y}}_0 = (\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-1} \mathbf{D}_0^{-1} \tilde{\mathbf{y}}_0$. To simplify the analysis, we assume that the powers of all users' information-bearing signals are identical, which means that the transmit signal-to-noise ratio (SNR) can be defined as $\rho = \mathcal{E}\{|x_0[k, l]|^2\} = \mathcal{E}\{|x_i(n)|^2\}$, since the noise power is assumed to be normalized.⁶ The following lemma provides the signal-to-interference-plus-noise ratio (SINR) achieved by FD-LE.

⁶Following steps in the proof for Proposition 1 to show the statistical property of $\tilde{\mathbf{z}}_0$ in (66), we can also show that $W_i[n, m] \sim \mathcal{CN}(0, 1)$ if $z_i[k, l] \sim \mathcal{CN}(0, 1)$.

Lemma 1: Assume that $\gamma_i = \gamma_1$, for $1 \leq i \leq N$. By using FD-LE, the SINRs for detecting all $x_0[k, l]$, $0 \leq k \leq N-1$ and $0 \leq l \leq M-1$, are identical and given by

$$\text{SINR}_{0,kl}^{\text{LE}} = \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \frac{1}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} |D_0^{\tilde{k}, \tilde{l}}|^{-2}}. \quad (19)$$

Proof: Please refer to Appendix B. \square

Remark 2: The proof of Lemma 1 shows that $\sum_{i=0}^M \gamma_i^2 = 1$ can be simplified as $\gamma_0^2 + \gamma_i^2 = 1$ for $1 \leq i \leq M$, which is the motivation for assuming $\gamma_i = \gamma_1$. Following steps similar to those in the proofs for Proposition 1 and Lemma 1, one can show that directly applying \mathbf{H}_0^{-1} to the observation vector yields the same SINR. However, the proposed FD-LE scheme can be implemented more efficiently since $(\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-1} = \mathbf{F}_N^H \otimes \mathbf{F}_M$ and \mathbf{D}_0 is a diagonal matrix. Hence, the inversion of a full $NM \times NM$ matrix is avoided.

In this paper, the outage probability and the outage rate are used as performance criteria, since the outage probability can provide a tight bound on the probability of erroneous detection and is general in the sense that it does not depend on particular channel coding and modulation schemes used [36]. The outage probability achieved by FD-LE is given by $P(\log(1 + \text{SINR}_{0,kl}^{\text{LE}}) < R_0)$, where R_i , $0 \leq i \leq M$, denotes \mathbf{U}_i 's target data rate. It is difficult to analyze the outage probability for the following two reasons. First, the $D_0^{k,l}$, $k \in \{0, \dots, N-1\}$, $l \in \{0, \dots, M-1\}$, are not statistically independent, and second, the distribution of a sum of the inverse of exponentially distributed random variables is difficult to characterize. The following lemma provides an asymptotic result for the outage probability based on the SINR provided in Lemma 1.

Lemma 2: If $\gamma_0^2 > \gamma_1^2 \epsilon_0$, the diversity order achieved by FD-LE is one, where $\epsilon_0 = 2^{R_0} - 1$. Otherwise, the outage probability is always one.

Proof: Please refer to Appendix C. \square

Remark 3: Recall that the diversity order achieved by OTFS-OMA, where the high-mobility user, \mathbf{U}_0 , solely occupies the bandwidth resources, is also one. Therefore, the use of OTFS-NOMA ensures that the additional M low-mobility users are served without compromising \mathbf{U}_0 's diversity order, which improves the spectral efficiency compared to OTFS-OMA.

B. Design and Performance of FD-DFE

Different from FD-LE, which is a linear equalizer, FD-DFE is based on the idea of feeding back previously detected symbols. Since both \mathbf{x}_0 and \mathbf{x}_q , $q \geq 1$, experience the same fading channel, we first define $\mathbf{x} = \gamma_0 \mathbf{x}_0 + \sum_{q=1}^M \gamma_q \mathbf{x}_q$, which are the signals to be recovered by FD-DFE. Given the received signal vector shown in (12), the outputs of FD-DFE are given by

$$\hat{\mathbf{x}} = \mathbf{P}_0 \mathbf{y}_0 - \mathbf{G}_0 \check{\mathbf{x}}, \quad (20)$$

where $\check{\mathbf{x}}$ contains the decisions made on the symbols \mathbf{x} , \mathbf{P}_0 is the feed-forward part of the equalizer, and \mathbf{G}_0 is the feedback part of the equalizer. Similar to [34], [35], we use the following choices for \mathbf{P}_0 and \mathbf{G}_0 : $\mathbf{P}_0 = \mathbf{L}_0 (\mathbf{H}_0^H \mathbf{H}_0)^{-1} \mathbf{H}_0^H$, $\mathbf{G}_0 = \mathbf{L}_0 - \mathbf{I}_{NM}$, where \mathbf{L}_0 is a lower triangular matrix

with its main diagonal elements being ones in order to ensure causality of the feedback signals. With the above choices for \mathbf{P}_0 and \mathbf{G}_0 , U_0 's signals can be detected as follows:

$$\hat{\mathbf{x}} = \mathbf{L}_0(\mathbf{H}_0^H \mathbf{H}_0)^{-1} \mathbf{H}_0^H \mathbf{y}_0 - (\mathbf{L}_0 - \mathbf{I}_{NM}) \hat{\mathbf{x}}. \quad (21)$$

For FD-DFE, \mathbf{L}_0 is obtained from the Cholesky decomposition of \mathbf{H}_0 , i.e., $\mathbf{H}_0^H \mathbf{H}_0 = \mathbf{L}_0^H \mathbf{\Lambda}_0 \mathbf{L}_0$, where \mathbf{L}_0 is the desirable lower triangular matrix, and $\mathbf{\Lambda}_0$ is a diagonal matrix. Therefore, the estimates of \mathbf{x}_0 can be rewritten as follows:

$$\hat{\mathbf{x}} = \mathbf{x} + \mathbf{L}_0(\mathbf{H}_0^H \mathbf{H}_0)^{-1} \mathbf{H}_0^H \mathbf{z}_0 \quad (22)$$

$$= \gamma_0 \mathbf{x}_0 + \underbrace{\sum_{q=1}^M \gamma_q \mathbf{x}_q + \mathbf{L}_0(\mathbf{H}_0^H \mathbf{H}_0)^{-1} \mathbf{H}_0^H \mathbf{z}_0}_{\text{Interference and noise terms}}, \quad (23)$$

where perfect decision-making is assumed, i.e., $\check{\mathbf{x}} = \mathbf{x}$, and there is no error propagation [35], [37], [38]. We note that (23) yields an upper bound on the reception reliability of FD-DFE when error propagation cannot be completely avoided.

Following steps similar to those in the proof of Lemma 1, the covariance matrix for the interference-plus-noise term can be found as follows:

$$\mathbf{C}_{\text{cov}} = \rho \gamma_1^2 \mathbf{I}_{MN} + \mathbf{L}_0(\mathbf{H}_0^H \mathbf{H}_0)^{-1} \mathbf{L}_0^H = \rho \gamma_1^2 \mathbf{I}_{MN} + \mathbf{\Lambda}_0^{-1}, \quad (24)$$

where the last step follows from the fact that \mathbf{L}_0 is obtained from the Cholesky decomposition of \mathbf{H}_0 . Therefore, the SINR for detecting $x_0[k, l]$ can be expressed as follows:

$$\text{SINR}_{0,kl} = \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \lambda_{0,kl}^{-1}}, \quad (25)$$

where $\lambda_{0,kl}$ is the $(kM+l+1)$ -th element on the main diagonal of $\mathbf{\Lambda}_0$.

Remark 4: We note that there is a fundamental difference between the two equalization schemes. One can observe from (19) that the SINRs achieved by FD-LE for different $x_0[k, l]$ are identical. However, for FD-DFE, different symbols experience different effective fading gains, $\lambda_{0,kl}$. Therefore, FD-DFE can realize unequal error protection for data streams with different priorities. This comes at the price of a higher computational complexity.

We further note that the use of FD-DFE also ensures that multi-path diversity can be harvested, as shown in the following. The outage performance analysis for FD-DFE requires knowledge of the distribution of the effective channel gains, $\lambda_{0,kl}$. Because of the implicit relationship between $\mathbf{\Lambda}_0$ and \mathbf{H}_0 , a general expression for the outage probability achieved by FD-DFE is difficult to obtain. However, analytical results can be developed for special cases to show that the use of FD-DFE can realize the maximal multi-path diversity.

In particular, the SINR for $x_0[N-1, M-1]$ is a function of $\lambda_{0,(N-1)(M-1)}$ which is the last element on the main diagonal of $\mathbf{\Lambda}_0$. Recall that $\mathbf{\Lambda}_0$ is obtained via Cholesky decomposition, i.e., $\mathbf{H}_0^H \mathbf{H}_0 = \mathbf{L}_0^H \mathbf{\Lambda}_0 \mathbf{L}_0$. Because \mathbf{L}_0 is a lower triangular matrix, $\lambda_{0,(N-1)(M-1)}$ is equal to the element of $\mathbf{H}_0^H \mathbf{H}_0$ located in the NM -th column and the NM -th row,

which means

$$\lambda_{0,(N-1)(M-1)} = \sum_{p=0}^{P_0} |h_{0,p}|^2. \quad (26)$$

Since the channel gains are i.i.d. and follow $h_{0,p} \sim \mathcal{CN}(0, \frac{1}{P_0+1})$, the probability density function (pdf) of $\sqrt{P_0+1} \lambda_{0,(N-1)(M-1)}$ is given by

$$f(x) = \frac{1}{P_0!} e^{-x} x^{P_0}. \quad (27)$$

By using the above pdf, the outage probability and the diversity order can be obtained by some algebraic manipulations, as shown in the following corollary.

Corollary 1: Assume $\gamma_0^2 > \gamma_1^2 \epsilon_0$. The use of FD-DFE realizes the following outage probability for detection of $x_0[N-1, M-1]$:

$$P_{N-1, M-1}^0 = \frac{1}{P_0!} g\left(P_0+1, \frac{\epsilon_0(P_0+1)}{\rho(\gamma_0^2 - \gamma_1^2 \epsilon_0)}\right), \quad (28)$$

where $g(\cdot)$ denotes the incomplete Gamma function. The full multi-path diversity order, P_0+1 , is achievable for $x_0[N-1, M-1]$.

Remark 5: The results in Corollary 1 can be extended to OTFS-OMA with FD-DFE straightforwardly. We also note that diversity gains larger than one are not achievable with FD-LE as shown in Lemma 2, which is one of the disadvantages of FD-LE compared to FD-DFE.

Remark 6: We note that not all NM data streams can benefit from the full diversity gain. The simulation results provided in Section VII (Fig. 2) show that the diversity orders achievable for $x_0[k, l]$, $k < N-1$ and $l < M-1$, are smaller than that for $x_0[N-1, M-1]$, and the diversity order for $x_0[0, 0]$ is one, i.e., the same value as for FD-LE. We further note that the diversity result in Corollary 1 is obtained by assuming that there is no error propagation, i.e., it is assumed that when detecting the i -th element of \mathbf{x} in (21), the first $(i-1)$ elements of \mathbf{x} have already been correctly detected. Because of this assumption, the diversity gain developed in Corollary 1 is an upper bound on the diversity gain achieved by FD-DFE. If the assumption does not hold, the diversity orders for $x_0[k, l]$ will be capped by the worst case, i.e., the diversity gain for $x_0[0, 0]$ which is one.

