

Citation for the published version:

Papazafeiropoulos, A., & Ratnarajah, T. (2018). Modeling and Performance of Uplink Cache-Enabled Massive MIMO Heterogeneous Networks. IEEE Transactions on Wireless Communications, 17(12), 8136-8149. [8490665].
<https://doi.org/10.1109/TWC.2018.2874229>

Document Version: Accepted Version

Link to the final published version available at the publisher:

<https://doi.org/10.1109/TWC.2018.2874229>

© 2018 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

General rights

Copyright© and Moral Rights for the publications made accessible on this site are retained by the individual authors and/or other copyright owners.

Please check the manuscript for details of any other licences that may have been applied and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights. You may not engage in further distribution of the material for any profitmaking activities or any commercial gain. You may freely distribute both the url (<http://uhra.herts.ac.uk/>) and the content of this paper for research or private study, educational, or not-for-profit purposes without prior permission or charge.

Take down policy

If you believe that this document breaches copyright please contact us providing details, any such items will be temporarily removed from the repository pending investigation.

Enquiries

Please contact University of Hertfordshire Research & Scholarly Communications for any enquiries at rsc@herts.ac.uk

Modeling and Performance of Uplink Cache-Enabled Massive MIMO Heterogeneous Networks

Anastasios Papazafeiropoulos and Tharmalingam Ratnarajah

Abstract—A significant burden on wireless networks is brought by the uploading of user-generated contents to the Internet by means of applications such as the social media. To cope with this mobile data tsunami, we develop a novel multiple-input multiple-output (MIMO) network architecture with randomly located base stations (BSs) a large number of antennas employing cache-enabled *uplink* transmission. In particular, we formulate a scenario, where the users upload their content to their strongest base stations (BSs), which are Poisson point process (PPP) distributed. In addition, the BSs, exploiting the benefits of massive MIMO, upload their contents to the core network by means of a finite-rate backhaul. After proposing the caching policies, where we propose the modified von Mises distribution as the popularity distribution function, we derive the outage probability and the average delivery rate by taking advantage of tools from the deterministic equivalent (DE) and stochastic geometry analyses. Numerical results investigate the realistic performance gains of the proposed heterogeneous cache-enabled uplink on the network in terms of cardinal operating parameters. For example, insights regarding the BSs storage size are exposed. Moreover, the impacts of the key parameters such the file popularity distribution, and the target bitrate are investigated. Specifically, the outage probability decreases if the storage size is increased, while the average delivery rate increases. In addition, the concentration parameter, defining the number of files stored at the intermediate nodes (popularity), affects directly the proposed metrics. Furthermore, a higher target rate results in higher outage because fewer users obey this constraint. Also, we demonstrate that a denser network decreases the outage and increases the delivery rate. Hence, the introduction of caching at the uplink of the system design ameliorates the network performance.

Index Terms—Caching, channel aging, heterogeneous networks, massive MIMO, stochastic geometry

I. INTRODUCTION

A vast majority of new wireless services such as social networks, web-browsing applications, and multimedia streaming has fueled the mobile data traffic with an imminent 500-fold boost over the next 10 years [1]. As a result, mobile operators need to redesign their current networks and delve into more

sophisticated techniques to surge system capacity forward and expand coverage in fifth generation (5G) networks [2].

A promising solution towards this direction relies on the deployment of low-power, short-range, and cost-efficient small cell networks (SCNs) or else heterogeneous networks (HetNets). In fact, irregular cellular networks, deployed opportunistically and in hot spots have been researched for a fairly long time now [3]. Specifically, downlink single-input single-output (SISO) HetNets have already presented sufficient progress [4]–[7]. Literally, relevant standardization activities have started in 3GPP release a long time ago [8]. Having assumed that SISO studies are anachronistic, efforts have been devoted to the challenge of modeling multi-antenna HetNets by assuming perfect channel state information (CSI) [9]–[11]. In addition, the practical consideration of imperfect CSI has taken place in several works [12], [13]. Moreover, research has been devoted to the analysis of stochastically geometric uplink models [14], [15].

In a parallel avenue, massive multiple-input multiple-output (MIMO) has emerged as another technology supporting the backbone of 5G networks. Remarkably, massive MIMO point to the increase of spectral and energy efficiencies [16], [17]. The achievement of high cell-throughput along with simple signal processing has been contrived by deploying large-scale antenna arrays at the BS and multi-user (MU) transmission. Starting from the strong assumption of perfect CSI, research has faced realistic impediments such as the presence of pilot contamination [18] and [16], the inevitable hardware impairments [19] and [20], as well as the channel aging [21]–[23]. Especially, by applying the theory of deterministic equivalent (DE) analysis, massive MIMO systems were studied under the presence of channel aging [21]. The key assumption of the DE analysis is that $M \rightarrow \infty$ and $K \rightarrow \infty$ with a given ratio, where M and K are the numbers of BS antennas and users, respectively. In particular, channel aging refers to the channel variation between the time instance the channel is estimated and the time instance it is used for precoding or detection. The sources of channel aging are mainly the relative movement of the users with the BS antennas, the phase noise, and any processing delays [23]. Despite its significant implications, few works have scrutinized its impact on massive MIMO systems [13], [21]–[23].

The growing trend of user-generated content such as the sharing of real-time events by means of smartphones inflicts a great uploading strain to the wireless networks. Despite the importance of uplink on mobile cellular networks, most efforts

A. Papazafeiropoulos was with the Institute for Digital Communications (IDCOM), University of Edinburgh, Edinburgh, EH9 3JL, U.K. He is now with the Communications and Intelligent Systems Group, University of Hertfordshire, AL10 9AB Hatfield, U. K. and with SnT at the University of Luxembourg, L-1855 Luxembourg. T. Ratnarajah is with IDCOM, University of Edinburgh, Edinburgh. Email: tapapazaf@gmail.com; t.ratnarajah@ed.ac.uk

This work was supported by the U.K. Engineering and Physical Sciences Research Council (EPSRC) under Grant EP/N014073/1 and the UK-India Education and Research Initiative Thematic Partnerships under grant number DST-UKIERI-2016-17-0060.

have been focused on the downlink scenario [24]–[26]. In fact, few attempts have been dedicated to the expanding demands of users’ transmission (uploading). Notably, differences appear between the procedures of uploading and downloading. Specifically, an asymmetry regarding the downlink and uplink bandwidths takes place, since the downlink bandwidth can reach 10 – 1000 times the corresponding uplink bandwidth. As a consequence, the uploading time is longer, and the throughput is lower. Hence, lower quality of experience is met. Another difference concerns the limited resources of mobile devices such as the battery capacity and the transmit power. Obviously, it is a critical solution to moderate the uplink traffic pressure. An efficient remedy that is brought to the play is the employment of caching in SCNs by exploiting the content popularity appearing as redundancy [27]. In fact, caching the users’ content in the edge of the network brings gains regarding the user satisfaction and the traffic load [24]. However, most works address the problem of caching in the downlink direction.

A. Motivation-Central Idea

This paper is motivated by the following observations: 1) Content providers relocate their users’ contents from the core network to the intermediate nodes in the network (caching). Hence, a key question to answer is the design and investigation of the converse scenario, where the users move their content to the intermediate nodes to alleviate the upload traffic. 2) We employ a large number of antennas at each BS to take advantage of the benefits of massive MIMO such as the elimination of intra-cell interference. 3) Networks have the tendency to become denser resulting in their irregularity, i.e., it is required to introduce the concept of SCNs. 4) User mobility and its resultant channel aging is a common phenomenon in wireless communications. The main contributions are summarized as follows.

- Contrary to existing works [24], [27], which have studied the downlink caching, we focus on a novel strategy described as *uplink caching*. Hence, instead of having the users requesting contents from their associated BS, the users upload their contents to their strongest BS. More concretely, we formulate the caching problem in the uplink.
- In parallel of considering uplink caching, we take advantage of the gained performance benefits when the massive MIMO technology is employed and HetNet design is encountered. As far as the authors are aware, it is worthwhile to mention that this work is unique regarding the study of the notion of massive MIMO in caching.
- We introduce the modified von Mises distribution instead of a power-law or the Zipf distribution to describe the content popularity distribution, in order to represent the locality and the concentration of the content popularities in a specific region.
- It is the first work introducing channel aging in an architecture including caching. Although we do not show the degradation of the system performance in terms of plots by varying the user mobility due to users’ relative movement, the analytical results allow the observation of its dependence and its loss quantification.

- We derive the outage probability and the average delivery rate of an uplink massive MIMO HetNet, where the intermediate nodes are enriched with caching resources. In particular, after having obtained the deterministic signal-to-interference-plus-noise ratio (SINR) by means of the theory of DEs, we achieve to obtain a statistical expression. Notably, the main benefit of the DEs is the provision of deterministic expressions allowing to avoid any Monte Carlo simulations.
- Relied on the numerical results, we elaborate on the impact of various parameters such as the BSs density and the storage size. For example, a high storage size induces improvement of the system, since the outage probability decreases and the average delivery rate increases. For the sake of comparison, we also present the results corresponding to the absence of caching, where applicable.

B. Paper Outline

The remainder of this paper has the following structure. Section II develops the system model of the uplink of a massive MIMO HetNet with channel aging. Section III presents the caching model, while Section IV provides the estimated channel including the effects of pilot contamination and channel aging. Next, Section V presents the uplink transmission under the presence of channel aging and introduction of the caching concept. Section VI provides the main results of this study. Especially, Subsection VI-A includes the derivation and investigation of the outage probability, while Subsection VI-B, provides the presentation of the average delivery rate of this general model. The numerical results are placed in Section VII, and Section VIII concludes the paper.

C. Notation

Vectors and matrices are denoted by boldface lower and upper case symbols. The notations $\mathcal{C}^{M \times 1}$ and $\mathcal{C}^{M \times N}$ refer to complex M -dimensional vectors and $M \times N$ matrices, respectively. The symbols $(\cdot)^T$, $(\cdot)^H$, and $\text{tr}(\cdot)$ express the transpose, Hermitian transpose, and trace operators, respectively. The expectation operator is denoted by $\mathbb{E}[\cdot]$, and the symbol \triangleq declares definition. Moreover, $J_0(\cdot)$ is the zeroth-order Bessel function of the first kind, and $\Gamma(x, y)$ denotes the Gamma distribution with shape and scale parameters x and y , respectively. Finally, $\mathbf{b} \sim \mathcal{CN}(\mathbf{0}, \Sigma)$ represents a circularly symmetric complex Gaussian vector with zero-mean and covariance matrix Σ .

