# How Many Small Cells Can Be Turned off via Vertical Offloading under a Separation Architecture?

Shan Zhang, *Student Member, IEEE,* Jie Gong, *Member, IEEE,* Sheng Zhou, *Member, IEEE,* and Zhisheng Niu, *Fellow, IEEE*

*Abstract*—To further improve the energy efficiency of heterogeneous networks, a separation architecture called hyper-cellular network (HCN) has been proposed, which decouples the control signaling and data transmission functions. Specifically, the control coverage is guaranteed by macro base stations (MBSs), whereas small cells (SCs) are only utilized for data transmission. Under HCN, SCs can be dynamically turned off when traffic load decreases for energy saving. A fundamental problem then arises: how many SCs can be turned off as traffic varies? In this paper, we address this problem in a theoretical way, where two sleeping schemes (i.e., random and repulsive schemes) with vertical inter-layer offloading are considered. Analytical results indicate the following facts: (1) Under the random scheme where SCs are turned off with certain probability, the expected ratio of sleeping SCs is inversely proportional to the traffic load of SC-layer and decreases linearly with the traffic load of MBS-layer; (2) The repulsive scheme, which only turns off the SCs close to MBSs, is less sensitive to the traffic variations; (3) deploying denser MBSs enables turning off more SCs, which may help to improve network energy-efficiency. Numerical results show that about 50% SCs can be turned off on average under the predefined daily traffic profiles, and 10% more SCs can be further turned off with inter-layer channel borrowing.

## I. INTRODUCTION

The multi-tier heterogeneous networks (HetNets), which consist of different types of base stations (BSs) (such as macro BSs, micro BSs, pico BSs and femto BSs), can effectively improve network capacity and thus are expected to be the dominant scenarios in the 5G era [1] [2] [3]. However, the huge energy consumption of HetNets brings heavy burdens to the network operators [4] [5]. Meanwhile, due to the dynamics of wireless traffic load, many BSs are lightly-loaded but still consume almost their peak energy on account of elements like airconditioner and power amplifier. Unfortunately, these low-efficient BSs can not be turned off for coverage guarantee, which makes the existing network energy inefficient [6].

To solve this problem, we have proposed a new separation architecture called *Hyper-Cellular Network* (HCN), whose
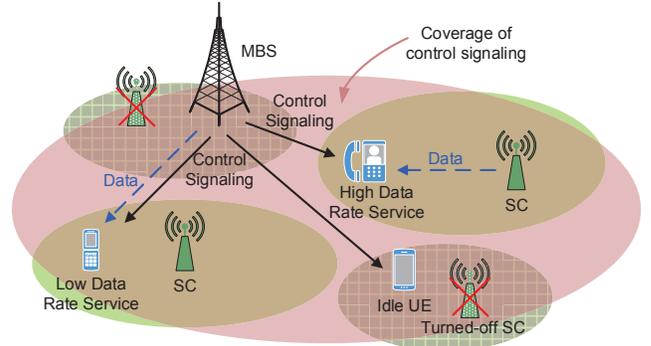
Fig. 1: Network architecture of HCN.

main idea is to decouple the coverage of control signaling from the coverage of data transmission such that the data coverage can be more elastic in accordance with the traffic dynamics [7]. Under HCN, macro base stations (MBSs) and small cells[1] (SCs) play different roles as shown in Fig. 1. Generally, SCs are only utilized for high data rate transmission, whereas MBSs guarantee the network coverage and provide low data rate service. Therefore, each UE is always connected with the MBS-layer for control signaling, while their data traffics are served by MBSs or SCs depending on their service demands. For example, the UEs making video calls should keep dual connections with both MBSs and SCs, while UEs making voice calls are only served by MBSs. The HCN architecture can be realized by separating C/U planes [8]-[11], which is a key technology for future 5G networks [12][13].

With network coverage well-guaranteed by the MBS-layer, SCs can be turned off flexibly without causing coverage holes. At the same time, the quality of service (QoS) of the UEs within sleeping SCs can be satisfied by being offloaded either horizontally to neighboring or vertically to high-layer cells. Although offloading traffic from SCs to MBS increases the transmit power, turning off small cells saves energy as the constant power of SCs is much larger than the power consumed by data transmission [14]. Then a fundamental problem arises: how many SCs can be turned off via traffic offloading for given QoS requirements under HCN.

As a starting point, we analyze the maximum ratio of sleeping SCs due to the traffic dynamics in time domain

---

[1]Small cells are the coverage of the BSs with relatively low transmit power (such as micro and pico BSs).

from the perspective of the whole network, where the UEs within sleeping SCs are only offloaded to MBSs (vertical offloading). In fact, the performance of horizontal offloading is limited due to the low transmit power of SCs, and accordingly a SC can only help to offload the traffic of neighboring cells within certain range. Therefore, horizontal offloading between SCs can not always be realized, and vertical offloading is sometimes the only choice for some SC to go into sleep if its neighboring SCs are not close enough. Besides, due to the dual connectivity of UEs in HCN, vertical offloading can be easily implemented without handovers. Two sleeping schemes are considered in this paper:

1) Random scheme: each SC is turned off with equivalent probability $p_s$;
2) Repulsive scheme: the SCs whose distance to the nearest MBSs smaller than $R_s$ go into sleep, whereas the other SCs remain active.

Then, our problem is to find the maximum $p_s$ and $R_s$ which satisfy the given outage constraints. To conduct theoretical analysis, a two-layer HCN is considered, where MBSs are assumed to be regularly deployed as hexagonal cells whereas the distribution of SCs is modeled as Voronoi tessellation of a homogeneous Poisson Point Process (PPP). The main contributions of this paper include:

- The approximated closed-form expressions of the outage probability are derived under the two sleeping schemes, which are validated through extensive simulations.
- The maximum sleeping ratio of SCs under the two scheme is derived based on the analytical outage probability, and the influences of system parameters (such as traffic load and BS density) are analyzed.
- We also consider the case when the redundant bandwidth of the SC-layer can be released and re-allocated to the MBS-layer to serve more offloaded UEs (channel borrowing).
- Under a typical scenario and two different daily traffic profiles, we show that about 50% SCs can go into sleep on average, and 10% more SCs can be further turned off if channel borrowing is conducted.

In summary, the analytical results of the maximum sleeping ratio of SCs are obtained under the random and repulsive schemes, which offers a guideline for real network planning and operations. In addition to the implementation considerations, the sleeping algorithm design of the separation architecture is consistent with that of the conventional networks. Therefore, our method also applies to the conventional two-tier heterogeneous networks.