Remark 7: FD-DFE entails a higher implementation complexity than FD-LE, as explained in the following. The complexity of FD-LE is mainly caused by computing the inversion of $\mathbf{H}_0^H \mathbf{H}_0$. However, for FD-DFE, \mathbf{L}_0 needs to be computed, in addition to $(\mathbf{H}_0^H \mathbf{H}_0)^{-1}$, as shown in (21). Recall that \mathbf{L}_0 is obtained from the Cholesky decomposition of the $NM \times NM$ matrix \mathbf{H}_0 , which entails a computational complexity of $\mathcal{O}(N^3 M^3)$. Therefore, the computational complexity of FD-DFE is higher than that of FD-LE, but FD-DFE offers a performance gain in terms of reception reliability compared to FD-LE, as shown in Section VII.

V. DOWNLINK OTFS-NOMA - DETECTING THE NOMA USERS' SIGNALS

Successive interference cancellation (SIC) will be carried out by the NOMA users, where each NOMA user first decodes

540 the high mobility user's signal in the delay-Doppler plane
541 and then decodes its own signal in the time-frequency plane.
542 The two stages of SIC are discussed in the following two
543 subsections, respectively.

544 A. Stage I of SIC

545 Following steps similar to the ones in the previous section,
546 each NOMA user also observes the mixture of the $(M + 1)$
547 users' signals in the delay-Doppler plane as follows:

$$548 \quad \mathbf{y}_i = \gamma_0 \mathbf{H}_i \mathbf{x}_0 + \underbrace{\sum_{q=1}^M \gamma_q \mathbf{H}_i \mathbf{x}_q}_{\text{Interference and noise terms}} + \mathbf{z}_i, \quad (29)$$

549 where \mathbf{H}_i and \mathbf{z}_i are defined similar to \mathbf{H}_0 and \mathbf{z}_0 , respectively.

550 We assume that the low-mobility NOMA users do not
551 experience Doppler shift, and therefore, their channels can be
552 simplified as follows:

$$553 \quad h_i(\tau) = \sum_{p=0}^{P_i} h_{i,p} \delta(\tau - \tau_{i,p}), \quad (30)$$

554 for $1 \leq i \leq K$, which means that each NOMA user's
555 channel matrix, \mathbf{H}_i , $1 \leq i \leq N$, is a block-diagonal matrix,
556 i.e., $\mathbf{A}_{i,0}$ is a non-zero circulant matrix and $\mathbf{A}_{i,n} = \mathbf{0}_{M \times M}$,
557 for $1 \leq n \leq N - 1$. Therefore, each NOMA user can
558 divide its observation vector into N equal-length sub-vectors,
559 i.e., $\mathbf{y}_i = [\mathbf{y}_{i,0}^T \cdots \mathbf{y}_{i,N-1}^T]^T$, which yields the following
560 simplified system model:

$$561 \quad \mathbf{y}_{i,n} = \gamma_0 \mathbf{A}_{i,0} \mathbf{x}_{0,n} + \sum_{q=1}^M \gamma_q \mathbf{A}_{i,0} \mathbf{x}_{q,n} + \mathbf{z}_{i,n}, \quad (31)$$

562 where, similar to $\mathbf{y}_{i,n}$, $\mathbf{x}_{i,n}$ and $\mathbf{z}_{i,n}$ are obtained from \mathbf{x}_i
563 and \mathbf{z}_i , respectively. Therefore, unlike the high-mobility user,
564 the NOMA users can perform their signal detection based on
565 reduced-size observation vectors, which reduces the computa-
566 tional complexity.

567 Since $\mathbf{A}_{i,0}$ is a circulant matrix, the two equalization
568 approaches used in the previous section are still applicable.
569 First, we consider the use of FD-LE. Following the same steps
570 as in the proof for Proposition 1, in the first step of FD-LE, the
571 DFT matrix is applied to the reduced-size observation vector,
572 which yields the following:

$$573 \quad \tilde{\mathbf{y}}_{i,n} = \tilde{\mathbf{D}}_i \mathbf{F}_M^H \left(\gamma_0 \mathbf{x}_{0,n} + \sum_{q=1}^M \gamma_q \mathbf{x}_{q,n} \right) + \tilde{\mathbf{z}}_{i,n}, \quad (32)$$

574 where $\tilde{\mathbf{y}}_{i,n} = \mathbf{F}_M^H \mathbf{y}_{i,n}$ and $\tilde{\mathbf{z}}_{i,n} = \mathbf{F}_M^H \mathbf{z}_{i,n}$. Compared to \mathbf{D}_i
575 in Proposition 1 which is an $NM \times NM$ matrix, $\tilde{\mathbf{D}}_i$ is an
576 $M \times M$ diagonal matrix, and its $(l + 1)$ -th main diagonal
577 element is given by $\tilde{D}_i^l = \sum_{m=0}^{M-1} a_{i,0}^{m,1} e^{j2\pi \frac{lm}{M}}$, for $0 \leq l \leq$
578 $M - 1$, where $a_{i,0}^{m,1}$ is the element located in the $(m + 1)$ -th
579 row and the first column of $\mathbf{A}_{i,0}$. Unlike conventional OFDM,
580 which uses \mathbf{F}_M at the receiver, \mathbf{F}_M^H is used here. Because
581 $\mathbf{F}_M^H \mathbf{A}_{i,0} \mathbf{F}_M = [\mathbf{F}_M \mathbf{A}_{i,0}^* \mathbf{F}_M^H]^*$, the sign of the exponent of
582 the exponential component of \tilde{D}_i^l is different from that in the
583 conventional case.

584 In the second step of FD-LE, $\mathbf{F}_M \tilde{\mathbf{D}}_i^{-1}$ is applied to $\tilde{\mathbf{y}}_{i,n}$.
585 Following steps similar to the ones in the proof for Lemma 1,
586 the SINR for detecting $x_0[k, l]$ can be obtained as follows:

$$587 \quad \text{SINR}_{0,kl}^{i,\text{LE}} = \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \frac{1}{M} \sum_{l=0}^{M-1} |\tilde{D}_i^l|^{-2}}. \quad (33)$$

588 We note that $\text{SINR}_{0,k_1 l}^{i,\text{LE}} = \text{SINR}_{0,k_2 l}^{i,\text{LE}}$, for $k_1 \neq k_2$, due to the
589 time invariant nature of the channels.

590 If FD-DFE is used, the corresponding SINR for detecting
591 $x_0[k, l]$ is given by

$$592 \quad \text{SINR}_{0,kl}^{i,\text{DFE}} = \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \tilde{\lambda}_{0,l}^{-1}}, \quad (34)$$

593 where $\tilde{\lambda}_{0,l}$ is obtained from the Cholesky decomposition
594 of $\mathbf{A}_{i,0}$. The details for the derivation of (34) are omitted
595 here due to space limitations.

596 B. Stage II of SIC

597 Assume that \mathbf{U}_0 's NM signals can be decoded and removed
598 successfully, which means that, in the time-frequency plane,
599 the NOMA users observe the following:

$$600 \quad Y_i[n, m] = \sum_{q=1}^M \gamma_q H_i[n, m] X_q[n, m] + W_i[n, m] \\ 601 \quad = \gamma_1 H_i[n, m] x_{m+1}(n) + W_i[n, m], \quad (35)$$

602 where the last step follows from the mapping scheme used
603 in (6) and it is assumed that all NOMA users employ the
604 same power allocation coefficient. We note that \mathbf{U}_i is only
605 interested in $Y_i[n, i - 1]$, $0 \leq n \leq N - 1$. Therefore, \mathbf{U}_i 's n -th
606 information bearing signal, $x_i(n)$, can be detected by applying
607 a one-tap equalizer as follows:

$$608 \quad \hat{x}_i(n) = \frac{Y_i[n, i - 1]}{\gamma_1 H_i[n, i - 1]}, \quad (36)$$

609 which means that the SNR for detecting $x_i(n)$ is given by

$$610 \quad \text{SNR}_{i,n} = \rho \gamma_1^2 |\tilde{D}_i^{i-1}|^2, \quad (37)$$

611 since $W_i[n, i - 1]$ is white Gaussian noise and $H_i[n, i - 1] =$
612 \tilde{D}_i^{i-1} . We note that $\text{SNR}_{i,n_1} = \text{SNR}_{i,n_2}$, for $n_1 \neq n_2$, which
613 is due to the time-invariant nature of the channel.

614 Without loss of generality, assume that the same target data
615 rate R_i is used for $x_i(n)$, $0 \leq n \leq N - 1$. Therefore, the outage
616 probability for $x_i(n)$ is given by

$$617 \quad \text{P}_{i,n}^{\text{LE}} \\ 618 \quad = 1 - \text{P} \left(\text{SNR}_{i,n} > \epsilon_i, \text{SINR}_{0,kl}^{i,\text{LE}} > \epsilon_0, \forall l \right) \\ 619 \quad = 1 - \text{P} \left(\rho \gamma_1^2 |\tilde{D}_i^{i-1}|^2 > \epsilon_i, \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \frac{1}{M} \sum_{l=0}^{M-1} |\tilde{D}_i^l|^{-2}} > \epsilon_0 \right), \quad (38)$$

620 if FD-LE is used in the first stage of SIC. If FD-DFE is used
621 in the first stage of SIC, the outage probability for $x_i(n)$ is
622 given by

$$623 \quad \text{P}_{i,n}^{\text{DFE}} = 1 - \text{P} \left(\text{SNR}_{i,n} > \epsilon_i, \text{SINR}_{0,kl}^{i,\text{DFE}} > \epsilon_0, \forall l \right) \\ 624 \quad = 1 - \text{P} \left(\rho \gamma_1^2 |\tilde{D}_i^{i-1}|^2 > \epsilon_i, \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \tilde{\lambda}_{0,l}^{-1}} > \epsilon_0, \forall l \right), \quad (39)$$

where $\epsilon_i = 2^{R_i} - 1$. Again because of the correlation between the random variables $|\tilde{D}_i^l|^{-2}$ and $\tilde{\lambda}_{0,l}$, the exact expressions for the outage probabilities are difficult to obtain. Alternatively, the achievable diversity order is analyzed in the following subsections.

1) *Random User Scheduling*: If the M users are randomly selected from the K available users, which means that each $|\tilde{D}_i^l|^2$ is complex Gaussian distributed. For the FD-LE case, the outage probability, $P_{i,n}^{\text{LE}}$, can be upper bounded as follows:

$$P_{i,n}^{\text{LE}} \leq 1 - \text{P} \left(\rho\gamma_1^2 |\tilde{D}_i^{\min}|^2 > \epsilon_i, \frac{\rho\gamma_0^2}{\rho\gamma_1^2 + |\tilde{D}_i^{\min}|^{-2}} > \epsilon_0 \right), \quad (40)$$

where $|\tilde{D}_i^{\min}|^2 = \min\{|\tilde{D}_i^m|^2, 0 \leq m \leq M-1\}$. The upper bound on the outage probability in (40) can be rewritten as follows:

$$P_{i,n}^{\text{LE}} \leq 1 - \text{P} \left(|\tilde{D}_i^{\min}|^2 > \bar{\epsilon} \right), \quad (41)$$

where $\bar{\epsilon} = \max \left\{ \frac{\epsilon_0}{\rho(\gamma_0^2 - \gamma_1^2 \epsilon_0)}, \frac{\epsilon_i}{\rho\gamma_1^2} \right\}$. As a result, an upper bound on the outage probability can be obtained as follows:

$$P_{i,n}^{\text{LE}} \leq \text{P} \left(|\tilde{D}_i^{\min}|^2 < \bar{\epsilon} \right) \leq MP \left(|\tilde{D}_i^0|^2 < \bar{\epsilon} \right) \doteq \frac{1}{\rho}, \quad (42)$$

where $\text{P}^o \doteq \rho^{-d}$ denotes exponential equality, i.e., $d = -\lim_{\rho \rightarrow \infty} \frac{\log \text{P}^o}{\log \rho}$ [36]. Therefore, the following corollary can be obtained.