II. SYSTEM MODEL

This section considers the setup for the uplink of a massive MIMO cellular network consisted of BSs with locations drawn according to an independent PPP Φ_B with density λ_B , i.e., this formulation corresponds to the generalized and quite interesting design of uplink massive MIMO HetNets. For the sake of better description, a graphical representation of the system layout is shown in Fig 1. Specifically, let the BS of the l th cell having a large number of antennas, denoted by M_l . The users are assumed to be distributed as an independent PPP

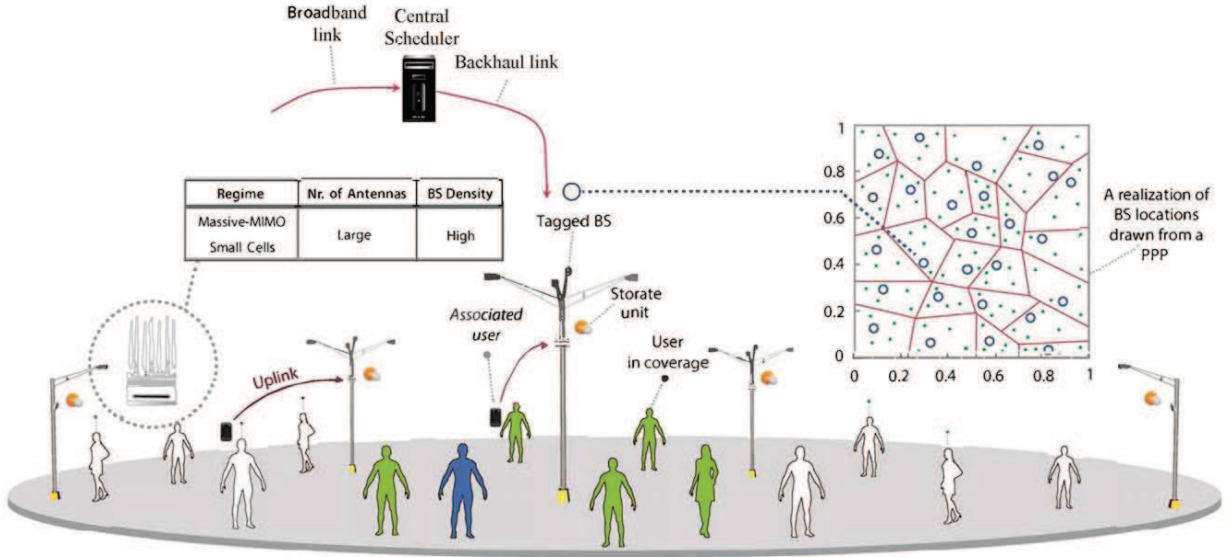


Figure 1: An illustration of the massive MIMO heterogeneous network model with caching. The main figure shows the configuration of BSs, users, storage units, the backhaul, and their interconnection. The rectangular at the top right side illustrates a possible HetNet. The green solid disks and the grey rings represent the users and the BSs, respectively.

with a sufficiently high density λ_b . In fact, in any resource block, we assume that the l th BS randomly schedules K_l users according to a distance-based criterion. Actually, the K_l users are connected with the nearest BS constituting its Voronoi cell, while a Voronoi tessellation is structured by the set of all these cells [28]¹. In other words, these users are connected with the strongest BS constituting its Voronoi cell, while a Voronoi tessellation is structured by the set of all these cells². Also, we assume that $M_l \gg K_l$, as stated by the basic principle of massive MIMO technology³. Moreover, we consider that each user is equipped with just a single-antenna mobile terminal. Evidently, since the massive MIMO concept is employed, many degrees of freedom are shared across each cell. A central scheduler provides a fixed broadband connection to these BSs by means of wired backhaul links. Obviously, the capacity of the link between the backhaul and each BS is a decreasing function of λ_B , since the deployment of more BSs per given area results in less capacity per backhaul link. A capacity expression, obeying to this property, will be introduced in Subsection VI-B.

Further to the network topology, by exploiting Slivnyak's theorem, we are able to conduct the analysis just by focusing on a randomly chosen BS found at the origin [30]. Hereafter, we refer to this BS as the tagged BS. Hence, we assume that at time n , the location of the associated scheduled user is at

¹In this work, we assume only a non-line-of-sight (NLoS) transmission, while the consideration of an LoS component is left for future work. Its introduction in the analysis could be made by means of a multi-slope path-loss model [29].

²The users are assumed to be distributed as an independent point process, but each cell is large enough (the density of the users' PPP is sufficiently large) to shelter K_l users. The users in each cell are independently and uniformly distributed.

³Although, in practice the number of BS antennas M_l , and the number of associated users K_l differ across cells, henceforth, we assume $M_l = M$ and $K_l = K$ for the sake of simplicity.

$x_{lk,n}$, while the location of the k th user found in the j th cell is denoted by $x_{ljk,n}$. Similarly, $\mathbf{h}_{ljk,n} \in \mathbb{C}^{M \times 1}$ is the channel vector from the associated k th user in the cell located at $x_{ljk,n}$ to the tagged BS (located at the origin), while the interference term $\mathbf{h}_{ljk,n} \in \mathbb{C}^{M \times 1}$ is the channel vector corresponding to the link from the k th user of the j th cell located at $x_{ljk,n} \in \mathbb{R}^2$. The locations of the k th scheduled users from all the cells are formed by a non-stationary point process Φ_k , which is not a PPP because of the correlation of their locations with the BS process. The explanation relies on the prohibition of the presence of all other users in Φ_k in the tagged cell [14], [15]. Although this kind of correlations regarding the scheduled users' locations make the exact analysis intractable, we endorse the uplink model, accounting for the pairwise correlations, proposed in [31] and followed in [15]. Moreover, we consider an exclusion ball approximation on the distribution of the scheduled user process Φ_k , being a first-order approximation of the model in [31]. On the top of the ball approximation, we assume that the random variable, expressing the distance $\|x_{lk,n}\|^{-\alpha}$ from the scheduled user to its tagged BS at the origin, is assumed to be a Rayleigh variable with a mean of $0.5\sqrt{1/\lambda_B}$ [14]. We denote $R_e = \sqrt{1/(\pi\lambda_B)}$ the radius of the ball, in order to let the size of the surface of the exclusion ball equal the average cell size, which is $1/\lambda_B$ [32]. Furthermore, the scheduled user process Φ_k , describing the locations of the other scheduled users, is formed as an inhomogeneous PPP of density λ_k where the users are found outside an exclusion ball having as center the tagged BS. Especially, the locations of the k th users of the other cells, belonging to Φ_k , are modeled by means of an inhomogeneous PPP with a density function of

$$\lambda_k(r) = \lambda_B \left(1 - e^{-\lambda_B \pi r^2}\right), \quad (1)$$

where $r = \|x_{ljk,n}\|$.

Both the uplink training and data transmission phases necessitate the introduction of fractional power control in our analysis. Thus, the k user in the l th cell transmits with power

$$P_{lk,n} = P_t \beta_{lk,n}^{-\epsilon}, \quad (2)$$

where $\epsilon \in [0, 1]$ expresses the fraction of the compensation of the path-loss given by $\beta_{lk,n}^{-\epsilon}$, while P_t is the open loop transmit power assuming no power control.

As far as the channel model is concerned, let the point-to-point channels be characterized by independent and identically distributed (i.i.d.) Rayleigh block fading with unit mean, while we assume a block fading model, where the channel is assumed constant during one block, but varies independently from block to block. Note that although the assumption of both line and non-line of sight signals appear in small cells, we consider only Rayleigh fading for the sake of simplicity. Relaxation of the Rayleigh fading assumption as well as the introduction and study of other fading models can be done with techniques found in [33], and is left for future work. Hence, in the proposed model, the channel vector $\mathbf{h}_{ljk,n}$ from user k in the j th cell at the n -th time slot is modelled as

$$\mathbf{h}_{ljk,n} = \beta_{ljk}^{1/2} \mathbf{w}_{ljk,n}, \quad (3)$$

where β_{ljk} is the large-scale path-loss and $\mathbf{w}_{ljk,n} \in \mathbb{C}^{M \times 1}$ is an uncorrelated fast fading Gaussian channel vector with elements having zero mean and unit variance, i.e., $\mathbf{w}_{ljk} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_M)$. Note that the incorporation of spatial correlation due to lack of limited antenna spacing, and different antenna patterns is left for future work. The path-loss at the tagged cell is described by

$$\beta_{ljk} = Cr_{ljk}^{-\alpha}, \quad (4)$$

where $\alpha > 2$ is the path-loss exponent, and C expresses a constant determined by the carrier frequency and the reference distance. Given that we have assumed an NLoS component, the distance-based criterion is translated to path-loss based, i.e., the minimum distance corresponds to the minimum path loss signal. For the sake of exposition, we have assumed a simplistic single-slope path-loss model, in order not to distract the reader from the main contributions. The application of a more complex model, such as the multi-slope path-loss model presented in [34], is outside of the scope of this paper and left to future work.

The transmission scheme includes an uplink channel estimation phase, and continues with an uplink data transmission phase, allowing the derivation of the outage probability and the average delivery rate, in order to shed light on their behavior. However, we need first to introduce the caching model.

III. CACHING MODEL

There are definitely certain cases that BSs need to cache the files of users in the uplink and backhaul for further cost savings, latency reduction, etc. For example, imagine a crowded scenario where many users are willing to upload their video recordings to the Internet/network. If uploaded files could be proactively cached at the BSs, that could be beneficial to the network if

nearby users have later the interest to download/watch those uploaded videos. Proactive caching at the uplink could alleviate the upload traffic. This can be also a criterion to select which file is of high interest. Specifically, assuming that the nodes/users upload files via uplink, there is a chance that the user who is uploading the file will have many downloads (i.e., Justin Bieber sharing a video, most likely will be viewed by his followers). Therefore, the file of interest could be inferred based on proactive prediction methods of how the uploading device is "influential", relying on machine learning, social networks, etc. Moreover, although there is no not too much study about the role of caching in the uplink, uplink caching will be very beneficial when the Internet of Things (IoT) devices will be introduced on cellular networks [35, Fig. 37].

Let us consider that the network has a *content catalog* of F contents represented by the set of $\mathcal{F} = \{f_1, \dots, f_F\}$. User k in j th cell at the n -th time slot (located at $x_{ljk,n}$) demands a content from a sub-catalog $\mathcal{F}_i \subseteq \mathcal{F}$ according to a content popularity distribution f_{pop} . In particular, the BS at the tagged cell has a content popularity distribution f_{pop} and is modelled by a *modified* von Mises distribution [36], which is a symmetric circular distribution defined as

$$f_{\text{pop}}(f, \mu_j, \gamma_j) = \frac{e^{\gamma_j \cos((2\pi(f-\mu_j)/F) - \pi)}}{2\pi J_0(\gamma_j)}, \quad (5)$$

where f is a point in the support such as $0 \leq f \leq F$, the parameter μ_j is a measure of location such as $0 \leq \mu_j \leq F$, the parameter γ_j is a measure of concentration with $0 \leq \gamma_j \leq +\infty$, and the function $J_0(\gamma_j)$ is the zeroth-order Bessel function of the first kind. The distribution becomes uniform when $\gamma_j = 0$ and highly concentrated on the point μ_j when $\gamma_j \rightarrow \infty$. The parameters μ_j and $1/\gamma_j$ are analogous to the mean and variance in Gaussian distribution. In fact, when $\gamma_j \rightarrow \infty$, we obtain the Dirac delta function. The intuition behind such a distribution and modelling is due to the observation that the content catalog is finite $[0, F]$ and the content popularities might be concentrated on specific region, where parameters μ_j and $1/\gamma_j$ are used for its description. Suppose that the j -th BS has a storage capacity of S_j nats with $S_j \leq F$ (1 bit = $\ln(2) = 0.693$ nats), and caches files according to the policy provided below. Henceforth, all the parameters concerning caching are the same across all cells, e.g., we assume that $\mu_j = \mu$ and $\gamma_j = \gamma$. The length of each file in the catalog is L nats, while T expresses its bitrate in nats/s/Hz. Note that the uplink rate of each user has to be equal or higher than the file bitrate T , in order to avoid any interruption during its experience.

We assume that we store the most popular files from the catalog in advance offline. Storing most popular files requires perfect knowledge of the content popularity, which might not be possible to be constructed locally. In order to make the things local/geographical, an interesting caching policy is to store the S th closest different files mentioned above.