The rest of this paper is organized as follows. Related work is introduced in Section II. System model is described in Section III, and the outage probability is derived in Section IV. Then, the random and repulsive schemes are analyzed in Section V and the two schemes are evaluated under daily traffic profiles in Section VI. At last, Section VII concludes this paper.

## II. RELATED WORK

Under the traditional network architecture, BSs can not be turned off flexibly due to the requirement of coverage guarantee. To tackle this problem, *cell zooming* is proposed to compensate for the sleeping cells by enlarging the coverage of neighboring BSs, which can be realized through power control and antenna titling [15]. Meanwhile, *horizontal offloading* is widely adopted for QoS guarantee [16]-[21], under which the UEs of sleeping cells are offloaded to neighboring cells regardless of their cell types. SC sleeping through horizontal offloading under HetNets is investigated theoretically in [23]-[24]. In [23], the upper and lower bounds of the optimal density of active SCs are derived under time-varying traffic load, where SCs are turned off randomly. Instead, the SCs close to MBSs are turned off in [24], after which the corresponding UEs will be re-associated to active cells based on the received signal strength. However, coverage holes may still exist via these horizontal offloading schemes in real systems due to the complex wireless environment. Therefore, *vertical offloading* is adopted as an alternative solution to guarantee QoS in the multi-tier HetNets [25]-[28], where the UEs of sleeping SCs are offloaded to MBSs instead of neighboring SCs. In [25]-[27], SCs are dynamically turned on/off based on the predicted future traffic load. Besides, a dynamic energy-optimal SC sleeping algorithm is proposed based on Markov Decision Process [28]. Whereas these studies have high complexity and are limited to small scale networks.

Recently, separation architecture has received more attention. A SC sleeping scheme under HCN is proposed and optimized in [29], where the SCs are turned off probabilistically according to their traffic load via vertical offloading. Furthermore, the energy saving gain of HCN is firstly analyzed in [30], where random sleeping scheme via horizontal offloading are considered. The main differences between [30] and our work are that we focus on vertical offloading, which has been rarely addressed under the separation architecture.

## III. SYSTEM MODEL

We consider a two-layer HCN, under which the service procedure for a UE is described in Fig. 2. When a detached UE arrives, it will firstly connect to the MBS-layer for basic control signaling, and the MBS which offers maximal average Signal to Interference and Noise Ratio (SINR) will be selected. Then, the UE will always connect the MBS-layer until it leaves the HCN. During this period, the UE initiated sessions will be served according to the data rate demand and mobility. The SCs are utilized for high data rate transmission (like real-time video, online game), whereas MBSs mainly guarantee lower data rate services (such as voice). Besides, high mobile UEs will be served by MBSs to avoid frequent handovers. Therefore, only the UEs with low mobility and high data rate requirement will choose a SC for data service, during which they maintain dual connections with both layers for signaling and data
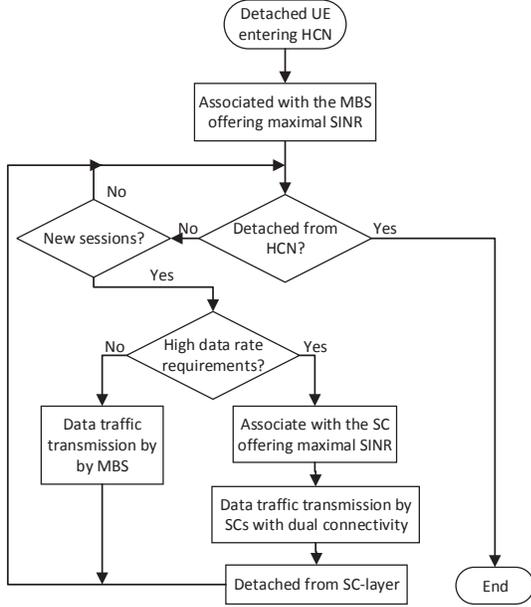
Fig. 2: Service process for a typical UE under HCN.



MBS    • Active SC

Fig. 3: Network topology.

respectively. Notice that their connections with the SC-layer will be dropped once their sessions are completed.

According to the service status, all active UEs in HCN can be classified into 2 classes:

1) *MBS UEs*: UEs served by MBSs with low data rate;
2) *SC UEs*: UEs served by SCs with high data rate.

Because we focus on the time dynamics of traffic load instead of the spatial non-uniformity, we assume the distributions of MBS UEs and SC UEs both follow homogeneous PPPs with different densities.

As for the network topology, the MBSs are assumed to be regularly deployed as hexagonal cells whereas SCs are assumed to be the Voronoi tessellation of a homogeneous PPP process as shown in Fig. 3, due to the different roles of MBSs and SCs. Recall that MBSs are expected to guarantee the network coverage, therefore their locations should be carefully designed. On the contrary, SCs are deployed only to boost the network capacity, whose locations can be quite random. In fact, the traditional regular hexagonal cells and the PPP have similar accuracy when used to model the distribution of BSs in real systems, whereas the former model is the ideal case and the latter offers a performance lower bound [32].

### A. Bandwidth Allocation

Orthogonal bandwidth is used by different layers to avoid the severe inter-layer interference, especially for protecting the signaling coverage. Furthermore, the spectrum reuse factor within each layer is set to 1. When some SCs are turned off with their UEs offloaded to MBSs, the spectrum resource of the MBS-layer is further divided into two parts to serve their associated MBS UEs and the offloaded SC UEs respectively. In addition, more SCs can be turned off if
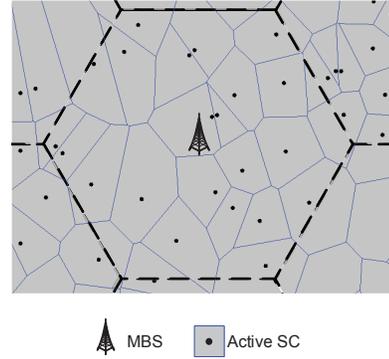
the redundant bandwidth of the SC-layer can be released and reallocated to the MBS-layer for traffic offloading (named as *Channel Borrowing (CB)*). Note that the borrowed bandwidth cannot be used by the SC-layer. In other words, the bandwidth used by the offloaded UEs is given by

$$w_{\mathrm{o}} = \begin{cases} W_{\mathrm{m}} + W_{\mathrm{s}} - w_{\mathrm{m}} - w_{\mathrm{s}}, & \text{with CB} \\ W_{\mathrm{m}} - w_{\mathrm{m}}, & \text{no CB} \end{cases}, \quad (1)$$

where $W_{\mathrm{m}}$ and $W_{\mathrm{s}}$ are the total bandwidth initially pre-allocated to the MBS-layer and SC-layer, $w_{\mathrm{m}}$ and $w_{\mathrm{s}}$ are the bandwidth used by the non-offloaded UEs at the two layers respectively. Obviously, CB offers more opportunity for SC sleeping by making better use of the spectrum resource.