Corollary 2: For random user scheduling and FD-LE, a diversity order of 1 is achievable at the NOMA users.

Our simulation results in Section VII show that a diversity order of 1 is also achievable for FD-DFE, although we do not have a formal proof for this conclusion, yet.

2) *Realizing Multi-User Diversity*: The diversity order of OTFS-NOMA can be improved by carrying out opportunistic user scheduling, which yields multi-user diversity gains. For illustration purpose, we propose a greedy user scheduling policy, where a single NOMA user is scheduled to transmit in all resource blocks of the time-frequency plane. From the analysis of the random scheduling case we deduce that $|\tilde{D}_i^{\min}|^2$ is critical to the outage performance. Therefore, the scheduled NOMA user, denoted by U_{i^*} , is selected based on the following criterion:

$$i^* = \arg \max_{i \in \{1, \dots, K\}} \left\{ |\tilde{D}_i^{\min}|^2 \right\}. \quad (43)$$

By using the assumption that the users' channel gains are independent and following steps similar to the ones in the proof for Lemma 2, the following corollary can be obtained in a straightforward manner.

Corollary 3: For FD-LE, the user scheduling strategy shown in (43) realizes the maximal multi-user diversity gain, K .

Remark 8: The reason why a multi-user diversity gain of K can be realized by the proposed scheduling strategy is explained in the following. Recall that the SINR for FD-LE to detect $x_0[k, l]$ is $\text{SINR}_{0,kl}^{i,\text{LE}} = \frac{\rho\gamma_0^2}{\rho\gamma_1^2 + \frac{1}{M} \sum_{i=0}^{M-1} |\tilde{D}_i^l|^{-2}}$. If this SINR is too small, the first stage of SIC will fail and an

outage event will occur. To improve the SINR, it is important to ensure that for a scheduled user, its weakest channel gain, $|\tilde{D}_i^{\min}|^2 = \min\{|\tilde{D}_i^m|^2, 0 \leq m \leq M-1\}$, is not too small. The used scheduling strategy shown in (43) is essentially a max-min strategy and ensures that the user with the strongest $|\tilde{D}_i^{\min}|^2$ is selected from the K candidates, which effectively exploits multi-user diversity.

We note that the user scheduling strategy shown in (43) is also useful for improving the performance of FD-DFE, as shown in Section VII.

VI. UPLINK OTFS-NOMA TRANSMISSION

The design of uplink OTFS-NOMA is similar to that of downlink OTFS-NOMA, and due to space limitations, we mainly focus on the difference between the two cases in this section. Again, we assume that U_0 is grouped with M NOMA users, selected from the K available users. U_0 's NM signals are placed in the delay-Doppler plane, and are denoted by $x_0[k, l]$, where $0 \leq k \leq N-1$ and $0 \leq l \leq M-1$. The corresponding time-frequency signals, $X_0[n, m]$, are obtained by applying ISFFT to $x_0[k, l]$. On the other hand, the NOMA users' signals, $x_i(n)$, are mapped to time-frequency signals, $X_i[n, m]$, according to (6).

Following steps similar to the ones for the downlink case, the base station's observations in the time-frequency plane are given by

$$\begin{aligned} Y[n, m] &= \sum_{q=0}^M H_q[n, m] X_q[n, m] + W[n, m] \\ &= \frac{H_0(n, m)}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} x_0[k, l] e^{j2\pi(\frac{kn}{N} - \frac{ml}{M})} \\ &\quad + \sum_{q=1}^M H_q[n, m] X_q[n, m] + W[n, m], \end{aligned} \quad (44)$$

where $W[n, m]$ is the Gaussian noise at the base station in the time-frequency plane. We assume that all users employ the same transmit pulse as well as the same transmit power. The base station applies SIC to first detect the NOMA users' signals in the time-frequency plane, and then tries to detect the high-mobility user's signals in the delay-Doppler plane, as shown in the following two subsections.

A. Stage I of SIC

The base station will first try to detect the NOMA users' signals in the time-frequency plane by treating the signals from U_0 as noise, which is the first stage of SIC.

By using (6), $x_i(n)$ can be estimated as follows:

$$\begin{aligned} \hat{x}_i(n) &= \frac{Y[n, i-1]}{H_i[n, i-1]} \\ &= x_i[n] + \frac{H_0[n, i-1] X_0[n, i-1] + W[n, i-1]}{H_i[n, i-1]}. \end{aligned} \quad (45)$$

Define an $NM \times 1$ vector, $\bar{\mathbf{x}}_0$, whose $(nM+m+1)$ -th element is $X_0[n, m]$. Recall that $X_0[n, m]$ is obtained from the ISFFT of $x_0[k, l]$, i.e.,

$$\bar{\mathbf{x}}_0 = (\mathbf{F}_N^H \otimes \mathbf{F}_M) \mathbf{x}_0, \quad (46)$$

721 which means $X_0[n, m]$ follows the same distribution as
 722 $x_0[k, l]$. By applying steps similar to those in the proof for
 723 Lemma 1, the SINR for detecting $x_i(n)$ is given by

$$724 \quad \text{SINR}_{i,n} = \frac{\rho |H_i[n, i-1]|^2}{\rho |H_0[n, i-1]|^2 + 1}. \quad (47)$$

725 Unlike downlink OTFS-NOMA, there are two possible
 726 strategies for uplink OTFS-NOMA to combat multiple access
 727 interference, as shown in the following two subsections.

728 1) *Adaptive-Rate Transmission*: One strategy to combat
 729 multiple access interference is to impose the following con-
 730 straint on $x_i(n)$:

$$731 \quad R_{i,n} \leq \log \left(1 + \frac{\rho |H_i[n, i-1]|^2}{\rho |H_0[n, i-1]|^2 + 1} \right), \quad (48)$$

732 which means that the first stage of SIC is guaranteed to be
 733 successful. Therefore, the M low-mobility users are served
 734 without affecting U_0 's outage probability, i.e., the use of
 735 NOMA is transparent to U_0 .

736 Because U_i 's data rate is adaptive, outage events when
 737 decoding $x_i(n)$ do not happen, which means that an appropri-
 738 ate criterion for the performance evaluation is the ergodic rate.
 739 Recall that $H_i[n, i-1] = \tilde{D}_i^{i-1}$ and $H_0[n, i-1] = D_0^{n, i-1}$.
 740 Therefore, U_i 's ergodic rate is given by

$$741 \quad \mathcal{E}\{R_{i,n}\} \leq \mathcal{E} \left\{ \log \left(1 + \frac{\rho |\tilde{D}_i^{i-1}|^2}{\rho |D_0^{n, i-1}|^2 + 1} \right) \right\}. \quad (49)$$

742 We note that the ergodic rate of uplink OTFS-NOMA can
 743 be further improved by modifying the user scheduling strategy
 744 proposed in (43), as shown in the following. Particularly,
 745 denote the NOMA user which is scheduled to transmit in the
 746 m -th frequency subchannel by $U_{i_m^*}$, and this user is selected
 747 by using the following criterion:

$$748 \quad i_m^* = \arg \max_{i \in \{1, \dots, K\}} \left\{ |\tilde{D}_i^m|^2 \right\}. \quad (50)$$

749 We note that a single user might be scheduled on multiple
 750 frequency channels, which reduces user fairness.

751 Because the integration of the logarithm function appearing
 752 in (49) leads to non-insightful special functions, we will use
 753 simulations to evaluate the ergodic rate of OTFS-NOMA in
 754 Section VII.

755 2) *Fixed-Rate Transmission*: If the NOMA users do not
 756 have the capabilities to adapt their transmission rates, they
 757 have to use fixed data rates R_i for transmission, which means
 758 that outage events can happen and the achieved outage perfor-
 759 mance is analyzed in the following. For illustration purposes,
 760 we focus on the case when the user scheduling strategy shown
 761 in (50) is used.

762 The outage probability for detecting $x_{i_m^*}(n)$ is given by

$$763 \quad P_{i_m^*, n} = \text{P} \left(\log \left(1 + \frac{\rho |\tilde{D}_{i_m^*}^{i_m^*-1}|^2}{\rho |D_0^{n, i_m^*-1}|^2 + 1} \right) < R_{i_m^*} \right). \quad (51)$$

764 Following steps similar to the ones in the proof for Lemma 2,
 765 we can show that $|\tilde{D}_{i_m^*}^{i_m^*-1}|^2$ and $|D_0^{n, i_m^*-1}|^2$ are independent,

and the use of the user scheduling scheme in (50) simplifies
 the outage probability as follows:

$$766 \quad P_{i_m^*, n} = \text{P} \left(\log \left(1 + \frac{\rho |\tilde{D}_{i_m^*}^{i_m^*-1}|^2}{\rho |D_0^{n, i_m^*-1}|^2 + 1} \right) < R_{i_m^*} \right) \quad 767$$

$$768 \quad = \int_0^\infty \left(1 - e^{-\frac{\epsilon_{i_m^*}(1+\rho y)}{\rho}} \right)^K e^{-y} dy, \quad (52) \quad 769$$

770 where we use the fact that the cumulative distribution function
 771 of $|\tilde{D}_{i_m^*}^{i_m^*-1}|^2$ is $(1 - e^{-x})^K$ because of the adopted user
 772 scheduling strategy.

773 The outage probability can be further simplified as follows:

$$774 \quad P_{i_m^*, n} = \sum_{k=0}^K \binom{K}{k} (-1)^k \int_0^\infty e^{-\frac{k\epsilon_{i_m^*}(1+\rho y)}{\rho}} e^{-y} dy \quad 775$$

$$776 \quad = \sum_{k=0}^K \binom{K}{k} (-1)^k e^{-\frac{k\epsilon_{i_m^*}}{\rho}} \frac{1}{k\epsilon_{i_m^*} + 1}. \quad (53) \quad 777$$

778 At high SNR, the outage probability can be approximated
 779 as follows:

$$780 \quad P_{i_m^*, n} \approx \sum_{k=0}^K \binom{K}{k} (-1)^k \frac{1}{k\epsilon_{i_m^*} + 1}, \quad (54) \quad 781$$

782 which is no longer a function of ρ , i.e., the outage probability
 783 has an error floor at high SNR. This is due to the fact that
 784 $U_{i_m^*}$ is subject to strong interference from U_0 .

785 However, we can show that the error floor experienced by
 786 $U_{i_m^*}$ can be reduced by increasing K , i.e., inviting more oppor-
 787 tunistic users for NOMA transmission. In particular, assuming
 788 $K\epsilon_{i_m^*} \rightarrow 0$, the outage probability can be approximated as
 789 follows:

$$790 \quad P_{i_m^*, n} \approx \sum_{k=0}^K \binom{K}{k} (-1)^k (1 + k\epsilon_{i_m^*})^{-1} \quad 791$$

$$792 \quad \approx \sum_{k=0}^K \binom{K}{k} (-1)^k \sum_{l=0}^\infty (-1)^l k^l \epsilon_{i_m^*}^l, \quad (55) \quad 793$$

794 where we use the fact that $(1+x)^{-1} = \sum_{l=0}^\infty (-1)^l x^l$, $|x| < 1$.
 795 Therefore, the error floor at high SNR can be approximated
 796 as follows:

$$797 \quad P_{i_m^*, n} \approx \sum_{l=0}^\infty (-1)^l \epsilon_{i_m^*}^l \sum_{k=0}^K \binom{K}{k} (-1)^k k^l \quad 798$$

$$799 \quad \approx (-1)^K \epsilon_{i_m^*}^K (-1)^K K! = K! \epsilon_{i_m^*}^K, \quad (56) \quad 800$$

801 where we use the identities $\sum_{k=0}^K \binom{K}{k} (-1)^k k^l = 0$, for $l < K$
 802 and $\sum_{k=0}^K \binom{K}{k} (-1)^k k^K = (-1)^K K!$.