If the file is of high interest for the BS (if it will be a popular file in the future and is not cached yet), then the BS should cache it, i.e., uplink transmission incurs. If the file is of high interest for the BS (if it will be a popular file in the future) and is already cached, then the BS should do nothing (cache miss).

In other words, in such case, cache hit means that the user is in coverage and the file is not included in the BS. Hence, uploading to the BS is meaningful; otherwise, if the user is not in coverage or the file is included at the BS, the file will not be uploaded or it will be uploaded to the core network from the BS.

Notably, the uplink caching process is dynamic, since the users are likely to have a popular content at any time, which should be uploaded for the better performance of the network. Hence, the proposed model is constantly vital. However, even when all the users have uploaded their contents, they are able to exploit the model and focus on downloading a content that already has been uploaded. The latter scenario is very rare because it is out of chance that at some point all the users will have uploaded their contents. At any time, there will be at least one user that will have some popular content to upload.

IV. CHANNEL ESTIMATION

Let us denote the channel coherence time is T_c . We assume that the same time-frequency resources are shared by the users across all cells. Aiming to the characterization of realistic systems, we account for imperfect CSIT due to pilot contamination and channel aging. Let τ denote the length of the training period. Obeying to time-division duplex (TDD) design, during the uplink training phase, having duration τ symbols, the tagged BS obtains the estimated channel. The uplink data transmission phase consists of $T_c - \tau$ symbols.

Having in mind that the signal from each user is attenuated with distance because of the path-loss, we present the pilot contamination occurred due to the re-use of the pilot sequences during the training phase.

A. Pilot Contamination

According to TDD, estimation of the local CSI takes place during the uplink training phase, where the same band of frequencies is shared across all cells. Moreover, the k th user in each cell is assigned with the same pilot sequence. As a result, pilot contamination occurs and the degradation of the system performance is inevitable. Let the superscript tr describe the training stage. Furthermore, the scheduled user processes with different pilots, i.e., $\Phi_k, \Phi_{k'}$, are assumed to be independent. The tagged BS receives a noisy observation of the channel vector from the associated scheduled user at time instance n . The average power of each transmitted pilot symbol from the scheduled user is $P_{lk,n}$. Hence, the associated BS observes the channel $\mathbf{h}_{lk,n}$ as

$$\mathbf{y}_{lk,n}^{\text{tr}} = \underbrace{\sqrt{P_{lk,n}}\mathbf{h}_{lk,n}}_{\text{Desired signal}} + \underbrace{\sum_{j \neq l} \sqrt{P_{jk,n}}\mathbf{h}_{lj,n}}_{\text{Interference part}} + \underbrace{\mathbf{N}_{lk,n}^{\text{tr}}\boldsymbol{\psi}_k^{\text{H}}}_{\text{noise}}, \quad (6)$$

where the vector $\boldsymbol{\psi}_k \in \mathbb{C}^{\tau \times 1}$ denotes the training sequence of the k th user with $\boldsymbol{\psi}_k\boldsymbol{\psi}_k^{\text{H}}=1$, and $\mathbf{N}_{lk,n}^{\text{tr}} \in \mathbb{C}^{M \times \tau}$ is the spatially white additive Gaussian noise matrix with i.i.d. entries distributed as $\mathcal{CN}\left(0, \frac{\sigma^2}{K}\right)$. Note that the channel vectors $\mathbf{h}_{lj,n}$,

being independent across cells and user distances, are Gaussian distributed as $\mathcal{CN}\left(\mathbf{0}, \mathbf{I}_M\right)$.

The tagged BS estimates $\mathbf{h}_{lk,n}$ by applying minimum mean square error (MMSE) estimation to (6), and by assuming that the tagged BS knows perfectly the large-scale path-losses β_{ljk} for $j \neq l$. Thus, the estimated channel is

$$\begin{aligned} \hat{\mathbf{h}}_{lk,n} \Big|_{\|r_{lj,k,n}\|} &= \mathbb{E}\left[\mathbf{h}_{lk,n}\mathbf{y}_{lk,n}^{\text{tr,H}}\right]\mathbb{E}^{-1}\left[\mathbf{y}_{lk,n}^{\text{tr}}\mathbf{y}_{lk,n}^{\text{tr,H}}\right]\mathbf{y}_{lk,n}^{\text{tr}} \\ &= \frac{\sqrt{P_{lk,n}}\beta_{llk}}{\sum_j P_{jk,n}\beta_{ljk} + \frac{\sigma^2}{K}}\mathbf{y}_{lk,n}^{\text{tr}}, \end{aligned} \quad (7)$$

and it is distributed as $\hat{\mathbf{h}}_{lk,n} \sim \mathcal{CN}\left(\mathbf{0}, \sigma_{\hat{\mathbf{h}}_k}^2 \mathbf{I}_M\right)$ with variance given as

$$\sigma_{\hat{\mathbf{h}}_k}^2 = \frac{P_{lk,n}\beta_{llk}^2}{\sum_j P_{jk,n}\beta_{ljk} + \frac{\sigma^2}{K}}. \quad (8)$$

Based on the orthogonality principle of the MMSE estimation, the uncorrelated estimation error vector at time instance n is $\hat{\mathbf{e}}_{lk,n} = \mathbf{h}_{lk,n} - \hat{\mathbf{h}}_{lk,n}$, being distributed as $\hat{\mathbf{e}}_{lk,n} \sim \mathcal{CN}\left(\mathbf{0}, \sigma_{\hat{\mathbf{e}}_{lk}}^2 \mathbf{I}_M\right)$ with

$$\sigma_{\hat{\mathbf{e}}_k}^2 = \beta_{llk} \left(1 - \frac{P_{lk,n}\beta_{llk}}{\sum_j P_{jk,n}\beta_{ljk} + \frac{\sigma^2}{K}}\right). \quad (9)$$

B. Channel Aging

The relative movement of the k th associated user with a comparison to the tagged BS antennas results in the variation of the channel. Hence, this source of imperfection contributes further to the need for estimation of the channel. Mathematically, we are able to relate the current sample of the channel with its past samples by means of an autoregressive model of order s [37]. Herein, for the sake of computational complexity and tractability, we choose an autoregressive model of order 1, which is a common approach in the literature [38]. Thus, the current channel at the tagged BS is modeled as

$$\mathbf{h}_{lk,n} = \delta\mathbf{h}_{lk,n-1} + \mathbf{e}_{lk,n}, \quad (10)$$

where $\mathbf{h}_{lk,n-1}$ is the channel in the previous symbol duration, and $\mathbf{e}_{lk,n} \in \mathbb{C}^N$, modelled as a stationary Gaussian random process with i.i.d. entries and distribution $\mathcal{CN}\left(\mathbf{0}, (1 - \delta^2)\mathbf{I}_M\right)$, is the uncorrelated channel error because of the channel variation [38]. Regarding δ , it is related to the second-order statistics. Specifically, an appropriate measure for modeling the variation of the channel is its second-order statistics, which can be described by means of the autocorrelation function of the channel. For this role, a widely accepted model is the Jakes model due to its generality and simplicity [37]. The Jakes model describes a propagation medium with two-dimensional isotropic scattering and a monopole antenna at the receiver [39]. In such case, the normalized discrete-time autocorrelation function of the fading channel is expressed by

$$r(s) = J_0(2\pi f_D T_s |s|), \quad (11)$$

where f_D and T_s are the maximum Doppler shift and the channel sampling period. Especially, the maximum Doppler shift f_D can be expressed by means of the relative velocity of

the scheduled user v , i.e., $f_D = \frac{vf_c}{c}$, where $c = 3 \times 10^8$ m/s is the speed of light and f_c is the carrier frequency. Also, s denotes the delay. Increasing the argument of the Bessel function results in a decrease of the magnitude to zero but with some ripples in the meanwhile. We set $\delta = r[1]$, i.e., we consider a single symbol delay. To this end, we assume that the BS has perfect knowledge of δ .

Remarkably, following the procedure in [21], we are able to write both pilot contamination and time-variation of the channel as a combination. More concretely, the channel at time slot n can be written as

$$\begin{aligned} \mathbf{h}_{lk,n} &= \delta \mathbf{h}_{lk,n-1} + \mathbf{e}_{lk,n} \\ &= \delta \hat{\mathbf{h}}_{lk,n-1} + \tilde{\mathbf{e}}_{lk,n}, \end{aligned} \quad (12)$$

where $\hat{\mathbf{h}}_{lk,n-1}$ and $\tilde{\mathbf{e}}_{lk,n} = \delta \tilde{\mathbf{h}}_{lk,n-1} + \mathbf{e}_{lk,n} \sim \mathcal{CN}(\mathbf{0}, \sigma_{\tilde{\mathbf{e}}_k}^2 \mathbf{I}_M)$ with $\sigma_{\tilde{\mathbf{e}}_k}^2 = \left(1 - \delta^2 \sigma_{\mathbf{h}_k}^2\right)$ are mutually independent. Hence, the estimated channel of the k scheduled user at time n is provided by $\hat{\mathbf{h}}_{lk,n} = \delta \hat{\mathbf{h}}_{lk,n-1}$. Note that in the special case, where $\delta = 1$, we obtain a static environment with no user mobility.

V. UPLINK TRANSMISSION

In general, the physical representation of a link defines the probability distribution function (PDF) of this link. Specifically, we face different distributions depending if we model the desired or the interference part of the received signal. Another example, affecting the PDF, concerns the choices between multi-antenna and single-antenna BS architecture, and between single or multi-user transmission. Notable, herein, we employ the general setting of a large number of antennas deployed by the tagged BS serving multiple users simultaneously. The first step towards the statistical characterization of the powers of the received signal's parts is to model the uplink transmission.

Thus, accounting for a quasi-static block fading model with frequency-flat fading channels varying for symbol to symbol, the received signal from the associated scheduled user at $x_{lk,n}$ to the tagged BS during the n th time-slot after applying a general decoder $\mathbf{q}_{lk,n}$ can be expressed as

$$\begin{aligned} y_{lk,n} &= \mathbf{q}_{lk,n}^H \mathbf{h}_{lk,n} s_{lk,n} + \sum_{(j,k') \neq (l,k)} \mathbf{q}_{lk,n}^H \mathbf{h}_{lj'k',n} s_{j'k',n} \\ &\quad + \mathbf{q}_{lk,n}^H \mathbf{n}_{lk}, \end{aligned} \quad (13)$$

where $s_{lk,n}$ is the uplink data symbol of the k th scheduled user with $\mathbb{E}[|s_{lk,n}|^2] = P_{lk,n}$. The channel vector $\mathbf{h}_{lk,n} \in \mathbb{C}^{M \times 1}$ denotes the desired channel vector between the tagged BS and the associated k th scheduled user located at $r_{lk,n} \in \mathbb{R}^2$ at time-instant n . Similarly, $\mathbf{h}_{lj'k',n} \in \mathbb{C}^{M \times 1}$ expresses the interference channel vector from the other users found at $r_{lj'k',n} \in \mathbb{R}^2$ far from the typical BS at time-instant n . Also, $\mathbf{n}_{lk} \in \mathbb{C}^{M \times 1} \sim \mathcal{CN}(0, \sigma^2 \mathbf{I}_M)$ is the Gaussian thermal noise vector in the uplink data transmission.