### B. Link Layer model

Assume all MBSs transmit at a fixed power $P^{\mathrm{m}}$. For a typical $\mathrm{UE}_u$, the received SINR is given as follows if it is served by $\mathrm{MBS}_i$.

$$\gamma_{iu}^{\mathrm{m}} = \frac{P^{\mathrm{m}}(d_{iu}^{\mathrm{m}})^{-\alpha_{\mathrm{m}}} h_{iu}^{\mathrm{m}}}{\sum\limits_{j \in \mathcal{B}^{\mathrm{m}}, j \neq i} P^{\mathrm{m}}(d_{ju}^{\mathrm{m}})^{-\alpha_{\mathrm{m}}} h_{ju}^{\mathrm{m}} + \sigma^2}, \quad (2)$$

where $\sigma^2$ is the noise power, $\mathcal{B}^{\mathrm{m}}$ is the set of active MBSs in the network, $d_{iu}^{\mathrm{m}}$ is the distance between $\mathrm{UE}_u$ and $\mathrm{MBS}_i$, $\alpha_{\mathrm{m}}$ is the path loss factor of MBS-layer, and $h_{iu}^{\mathrm{m}}$ is an exponential random variable with mean 1 incorporating the effect of Rayleigh fading. Assume each BS allocates resource (e.g., time slots or wireless spectrum) equally to its associated active UEs, then the achievable rate of $\mathrm{UE}_u$ is given by:

$$r_{iu}^{\mathrm{m}} = \frac{w_{\mathrm{m}}}{N_i^{\mathrm{m}} + 1} \log_2(1 + \gamma_{iu}^{\mathrm{m}}), \quad (3)$$

where $w_{\mathrm{m}}$ is the bandwidth used by each MBS for its associated UEs, and $N_i^{\mathrm{m}}$ is a random variable denoting the number of active residual UEs served by $\mathrm{MBS}_i$ except $\mathrm{UE}_u$. Then the outage probability that the rate of $\mathrm{UE}_u$ is less than the predefined threshold $U_{\mathrm{m}}$ is given by $\mathbb{P}\{r_{iu}^{\mathrm{m}} < U_{\mathrm{m}}\}$. By averaging this probability over the possible position of $\mathrm{UE}_u$, and $N_i^{\mathrm{m}}$, we can get the service outage constraint of MBS UEs as[2]:

$$G_{\mathrm{m}} = \mathbb{E}_{\{N_{\mathrm{m}}, d_{\mathrm{m}}\}} \left\{ \mathbb{P}\left( \frac{w_{\mathrm{m}}}{N_{\mathrm{m}} + 1} \log_2(1 + \gamma_{\mathrm{m}}) < U_{\mathrm{m}} \middle| N_{\mathrm{m}}, d_{\mathrm{m}} \right) \right\} < \eta_{\mathrm{m}}. \quad (4)$$

---

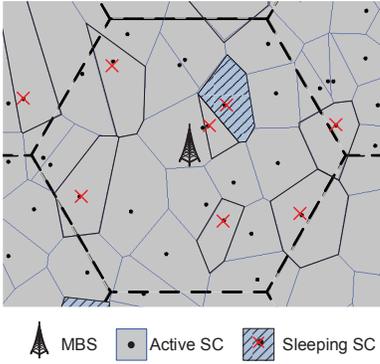[2]The subscripts $i$ and $u$ are omitted here for simplicity.
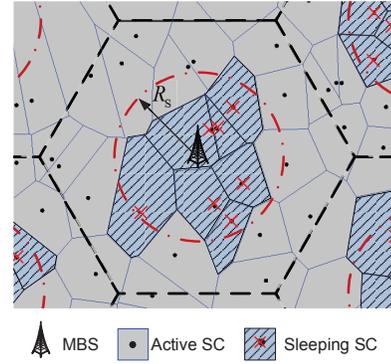
Fig. 4: Illustration of the random scheme.



Fig. 5: Illustration of the repulsive scheme.

Assume all SBSs also adopt the constant transmit power $P^s$. If $UE_u$ is served by $SC_k$, its received SINR is given by:

$$\gamma_{ku}^s = \frac{P^s(d_{ku}^s)^{-\alpha_s} h_{ku}^s}{\sum\limits_{l \in \mathcal{B}^s, l \neq k} P^s(d_{lu}^s)^{-\alpha_s} h_{ku}^s + \sigma^2}, \tag{5}$$

where $\mathcal{B}^s$ is the set of active SCs in the network, $d_{ku}^s$ is the distance between $UE_u$ and $SC_k$, $\alpha_s$ is the path loss factor of SC-layer, and $h_{ku}^s$ is an exponential random variable with mean 1 incorporating the effect of Rayleigh fading. Similarly, the outage probability constraint of the SC UEs can be obtained[3]:

$$G_s = \mathbb{E}_{\{A_s, N_s, d_s\}} \left\{ \mathbb{P}\left( \frac{w_s}{N_s + 1} \log_2(1 + \gamma_s) < U_s \Big| N_s, d_s, A_s \right) \right\} < \eta_s, \tag{6}$$

where $w_s$ is the bandwidth used by each SC, $N_s$ is the number of residual UEs in the target SC, $\gamma_s$ is the spectrum efficiency of the considered SC UE, whose distance to the target SC is denoted as $d_s$. In addition, the cell size of the target SC is also random, denoted as $A_s$.

As for a typical offloaded UE, its outage probability can be derived in the same way:

$$G_o = \mathbb{E}_{\{N_o, d_o\}} \left\{ \mathbb{P}\left( \frac{w_o}{N_o + 1} \log_2(1 + \gamma_o) < U_o \Big| N_o, d_o \right) \right\} < \eta_o, \tag{7}$$

where $d_o$ is the distance between the typical offloaded UE to its associated MBS, $N_o$ denotes the number of remaining offloaded UEs in the same MBS cell, and $\gamma_o$ means the received SINR varying with channel fading. Notice that the QoS of the offloaded UEs will not degrade if $U_o = U_s$ and $\eta_o = \eta_s$.