803 The conclusion that increasing K reduces the error floor
 804 can be confirmed by defining $f(k) = k! \epsilon_{i_m^*}^k$ and using the
 805 following fact:

$$806 \quad f(k) - f(k+1) = k! \epsilon_{i_m^*}^k (1 - (k+1)\epsilon_{i_m^*}) > 0, \quad (57)$$

807 where it is assumed that $k\epsilon_{i_m^*} \rightarrow 0$.

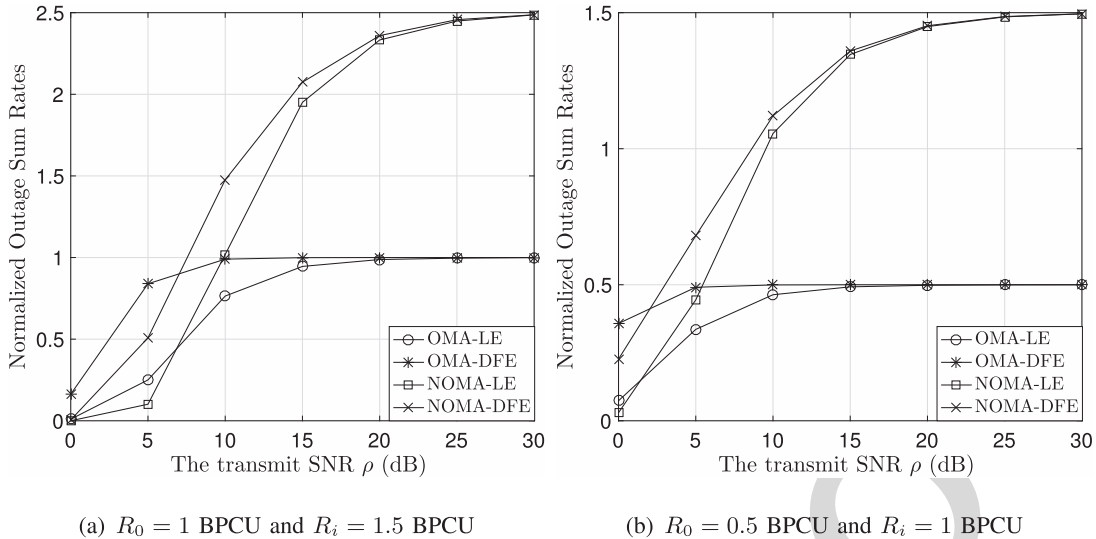


Fig. 1. Impact of OTFS-NOMA on the downlink sum rates. $M = N = K = 16$. $P_0 = P_i = 3$. BPCU denotes bit per channel use. $\gamma_0^2 = \frac{3}{4}$ and $\gamma_i^2 = \frac{1}{4}$ for $i > 0$. Random user scheduling is used.

B. Stage II of SIC

If adaptive transmission is used, the NOMA users' signals can be detected successfully during the first stage of SIC. Therefore, they can be removed from the observations at the base station, i.e., $\bar{Y}[n, m] = Y[n, m] - \sum_{q=1}^N H_q(n, m)X_q[n, m]$, and SFFT is applied to obtain the delay-Doppler observations as follows:

$$y_0[k, l] = \frac{1}{NM} \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} \bar{Y}[n, m] e^{-j2\pi(\frac{nk}{N} - \frac{ml}{M})}$$

$$= \sum_{p=0}^{P_0} h_{0,p} x_0[(k - k_{\mu_{0,p}})_N, (l - l_{\tau_{0,p}})_M] + z[k, l], \quad (58)$$

where $z[k, l]$ denote additive noise. U_0 's signals can be detected by applying either of the two considered equalization approaches, and the same performance as for OTFS-OMA can be realized. The analytical development is similar to the downlink case, and hence is omitted due to space limitations.

However, if fixed-rate transmission is used, the uplink outage events for decoding $x_0[k, l]$ are different from the downlink ones, as shown in the following. Particularly, the use of FD-LE yields the following SINR expression for decoding $x_0[k, l]$:

$$\text{SINR}_{0,kl}^{\text{LE}} = \frac{\rho}{\frac{1}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} |D_0^{k,l}|^{-2}}. \quad (59)$$

If FD-DFE is used, the SNR for detection of $x_0[k, l]$ is given by

$$\text{SINR}_{0,kl}^{\text{DFE}} = \rho \lambda_{0,kl}. \quad (60)$$

Therefore, the outage probability for detecting $x_0[k, l]$ is given by

$$P_{kl} = 1 - \text{P}(\text{SINR}_{0,kl}^{\text{DFE/LE}} > \epsilon_0, \text{SNR}_{i,n} > \epsilon_i \forall i, n)$$

$$\geq 1 - \text{P}(\text{SNR}_{i,n} > \epsilon_i \forall i, n) \geq \text{P}(\text{SNR}_{1,0} < \epsilon_i).$$

Since $\text{P}(\text{SNR}_{1,0} < \epsilon_i)$ has an error floor as shown in the previous subsection, the uplink outage probability for detection

TABLE I

DELAY-DOPPLER PROFILE FOR U_0 'S CHANNEL

Propagation path index (p)	0	1	2	3
Delay ($\tau_{0,p}$) μs	8.33	25	41.67	58.33
Delay tap index ($l_{\tau_{0,p}}$)	2	6	10	14
Doppler ($\nu_{0,p}$) Hz	0	0	468.8	468.8
Doppler tap index ($k_{\nu_{0,p}}$)	0	0	1	1

of U_0 's signals does not go to zero even if $\rho \rightarrow \infty$, which is different from the downlink case. Therefore, if fixed-rate transmission is used, adding the M low-mobility users into the bandwidth, which would be solely occupied by U_0 in OTFS-OMA, improves connectivity but degrades U_0 's performance.

VII. NUMERICAL STUDIES

In this section, the performance of OTFS-NOMA is evaluated via computer simulations. Similar to [26]–[28], we first define the delay-Doppler profile for U_0 's channel as shown in Table I, where $P_0 = 3$ and the subchannel spacing is $\Delta f = 7.5$ kHz. Therefore, the maximal speed corresponding to the largest Doppler shift $\nu_{0,3} = 468.8$ Hz is 126.6 km/h if the carrier frequency is $f_c = 4$ GHz. On the other hand, the NOMA users' channels are assumed to be time invariant with $P_i = 3$ propagation paths, i.e., $\tau_{i,p} = 0$ for $p \geq 4, i \geq 1$. For all the users' channels, we assume that $\sum_{p=0}^{P_i} \mathcal{E}\{|h_{i,p}|^2\} = 1$ and $|h_{i,p}|^2 \sim CN(0, \frac{1}{P_i+1})$. For the fixed rate transmission scheme, a simple choice for power allocation ($\gamma_0^2 = \frac{3}{4}$ and $\gamma_i^2 = \frac{1}{4}$ for $i > 0$) is considered. The performance of OTFS-NOMA could be further improved by optimizing γ_i according to the users' channel conditions and QoS requirements.

In Fig. 1, downlink OTFS-NOMA transmission is evaluated by using the normalized outage sum rate as the performance criterion which is defined as $\frac{1}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} (1 - P_{0,kl}) R_0$

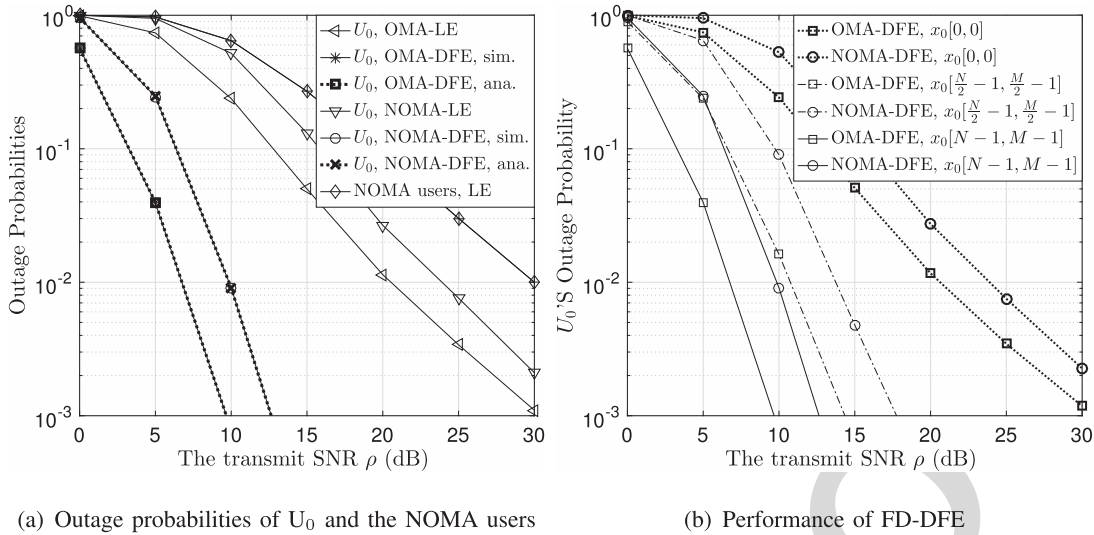


Fig. 2. The outage performance of downlink OTFS-OMA and OTFS-NOMA. $M = N = K = 16$. $P_0 = P_i = 3$. $\gamma_0^2 = \frac{3}{4}$ and $\gamma_i^2 = \frac{1}{4}$ for $i > 0$. $R_0 = 0.5$ BPCU and $R_i = 1$ BPCU. In Fig. 2(a), for FD-DFE, the performance of $x_0[N-1, M-1]$ is shown. Random user scheduling is used.

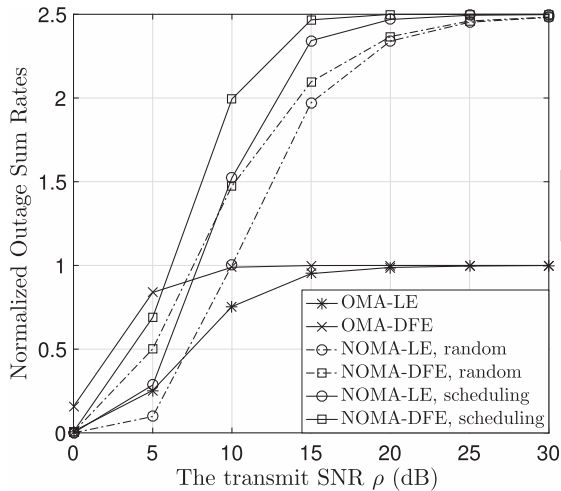


Fig. 3. Impact of user scheduling on the downlink outage sum rates. $P_0 = P_i = 3$. $R_0 = 1$ BPCU and $R_i = 1.5$ BPCU. $M = N = K = 16$, $\gamma_0^2 = \frac{3}{4}$ and $\gamma_i^2 = \frac{1}{4}$ for $i > 0$.