Taking into account for the realistic case, where imperfect CSI due to pilot contamination and time-variation of the channel (see (12)), is considered, the received signal by the tagged BS

can be written as

$$\begin{aligned} y_{lk,n} &= \delta \mathbf{q}_{lk,n}^H \hat{\mathbf{h}}_{lk,n-1} s_{lk,n} + \mathbf{q}_{lk,n}^H \tilde{\mathbf{e}}_{lk,n} s_{lk,n} \\ &\quad + \sum_{(j,k') \neq (l,k)} \mathbf{q}_{lk,n}^H \mathbf{h}_{lj'k',n} s_{j'k',n} + \mathbf{q}_{lk,n}^H \mathbf{n}_{lk}, \end{aligned} \quad (14)$$

where we have replaced the current channel by means of (12) with its estimated version⁴. In (14), the first term expresses the desired signal received by the tagged BS. The second term describes the estimation error effect. Furthermore, the third term represents the other users interference, while the last term denotes the post-processed noise.

The achievable uplink SINR from the k scheduled user to the tagged BS, denoted by SINR_k , is shown in (16), where we have treated the unknown terms at the tagged BS as uncorrelated additive noise. The encoding of the message takes place over many realizations of certain sources of randomness in the model. Specifically, the expectation operators are taken over the channel estimation error, the small-scale fading in the interference links, and the thermal noise. Notably, the resultant SINR, provided by (17), is a random variable because of the randomness accompanying the large-scale path-losses. Thus, we result in the uplink of a multi-user large MIMO HetNet, where the SINR expression will be investigated by using tools from stochastic geometry and large random matrix theory.

Herein, we present the derivation of an approximation of the SINR, when maximal-ratio combining (MRC) receivers are employed. As the number of BS antennas grows large, the approximation becomes tighter according to the theory of DEs [40]. Starting from the DE SINR distribution, we derive below the outage probability and the average delivery rate.

Let the tagged BS apply the decoder $\mathbf{q}_{lk,n}$ to the received signal, being a scaled version of the channel estimate $\delta \hat{\mathbf{h}}_{lk,n-1}$. The mathematical expression of the MRC decoder is

$$\mathbf{q}_{lk,n} = \frac{\sum_j P_{jk,n} \beta_{ljk} + \frac{\sigma^2}{K}}{\sqrt{P_{lk,n} \beta_{llk}}} \delta \hat{\mathbf{h}}_{lk,n-1}, \quad (15)$$

where the scaling of the decoder is applied for the sake of simplicity, but it will not affect the SINR distribution. Also, note that the decoder depends on the estimated channel, obtained during the training phase. Hereafter, we omit the time index from the expressions, while note that the DE expressions are calculated over the channel distributions, i.e., they are conditioned on the BSs positions.

Let $\overline{\text{SINR}}_{lk}$ be the deterministic SINR, obtained such that $\text{SINR}_{lk} - \overline{\text{SINR}}_{lk} \xrightarrow[M \rightarrow \infty]{\text{a.s.}} 0^5$.

Proposition 1. *The uplink achievable DE SINR with MRC decoding under the presence of pilot contamination and channel aging is given by (17) with $W_{jk} = \beta_{jjk}^{-\epsilon} \beta_{ljk}$.*

Proof: See Appendix B. ■

Notably, the terms W_{jk} correspond to the interference terms from other cells.

⁴Note that the replacement concerns only the current desired channel because the interference part is not of direct interest and can be seen as additive noise.

⁵The notation $\xrightarrow[M \rightarrow \infty]{\text{a.s.}}$ denotes almost sure convergence, while the definition of the term "deterministic equivalent" is given by [40, Def. 6.1].

$$\text{SINR}_{llk,n} = \frac{\delta^2 P_{lk,n} |\mathbf{q}_{llk,n}^H \hat{\mathbf{h}}_{llk,n}|^2}{P_{lk,n} \mathbb{E} \left[|\mathbf{q}_{llk,n}^H \tilde{\mathbf{e}}_{llk,n}|^2 \right] + \sum_{(j,k') \neq (l,k)} P_{jk',n} \mathbb{E} \left[|\mathbf{q}_{llk,n}^H \mathbf{h}_{ljk',n}|^2 \right] + \|\mathbf{q}_{llk,n}\|^2 \sigma^2}. \quad (16)$$

$$\overline{\text{SINR}}_{llk} = \frac{P_t \delta^2 \beta_{llk}^{2(1-\epsilon)}}{P_t \beta_{llk}^{(1-\epsilon)} \left(\sum_{j \neq l} W_{jk} + \frac{\sigma^2}{P_t K} \right) - P_t \delta^2 \beta_{llk}^{2(1-\epsilon)} + \sum_{(j,k') \neq (l,k)} P_t \left(\left(\frac{\sigma^2}{P_t K} + W_{jk} \right) W_{jk'} + W_{jk}^2 \right) + \sum_j W_{jk} + \frac{\sigma^2}{P_t K}}. \quad (17)$$

VI. PERFORMANCE ANALYSIS

The proposed realistic system depends on several practical factors, e.g., the pilot contamination and the channel aging as well as caching parameters such as the storage size. As already known, the quality-of-experience constraints specify that the uplink rate of the scheduled user should be equal or higher than the file bitrate T so that the user does not observe any interruption during its experience. In addition, another impediment is not been taken into account in most cases. It concerns the rate of backhaul, which becomes quite important when the cache misses.

The quantification and the assessment of the system necessitate the definition of certain metrics, namely the outage probability and the average delivery rate.

A. Outage probability

In this section, we present the uplink outage probability of the associated user in a large antenna MU HetNet with imperfect CSIT due to pilot contamination and channel aging, while caching is employed. The technical derivation is given in Appendix C. As a performance metric, the outage probability is given as the complementary of the success (coverage) probability, expressing the joint probabilities of the uplink rate exceeding the file bitrate T and the received file missing from the local cache. It is worthwhile to mention that only the BSs connect with the backhaul. In other words, the associated user is not able to upload its content if the BS already has it. Actually, there is no reason to do it, since the content is stored at the BS and it is BS's task to upload it through its wired link. Hence, we have

$$\mathbb{P}_{\text{out}} \triangleq 1 - \tilde{\mathbb{P}} \left(\overline{\text{SINR}}_{llk} > \tilde{T}, f \notin \Delta_l \right), \quad (18)$$

where $\tilde{T} = e^T - 1$, f is the received file by the typical BS, and Δ_l is the local cache of the served BS at the l th cell. Differently to [27], this definition follows another line of reasoning. In particular, if the requested file is not in the cache of the served BS, and if the uplink rate is higher than the file bitrate T , then, the user uploads its content and does not observe any interruption during its communication. Hence, we expect the outage probability to be close to zero. Formally, the outage probability is given by the following theorem.

Theorem 1. *The approximated uplink outage probability \mathbb{P}_{out} in a large MU-MIMO HetNet with caching attributes, accounting for imperfect CSIT due to pilot contamination and*

channel aging, is given by

$$\mathbb{P}_{\text{out}} \approx 1 - \left(1 - \int_0^{\frac{\tilde{T}}{T}} f_{\text{pop}}(f, \mu, \gamma) df \right) \tilde{\mathbb{P}} \left(\overline{\text{SINR}} > \tilde{T} \right) \quad (19)$$

with $f_{\text{pop}}(f, \mu, \gamma)$ given by (5) and the coverage probability $\tilde{\mathbb{P}} \left(\overline{\text{SINR}} > \tilde{T} \right)$ given by (21) with $t = \pi \lambda_B x^2$. The variable N represents the number of terms used in the calculation, $\eta = N(N!)^{-\frac{1}{N}}$, while $D_{\sigma^2} = \frac{\sigma^2}{K P_t C^{1-\epsilon} (\lambda_b \pi)^{\frac{\alpha(1-\epsilon)}{2}}}$,

$$D_1 = \frac{2\Gamma^\alpha \left(\frac{\epsilon}{2} + 1 \right)}{\alpha - 2} + D_{\sigma^2}, \quad D_2 = (K - 1) D_1, \quad D_3 = (K - 1) \frac{\Gamma^\alpha(\epsilon + 1)}{\alpha - 1}, \quad C_4 = D_1 + 2C^{1-\epsilon} (\pi \lambda_b)^{\frac{1-\alpha\epsilon}{2}} \Gamma^\alpha \left(\frac{\epsilon}{2} + 1 \right),$$

$$D_5(t) = \int_0^\infty \frac{e^{-u} du}{1 + D_2 t^{2\alpha(1-\epsilon)} u^{-\frac{\alpha}{2}(1-\epsilon)}}, \quad \text{and } D_6 = -\frac{1}{P_t}.$$

Proof: See Appendix C. ■

B. Average Delivery Rate

This section presents the derivation of the average delivery rate, defined as

$$R \triangleq \begin{cases} T, & \text{if } \psi \ln(1 + \overline{\text{SINR}}_{llk}) > T \text{ and } f \notin \Delta_{b_0} \\ C(\lambda_B), & \text{if } \psi \ln(1 + \overline{\text{SINR}}_{llk}) > T \text{ and } f \in \Delta_{b_0} \\ 0, & \text{otherwise} \end{cases} \quad (20)$$

where $C(\lambda_B) = \frac{C_1}{\lambda_B} + C_2$ with C_1 and C_2 being arbitrary coefficients under the constraint the ceiling of the delivery rate is $C(\lambda_B)$ with $C(\lambda_B) < T$. $C(\lambda_B)$ denotes the backhaul capacity being available to the intermediate nodes. Also, $\psi = \frac{T_c - \tau}{T_c}$ is the fraction of time expressing the training overhead which occurs during the estimation channel. However, (20) refers to the opposite direction (uplink) and then it includes an interesting and insightful explanation, especially, because the conditions are different. In the case that the uplink rate is higher than the target file rate (bitrate) T and the file is not found in the BSs, the user uploads at full rate T . On the contrary, if the rate is greater than T and the file is already to the local cache, the associated user does not upload its content, but the tagged BS does. The latter constraint relies on the assumption that a high-speed backhaul is not cost-efficient in dense networks.

Theorem 2. *The approximated uplink average delivery rate of the typical BS in a large MU-MIMO HetNet with caching attributes, accounting for imperfect CSIT due to pilot contamination and channel aging, is given by*

$$R = \left(T - (C(\lambda_B) - T) \int_0^{\frac{\tilde{T}}{T}} f_{\text{pop}}(f, \mu, \gamma) df \right) \tilde{\mathbb{P}} \left(\overline{\text{SINR}} > T \right) \quad (22)$$

$$\tilde{\mathbb{P}}(\overline{\text{SINR}} > \tilde{T}) \approx \sum_{n=1}^N \binom{N}{n} (-1)^{n+1} \int_0^\infty e^{-t-n\eta\tilde{T}\delta^{-2}} \left(D_1 t^{\alpha(1-\epsilon)} + D_3 t^{2\alpha(1-\epsilon)} + D_4 t^{2\alpha(1-\epsilon)} + D_6 \right) \left(1 - D_2 t^{2\alpha(1-\epsilon)} D_5(t) \right)^{K-1} dt. \quad (21)$$

where $f_{\text{pop}}(f, \mu, \gamma)$ is given by (5), and the coverage probability is given by (21) with $t = \pi\lambda_B x^2$ if we substitute \tilde{T} with T .