### C. SC Sleeping Schemes

Our problem can be formulated to maximize the ratio of SCs turned off under the outage probability constrains of Eqs. (4), (6) and (7). However, this problem is hard to solve as there exist no closed-form expressions of the outage probability. In fact, even the distributions of $N_s$, $N_o$, $A_s$, and $d_o$ do not have general expressions.

To conduct theoretical analysis, we consider two basic sleeping schemes: random scheme and repulsive scheme. Under the random scheme, SCs are treated equivalently and

[3] The subscripts $k$ and $u$ are omitted here for simplicity.

go into sleep independently with probability $p_s$. Whereas, SCs are differentiated by their distance to the MBSs under the repulsive scheme, and only the SCs whose distance to the nearest MBSs is less than sleeping radius $R_s$ can be turned off. The two schemes are illustrated in Fig. 4 and Fig. 5 respectively.

The random scheme has been adopted in many studies, which is considered as a a baseline for cell sleeping [21], [23], and [30]. Besides, the random sleeping can be easily implemented as it requires no extra information like traffic load and locations of each cell, and the average ratio of sleeping SCs is equal to the sleeping probability $p_s$. Under the repulsive scheme, the offloaded UEs generally enjoy smaller path loss and larger spectrum efficiency for shorter access distance to the MBSs, and thus the QoS of the offloaded UEs can be more easily guaranteed. In this case, the average sleeping ratio of SCs is given by $\frac{\pi R_s^2}{\frac{3\sqrt{3}}{2} D^2}$, where $D$ is the coverage radius of each MBS cell.

When power control is not implemented, the power consumption of a BS becomes a constant [14]. Thus, the network power consumption can be treated as the weighted sum of the MBS and SC densities. Therefore, the ratio of sleeping SCs reflects the energy saving gain.

### IV. OUTAGE PROBABILITY ANALYSIS

#### A. Outage Constraint of MBS UEs

For the outage probability of MBS UEs given by Eq. (4), we have

$$G_m = \mathbb{E}_{\{N_m, d_m\}} \left\{ \mathbb{P}\left\{ \frac{w_m}{N_m + 1} \log_2(1 + \gamma_m) < U_m \Big| N_m, d_m \right\} \right\}$$
$$= \int_0^D \sum_{n=0}^{\infty} \mathbb{P}\left\{ \gamma_m < 2^{\frac{(n+1)U_m}{w_m}} - 1 \Big| d \right\} p_{N_m}(n) f_{d_m}(d) \mathrm{d}d, \tag{8}$$

where $p_{N_m}(n)$ is the probability that the target MBS has $n$ residential MBS UEs except $UE_u$, and $f_{d_m}(d)$ is the probability density function of $d_m$. As the distribution of MBS UEs follows homogeneous PPP process, $N_m$ follows the Poisson distribution of parameter $\frac{3\sqrt{3}}{2}\lambda_m D^2$ according to Slivnyak-Mecke theorem [33], where $\lambda_m$ is the density of MBS UEs. In addition, $f_{d_m}(d) = \frac{2\pi}{D} d$ by assuming UEs to be uniformly distributed within a circle cell of radius $D$. The closed-form expression of $G_m$ can be derived when the

received SNR is high, which generally holds for current cellular systems.

**Theorem 1.** As $\frac{\sigma^2}{P^{\mathrm{m}}} \to 0$, the outage probability of MBS UEs is given by

$$G_{\mathrm{m}} = \frac{2D^{\alpha_{\mathrm{m}}}(I+1)\sigma^2}{P^{\mathrm{m}}(\alpha_{\mathrm{m}}+2)} \left( 2^{\frac{U_{\mathrm{m}}}{w_{\mathrm{m}}}} \exp\left( \frac{3\sqrt{3}}{2} D^2 \lambda_{\mathrm{m}} \left( 2^{\frac{U_{\mathrm{m}}}{w_{\mathrm{m}}}} - 1 \right) \right) - 1 \right),$$
(9)

where $I$ denotes the ratio of inter-cell interference to noise. *Proof*: See Appendix A. ∎

Furthermore, when $\frac{U_{\mathrm{m}}}{w_{\mathrm{m}}} \to 0$, Eq. (9) can be further simplified:

$$G_{\mathrm{m}} = \frac{2D^{\alpha_{\mathrm{m}}}(I+1)\sigma^2}{P^{\mathrm{m}}(\alpha_{\mathrm{m}}+2)} \left( 2^{\frac{U_{\mathrm{m}}}{w_{\mathrm{m}}} \left( 1 + \frac{3\sqrt{3}}{2} D^2 \lambda_{\mathrm{m}} \frac{U_{\mathrm{m}}}{w_{\mathrm{m}}} \right)} - 1 \right),$$
(10)

and the service outage constraint of MBS UEs Eq. (4) is equivalent to

$$\bar{w}_{\mathrm{m}} \log_2(1 + \tau_{\mathrm{m}}) \geq U_{\mathrm{m}},$$
(11)

where $\bar{w}_{\mathrm{m}} = \frac{w_{\mathrm{m}}}{1 + \frac{3\sqrt{3}}{2}\lambda_{\mathrm{m}} D^2}$ is the expected bandwidth allocated to each MBS UE, and $\tau_{\mathrm{m}}$ denotes the received SINR of cell edge UEs depending on $\eta_{\mathrm{m}}$ given by

$$\tau_{\mathrm{m}} = \frac{P^{\mathrm{m}}}{\sigma^2(I+1)} \frac{\alpha_{\mathrm{m}}+2}{2} \frac{\eta_{\mathrm{m}}}{D^{\alpha_{\mathrm{m}}}}.$$
(12)

Notice that Eq. (11) is a linear constraint on $w_{\mathrm{m}}$, and its physical meaning is that the average data rate of the non-cell-edge UEs (received SINR above $\tau_{\mathrm{m}}$) should be no smaller than the $U_{\mathrm{m}}$.

The physical meaning of $\frac{U_{\mathrm{m}}}{w_{\mathrm{m}}} \to 0$ is that the data rate requirement is relatively low compared with the spectrum resource, which is reasonable as each MBS usually supports large number of UEs simultaneously in real systems. Therefore, Eq. (11) can be applied to simplify the service outage constraint of MBS UEs.