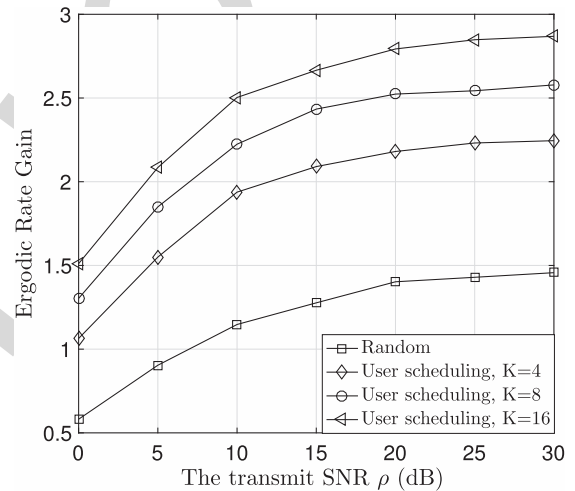


Fig. 4. The ergodic rate gain of OTFS-NOMA over OTFS-OMA. The NOMA users adapt their data rates according to (48). $P_0 = P_i = 3$. $M = N = 16$.

and $\frac{1}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} (1 - P_{0,kl}) R_0 + \frac{1}{NM} \sum_{i=1}^M \sum_{n=0}^{N-1} (1 - P_{i,n}) R_i$ for OTFS-OMA and OTFS-NOMA, respectively. Fig. 1 shows that the use of OTFS-NOMA can significantly improve the sum rate at high SNR for both considered choices of R_0 and R_i . The reason for this performance gain is the fact that the maximal sum rate achieved by OTFS-OMA is capped by R_0 , whereas OTFS-NOMA can provide sum rates up to $R_0 + R_i$. Comparing Fig. 1(a) to Fig. 1(b), one can observe that the performance loss of OTFS-NOMA at low SNR can be mitigated by reducing the target data rates, since reducing the target rates improves the probability of successful SIC. Furthermore, both figures show that FD-DFE outperforms FD-LE in the entire considered range of SNRs; however, we note that the performance gain of FD-DFE over FD-LE is achieved at the expense of increased computational complexity.

In Fig. 2, the outage probabilities achieved by downlink OTFS-OMA and OTFS-NOMA are shown. As can be seen

from Fig. 2(a), the diversity order achieved with FD-LE for detection of $x_0[k, l]$ is one, as expected from Lemma 2. As discussed in Section IV-B, one advantage of FD-DFE over FD-LE is that FD-DFE facilitates multi-path fading diversity gains, whereas FD-LE is limited to a diversity gain of one. This conclusion is confirmed by Fig. 2(a), where the analytical results developed in Corollary 1 are also verified. Fig. 2(b) shows the outage probabilities achieved by FD-DFE for different $x_0[k, l]$. As shown in the figure, the lowest outage probability is obtained for $x_0[N-1, M-1]$, whereas the outage probability of $x_0[0, 0]$ is the largest, which is due to the fact that, in FD-DFE, different signals $x_0[k, l]$ are affected by different effective channel gains, $\lambda_{0,kl}$. Another important observation from the figures is that the FD-LE outage probability is the same as the FD-DFE outage probability for detection of $x_0[0, 0]$, which fits the intuition that for FD-DFE the reliability of the first decision ($x_0[0, 0]$) is the same as that of FD-LE. For the same reason, FD-LE and FD-DFE yield similar performance for detection of the NOMA users'

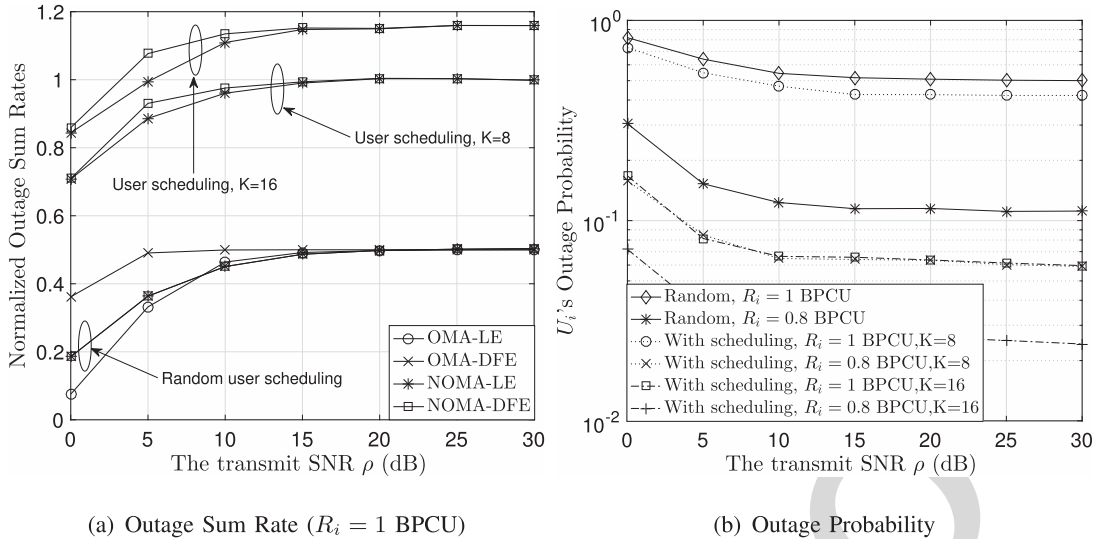


Fig. 5. The performance of uplink OTFS-NOMA. Fixed-rate transmission is used by the NOMA users. $M = N = 16$. $P_0 = P_i = 3$. $R_0 = 0.5$ BPCU. $\gamma_0^2 = \frac{3}{4}$ and $\gamma_i^2 = \frac{1}{4}$ for $i > 0$.

993 signals, since the FD-DFE outage performance is dominated
 994 by the reliability for detection of $x_0[0, 0]$, and hence is the
 995 same as that of FD-LE.

996 In addition to multi-path diversity, another degree of free-
 997 dom available in the considered OTFS-NOMA downlink
 998 scenario is multi-user diversity, which can be harvested by
 999 applying user scheduling as discussed in Section V-B. Fig. 3
 1000 demonstrates the benefits of exploiting multi-user diversity.
 1001 With random user scheduling, at low SNR, the performance
 1002 of OTFS-NOMA is worse than that of OTFS-OMA, which
 1003 is also consistent with Fig. 1. By increasing the number
 1004 of users participating in OTFS-NOMA, the performance of
 1005 OTFS-NOMA can be improved, particularly at low and moder-
 1006 ate SNR. For example, for FD-LE, the performance of
 1007 OTFS-NOMA approaches that of OTFS-OMA at low SNR
 1008 by exploiting multi-user diversity, and for FD-DFE, an extra
 1009 gain of 0.5 BPCU can be achieved at moderate SNR.

1010 In Figs. 4 and 5, the performance of uplink OTFS-NOMA is
 1011 evaluated. As discussed in Section VI, the NOMA users have
 1012 two choices for their transmission rates, namely adaptive and
 1013 fixed rate transmission. The use of adaptive rate transmission
 1014 can ensure that the implementation of NOMA is transparent
 1015 to U_0 , which means that U_0 's QoS requirements are strictly
 1016 guaranteed. Since U_0 achieves the same performance for
 1017 OTFS-NOMA and OTFS-OMA when adaptive rate transmis-
 1018 sion is used, we only focus on the NOMA users' performance,
 1019 where the ergodic rate in (49) is used as the criterion.
 1020 We note that this ergodic rate is the net performance gain of
 1021 OTFS-NOMA over OTFS-OMA, which is the reason why the
 1022 vertical axis in Fig. 4 is labeled 'Ergodic Rate Gain'. When
 1023 the M users are randomly selected from the K NOMA users,
 1024 the ergodic rate gain is moderate, e.g., 1.5 bit per channel
 1025 use (BPCU) at $\rho = 30$ dB. By applying the scheduling strategy
 1026 proposed in (50), the ergodic rate gain can be significantly
 1027 improved, e.g., nearly by a factor of two compared to the
 1028 random case with $K = 16$ and $\rho = 30$ dB.

Fig. 5 focuses on the case with fixed rate transmission, and
 similar to Fig. 1, the normalized outage sum rate is used as
 performance criterion in Fig. 5(a). One can observe that with
 random user scheduling, the sum rate of OTFS-NOMA is similar
 to that of OTFS-OMA. This is due to the fact that no inter-
 ference mitigation strategy, such as power or rate allocation,
 is used for NOMA uplink transmission, which means that U_0
 and the NOMA users cause strong interference to each other
 and SIC failure may happen frequently. By applying the user
 scheduling strategy proposed in (50), the channel conditions of
 the scheduled users become quite different, which facilitates
 the implementation of SIC. This benefit of user scheduling
 can be clearly observed in Fig. 5(a), where NOMA achieves
 a significant gain over OMA although advanced power or rate
 allocation strategies are not used. Fig. 5(a) also shows that the
 difference between the performance of FD-LE and FD-DFE is
 insignificant for the uplink case. This is due to the fact that the
 outage events during the first stage of SIC dominate the outage
 performance, and they are not affected by whether FD-LE or
 FD-DFE is employed. Another important observation from
 Fig. 5(a) is that the maximal sum rate $R_0 + R_i$ cannot be
 realized, even at high SNR. The reason for this behaviour is
 the existence of the error floor for the NOMA users' outage
 probabilities, as shown in Fig. 5(b). The analytical results
 provided in Section V-B show that increasing K can reduce
 the error floor, which is confirmed by Fig. 5(b).

VIII. CONCLUSION

In this paper, we have proposed OTFS-NOMA uplink and
 downlink transmission schemes, where users with different
 mobility profiles are grouped together for the implementa-
 tion of NOMA. The analytical results developed in the
 paper demonstrate that both the high-mobility and the low-
 mobility users benefit from the application of OTFS-NOMA.
 In particular, the use of NOMA enables the spreading of
 the signals of a high-mobility user over a large amount

of time-frequency resources, which enhances the OTFS resolution and improves the detection reliability. In addition, OTFS-NOMA ensures that the low-mobility users have access to the bandwidth resources which would be solely occupied by the high-mobility users in OTFS-OMA. Hence, OTFS-NOMA improves the spectral efficiency and reduces latency. An interesting topic for future works is studying the impact of non-zero fractional delays and fractional Doppler shifts on the performance of the developed OTFS-NOMA protocol. Furthermore, in this paper, the users' channel gains (the taps of the delay-Doppler impulse response) have been assumed to be Gaussian distributed, and an important direction for future research is to investigate the impact of other types of channel distributions on the performance of OTFS-NOMA. Moreover, the combination of emerging spectrally efficient 5G solutions, such as 5G New Radio Bandwidth Part (5G-NR-BWP) [39], [40] and software-controlled metasurfaces [41], with OTFS-NOMA is also a promising topic for future research.