Proof: See Appendix D. ■

VII. NUMERICAL RESULTS

In this section, we illustrate the behavior of the analytical expressions concerning the outage probability \mathbb{P}_{out} and the average delivery rate R , which are provided by means of (19) and (22)⁶. In fact, we investigate the impact of various design parameters such as the BS density λ_B , the storage size of BSs S in nats, and the target bit-rate T in nats/sec/Hz. Also, the analytical expressions are verified by Monte Carlo simulations. The simulated curves were obtained by averaging the corresponding expressions over 10^3 random instances. Actually, the simulated results of the outage probability \mathbb{P}_{out} and the average delivery user rate R are depicted along with the proposed analytical expressions. Specifically, the bullets correspond to the simulation results, while the “solid” lines represent the proposed analytical results by varying their parameters. The discrimination between “solid” and “dot” lines, where applicable, designates the results with “caching” and “no caching”, respectively. The “no caching” scenario is obtained by assuming that the content popularity distribution coincides with the Dirac delta function, i.e., $\gamma_j \rightarrow \infty$.

The simulations are conducted by following a specific procedure. Specifically, we choose a sufficiently large area of $5 \text{ km} \times 5 \text{ km}$, where the locations of the BSs are simulated as a realization of a PPP with given density $\lambda_B = 0.2 \text{ m}^{-2}$. Next, the users’ PPP density is considered to be $\lambda_K = 60\lambda_B$ ⁷. The association relies on the minimum path-loss (distance-based) rule, while K users from each cell are randomly scheduled. Hence, we select the strongest user to the tagged BS, found at the origin, as the associated scheduled user at $x_{lk,n}$. It is worthwhile to mention that the users could employ other schemes to upload their contents to the BSs. For example, they could select the serving BSs rather than the closest BSs. The relevant comparison with other approaches regarding the selection of the appropriate BSs is interesting and is left for future work. Furthermore, the setup includes BSs of $M = 25$ number of antennas, while we pick $K = 5$ users per BS. The system under study, embodying a such number of BS antennas, is considered to describe a massive MIMO model, since the simulations coincide with the DEs. In other words, the DEs are tight approximations even for this number of antennas. Hence,

⁶Remarkably, there is no known result in the literature studying caching in the uplink of a HetNet employing a large number of antennas (massive MIMO). In addition, there is no known reference investigating channel aging in the case of cached-enabled BSs.

⁷Although the analytical expressions rely on the assumption of an infinite plane, the simulation takes place over a finite window.

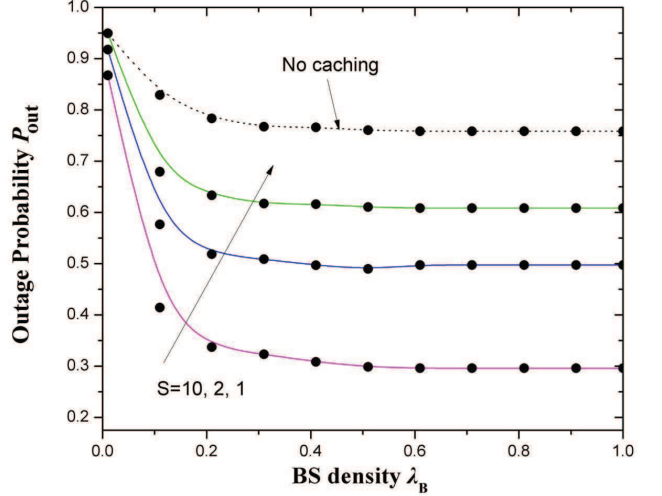


Figure 2: Outage probability versus the BS density λ_B for varying storage size S . Solid lines and bullets correspond to the theoretical and simulated results, respectively, while the dotted line refers to the “No caching” scenario.

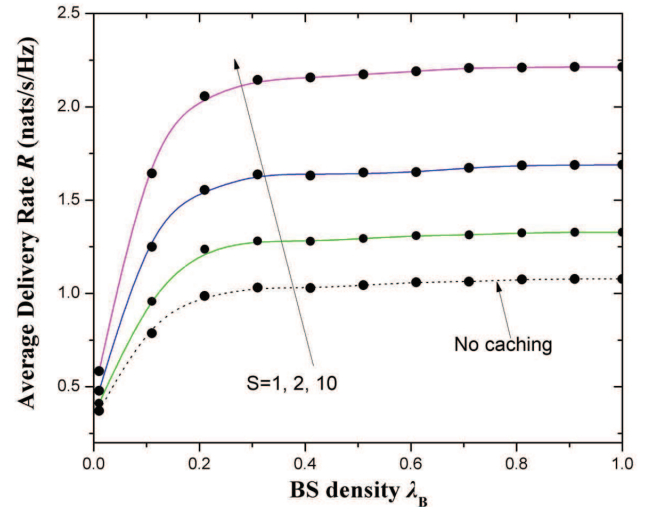


Figure 3: Average delivery rate versus the BS density λ_B for varying storage size S . Solid lines and bullets correspond to the theoretical and simulated results, respectively, while the dotted line refers to the “No caching” scenario.

such a number of BS antennas can represent a massive MIMO model. However, this is not a new observation. According to the literature, similar observations have been made in the literature even for an 8×8 system [40]–[43]. The average uplink transmit power for both training and transmission phases is $P_{l_{k,n}} = 2 \text{ dBW}$, and the bandwidth allocated for each user is

20 MHz. Also, regarding the rest parameters, we set $L = 1$ nats, $\alpha = 3$, $P_t = 1$, $\epsilon = 0.7$, $C_1 = 0.005$, $C_2 = 0$, $\delta = 0.8$, $\mu = 0$, $\gamma = 0.5$ unless otherwise stated. Due to limited space, in this work, we do not focus on channel aging, studied in other works such as [13], but the cynosure is the impact of caching in the uplink.

A. Impact of BS Density

In Fig. 2, we illustrate the behavior of the outage probability \mathbb{P}_{out} with respect to the BS density λ_B for different values of the storage size S . We observe a decrement of the outage probability as the BS density increases. In other words, a denser HetNet provides better coverage. At the same time, an increase of the storage size of the intermediate nodes brings a decrease in an outage, since the users do not have to upload their content to the core network because the BSs have plenty of space to save the receiver information.

Regarding the average delivery rate R , it increases with the BS density λ_B as can be seen in Fig. 3. However, it saturates soon due to the increasing intra-cell interference. Moreover, higher storage size contributes to the increase of the rate because of the traffic load towards the backhaul is alleviated.

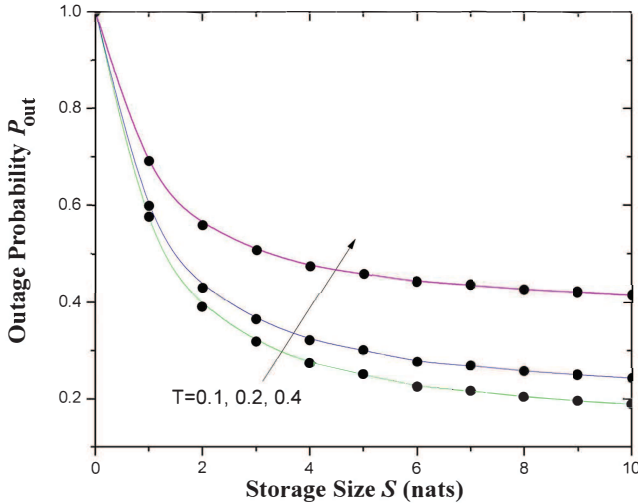


Figure 4: Outage probability versus the storage size S for varying target file bit-rate T . Solid lines and bullets correspond to the theoretical and simulated results, respectively.

B. Impact of Storage Size

Fig. 4 shows the relationship of the outage probability \mathbb{P}_{out} with the storage size S of the BSs. Notably, the storage capability of networks with caching is one of the most crucial parameters during the design. Obviously, the outage probability increases with the storage size, but decreases with the target bit-rate. In other words, the larger the target bit rate is set, the larger the outage probability will be.

In the same direction, in Fig. 5, the average delivery rate becomes higher with increasing storage size, but after a value of S , further increment is not beneficial, since all users content

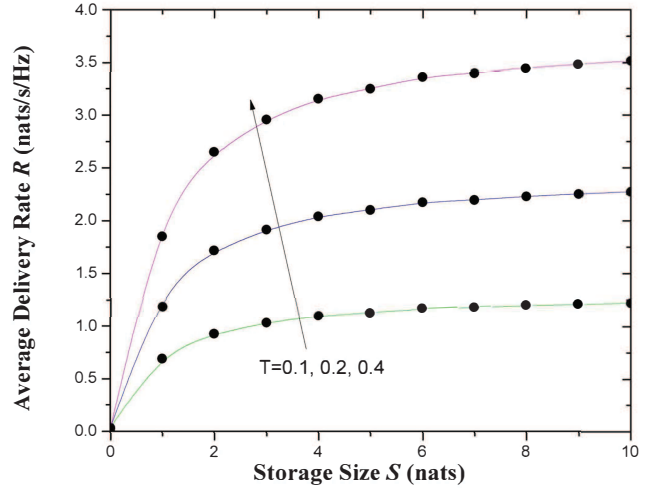


Figure 5: Average delivery rate versus the storage size S for varying target file bit-rate T . Solid lines and bullets correspond to the theoretical and simulated results, respectively.

will be already available to the corresponding BSs. Especially, less target rate allows better coverage.

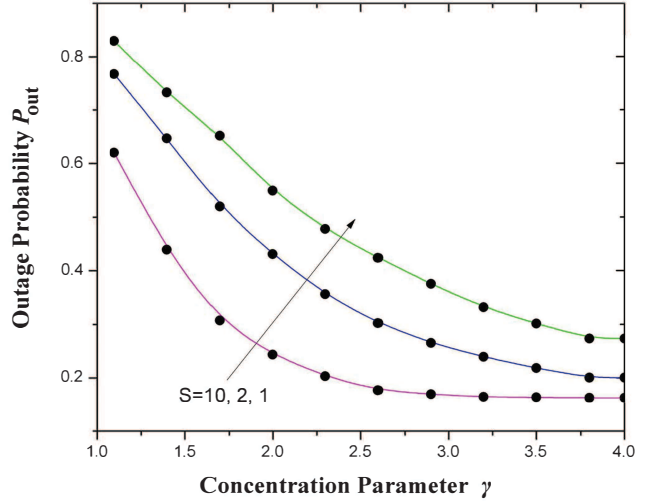


Figure 6: Outage probability versus the concentration parameter γ for varying storage size S . Solid lines and bullets correspond to the theoretical and simulated results, respectively.

C. Impact of Concentration Parameter

The variation of the concentration parameter, described by γ , is depicted in Fig. 6. Small γ means that a high quota of files is already at the intermediate nodes, i.e., many files are popular. Hence, the users do not need to upload their files and significant outage is observed. Moreover, higher storage size allows more files to be uploaded in the BSs. As a result, it is likely that the contents of the users are already at the BSs and the users are inactive since they do not need to upload their contents.

The dependence of the average delivery rate R with the concentration parameter is provided by Fig. 7. Specifically, a

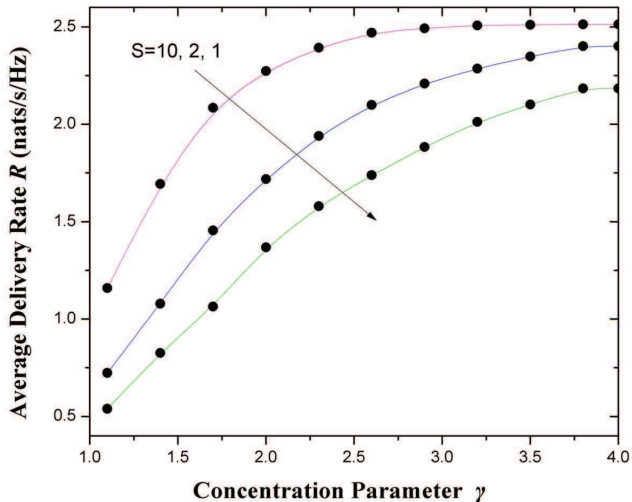


Figure 7: Average delivery rate versus the concentration parameter γ for varying storage size S . Solid lines and bullets correspond to the theoretical and simulated results, respectively.

high concentration parameter means that many files will be uploaded, and thus, the average rate increases.