### B. Outage Constraint of SC UEs

The outage probability for a typical SC UE (Eq. (6)) is given by

$$G_{\mathrm{s}} = \mathbb{E}_{\{A_{\mathrm{s}}, N_{\mathrm{s}}, d_{\mathrm{s}}\}} \left\{ \mathbb{P}\left\{ \frac{w_{\mathrm{s}}}{N_{\mathrm{s}}+1} \log_2(1 + \gamma_{\mathrm{s}}) < U_{\mathrm{s}} \Big| A_{\mathrm{s}}, N_{\mathrm{s}}, d_{\mathrm{s}}, \right\} \right\}$$
$$= \int_0^\infty \int_0^\infty \sum_{n=0}^\infty \mathbb{P}\left\{ \gamma_{\mathrm{s}} < 2^{\frac{(n+1)U_{\mathrm{s}}}{w_{\mathrm{s}}}} - 1 \right\} p_{N_{\mathrm{s}}}(n) f_{A_{\mathrm{s}}}(a) f_{d_{\mathrm{s}}}(d) \mathrm{d}a \mathrm{d}d,$$
(13)

where $p_{N_{\mathrm{s}}}(n)$, $f_{A_{\mathrm{s}}}(a)$ and $f_{d_{\mathrm{s}}}(d)$ denote the probability distribution functions of $N_{\mathrm{s}}$, $A_{\mathrm{s}}$ and $d_{\mathrm{s}}$ respectively.

### (1) No SC Sleeping

Firstly, we analyze the outage probability when all SCs are active. In this case, $N_{\mathrm{s}}$ follows Poisson distribution of $A_{\mathrm{s}}\lambda_{\mathrm{s}}$, where $\lambda_{\mathrm{s}}$ is the density of SC UEs. Besides, $A_{\mathrm{s}}$ follows Gamma distribution with shape $K = 3.575$ and scale $\frac{1}{K\rho_{\mathrm{s}}}$, where $\rho_{\mathrm{s}}$ is the density of SCs [34].

Notice that SCs will be more densely deployed in the future to boost network capacity, and the inter-cell interference, instead of noise, will be the main factor influencing the received SINR of SC UEs. In this case, we

can derive the approximated outage probability as Theorem 2.

**Theorem 2.** If $U_{\mathrm{s}}/w_{\mathrm{s}} \to 0$ and the SC-layer is interference-limited, the outage probability of the SC UEs without SC sleeping is given by:

$$G_{\mathrm{s}} = 1 - \frac{\frac{\alpha_{\mathrm{s}}-2}{2} 2^{-\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}}\left(1 + \frac{\lambda_{\mathrm{s}}}{\rho_{\mathrm{s}}}\right)}}{1 - \frac{4-\alpha_{\mathrm{s}}}{2} 2^{-\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}}\left(1 + \frac{\lambda_{\mathrm{s}}}{\rho_{\mathrm{s}}}\right)}}.$$
(14)

*Proof*: See Appendix B. ∎

Similarly, $U_{\mathrm{s}}/w_{\mathrm{s}} \to 0$ means the bandwidth is relatively large compared with the data rate requirement. Based on Theorem 2, the outage constraint for SC UEs without SC sleeping is equivalent to

$$\bar{w}_{\mathrm{s}} \log_2(1 + \tau_{\mathrm{s}}) \geq U_{\mathrm{s}},$$
(15)

where $\bar{w}_{\mathrm{s}} = \frac{w_{\mathrm{s}}}{1+N_{\mathrm{s}}}$ is the average bandwidth allocated to each SC UE, and $\tau_{\mathrm{s}}$ is defined as

$$\tau_{\mathrm{s}} = \frac{\alpha_{\mathrm{s}}-2}{2} \frac{\eta_{\mathrm{s}}}{1-\eta_{\mathrm{s}}},$$
(16)

denoting the received SINR of cell edge UEs of the SC-layer.

### (2) Repulsive Scheme

Under the repulsive scheme, the distributions of $N_{\mathrm{s}}$, $A_{\mathrm{s}}$ and $d_{\mathrm{s}}$ remain the same after some SCs are turned off. As for the received SINR, only the UEs which locate around the sleeping area enjoy reduced inter-cell interference, while the received SINR of the other UEs is barely influenced. Therefore, we use Eq. (15) to approximate the outage constraint of SC UEs by ignoring the benefit brought by SC sleeping. This is a conservative approximation of the real case.

### (3) Random Scheme

Similarly to the repulsive scheme, the distributions of $N_{\mathrm{s}}$, $A_{\mathrm{s}}$ and $d_{\mathrm{s}}$ are not influenced by SC sleeping under the random scheme. However, the inter-cell interference received by a typical SC UE decreases by $p_{\mathrm{s}}$ on average. If the network is still interference-limited after SC sleeping, then noise can be ignored and the received SINR of the SC UEs will increase by $1/(1 - p_{\mathrm{s}}) - 1$. Nonetheless, the received SINR will finally level off as the network becomes noise-limited.

Inspired by Theorem 2, we use the following inequality to approximate the outage constraint of SC UEs under the random scheme:

$$\bar{w}_{\mathrm{s}} \log_2\left(1 + \tau_{\mathrm{s}}'(p_{\mathrm{s}})\right) \geq U_{\mathrm{s}},$$
(17)

where $\tau_{\mathrm{s}}'(p_{\mathrm{s}})$ is defined as

$$\tau_{\mathrm{s}}'(p_{\mathrm{s}}) = \frac{\alpha_{\mathrm{s}}-2}{2(1 - \min(\hat{p}_{\mathrm{s}}, p_{\mathrm{s}}))} \frac{\eta_{\mathrm{s}}}{1-\eta_{\mathrm{s}}},$$
(18)

and $\hat{p}_{\mathrm{s}}$ is an experimental threshold of $p_{\mathrm{s}}$, indicating whether the network is interference-limited or not. Specifically, the network is considered as noise-limited if $\hat{p}_{\mathrm{s}} < p_{\mathrm{s}}$, in which case turning off SCs no longer improve the received SINR. In fact, the additional noise is approximated by the inter-cell interference of $p_{\mathrm{s}} = \hat{p}_{\mathrm{s}}$ in Eq. (18). Thus the approximated QoS constraint Eq. (17) is more strict for smaller $\hat{p}_{\mathrm{s}}$. Notice that $\tau_{\mathrm{s}}'(p_{\mathrm{s}})$ reflects both inter-cell interference and noise.