APPENDIX A PROOF FOR PROPOSITION 1

Intuitively, the use of $\mathbf{F}_N \otimes \mathbf{F}_M^H$ is analogous to the application of the ISFFT which transforms signals from the delay-Doppler plane to the time-frequency plane, where inter-symbol interference is removed, i.e., the user's channel matrix is diagonalized. The following proof confirms this intuition and reveals how the diagonalized channel matrix is related to the original block circulant matrix. We first apply $\mathbf{F}_N \otimes \mathbf{I}_M$ to \mathbf{y}_0 , which yields the following:

$$\begin{aligned} & (\mathbf{F}_N \otimes \mathbf{I}_M) \mathbf{y}_0 \\ &= (\mathbf{F}_N \otimes \mathbf{I}_M) \mathbf{H}_0 \left(\gamma_0 \mathbf{x}_0 + \sum_{q=1}^M \gamma_q \mathbf{x}_q \right) + (\mathbf{F}_N \otimes \mathbf{I}_M) \mathbf{z}_0 \\ &= \text{diag} \left\{ \sum_{n=0}^{N-1} \mathbf{A}_{0,n} e^{-j \frac{2\pi l n}{N}}, 0 \leq l \leq N-1 \right\} (\mathbf{F}_N \otimes \mathbf{I}_M) \\ & \quad \times \left(\gamma_0 \mathbf{x}_0 + \sum_{q=1}^M \gamma_q \mathbf{x}_q \right) + (\mathbf{F}_N \otimes \mathbf{I}_M) \mathbf{z}_0, \end{aligned} \quad (61)$$

where $\text{diag}\{\mathbf{B}_1, \dots, \mathbf{B}_N\}$ denotes a block-diagonal matrix with \mathbf{B}_n , $1 \leq n \leq N$, on its main diagonal. Note that $\sum_{n=0}^{N-1} \mathbf{A}_{0,n} e^{-j \frac{2\pi l n}{N}}$, $0 \leq l \leq N-1$, is a sum of $N M \times M$ circulant matrices, each of which can be further diagonalized by \mathbf{F}_M . Therefore, we can apply $\mathbf{I}_N \otimes \mathbf{F}_M^H$ to $(\mathbf{F}_N \otimes \mathbf{I}_M) \mathbf{y}_0$, which yields the following:

$$\begin{aligned} & (\mathbf{I}_N \otimes \mathbf{F}_M^H) (\mathbf{F}_N \otimes \mathbf{I}_M) \mathbf{y}_0 \\ &= \text{diag} \left\{ \sum_{n=0}^{N-1} \mathbf{A}_{0,n} e^{-j \frac{2\pi l n}{N}}, 0 \leq l \leq N-1 \right\} \\ & \quad \times (\mathbf{F}_N \otimes \mathbf{I}_M) (\mathbf{I}_N \otimes \mathbf{F}_M^H) \left(\gamma_0 \mathbf{x}_0 + \sum_{q=1}^M \gamma_q \mathbf{x}_q \right) \\ & \quad + (\mathbf{I}_N \otimes \mathbf{F}_M^H) (\mathbf{F}_N \otimes \mathbf{I}_M) \mathbf{z}_0, \end{aligned} \quad (62)$$

where $\mathbf{A}_{0,n}$ is a diagonal matrix, $\mathbf{A}_{0,n} = \text{diag} \left\{ \sum_{m=0}^{M-1} a_{0,n}^{m,1} e^{j \frac{2\pi t m}{M}}, 0 \leq t \leq M-1 \right\}$, and $a_{0,n}^{m,1}$ is the element located in the m -th row and first column of $\mathbf{A}_{0,n}$.

By applying a property of the Kronecker product, $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{A}\mathbf{C}) \otimes (\mathbf{B}\mathbf{D})$, the received signals can be simplified as follows:

$$\begin{aligned} & (\mathbf{F}_N \otimes \mathbf{F}_M^H) \mathbf{y}_0 \\ &= \text{diag} \left\{ \sum_{n=0}^{N-1} \mathbf{A}_{0,n} e^{-j \frac{2\pi l n}{N}}, 0 \leq l \leq N-1 \right\} (\mathbf{F}_N \otimes \mathbf{F}_M^H) \\ & \quad \times \left(\gamma_0 \mathbf{x}_0 + \sum_{q=1}^M \gamma_q \mathbf{x}_q \right) + (\mathbf{F}_N \otimes \mathbf{F}_M^H) \mathbf{z}_0, \end{aligned} \quad (63)$$

where the $(kM + l + 1)$ -th element on the main diagonal of \mathbf{D}_0 is $D_0^{k,l}$ as defined in the proposition. The proof for the proposition is complete.

APPENDIX B PROOF FOR LEMMA 1

In order to facilitate the SINR analysis, the system model in (18) is further simplified. Define $\tilde{X}[n, m] = \sum_{i=1}^M X_i[n, m]$. With the mapping scheme used in (6), the NOMA users' signals are interleaved and orthogonally placed in the time-frequency plane, i.e., $\tilde{X}[n, m]$ is simply U_{m+1} 's n -th signal, $x_{m+1}(n)$. Denote the outcome of the SFFT of $\tilde{X}[n, m]$ by $\tilde{x}[k, l]$, which yields the following transform:

$$\tilde{x}[k, l] = \frac{1}{\sqrt{NM}} \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} \tilde{X}[n, m] e^{-j 2\pi \left(\frac{nk}{N} - \frac{ml}{M} \right)}. \quad (64)$$

Denote the $NM \times 1$ vector collecting the $\tilde{x}[k, l]$ by $\tilde{\mathbf{x}}$ and the $NM \times 1$ vector collecting the $\tilde{X}[n, m]$ by $\tilde{\mathbf{x}}$, which means that (64) can be rewritten as follows:

$$\tilde{\mathbf{x}} = (\mathbf{F}_N \otimes \mathbf{F}_M^H) \tilde{\mathbf{x}}. \quad (65)$$

Therefore, the model for the received signals in (18) can be re-written as follows:

$$\begin{aligned} \check{\mathbf{y}}_0 &= \gamma_0 \mathbf{x}_0 + \gamma_1 \tilde{\mathbf{x}} + (\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-1} \mathbf{D}_0^{-1} \tilde{\mathbf{z}}_i \\ &= \gamma_0 \mathbf{x}_0 + \gamma_1 \underbrace{(\mathbf{F}_N \otimes \mathbf{F}_M^H) \tilde{\mathbf{x}}}_{\text{Interference and noise terms}} + (\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-1} \mathbf{D}_0^{-1} \tilde{\mathbf{z}}_0, \end{aligned} \quad (66)$$

where we have used the assumption that $\gamma_i = \gamma_1$, for $1 \leq i \leq N$. Note that the power of the information-bearing signals is simply $\gamma_0^2 \rho$, and therefore, the key step to obtain the SINR is to find the covariance matrix of the interference-plus-noise term.

We first show that $\tilde{\mathbf{z}}_0 \triangleq (\mathbf{F}_N \otimes \mathbf{F}_M^H) \mathbf{z}_0$ is still a complex Gaussian vector, i.e., $\tilde{\mathbf{z}}_i \sim CN(0, \mathbf{I}_{NM})$. Recall that \mathbf{z}_0 contains NM i.i.d. complex Gaussian random variables. Furthermore, $\mathbf{F}_N \otimes \mathbf{F}_M^H$ is a unitary matrix as shown in the following:

$$\begin{aligned} & (\mathbf{F}_N \otimes \mathbf{F}_M^H) (\mathbf{F}_N \otimes \mathbf{F}_M^H)^H \stackrel{(a)}{=} (\mathbf{F}_N \otimes \mathbf{F}_M^H) (\mathbf{F}_N^H \otimes \mathbf{F}_M) \\ & \stackrel{(b)}{=} (\mathbf{F}_N \mathbf{F}_N^H) \otimes (\mathbf{F}_M^H \mathbf{F}_M) \\ & = \mathbf{I}_{NM}, \end{aligned} \quad (67)$$

where step (a) follows from the fact that $(\mathbf{A} \otimes \mathbf{B})^H = \mathbf{A}^H \otimes \mathbf{B}^H$ and step (b) follows from the fact that

1051 $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD})$. Therefore, $(\mathbf{F}_N \otimes \mathbf{F}_M^H) \mathbf{z}_0 \sim$
 1052 $CN(0, \mathbf{I}_{NM})$ given the fact that $\mathbf{z}_0 \sim CN(0, \mathbf{I}_{NM})$ and a
 1053 unitary transformation of a Gaussian vector is still a Gaussian
 1054 vector.

1055 Therefore, the covariance matrix of the interference-plus-
 1056 noise term is given by

$$1057 \mathbf{C}_{\text{cov}} = \gamma_1^2 \mathcal{E} \left\{ (\mathbf{F}_N \otimes \mathbf{F}_M^H) \check{\mathbf{x}} \check{\mathbf{x}}^H (\mathbf{F}_N \otimes \mathbf{F}_M^H)^H \right\} \\
 1058 + \mathcal{E} \left\{ (\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-1} \mathbf{D}_0^{-1} \check{\mathbf{z}}_0 \check{\mathbf{z}}_0^H \mathbf{D}_0^{-H} (\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-H} \right\}. \quad (68)$$

1061 Recall that the $(nM + m + 1)$ -th element of $\check{\mathbf{x}}$ is $\tilde{X}[n, m]$
 1062 which is equal to $x_{m+1}(n)$. Therefore, the covariance matrix
 1063 can be further simplified as follows:

$$1064 \mathbf{C}_{\text{cov}} = \gamma_1^2 \rho (\mathbf{F}_N \otimes \mathbf{F}_M^H) (\mathbf{F}_N \otimes \mathbf{F}_M^H)^H \\
 1065 + (\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-1} \mathbf{D}_0^{-1} \mathbf{D}_0^{-H} (\mathbf{F}_N \otimes \mathbf{F}_M^H)^{-H} \\
 1066 = \gamma_1^2 \rho \mathbf{I}_{MN} + (\mathbf{F}_N^H \otimes \mathbf{F}_M) \mathbf{D}_0^{-1} \mathbf{D}_0^{-H} (\mathbf{F}_N \otimes \mathbf{F}_M^H), \quad (69)$$

1068 where the noise power is assumed to be normalized.

1069 Following the same steps as in the proof of Proposi-
 1070 tion 1, we learn that, by construction, $(\mathbf{F}_N^H \otimes \mathbf{F}_M) \mathbf{D}_0^{-1} \mathbf{D}_0^{-H}$
 1071 $(\mathbf{F}_N \otimes \mathbf{F}_M^H)$ is also a block-circulant matrix, which means
 1072 that the elements on the main diagonal of $(\mathbf{F}_N^H \otimes \mathbf{F}_M)$
 1073 $\mathbf{D}_0^{-1} \mathbf{D}_0^{-H} (\mathbf{F}_N \otimes \mathbf{F}_M^H)$ are identical. Without loss of gener-
 1074 ality, denote the diagonal elements of $(\mathbf{F}_N^H \otimes \mathbf{F}_M) \mathbf{D}_0^{-1} \mathbf{D}_0^{-H}$
 1075 $(\mathbf{F}_N \otimes \mathbf{F}_M^H)$ by ϕ . Therefore, ϕ can be found by using the
 1076 trace of the matrix as follows:

$$1077 \phi = \frac{1}{NM} \text{Tr} \left\{ (\mathbf{F}_N^H \otimes \mathbf{F}_M) \mathbf{D}_0^{-1} \mathbf{D}_0^{-H} (\mathbf{F}_N \otimes \mathbf{F}_M^H) \right\} \\
 1078 = \frac{1}{NM} \text{Tr} \left\{ (\mathbf{F}_N \otimes \mathbf{F}_M^H) (\mathbf{F}_N^H \otimes \mathbf{F}_M) \mathbf{D}_0^{-1} \mathbf{D}_0^{-H} \right\} \\
 1079 = \frac{1}{NM} \text{Tr} \left\{ \mathbf{D}_0^{-1} \mathbf{D}_0^{-H} \right\} = \frac{1}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} |D_0^{k,l}|^{-2}. \quad (70)$$

1080 Therefore, the SINR for detection of $x_0[k, l]$ is given by

$$1081 \text{SINR}_{0,kl}^{LE} = \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \phi}, \quad (71)$$

1082 and the proof is complete.

1083 APPENDIX C 1084 PROOF FOR LEMMA 2

1085 The lemma is proved by first developing upper and lower
 1086 bounds on the outage probability, and then showing that both
 1087 bounds have the same diversity order.

1088 An upper bound on $\text{SINR}_{0,kl}$ is given by

$$1089 \text{SINR}_{0,kl} = \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \frac{1}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} |D_0^{k,l}|^{-2}} \\
 1090 \leq \frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \frac{1}{NM} |D_0^{0,0}|^{-2}}. \quad (72)$$

Therefore, the outage probability, denoted by $P_{0,kl}$, can be
 lower bounded as follows:

$$1091 P_{0,kl} \geq \mathbb{P} \left(\frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \frac{1}{NM} |D_0^{0,0}|^{-2}} < \epsilon_0 \right) \\
 1092 = \mathbb{P} \left(|D_0^{0,0}|^2 < \frac{\epsilon_0}{NM \rho (\gamma_0^2 - \gamma_1^2 \epsilon_0)} \right), \quad (73) \quad 1093$$

where we assume that $\gamma_0^2 > \gamma_1^2 \epsilon_0$. Otherwise, the outage
 probability is always one.