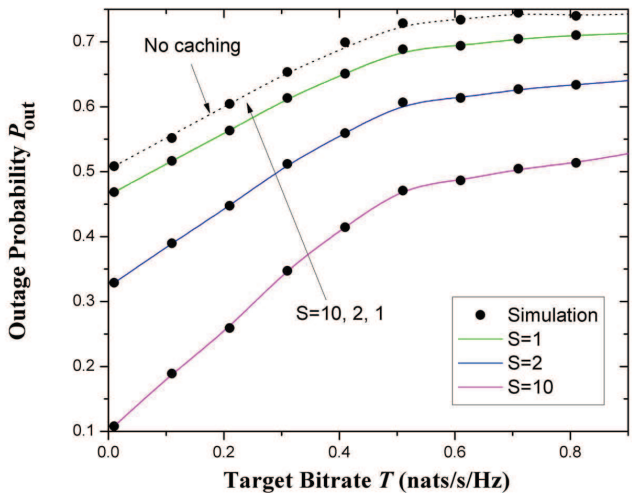


Figure 8: Outage probability versus the target bitrate T for varying storage size S . Solid lines and bullets correspond to the theoretical and simulated results, respectively, while the dotted line refers to the “No caching” scenario.

D. Impact of Target Bit-Rate

The target bit-rate T is another critical parameter that should be taken into account during the formation and study of the current architecture. In particular, Fig. 8 demonstrates the lines of the outage probability \mathbb{P}_{out} versus the target T for $S = 1, 2$, and 10 . Increasing the target rate, the outage probability increases, since less users are served. In addition, the performance is improved with increasing storage size because more content can be saved to the intermediate nodes without the need to upload it at the backhaul.

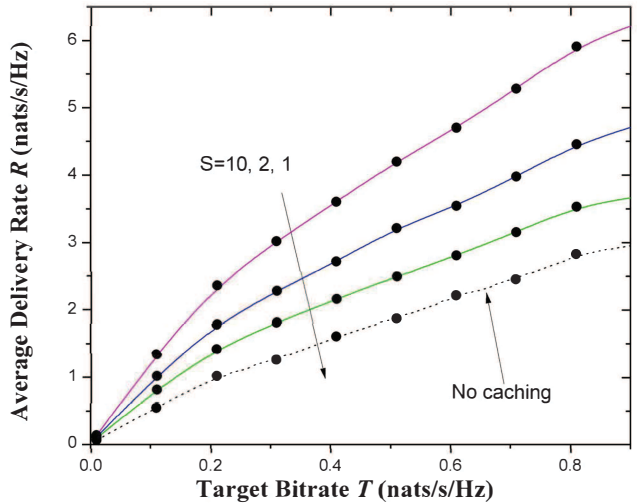


Figure 9: Average delivery rate versus the target bitrate T for varying storage size S . Solid lines and bullets correspond to the theoretical and simulated results, respectively, while the dotted line refers to the “No caching” scenario.

In a parallel avenue, Fig. 9 shows the increase in the performance with bigger storage capacities at the BSs, while it is apparent that a higher target rate results in a higher average delivery rate.

VIII. CONCLUSION

In this paper, we introduced the concept of caching in the uplink of a system with stochastically distributed massive MIMO BSs, where users upload their contents to servers through the BS by means of finite-rate backhaul links. In addition to significantly generalizing the state of the art cache-enabled PPP models to the uplink scenario, we enriched the uplink of the HetNet with the massive MIMO concept. Remarkably, it is the first work, where the caching nodes have a large number of antennas. Moreover, our approach considered imperfect CSI due to pilot contamination and channel aging. After deriving the DE of the SINR, we provided the outage probability and the average delivery rate. Our main purpose was to focus on fundamental parameters, being relevant to the caching design. Such parameters are the storage size of the serving BS and their target file rate. In particular, we demonstrated that by increasing the storage size, the performance of the system is improved, since the outage probability decreases and the average delivery rate increases. Furthermore, by increasing the target file bitrate, the majority of the users is not served. Hence, the outage probability increases. Overall, it was shown that the introduction of the notion of caching in the uplink enhances the system performance.

APPENDIX A USEFUL LEMMAS

Lemma 1 (Alzer’s inequality [44]). *Assuming that h is a normalized gamma random variable with parameter N and*

a constant $\gamma >$, then the probability $\mathbb{P}(h < \gamma)$ can be tightly upper bounded by

$$\mathbb{P}(h < \gamma) < (1 - e^{-\alpha\gamma})^N, \quad (23)$$

where $\alpha = N(N!)^{-\frac{1}{N}}$.

Lemma 2 ([43, Lem. B.26]). *Let $\mathbf{A} \in \mathbb{C}^{N \times N}$ with uniformly bounded spectral norm (with respect to N). Consider \mathbf{x} and \mathbf{y} , where $\mathbf{x}, \mathbf{y} \in \mathbb{C}^N$, $\mathbf{x} \sim \mathcal{CN}(\mathbf{0}, \Phi_x)$ and $\mathbf{y} \sim \mathcal{CN}(\mathbf{0}, \Phi_y)$, are mutually independent and independent of \mathbf{A} . Then, we have*

$$\frac{1}{N} \mathbf{x}^H \mathbf{A} \mathbf{x} - \frac{1}{N} \text{tr} \mathbf{A} \Phi_x \xrightarrow[N \rightarrow \infty]{a.s.} 0 \quad (24)$$

$$\frac{1}{N} \mathbf{x}^H \mathbf{A} \mathbf{y} \xrightarrow[N \rightarrow \infty]{a.s.} 0 \quad (25)$$

$$\mathbb{E} \left[\left| \left(\frac{1}{N} \mathbf{x}^H \mathbf{A} \mathbf{x} \right)^2 - \left(\frac{1}{N} \text{tr} \mathbf{A} \Phi_x \right)^2 \right| \right] \xrightarrow[N \rightarrow \infty]{a.s.} 0 \quad (26)$$

$$\frac{1}{N^2} |\mathbf{x}^H \mathbf{A} \mathbf{y}|^2 - \frac{1}{N^2} \text{tr} \mathbf{A} \Phi_x \text{tr} \mathbf{A}^H \Phi_y \xrightarrow[N \rightarrow \infty]{a.s.} 0. \quad (27)$$

APPENDIX B PROOF OF PROPOSITION 1

First, we divide both the numerator and the denominator of (16) by $\frac{1}{M^2}$. Then, we start with the numerator of the SINR. We insert in (16) the expression of the MRC decoder given by (15). We have⁸

$$\begin{aligned} \frac{1}{M^2} P_{lk} |\mathbf{q}_{llk}^H \hat{\mathbf{h}}_{llk}|^2 &\stackrel{(a)}{=} \frac{(P_{lk} \beta_{llk})^2}{\delta^2 \left(\sum_j P_{jk} \beta_{ljk} + \frac{\sigma^2}{K} \right)^2} \frac{1}{M^2} |\mathbf{q}_{llk}|^4 \\ &\stackrel{(b)}{\asymp} \frac{1}{M} P_t^2 \delta^4 \beta_{llk}^{2(1-\epsilon)}, \end{aligned} \quad (28)$$

where (a) follows after substituting the estimated channel given in (15), while (b) is obtained by means of Lemma 2, since the covariance of \mathbf{q}_{llk} is $\delta^2 \left(\sum_j P_{jk} \beta_{ljk} + \frac{\sigma^2}{K} \right) \mathbf{I}_M$. We continue with the first term in the denominator of the SINR including the estimation error. We have

$$\begin{aligned} \frac{1}{M^2} P_{lk} |\mathbf{q}_{llk}^H \tilde{\mathbf{e}}_{llk}|^2 &\asymp \frac{1}{M^2} P_{lk} \delta^2 M \beta_{llk} \left(\sum_j P_{jk} \beta_{ljk} + \frac{\sigma^2}{K} \right) \\ &\times \left(1 - \delta^2 \frac{P_{lk} \beta_{llk}}{\sum_j P_{jk} \beta_{ljk} + \frac{\sigma^2}{K}} \right) \\ &= P_t^2 \delta^2 \beta_{llk}^{(1-\epsilon)} \frac{1}{M} \left(\sum_{j \neq l} \beta_{jjk}^{-\epsilon} \beta_{ljk} + (1 - \delta^2) \beta_{llk}^{(1-\epsilon)} + \frac{\sigma^2}{P_t K} \right) \end{aligned}$$

because the covariance of the estimation error is $\beta_{llk} \left(1 - \delta^2 \frac{P_{lk} \beta_{llk}}{\sum_j P_{jk} \beta_{ljk} + \frac{\sigma^2}{K}} \right) \mathbf{I}_M$. The next step is to

⁸Let a_n and b_n two infinite sequences. $a_n \asymp b_n$ denotes the equivalence relation $a_n - b_n \xrightarrow[N \rightarrow \infty]{a.s.} 0$.

derive the second term in the denominator, which is written as

$$\begin{aligned} &\frac{1}{M^2} \sum_{(j,k') \neq (l,k)} |\sqrt{P_{jk'}} \mathbf{q}_{llk}^H \mathbf{h}_{lj'k'}|^2 \\ &= \frac{1}{M^2} \sum_{(j,k') \neq (l,k)} \delta^2 |\sqrt{P_{jk'}} \mathbf{n}_{llk}^{\text{tr}} \mathbf{h}_{lj'k'} + \sum_{j' \neq l} \sqrt{P_{j'k}} \mathbf{h}_{lj'k} \mathbf{h}_{lj'k'}|^2 \\ &\asymp \sum_{(j,k') \neq (l,k)} \delta^2 \left(\frac{\sigma^2}{MK} P_t \beta_{jjk'}^{-\epsilon} \beta_{ljk'} + \frac{1}{M^2} \sum_{j' \neq l} P_{j'k} P_{j'k'} |\mathbf{h}_{lj'k} \mathbf{h}_{lj'k'}|^2 \right). \end{aligned} \quad (29)$$

When $k' = k$ and $j' = j$, the second term of the previous expression simplifies to

$$\begin{aligned} \frac{1}{M^2} P_{jk'} P_{j'k} |\mathbf{h}_{lj'k} \mathbf{h}_{lj'k'}|^2 &= P_{jk}^2 \frac{1}{M^2} \|\mathbf{h}_{lj'k}\|^4 \\ &\asymp \frac{1}{M} P_t^2 \beta_{jjk}^{-2\epsilon} \beta_{ljk}^2, \end{aligned} \quad (30)$$

where we have used Lemma 2. In the case that $j' \neq j$ and $k' = k$ or $k' \neq k$, we have

$$\begin{aligned} \frac{1}{M^2} P_{jk'} P_{j'k} \mathbf{h}_{lj'k} \mathbf{h}_{lj'k'} &\asymp \frac{1}{M} P_{j'k} P_{jk'} \beta_{lj'k} \beta_{ljk'} \\ &= \frac{1}{M} P_t^2 \left(\beta_{j'j'k} \beta_{jjk'} \right)^{-\epsilon} \beta_{lj'k} \beta_{ljk'}. \end{aligned} \quad (31)$$