TABLE I: Simulation Parameters

| Parameter | Value | Parameter | Value |
|-----------|-------|-----------|-------|
| $D$ | 500m | $\rho_{\mathrm{s}}$ | 25/km² |
| $P^{\mathrm{m}}$ | 10W | $P^{\mathrm{s}}$ | 1W |
| $W_{\mathrm{m}}$ | 10MHz | $W_{\mathrm{s}}$ | 10MHz |
| $\sigma^2$ | -104dBm | $\alpha_{\mathrm{s}}$ | 4 |
| $U_{\mathrm{m}}$ | 64kbps | $\eta_{\mathrm{m}}$ | 0.05 |
| $U_{\mathrm{s}}$ | 100kbps | $\eta_{\mathrm{s}}$ | 0.05 |
| $U_{\mathrm{o}}$ | 100kbps | $\eta_{\mathrm{o}}$ | 0.05 |

### C. Outage Constraint of Offloaded UEs

#### (1) No CB

Under the random scheme, the outage probability of the offloaded UEs can be obtained in the same way as Theorem 1, and the outage QoS is equivalent to:

$$\frac{W_{\mathrm{m}} - w_{\mathrm{m}}}{1 + \frac{3\sqrt{3}}{2}\lambda_{\mathrm{s}} p_{\mathrm{s}} D^2} \log_2\left(1 + \frac{P^{\mathrm{m}}}{\sigma^2(1+I)} \frac{\alpha_{\mathrm{m}}+2}{2} \frac{\eta_{\mathrm{o}}}{D^{\alpha_{\mathrm{m}}}}\right) \geq U_{\mathrm{o}}. \quad (19)$$

Similarly, the outage constraint of the offloaded UEs under repulsive scheme can also be simplified:

$$\frac{W_{\mathrm{m}} - w_{\mathrm{m}}}{1 + \pi R_{\mathrm{s}}^2 \lambda_{\mathrm{s}}} \log_2\left(1 + \frac{P^{\mathrm{m}}}{\sigma^2(1+I)} \frac{\alpha_{\mathrm{m}}+2}{2} \frac{\eta_{\mathrm{o}}}{R_{\mathrm{s}}^{\alpha_{\mathrm{m}}}}\right) \geq U_{\mathrm{o}}. \quad (20)$$

#### (2) With CB

In this case, there are two bands serving the offloaded UEs which may have different path loss factors: $W_{\mathrm{m}} - w_{\mathrm{m}}$ and $W_{\mathrm{s}} - w_{\mathrm{s}}$. For simplicity, we assume the offloaded UEs are divided into two groups randomly, and each group shares one band.

For the random scheme, the outage constraint for the offloaded UEs is as follows:

$$\begin{cases} \dfrac{W_{\mathrm{m}} - w_{\mathrm{m}}}{1 + \frac{3\sqrt{3}}{2}\lambda_{\mathrm{s}} p_{\mathrm{s}} p_{\mathrm{m}} D^2} \log_2\left(1 + \tau_{\mathrm{o}}(\alpha_{\mathrm{m}}, D)\right) \geq U_{\mathrm{o}} \\ \dfrac{W_{\mathrm{s}} - w_{\mathrm{s}}}{1 + \frac{3\sqrt{3}}{2}\lambda_{\mathrm{s}} p_{\mathrm{s}}(1 - p_{\mathrm{m}}) D^2} \log_2\left(1 + \tau_{\mathrm{o}}(\alpha_{\mathrm{s}}, D)\right) \geq U_{\mathrm{o}}, \end{cases} \quad (21)$$

where $p_{\mathrm{m}}$ is the probability that an offloaded UE uses the bandwidth of the MBS-layer, $\tau_{\mathrm{o}}(\alpha, r)$ is defined as

$$\tau_{\mathrm{o}}(\alpha, r) = \frac{P^{\mathrm{m}}}{\sigma^2(1+I)} \frac{\alpha+2}{2} \frac{\eta_{\mathrm{o}}}{r^{\alpha}}. \quad (22)$$

$\tau_{\mathrm{o}}(\alpha, r)$ is a threshold of the received SINR of the offloaded UEs, when they are uniformly distributed within circles of radius $r$ centered at MBSs with the path loss factor $\alpha$.

Similarly, the outage constraint of the offloaded UEs under the repulsive scheme is given by:

$$\begin{cases} \dfrac{W_{\mathrm{m}} - w_{\mathrm{m}}}{1 + \pi R_{\mathrm{s}}^2 \lambda_{\mathrm{s}} p_{\mathrm{m}}} \log_2\left(1 + \tau_{\mathrm{o}}(\alpha_{\mathrm{m}}, R_{\mathrm{s}})\right) \geq U_{\mathrm{o}} \\ \dfrac{W_{\mathrm{s}} - w_{\mathrm{s}}}{1 + \pi R_{\mathrm{s}}^2 \lambda_{\mathrm{s}}(1 - p_{\mathrm{m}})} \log_2\left(1 + \tau_{\mathrm{o}}(\alpha_{\mathrm{s}}, R_{\mathrm{s}})\right) \geq U_{\mathrm{o}}. \end{cases} \quad (23)$$

To achieve load balance of the two bands and turn off more SCs, the probability $p_{\mathrm{m}}$ is also considered as a variable to be optimized.

### D. Evaluations of Analytical Outage Probabilities

Now we evaluate the approximation errors of the derived closed-form outage probabilities under the two schemes. A HCN with 19 MBSs is considered for simulation, where the number of SCs (and UEs) and their locations are set up by Monte Carlo simulation method. The simulation parameters are listed in Table I [30]. By calculating the data rate of each UE, the outage probability can be obtained.

When $w_{\mathrm{m}} = w_{\mathrm{s}} = w_{\mathrm{o}} = 10\text{MHz}$, the simulation and analytical results are compared in Fig. 6. The analytical results of the four sub-figures come from Eqs. (11), (17), (20), and (15) respectively. Note that the analytical results are quite close to the simulation results in Fig. 6(a-c), which validates the corresponding assumptions and approximations. However, the error of the analytical results in Fig. 6(d) increases with the sleeping radius. This is because of the conservative approximation that uses Eq. (15) to calculate the outage probability of SC UEs under the repulsive scheme, which exaggerates the inter-cell interference. Nevertheless, the analytical outage probability is a upper bound of the real case, based on which the sub-optimal solution can be derived.

## V. PROBLEM ANALYSIS AND SOLUTIONS

Based on the derived outage constraints in the last section, we analyze how many SCs can be turned off under the random and repulsive schemes respectively.