To evaluate the lower bound on the outage probability,
 the distribution of $D_0^{u,v}$ is required. Recall from (16) that $D_0^{u,v}$
 is the $((v-1)M + u)$ -th main diagonal element of \mathbf{D}_0 and
 can be expressed as follows:

$$1094 D_0^{u,v} = \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} a_{0,n}^{m,1} e^{j2\pi \frac{um}{M}} e^{-j2\pi \frac{vn}{N}}, \quad (74) \quad 1095$$

which is the ISFFT of $a_{0,n}^{m,1}$. Therefore, we have the following
 property:

$$1096 \tilde{\mathbf{D}}_0 = \sqrt{NM} \mathbf{F}_M^H \mathbf{A}_0 \mathbf{F}_N, \quad (75) \quad 1097$$

where the element in the u -th row and the v -th column of $\tilde{\mathbf{D}}_0$
 is $D_0^{u,v}$ and the element in the m -th row and the n -th column
 of \mathbf{A}_0 is $a_{0,n}^{m,1}$.

The matrix-based expression shown in (75) can be vector-
 ized as follows:

$$1098 \text{Diag}(\mathbf{D}_0) = \text{vec}(\tilde{\mathbf{D}}_0) = \sqrt{NM} \text{vec}(\mathbf{F}_M^H \mathbf{A}_0 \mathbf{F}_N) \\
 1099 = \sqrt{NM} (\mathbf{F}_N \otimes \mathbf{F}_M^H) \text{vec}(\mathbf{A}_0), \quad (76) \quad 1100$$

where $\text{Diag}(\mathbf{A})$ denotes a vector collecting all elements on
 the main diagonal of \mathbf{A} and we use the facts that $(\mathbf{C}^T \otimes$
 $\mathbf{A}) \text{vec}(\mathbf{B}) = \text{vec}(\mathbf{D})$ if $\mathbf{ABC} = \mathbf{D}$, and $\mathbf{F}_N^T = \mathbf{F}_N$.

We note that $\text{vec}(\mathbf{A}_0)$ contains only $(P_0 + 1)$ non-zero
 elements, where the remaining elements are zero. Therefore,
 each element on the main diagonal of \mathbf{D}_0 is a superposition
 of $(P_0 + 1)$ i.i.d. random variables, $h_{i,p} \sim CN\left(0, \frac{1}{P_0+1}\right)$.
 We further note that the coefficients for the superposition are
 complex exponential constants, i.e., the magnitude of each
 coefficient is one. Therefore, each element on the main di-
 agonal of \mathbf{D}_0 is still complex Gaussian distributed, i.e., $D_0^{u,v} \sim$
 $CN(0, 1)$, which means that the lower bound on the outage
 probability shown in (73) can be expressed as follows:

$$1101 P_{0,kl} \geq 1 - e^{-\frac{\epsilon_0}{NM \rho (\gamma_0^2 - \gamma_1^2 \epsilon_0)}} \doteq \frac{1}{\rho}. \quad (77) \quad 1102$$

On the other hand, an upper bound on the outage probability
 is given by

$$1103 P_{0,kl} \leq \mathbb{P} \left(\frac{\rho \gamma_0^2}{\rho \gamma_1^2 + \frac{1}{NM} \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} |D_0^{\min}|^{-2}} < \epsilon_0 \right), \quad (78) \quad 1104$$

where $|D_0^{\min}| = \min\{|D_0^{k,l}|, \forall l \in \{0, \dots, M-1\}, k \in$
 $\{0, \dots, N-1\}\}$.

Therefore, the outage probability can be upper bounded as
 follows:

$$1105 P_{0,kl} \leq \mathbb{P} \left(|D_0^{\min}|^2 < \frac{\epsilon_0}{\rho (\gamma_0^2 - \gamma_1^2 \epsilon_0)} \right). \quad (79) \quad 1106$$

It is important to point out that the $|D_0^{k,l}|^2$, $l \in \{0, \dots, M-1\}$, $k \in \{0, \dots, N-1\}$, are identically but not independently distributed. This correlation property is shown as follows. The covariance matrix of the effective channel gains, i.e., the elements on the main diagonal of \mathbf{D}_0 , is given by

$$\begin{aligned} & \mathcal{E} \{ \text{Diag}(\mathbf{D}_0) \text{Diag}(\mathbf{D}_0)^H \} \\ &= NM \mathcal{E} \{ (\mathbf{F}_N \otimes \mathbf{F}_M^H) \text{vec}(\mathbf{A}_0) \text{vec}(\mathbf{A}_0)^H (\mathbf{F}_N \otimes \mathbf{F}_M^H)^H \} \\ &= NM (\mathbf{F}_N \otimes \mathbf{F}_M^H) \mathcal{E} \{ \text{vec}(\mathbf{A}_0) \text{vec}(\mathbf{A}_0)^H \} (\mathbf{F}_N \otimes \mathbf{F}_M^H)^H. \end{aligned} \quad (80)$$

Because the channel gains, $h_{0,p}$, are i.i.d., $\mathcal{E} \{ \text{vec}(\mathbf{A}_0) \text{vec}(\mathbf{A}_0)^H \}$ is a diagonal matrix, where only (P_0+1) of its main diagonal elements are non-zero. Following the same steps as in the proof for Proposition 1, one can show that the product of $(\mathbf{F}_N \otimes \mathbf{F}_M^H)$, a diagonal matrix, and $(\mathbf{F}_N \otimes \mathbf{F}_M^H)^H$ yields a block circulant matrix, which means that $\mathcal{E} \{ \text{Diag}(\mathbf{D}_0) \text{Diag}(\mathbf{D}_0)^H \}$ is a block-circulant matrix, not a diagonal matrix. Therefore, the $|D_0^{k,l}|^2$, $l \in \{0, \dots, M-1\}$, $k \in \{0, \dots, N-1\}$, are correlated, and not independent.

Although the $|D_0^{k,l}|^2$ are not independent, an upper bound on $P_{0,kl}$ can be still found as follows:

$$\begin{aligned} P_{0,kl} &\leq \mathbb{P} \left(|D_0^{\min}|^2 < \frac{\epsilon_0}{\rho(\gamma_0^2 - \gamma_1^2 \epsilon_0)} \right) \\ &\leq \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} \mathbb{P} \left(|D_0^{k,l}|^2 < \frac{\epsilon_0}{\rho(\gamma_0^2 - \gamma_1^2 \epsilon_0)} \right) \\ &\leq MNP \left(|D_0^{0,0}|^2 < \frac{\epsilon_0}{\rho(\gamma_0^2 - \gamma_1^2 \epsilon_0)} \right) \\ &= MN \left(1 - e^{-\frac{\epsilon_0}{\rho(\gamma_0^2 - \gamma_1^2 \epsilon_0)}} \right) \doteq \frac{1}{\rho}. \end{aligned} \quad (81)$$

Since both the upper and lower bounds on the outage probability have the same diversity order, the proof of the lemma is complete.

REFERENCES

[1] Z. Ding *et al.*, "Application of non-orthogonal multiple access in LTE and 5G networks," *IEEE Commun. Mag.*, vol. 55, no. 2, pp. 185–191, Feb. 2017.

[2] W. Shin, M. Vaezi, B. Lee, D. J. Love, J. Lee, and H. V. Poor, "Non-orthogonal multiple access in multi-cell networks: Theory, performance, and practical challenges," *IEEE Commun. Mag.*, vol. 55, no. 10, pp. 176–183, Oct. 2017.

[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," *IEEE Commun. Surveys Tuts.*, vol. 19, no. 2, pp. 721–742, 2nd Quart., 2017.

[4] H. Sari, F. Vanhaverbeke, and M. Moeneclaey, "Multiple access using two sets of orthogonal signal waveforms," *IEEE Commun. Lett.*, vol. 4, no. 1, pp. 4–6, Jan. 2000.

[5] Y. Saito, A. Benjebbour, Y. Kishiyama, and T. Nakamura, "System-level performance evaluation of downlink non-orthogonal multiple access (NOMA)," in *Proc. IEEE 24th Annu. Int. Symp. Pers., Indoor, Mobile Radio Commun. (PIMRC)*, London, U.K., Sep. 2013, pp. 611–615.

[6] Z. Ding, Z. Yang, P. Fan, and H. V. Poor, "On the performance of non-orthogonal multiple access in 5G systems with randomly deployed users," *IEEE Signal Process. Lett.*, vol. 21, no. 12, pp. 1501–1505, Dec. 2014.

[7] Z. Ding, P. Fan, and H. V. Poor, "Impact of user pairing on 5G nonorthogonal multiple-access downlink transmissions," *IEEE Trans. Veh. Technol.*, vol. 65, no. 8, pp. 6010–6023, Aug. 2016.

[8] A. Brighente and S. Tomasin, "Power allocation for non-orthogonal millimeter wave systems with mixed traffic," *IEEE Trans. Wireless Commun.*, vol. 18, no. 1, pp. 432–443, Jan. 2019.

[9] Y. Wu, L. P. Qian, K. Ni, C. Zhang, and X. Shen, "Delay-minimization nonorthogonal multiple access enabled multi-user mobile edge computation offloading," *IEEE J. Sel. Topics Signal Process.*, vol. 13, no. 3, pp. 392–407, Jun. 2019.

[10] Z. Ding, P. Fan, and H. V. Poor, "Random beamforming in millimeter-wave NOMA networks," *IEEE Access*, vol. 5, pp. 7667–7681, 2017.

[11] Y. Zhou, V. W. S. Wong, and R. Schober, "Coverage and rate analysis of millimeter wave NOMA networks with beam misalignment," *IEEE Trans. Wireless Commun.*, vol. 17, no. 12, pp. 8211–8227, Dec. 2018.

[12] Z. Ding and H. V. Poor, "Design of massive-MIMO-NOMA with limited feedback," *IEEE Signal Process. Lett.*, vol. 23, no. 5, pp. 629–633, May 2016.

[13] R. Chopra, C. R. Murthy, H. A. Suraweera, and E. G. Larsson, "Analysis of nonorthogonal training in massive MIMO under channel aging with SIC receivers," *IEEE Signal Process. Lett.*, vol. 26, no. 2, pp. 282–286, Feb. 2019.

[14] A. Maatouk, E. Çalışkan, M. Koca, M. Assaad, G. Gui, and H. Sari, "Frequency-domain NOMA with two sets of orthogonal signal waveforms," *IEEE Commun. Lett.*, vol. 22, no. 5, pp. 906–909, May 2018.

[15] Y. Sun, D. W. K. Ng, Z. Ding, and R. Schober, "Optimal joint power and subcarrier allocation for full-duplex multicarrier non-orthogonal multiple access systems," *IEEE Trans. Commun.*, vol. 65, no. 3, pp. 1077–1091, Mar. 2017.

[16] L. Yin, W. O. Popoola, X. Wu, and H. Haas, "Performance evaluation of non-orthogonal multiple access in visible light communication," *IEEE Trans. Commun.*, vol. 64, no. 12, pp. 5162–5175, Dec. 2016.

[17] H. Marshoud, V. M. Kapinas, G. K. Karagiannidis, and S. Muhaidat, "Non-orthogonal multiple access for visible light communications," *IEEE Photon. Technol. Lett.*, vol. 28, no. 1, pp. 51–54, Jan. 1, 2016.

[18] Z. Ding, P. Fan, and H. V. Poor, "Impact of non-orthogonal multiple access on the offloading of mobile edge computing," *IEEE Trans. Commun.*, vol. 67, no. 1, pp. 375–390, Jan. 2019.

[19] *Study on Downlink Multiuser Superposition Transmission for LTE*, document TR 36.859, 3GPP, Mar. 2015.