Thus, (29) becomes

$$\begin{aligned} &\sum_{(j,k') \neq (l,k)} \delta^2 \left(\frac{\sigma^2}{MK} M P_t \beta_{jjk'}^{-\epsilon} \beta_{ljk'} \sum_{j' \neq l} P_{j'k} P_{j'k'} \frac{1}{M^2} |\mathbf{h}_{lj'k} \mathbf{h}_{lj'k'}|^2 \right) \\ &\asymp \sum_{(j,k') \neq (l,k)} \delta^2 \frac{1}{M} \left(\frac{\sigma^2}{K} P_t \beta_{jjk'}^{-\epsilon} \beta_{ljk'} \right. \\ &\quad \left. + P_t^2 \beta_{jjk}^{-2\epsilon} \beta_{ljk}^2 + P_t^2 \left(\beta_{j'j'k} \beta_{jjk'} \right)^{-\epsilon} \beta_{lj'k} \beta_{ljk'} \right). \end{aligned} \quad (32)$$

Regarding the term that includes the thermal noise, after applying Lemma 2 we have

$$\frac{1}{M^2} \|\mathbf{q}_{llk}\|^2 \sigma^2 \asymp \frac{1}{M} P_t \delta^2 \left(\sum_j \beta_{jjk}^{-\epsilon} \beta_{ljk} + \frac{\sigma^2}{P_t K} \right). \quad (33)$$

APPENDIX C PROOF OF THEOREM 1

The proof starts by finding first the conditional coverage probability on x as

$$\begin{aligned} \tilde{\mathbb{P}}(\overline{\text{SINR}}_{llk} > \tilde{T}, f \notin \Delta_{b_0}) &= \mathbb{E}_x \left[\tilde{\mathbb{P}}(\overline{\text{SINR}}_{llk} > \tilde{T} | x) \right] \\ &\times \mathbb{E}_x [\mathbb{P}(f \notin \Delta_{b_0} | x)]. \end{aligned} \quad (34)$$

Hence, we focus on the derivation of $\tilde{\mathbb{P}}(\overline{\text{SINR}}_{llk} > \tilde{T} | x)$. Specifically, we propose an approximation for the out-of-cell interference, described by W_{jk} for all users in each cell, i.e., $k \in [1, K]$. This approximation will allow the decoupling of the correlated terms and will result in a tractable evaluation of

\mathbb{P}_c . Specifically, by approximating the out-of-cell interference with its mean, we have [45]

$$\begin{aligned}
\sum_{j \neq l} W_{jk} &= \sum_{j \neq l} \beta_{jjk}^{-\epsilon} \beta_{ljk} \\
&\approx \mathbb{E} \left[\sum_{j \neq l} \beta_{jjk}^{-\epsilon} \beta_{ljk} \right] \\
&= \mathbb{E} \left[\sum_{j \neq l} \mathbb{E} \left[\beta_{jjk}^{-\epsilon} \right] \beta_{ljk} \right] \quad (35) \\
&= \lambda_b C^{-\epsilon} (\pi \lambda_b)^{-\frac{\alpha \epsilon}{2}} \Gamma^\alpha \left(\frac{\epsilon}{2} + 1 \right) \mathbb{E} \left[\sum_{j \neq l} \beta_{ljk} \right] \\
&= \frac{2C^{1-\epsilon} (\pi \lambda_b)^{\frac{\alpha(1-\epsilon)}{2}}}{\alpha - 2} \Gamma^\alpha \left(\frac{\epsilon}{2} + 1 \right), \quad (36)
\end{aligned}$$

where in (35), we made the following substitution

$$\begin{aligned}
\mathbb{E} \left[\beta_{jjk}^{-\epsilon} \right] &= C^{-\epsilon} (\mathbb{E} [r_{ljk}^\epsilon])^\alpha \\
&= (\pi \lambda_b)^{-\frac{\alpha \epsilon}{2}} \Gamma^\alpha \left(\frac{\epsilon}{2} + 1 \right).
\end{aligned}$$

Moreover, (36) is obtained by means of the Campbell's theorem [30] and the exclusion ball model, described in Sec. II as

$$\begin{aligned}
\mathbb{E} \left[\sum_{j \neq l} \beta_{ljk} \right] &= 2\pi \lambda_b C \int_{R_e}^\infty x^{-\alpha} x dx \\
&= \frac{2C (\lambda_b)^{\frac{\alpha}{2}} \pi \lambda_b}{\alpha - 2}. \quad (37)
\end{aligned}$$

In a similar way, we can write

$$\sum_{j \neq l} W_{jk}^2 = \frac{2C^{2(1-\epsilon)} (\pi \lambda_b)^{\alpha(1-\epsilon)}}{\alpha - 1} \Gamma^\alpha (\epsilon + 1). \quad (38)$$

The approximations (36) and (38) allow the simplification of the SINR, conditioned on $r_{jjk} = x$, by its approximate

$$\begin{aligned}
\overline{\text{SINR}} &\approx \left(D_1 (\pi \lambda_B x^2)^{\alpha(1-\epsilon)} + D_2 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} \sum_{k' \neq k} r_{jjk'}^{\alpha(1-\epsilon)} \right. \\
&\quad \left. + D_3 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} + D_4 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} + D_6 \right)^{-1}, \quad (39)
\end{aligned}$$

where in (39) the variables $D_i \in 1, \dots, 4, 6$ are defined in Theorem 1. Conditioned on $r_{jjk} = x$, we obtain the approximate distribution of the SINR, given by (40), after substituting its expression from (39). In (41), we have considered the dummy variable \tilde{g} , having unit mean and shape parameter N , in order to approximate the constant number one. Actually, this approximation becomes tighter as N goes to infinity [44], since $\lim_{y \rightarrow \infty} \frac{y^y x^{y-1} e^{-yx}}{\Gamma(y)} = \delta(x-1)$ with $\delta(x)$ being Dirac's delta function. In (42), we have applied Alzer's inequality (see Lemma 1), where $\eta = N(N!)^{-\frac{1}{N}}$, while afterwards, in (43), we have used the Binomial theorem. Next, (44) is obtained by assuming that y is a Rayleigh random variable. In (45), we set $u = \lambda_B \pi y^2$, and we take into account the approximation $\exp(-x) \approx \frac{1}{1+x}$, in order to make the numerical integration

faster. Finally, given that x is a Rayleigh random variable, we obtain the uplink SINR distribution as

$$\begin{aligned}
\mathbb{E}_x \left[\tilde{\mathbb{P}} \left(\overline{\text{SINR}}_{llk} > \tilde{T} \right) | x \right] &= \int_0^\infty \tilde{\mathbb{P}} \left(\overline{\text{SINR}} > \tilde{T} | r_{jjk} = x \right) \\
&\quad \times e^{-\pi \lambda_B x^2} 2\pi \lambda_B x dx. \quad (46)
\end{aligned}$$

The derivation of the second term of (34) is straightforward. In particular, we assume that all the BSs cache the same amount of files (storage size), while the cache hit probability is independent of the distance $r_{jjk} = x$. Thus, we have

$$\mathbb{E}_x [\mathbb{P}(f \notin \Delta_{b_0} | x)] = 1 - \int_0^{\frac{\tilde{T}}{L}} f_{\text{pop}}(f, \mu, \gamma) df. \quad (47)$$

Inserting (46) and (47) into (34) and after some algebraic manipulations, we obtain the coverage probability. The proof is concluded by substituting the coverage probability $\tilde{\mathbb{P}}(\overline{\text{SINR}}_{llk} > \tilde{T}, f \notin \Delta_{b_0})$ in (18).

APPENDIX D PROOF OF THEOREM 2

The average delivery rate $\bar{R} = \mathbb{E}(R)$ is obtained by applying the expectation operator over both the fading distribution and the PPP. In particular, we have

$$\begin{aligned}
\bar{R} &= \mathbb{E}[R] \\
&= \psi \mathbb{E} \left[\tilde{\mathbb{P}}(\ln(1 + \overline{\text{SINR}}_{llk}) > T) \right. \\
&\quad \left. \times (T\mathbb{P}(f \notin \Delta_{b_0}) + C(\lambda_B) \mathbb{P}(f \in \Delta_{b_0})) \right] \quad (48)
\end{aligned}$$

$$\begin{aligned}
&= \psi \mathbb{E} \left[\tilde{\mathbb{P}}(\ln(1 + \overline{\text{SINR}}_{llk}) > T | x) \right] \\
&\quad \times (\mathbb{E}[T\mathbb{P}(f \notin \Delta_{b_0} | x)] + \mathbb{E}[C(\lambda_B) \mathbb{P}(f \in \Delta_{b_0} | x)]) \quad (49)
\end{aligned}$$

$$= \psi \mathcal{I}_1 (\mathcal{I}_2 + \mathcal{I}_3), \quad (50)$$

where (48) follows by applying the definition described by (20), while (49) is obtained because of the independence between the different events and the property of linearity of the expectation operator. The derivation of \bar{R} continues with the substitution of \mathcal{I}_1 , which is basically given by (21) after substituting \tilde{T} with T . Moreover, given that the cache hit probability does not depend on x , we have

$$\mathcal{I}_2 = T \int_0^{\frac{\tilde{T}}{L}} f_{\text{pop}}(f, \mu, \gamma) df, \quad (51)$$

while \mathcal{I}_3 is obtained by multiplying (47) with $C(\lambda_B)$. After appropriate substitutions in (50), the proof is concluded. \square

ACKNOWLEDGEMENT

The authors would like to express their gratitude to Dr. E. Baştuğ for his help and support in making this work possible.

REFERENCES

- [1] I. Cisco Visual Networking, "Global mobile data traffic forecast update 2014–2019. White paper c11-520862," Available on http://www.cisco.com/c/en/us/solutions/collateral/serviceprovider/visual-networking-index-vni/white_paper_c11-520862.html.
- [2] J. G. Andrews et al., "What will 5G be?" *IEEE J. Sel. Areas Commun.*, vol. 32, no. 6, pp. 1065–1082, June 2014.