### A. Optimal Random Scheme

Under the random scheme, the problem is formulated as follows.

$$\max_{p_{\mathrm{s}}, p_{\mathrm{m}}} \quad p_{\mathrm{s}}$$

$$\text{s.t.} \quad \frac{w_{\mathrm{s}}}{1 + \frac{\lambda_{\mathrm{s}}}{\rho_{\mathrm{s}}}} \log_2\left(1 + \tau'(p_{\mathrm{s}})\right) \geq U_{\mathrm{s}}$$

$$\frac{w_{\mathrm{m}}}{1 + \frac{3\sqrt{3}}{2}\lambda_{\mathrm{m}} D^2} \log_2\left(1 + \tau_{\mathrm{m}}\right) \geq U_{\mathrm{m}}$$

$$\frac{W_{\mathrm{m}} - w_{\mathrm{m}}}{1 + \frac{3\sqrt{3}}{2}\lambda_{\mathrm{s}} p_{\mathrm{s}} p_{\mathrm{m}} D^2} \log_2\left(1 + \tau_{\mathrm{o}}(\alpha_{\mathrm{m}}, D)\right) \geq U_{\mathrm{o}} \quad (24)$$

$$\frac{W_{\mathrm{s}} - w_{\mathrm{s}}}{1 + \frac{3\sqrt{3}}{2}\lambda_{\mathrm{s}} p_{\mathrm{s}}(1 - p_{\mathrm{m}}) D^2} \log_2\left(1 + \tau_{\mathrm{o}}(\alpha_{\mathrm{s}}, D)\right) \geq U_{\mathrm{o}}$$

$$p_{\mathrm{m}} \begin{cases} \in (0, 1), & \text{with CB} \\ = 1, & \text{without CB} \end{cases},$$

where $\tau_{\mathrm{m}}$, $\tau'_{\mathrm{s}}(p_{\mathrm{s}})$, and $\tau_{\mathrm{o}}(\alpha, r)$ are given by Eqs. (12), (18), (22) respectively. The equality of the service outage constraints should hold under the optimal solution, and thus all the residual bandwidth can be utilized to turn off more SCs.

#### (1) No CB

When CB is not conducted at MBSs, $p_{\mathrm{m}} = 1$. In this case, the first and fourth conditions are invalid, and the optimal value of this problem is given by

$$p_{\mathrm{s}}^* = \frac{\rho_{\mathrm{m}}}{\lambda_{\mathrm{s}}} \left( \frac{\log_2(1 + \tau_{\mathrm{o}}(\alpha_{\mathrm{m}}, D))}{U_{\mathrm{o}}} \left( W_{\mathrm{m}} - \frac{U_{\mathrm{m}}(1 + \frac{\lambda_{\mathrm{m}}}{\rho_{\mathrm{m}}})}{\log_2(1 + \tau_{\mathrm{m}})} \right) - 1 \right), \quad (25)$$

where $\rho_{\mathrm{m}} = \frac{1}{\frac{3\sqrt{3}}{2} D^2}$.

(a) MBS UEs

(b) SC UEs (random scheme)

(c) offloaded UEs (repulsive scheme: $R_\mathrm{s} = 300$m)
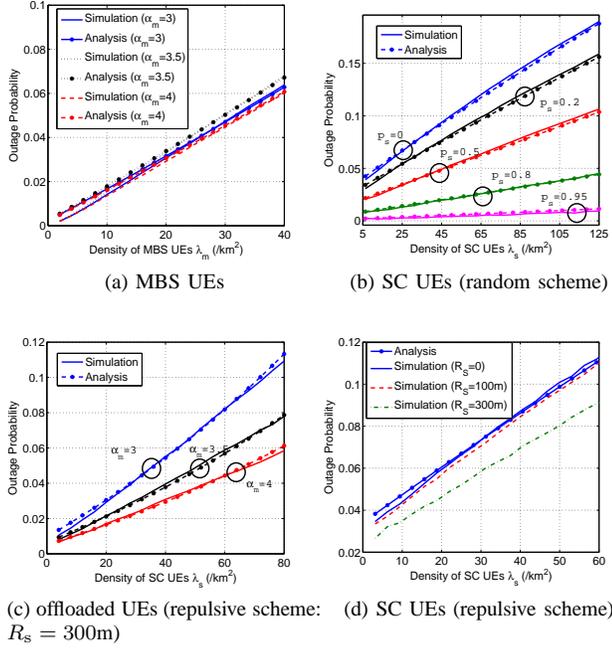
(d) SC UEs (repulsive scheme)

Fig. 6: Comparison of analytical results and simulation results.

***Proposition 1*** The ratio of sleeping SCs is inversely proportional to the density of SC UEs $\lambda_\mathrm{s}$, and decreases linearly with the density of MBS UEs $\lambda_\mathrm{m}$.

***Proposition 2*** Deploying denser MBSs can help to save energy, whereas the density of SCs should be minimized, i.e., just satisfying the peak traffic demand of high data rate service.

Notice that Proposition 1 roughly describes the amount of energy that can be saved due to the dynamics of traffic load in time domain. Recall that the bandwidth for the offloaded UEs decreases linearly with the density of MBS UEs based on Eq. (11), therefore, the capability of the MBS-layer to handle the offloaded UEs also decreases linearly with $\lambda_\mathrm{m}$. Furthermore, the density of the offloaded UE given by $p_\mathrm{s}\lambda_\mathrm{s}$ is limited by the offloading capability of the MBS-layer, thus the sleeping ratio is inversely proportional to the traffic load of the SC-layer.

Proposition 2 offers insights for energy-efficient network deployment. According to the traditional network planning method, the density of MBSs and SCs should both be minimized to meet the peak hour traffic demands. However, this conclusion is inaccurate if SC sleeping and vertical offloading are introduced. The density of active SCs given by $(1 - p_\mathrm{s})^*\rho_\mathrm{s}$ increases linearly with $\rho_\mathrm{s}$ as $p_\mathrm{s}^*$ is irrelevant with $\rho_\mathrm{s}$, which suggests that deploying more SCs only results in denser active SCs consuming more energy. Therefore, the traditional method is still energy-optimal for SC deployment. Whereas, deploying more MBSs helps to turn off more SCs and reduces the energy consumption of the SC-layer. Therefore, the energy-optimal density of MBSs may be larger than the one obtained by traditional method.