[20] *Technical Specification Group Radio Access Network; Evolved Universal Terrestrial Radio Access (E-UTRA); Physical Channels and Modulation (Release 15)*, document TS 36.211, 3GPP, Jan. 2019.

[21] *Study on Non-Orthogonal Multiple Access (NOMA) for NR (Release 16)*, document TR 38.812, 3GPP, Dec. 2018.

[22] B. Di, L. Song, Y. Li, and Z. Han, "V2X meets NOMA: Non-orthogonal multiple access for 5G-enabled vehicular networks," *IEEE Wireless Commun.*, vol. 24, no. 6, pp. 14–21, Dec. 2017.

[23] Y. Chen, L. Wang, Y. Ai, B. Jiao, and L. Hanzo, "Performance analysis of NOMA-SM in vehicle-to-vehicle massive MIMO channels," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 12, pp. 2653–2666, Dec. 2017.

[24] R. Hadani and A. Monk, "OTFS: A new generation of modulation addressing the challenges of 5G," 2018, *arXiv:1802.02623*. [Online]. Available: <https://arxiv.org/abs/1802.02623>

[25] R. Hadani *et al.*, "Orthogonal time frequency space modulation," in *Proc. IEEE Wireless Commun. Netw. Conf. (WCNC)*, San Francisco, CA, USA, Mar. 2017, pp. 1–6.

[26] K. R. Murali and A. Chockalingam, "On OTFS modulation for high-Doppler fading channels," in *Proc. Inf. Theory Appl. Workshop (ITA)*, San Diego, CA, USA, Feb. 2018, pp. 1–10.

[27] P. Raviteja, Y. Hong, E. Viterbo, and E. Biglieri, "Practical pulse-shaping waveforms for reduced-cyclic-prefix OTFS," *IEEE Trans. Veh. Tech.*, vol. 68, no. 1, pp. 957–961, Jan. 2019.

[28] P. Raviteja, K. T. Phan, Y. Hong, and E. Viterbo, "Interference cancellation and iterative detection for orthogonal time frequency space modulation," *IEEE Trans. Wireless Commun.*, vol. 17, no. 10, pp. 6501–6515, Oct. 2018.

[29] G. D. Surabhi, R. M. Augustine, and A. Chockalingam, "On the diversity of uncodet OTFS modulation in doubly-dispersive channels," 2018, *arXiv:1808.07747*. [Online]. Available: <https://arxiv.org/abs/1808.07747>

[30] V. Khammammetti and S. K. Mohammed, "OTFS-based multiple-access in high Doppler and delay spread wireless channels," *IEEE Wireless Commun. Lett.*, vol. 8, no. 2, pp. 528–531, Apr. 2019.

[31] A. RezaezadehReyhani, A. Farhang, M. Ji, R. R. Chen, and B. Farhang-Boroujeny, "Analysis of discrete-time MIMO OFDM-based orthogonal time frequency space modulation," in *Proc. IEEE Int. Conf. Communicat. (ICC)*, Kansas City, MO, USA, May 2018, pp. 1–6.

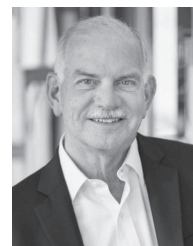
- 1264 [32] M. K. Ramachandran and A. Chockalingam, "MIMO-OTFS in high-
1265 Doppler fading channels: Signal detection and channel estimation,"
1266 in *Proc. IEEE GLOBECOM*, Kansas City, MO, USA, Dec. 2018,
1267 pp. 206–212.
- 1268 [33] P. Raviteja, E. Viterbo, and Y. Hong, "OTFS performance on static
1269 multipath channels," *IEEE Wireless Commun. Lett.*, vol. 8, no. 3,
1270 pp. 745–748, Jun. 2019.
- 1271 [34] D. Falconer, S. L. Ariyavisitakul, A. Benyamin-Seeyar, and B. Eidson,
1272 "Frequency domain equalization for single-carrier broadband wireless
1273 systems," *IEEE Commun. Mag.*, vol. 40, no. 4, pp. 58–66, Apr. 2002.
- 1274 [35] J. Louveaux, L. Vandendorpe, and T. Sartenar, "Cyclic prefixed single
1275 carrier and multicarrier transmission: Bit rate comparison," *IEEE*
1276 *Commun. Lett.*, vol. 7, no. 4, pp. 180–182, Apr. 2003.
- 1277 [36] L. Zheng and D. N. C. Tse, "Diversity and multiplexing: A fundamental
1278 tradeoff in multiple-antenna channels," *IEEE Trans. Inf. Theory*, vol. 49,
1279 no. 5, pp. 1073–1096, May 2003.
- 1280 [37] B. Devillers, "Cyclic prefixed block transmission for wireless commu-
1281 nications: Performance analysis and optimization," Ph.D. dissertation,
1282 Université catholique de Louvain, Ottignies-Louvain-la-Neuve, Belgium,
1283 2009.
- 1284 [38] B. Devillers, J. Louveaux, and L. Vandendorpe, "About the diversity in
1285 cyclic prefixed single-carrier systems," *Phys. Commun.*, vol. 1, no. 4,
1286 pp. 266–276, Dec. 2008.
- 1287 [39] J. Jeon, "NR wide bandwidth operations," *IEEE Commun. Mag.*, vol. 56,
1288 no. 3, pp. 42–46, Mar. 2018.
- 1289 [40] C. Sexton, N. Marchetti, and L. A. DaSilva, "Customization and
1290 trade-offs in 5G RAN slicing," *IEEE Commun. Mag.*, vol. 57, no. 4,
1291 pp. 116–122, Apr. 2019.
- 1292 [41] C. Liaskos, S. Nie, A. Tsioliaridou, A. Pitsillides, S. Ioannidis, and
1293 I. Akyildiz, "A new wireless communication paradigm through software-
1294 controlled metasurfaces," *IEEE Commun. Mag.*, vol. 56, no. 9,
1295 pp. 162–169, Sep. 2018.

the Vodafone Foundation for Research in Mobile Communications, the
2006 UBC Killam Research Prize, the 2007 Wilhelm Friedrich Bessel
Research Award of the Alexander von Humboldt Foundation, the 2008
Charles McDowell Award for Excellence in Research from UBC, the 2011
Alexander von Humboldt Professorship, the 2012 NSERC E.W.R. Steacie
Fellowship, and the 2017 Wireless Communications Recognition Award by
the IEEE Wireless Communications Technical Committee. He was listed as
a 2017 Highly Cited Researcher by the Web of Science. He is also the Chair
of the Steering Committee of IEEE TRANSACTIONS ON MOLECULAR, BIO-
LOGICAL AND MULTI-SCALE COMMUNICATION, a member of the Editorial
Board of PROCEEDINGS OF THE IEEE, a Member-at-Large of the Board
of Governors of ComSoc, and the ComSoc Director of journals. He is also a
Distinguished Lecturer of the IEEE Communications Society (ComSoc). From
2012 to 2015, he served as the Editor-in-Chief of IEEE TRANSACTIONS ON
COMMUNICATIONS.



Pingzhi Fan (M'93–SM'99–F'15) received the M.Sc. degree in computer science from Southwest Jiaotong University, China, in 1987, and the Ph.D. degree in electronic engineering from Hull University, U.K., in 1994.

He was the Chief Scientist of the National 973 Research Project (MoST) from 2012 to 2016. He is currently a Professor and the Director of the Institute of Mobile Communications, Southwest Jiaotong University. He has been a Visiting Professor with Leeds University, U.K., since 1997, and has been a Guest Professor with Shanghai Jiaotong University since 1999. He has over 280 research papers published in various international journals and eight books (including edited). He is the inventor of 22 granted patents. His research interests include vehicular communications, wireless networks for big data, and signal design and coding. He is also a fellow of IET, CIE, and CIC. He was a recipient of the U.K. ORS Award in 1992 and the Outstanding Young Scientist Award (NSFC) in 1998. He has served as the general chair or the TPC chair of a number of international conferences. He is also the Founding Chair of IEEE VTS BJ Chapter, IEEE ComSoc CD Chapter, and IEEE Chengdu Section. He is also the guest editor or editorial member of several international journals. He has also served as the Board Member of IEEE Region 10, IET (IEE) Council, and IET Asia–Pacific Region. He is also an IEEE VTS Distinguished Lecturer (2015–2019).



H. Vincent Poor (M'77–SM'82–F'87) received the Ph.D. degree in EECS from Princeton University in 1977.

From 1977 to 1990, he was on the faculty of the University of Illinois at Urbana–Champaign. Since 1990, he has been on the faculty at Princeton University, where he is currently the Michael Henry Strater University Professor of electrical engineering. From 2006 to 2016, he served as the Dean of the School of Engineering and Applied Science, Princeton University. He has also held visiting appointments at several other universities, including most recently at Berkeley and Cambridge. His research interests are in the areas of information theory and signal processing, and their applications in wireless networks, energy systems, and related fields. Among his publications in these areas is the recent book *Multiple Access Techniques for 5G Wireless Networks and Beyond*. (Springer, 2019). He is also a member of the National Academy of Engineering and the National Academy of Sciences. He is also a Foreign Member of the Chinese Academy of Sciences, the Royal Society, and other national and international academies. He was a recipient of the Marconi and Armstrong Awards of the IEEE Communications Society in 2007 and 2009, respectively. Recent recognition of his work includes the 2017 IEEE Alexander Graham Bell Medal, the 2019 ASEE Benjamin Garver Lamme Award, the D.Sc. (*honoris causa*) from Syracuse University awarded in 2017, and the D.Eng. (*honoris causa*) from the University of Waterloo awarded in 2019.



Zhiguo Ding (S'03–M'05–SM'15) received the B.Eng. degree in electrical engineering from the Beijing University of Posts and Telecommunications in 2000 and the Ph.D. degree in electrical engineering from Imperial College London in 2005.

From July 2005 to April 2018, he was with Queen's University Belfast, Imperial College, Newcastle University, and Lancaster University. Since April 2018, he has been a Professor of communications with The University of Manchester. From October 2012 to September 2018, he was an

Academic Visitor with Princeton University. His research interests are 5G networks, game theory, cooperative and energy harvesting networks, and statistical signal processing. He was a recipient of the Best Paper Award at IET ICWMC-2009 and IEEE WCSP-2014, the EU Marie Curie Fellowship (2012–2014), the Top IEEE TVT Editor 2017, the 2018 IEEE Communication Society Heinrich Hertz Award, the 2018 IEEE Vehicular Technology Society Jack Neubauer Memorial Award, and the 2018 IEEE Signal Processing Society Best Signal Processing Letter Award. He was an Editor of IEEE WIRELESS COMMUNICATIONS LETTERS and IEEE COMMUNICATIONS LETTERS from 2013 to 2016. He has been serving as an Editor for IEEE TRANSACTIONS ON COMMUNICATIONS, IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY, and *Journal of Wireless Communications and Mobile Computing*.



Robert Schober (M'01–SM'08–F'10) received the Diploma (Univ.) and Ph.D. degrees in electrical engineering from the Friedrich-Alexander University of Erlangen-Nuremberg (FAU), Germany, in 1997 and 2000, respectively.

From 2002 to 2011, he was a Professor and the Canada Research Chair with The University of British Columbia (UBC), Vancouver, Canada. Since January 2012, he has been an Alexander von Humboldt Professor and the Chair for Digital Communication with FAU. His research interests fall

into the broad areas of communication theory, wireless communications, and statistical signal processing. He is also a fellow of the Canadian Academy of Engineering and the Engineering Institute of Canada. He was a recipient of several awards for his work, including the 2002 Heinz Maier-Leibnitz Award of the German Science Foundation (DFG), the 2004 Innovations Award of

1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
13491350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
13731374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399

AQ:5

AQ:6