$$\begin{aligned} & \mathbb{P}(\overline{\text{SINR}} > \tilde{T} | r_{j j k} = x) \\ & \approx \mathbb{P}\left(1 > \tilde{T} \delta^{-2} \left(D_1 (\pi \lambda_B x^2)^{\alpha(1-\epsilon)} + D_2 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} \sum_{k' \neq k} r_{j j k'}^{\alpha(1-\epsilon)} + D_3 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} + D_4 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} + D_6 \right) \right) \end{aligned} \quad (40)$$

$$\approx \mathbb{P}\left(\tilde{g} > \tilde{T} \delta^{-2} \left(D_1 (\pi \lambda_B x^2)^{\alpha(1-\epsilon)} + D_2 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} \sum_{k' \neq k} r_{j j k'}^{\alpha(1-\epsilon)} + D_3 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} + D_4 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} + D_6 \right) \right) \quad (41)$$

$$\approx 1 - \mathbb{E} \left[\left(1 - \exp \left(- \eta \tilde{T} \delta^{-2} \left(D_1 (\pi \lambda_B x^2)^{\alpha(1-\epsilon)} + D_2 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} \sum_{k' \neq k} r_{j j k'}^{\alpha(1-\epsilon)} + D_3 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} + D_4 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} + D_6 \right) \right) \right)^N \right] \quad (42)$$

$$= \sum_{n=1}^N \binom{N}{n} (-1)^{n+1} \mathbb{E} \left[\exp \left(- n \eta \tilde{T} \delta^{-2} \left(D_1 (\pi \lambda_B x^2)^{\alpha(1-\epsilon)} + D_2 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} \sum_{k' \neq k} y_{k'}^{\alpha(1-\epsilon)} + D_3 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} + D_4 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} + D_6 \right) \right) \right] \quad (43)$$

$$= \sum_{n=1}^N \binom{N}{n} (-1)^{n+1} \exp \left(- n \eta \tilde{T} \delta^{-2} \left(D_1 (\pi \lambda_B x^2)^{\alpha(1-\epsilon)} + D_3 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} + D_4 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} + D_6 \right) \right) \times \int_0^\infty \exp \left(D_2 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} y^{\alpha(1-\epsilon) - \lambda_B \pi y^2} \right)^{K-1} \lambda_B \pi y dy \quad (44)$$

$$= \sum_{n=1}^N \binom{N}{n} (-1)^{n+1} \exp \left(- n \eta \tilde{T} \delta^{-2} \left(D_1 (\pi \lambda_B x^2)^{\alpha(1-\epsilon)} + D_3 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} + D_4 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} + D_6 \right) \right) \times \left(1 - D_2 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} \int_0^\infty \frac{e^{-u} du}{1 + D_2 (\pi \lambda_B x^2)^{2\alpha(1-\epsilon)} u^{-\frac{\alpha}{2}(1-\epsilon)}} \right)^{K-1} \quad (45)$$

- [3] J. G. Andrews, H. Claussen, M. Dohler, S. Rangan, and M. C. Reed, "Femtocells: Past, present, and future," *IEEE J. Sel. Areas Commun.*, vol. 30, no. 3, pp. 497–508, Apr. 2012.
- [4] J. G. Andrews, F. Baccelli, and R. K. Ganti, "A tractable approach to coverage and rate in cellular networks," *IEEE Trans. Commun.*, vol. 59, no. 11, pp. 3122–3134, Nov. 2011.
- [5] P. Madhusudhanan, J. G. Restrepo, Y. Liu, T. X. Brown, and K. R. Baker, "Multi-tier network performance analysis using a shotgun cellular system," in *IEEE Global Telecommunications Conference (GLOBECOM 2011)*, 2011, pp. 1–6.
- [6] H.-S. Jo, Y. J. Sang, P. Xia, and J. G. Andrews, "Heterogeneous cellular networks with flexible cell association: A comprehensive downlink SINR analysis," *IEEE Trans. Wireless Commun.*, vol. 11, no. 10, pp. 3484–3495, Oct. 2012.
- [7] S. Mukherjee, "Distribution of downlink SINR in heterogeneous cellular networks," *IEEE J. Sel. Areas Commun.*, vol. 30, no. 3, pp. 575–585, Apr. 2012.
- [8] C. B. Sankaran, "Data offloading techniques in 3GPP Rel-10 networks: A tutorial," *IEEE Commun. Mag.*, vol. 50, no. 6, pp. 46–53, June 2012.
- [9] R. W. Heath, M. Kountouris, and T. Bai, "Modeling heterogeneous network interference using Poisson point processes," *IEEE Trans. Signal Process.*, vol. 61, no. 16, pp. 4114–4126, Aug. 2013.
- [10] H. S. Dhillon, M. Kountouris, and J. G. Andrews, "Downlink MIMO HetNets: Modeling, ordering results and performance analysis," *IEEE Trans. Wireless Commun.*, vol. 12, no. 10, pp. 5208–5222, Oct. 2013.
- [11] A. K. Papazafeiropoulos and T. Ratnarajah, "Downlink MIMO HCNs with residual transceiver hardware impairments," *IEEE Commun. Lett.*, vol. 20, no. 10, pp. 2023–2026, Oct. 2016.
- [12] M. Kountouris and J. G. Andrews, "Downlink SDMA with limited feedback in interference-limited wireless networks," *IEEE Trans. Wireless Commun.*, vol. 11, no. 8, pp. 2730–2741, Aug. 2012.
- [13] A. Papazafeiropoulos and T. Ratnarajah, "Towards a realistic assessment of multiple antenna HCNs: Residual additive transceiver hardware impairments and channel aging," *IEEE Trans. Veh. Technol.*, vol. 66, no. 10, pp. 9061–9073, Oct. 2017.
- [14] T. D. Novlan, H. S. Dhillon, and J. G. Andrews, "Analytical modeling of uplink cellular networks," *IEEE Trans. Wireless Commun.*, vol. 12, no. 6, pp. 2669–2679, June 2013.
- [15] T. Bai and R. W. Heath, "Analyzing uplink SINR and rate in massive MIMO systems using stochastic geometry," *IEEE Trans. Commun.*, vol. 64, no. 11, pp. 4592–4606, Nov. 2016.
- [16] H. Q. Ngo, E. Larsson, and T. Marzetta, "Energy and spectral efficiency of very large multiuser MIMO systems," *IEEE Trans. Commun.*, vol. 61, no. 4, pp. 1436–1449, Apr. 2013.
- [17] L. Lu, G. Y. Li, A. L. Swindlehurst, A. Ashikhmin, and R. Zhang, "An overview of massive MIMO: Benefits and challenges," *IEEE J. Sel. Topics Signal Process.*, vol. 8, no. 5, pp. 742–758, Oct. 2014.
- [18] T. Marzetta, "Noncooperative cellular wireless with unlimited numbers of base station antennas," *IEEE Trans. Wireless Commun.*, vol. 9, no. 11, pp. 3590–3600, Nov. 2010.
- [19] E. Björnson, J. Hoydis, M. Kountouris, and M. Debbah, "Massive MIMO systems with non-ideal hardware: Energy efficiency, estimation, and capacity limits," *IEEE Trans. Inf. Theory*, vol. 60, no. 11, pp. 7112–7139, Nov. 2014.
- [20] E. Björnson, M. Matthaiou, and M. Debbah, "Massive MIMO with non-ideal arbitrary arrays: Hardware scaling laws and circuit-aware design," *IEEE Trans. Wireless Commun.*, vol. 14, no. 8, pp. 4353–4368, Aug. 2015.
- [21] A. K. Papazafeiropoulos and T. Ratnarajah, "Deterministic equivalent performance analysis of time-varying massive MIMO systems," *IEEE Trans. Wireless Commun.*, vol. 14, no. 10, pp. 5795–5809, Oct. 2015.
- [22] A. K. Papazafeiropoulos, "Impact of user mobility on optimal linear receivers in cellular networks," in *IEEE International Conference on Communications (ICC)*, 2015, pp. 2239–2244.
- [23] —, "Impact of general channel aging conditions on the downlink performance of massive MIMO," *IEEE Trans. Veh. Technol.*, vol. 66, no. 2, pp. 1428–1442, Feb. 2017.
- [24] E. Baştuğ, M. Bennis, and M. Debbah, "Living on the edge: The role of proactive caching in 5G wireless networks," *IEEE Commun. Mag.*, vol. 52, no. 8, pp. 82–89, May 2014.
- [25] H. Liu, Z. Chen, X. Tian, X. Wang, and M. Tao, "On content-centric wireless delivery networks," *IEEE Wireless Commun.*, vol. 21, no. 6, pp. 118–125, Dec. 2014.
- [26] Z. Zhao, M. Peng, Z. Ding, W. Wang, and H. V. Poor, "Cluster content caching: An energy-efficient approach to improve quality of service in cloud radio access networks," *IEEE J. Sel. Areas Commun.*, vol. 34, no. 5, pp. 1207–1221, May 2016.
- [27] E. Baştuğ, M. Bennis, M. Kountouris, and M. Debbah, "Cache-enabled

- small cell networks: Modeling and tradeoffs,” *EURASIP J. Wireless Commun. and Net.*, vol. 2015, no. 1, p. 41, Dec. 2015.
- [28] E. Björnson, L. Sanguinetti, and M. Kountouris, “Deploying dense networks for maximal energy efficiency: Small cells meet massive MIMO,” *IEEE J. Sel. Areas Commun.*, vol. 34, no. 4, pp. 832–847, Apr. 2016.
- [29] A. Tang, J. Sun, and K. Gong, “Mobile propagation loss with a low base station antenna for NLOS street microcells in urban area,” in *IEEE VTS 53rd Vehicular Technology Conference, Spring 2001.*, vol. 1. IEEE, 2001, pp. 333–336.
- [30] S. N. Chiu, D. Stoyan, W. S. Kendall, and J. Mecke, *Stochastic Geometry and its Applications*. John Wiley & Sons, 2013.
- [31] N. Liang, W. Zhang, and C. Shen, “An uplink interference analysis for massive MIMO systems with MRC and ZF receivers,” in *IEEE Wireless Commun. and Net. Conf. (WCNC), 2015*, pp. 310–315.
- [32] F. Baccelli and B. Błaszczyszyn, *Stochastic Geometry and Wireless Networks: Volume I Theory*. NoW Publishers, 2009.
- [33] H. ElSawy, E. Hossain, and M. Haenggi, “Stochastic geometry for modeling, analysis, and design of multi-tier and cognitive cellular wireless networks: A survey,” *IEEE Commun. Surveys Tuts.*, vol. 15, no. 3, pp. 996–1019, July 2013.
- [34] X. Zhang and J. G. Andrews, “Downlink cellular network analysis with multi-slope path loss models,” *IEEE Trans. Commun.*, vol. 63, no. 5, pp. 1881–1894, May 2015.
- [35] Cisco. Cisco visual networking index: Global mobile data traffic forecast update, 2016–2021 white paper. [Online]. Available: <https://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/mobile-white-paper-c11-520862.html>
- [36] K. Mardia and P. Zemerch, “Algorithm AS 86: The von Mises distribution function,” *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, vol. 24, no. 2, pp. 268–272, 1975.
- [37] K. E. Baddour and N. C. Beaulieu, “Autoregressive modeling for fading channel simulation,” *IEEE Trans. Wireless Commun.*, vol. 4, no. 4, pp. 1650–1662, July 2005.
- [38] M. Vu and A. Paulraj, “On the capacity of MIMO wireless channels with dynamic CSIT,” *IEEE J. Select. Areas Commun.*, vol. 25, no. 7, pp. 1269–1283, Sep. 2007.
- [39] W. C. Jakes and D. C. Cox, *Microwave mobile communications*. Wiley-IEEE Press, 1994.
- [40] R. Couillet and M. Debbah, *Random Matrix Methods for Wireless Communications*. Cambridge University Press, 2011.
- [41] A. W. van der Vaart, *Asymptotic statistics (Cambridge series in statistical and probabilistic mathematics)*. New York: Cambridge University Press, 2000.
- [42] S. Wagner, R. Couillet, M. Debbah, and D. Slock, “Large system analysis of linear precoding in correlated MISO broadcast channels under limited feedback,” *IEEE Trans. Inf. Theory*, vol. 58, no. 7, pp. 4509–4537, July 2012.
- [43] Z. Bai and J. W. Silverstein, *Spectral Analysis of Large Dimensional Random Matrices*. Springer, 2010, vol. 20.
- [44] H. Alzer, “On some inequalities for the incomplete gamma function,” *Mathematics of Computation of the American Mathematical Society*, vol. 66, no. 218, pp. 771–778, 1997.
- [45] R. K. Mungara, D. Morales-Jimenez, and A. Lozano, “System-level performance of interference alignment,” *IEEE Trans. Wireless Commun.*, vol. 14, no. 2, pp. 1060–1070, Feb. 2015.