*(2) With CB*

When CB is conducted, the optimal solution has no closed-form expression as the first, third and fourth constraints are coupled together by $p_\mathrm{s}$. Only numerical results can be obtained by methods like dichotomy.

To offer some insights, we consider a special case when $\alpha_\mathrm{m} = \alpha_\mathrm{s}$ (the same path loss factor of the two layers) and $p_\mathrm{s} \geq \hat{p}_\mathrm{s}$ (noise-limited case when many SCs are turned off). Under this condition, the three constraints can be decoupled and the optimal solution is

$$p_\mathrm{s}'^* = \frac{\rho_\mathrm{m}}{\lambda_\mathrm{s}} \left( \frac{\log_2(1 + \tau_\mathrm{o}(\alpha_\mathrm{m}, D))}{U_\mathrm{o}} \left( W_\mathrm{m} - \frac{U_\mathrm{m}(1 + \frac{\lambda_\mathrm{m}}{\rho_\mathrm{m}})}{\log_2(1 + \tau_\mathrm{m})} \right. \right.$$
$$\left. \left. + W_\mathrm{s} - \frac{U_\mathrm{s}(1 + \frac{\lambda_\mathrm{s}}{\rho_\mathrm{s}})}{\log_2(1 + \tau_\mathrm{s}'(\hat{p}_\mathrm{s}))} \right) - 1 \right). \tag{26}$$

Then, the performance gain by introducing CB is given by

$$p_\mathrm{s}'^* - p_\mathrm{s}^* = \frac{\rho_\mathrm{m}}{\lambda_\mathrm{s}} \left( \frac{\log_2(1 + \tau_\mathrm{o}(\alpha_\mathrm{m}, D))}{U_\mathrm{o}} \left( W_\mathrm{s} - \frac{U_\mathrm{s}(1 + \frac{\lambda_\mathrm{s}}{\rho_\mathrm{s}})}{\log_2(1 + \tau_\mathrm{s}'(\hat{p}_\mathrm{s}))} \right) \right), \tag{27}$$

which is inversely proportional to the density of SC UEs and increases linearly with the redundant bandwidth at SC-layer.

***Proposition 3*** CB is more beneficial when the traffic load of the SC-layer is lower.

***Proposition 4*** Denser networks (with larger $\rho_\mathrm{m}$ and $\rho_\mathrm{s}$) will benefit more from conducting CB.

Denser SCs and lower traffic load of SC-layer both means more redundant bandwidth at SC-layer, and thus conducting CB will bring higher capacity gains for the offloaded UEs. Similarly, networks with denser MBSs can make more use of the borrowed bandwidth to increase the network capacity.

### B. Optimal Repulsive Scheme

The problem under the repulsive scheme can be formulated as follows.

$$\max_{R_\mathrm{s}, p_\mathrm{m}} \quad \pi R_\mathrm{s}^2 \rho_\mathrm{m}$$
$$\text{s.t.} \quad \frac{w_\mathrm{s}}{1 + \frac{\lambda_\mathrm{s}}{\rho_\mathrm{s}}} \log_2(1 + \tau_\mathrm{s}) \geq U_\mathrm{s}$$
$$\frac{w_\mathrm{m}}{1 + \frac{3\sqrt{3}}{2}\lambda_\mathrm{m}D^2} \log_2(1 + \tau_\mathrm{m}) \geq U_\mathrm{m}$$
$$\frac{W_\mathrm{m} - w_\mathrm{m}}{1 + \pi R_\mathrm{s}^2 \lambda_\mathrm{s} p_\mathrm{m}} \log_2(1 + \tau_\mathrm{o}(\alpha_\mathrm{m}, R_\mathrm{s})) \geq U_\mathrm{o} \quad (28)$$
$$\frac{W_\mathrm{s} - w_\mathrm{s}}{1 + \pi R_\mathrm{s}^2 \lambda_\mathrm{s}(1 - p_\mathrm{m})} \log_2(1 + \tau_\mathrm{o}(\alpha_\mathrm{s}, R_\mathrm{s})) \geq U_\mathrm{o}$$
$$p_\mathrm{m} \begin{cases} \in (0, 1), & \text{with CB} \\ = 1, & \text{without CB} \end{cases}.$$

The four service outage constraints should take equality under the optimal solution, and we have

$$w_\mathrm{s} = \frac{U_\mathrm{s}(1 + \frac{\lambda_\mathrm{s}}{\rho_\mathrm{s}})}{\log_2(1 + \tau_\mathrm{s})}, w_\mathrm{m} = \frac{U_\mathrm{m}(1 + \frac{\lambda_\mathrm{m}}{\rho_\mathrm{m}})}{\log_2(1 + \tau_\mathrm{m})}. \tag{29}$$

*(1) Without CB*

When CB is not conducted, the optimal solution $R_s^*$ should satisfy

$$\frac{W_m - \frac{U_m(1+\frac{\lambda_m}{\rho_m})}{\log_2(1+\tau_m)}}{1+\pi R_s^{*2}\lambda_s p_m} \log_2\left(1+\tau_o(\alpha_m, R_s^{*2})\right) = U_o, \qquad (30)$$

based on which the numerical results of Eq. (28) can be obtained. Note that Proposition 2 still holds as $R_s^*$ is irrelevant with $\rho_s$ whereas increases with $\rho_m$.

The upper bound of the optimal value of problem Eq. (28) can be derived by rewriting the constraint of the offloaded UEs as

$$\frac{P^m}{\sigma^2(1+I)}\frac{\alpha_m+2}{2}\eta_c \geq \left(2^{\frac{(1+\pi R_s^2\lambda_s)U_o}{W_m - w_m}}-1\right)R_s^{\alpha_m}$$

$$\overset{(a)}{\geq} \frac{U_o\ln 2}{W_m - w_m}(1+\pi R_s^2\lambda_s)R_s^{\alpha_m} > \frac{U_o\ln 2}{W_m - w_m}\pi\lambda_s R_s^{\alpha_m+2}, \qquad (31)$$

where (a) applies the inequality $e^x - 1 \geq x (x \geq 0)$. Then the upper bound of SC sleeping ratio is $\pi\tilde{R}_s^2\rho_m$, with the maximum sleeping radius $\tilde{R}_s$ given by

$$\tilde{R}_s = \left(\frac{(\alpha_m+2)\eta_o P^m}{2\sigma^2(1+I)U_o\pi\lambda_s}\left(W_m - \frac{U_m(1+\frac{\lambda_m}{\rho_m})}{\log_2(1+\tau_m)}\right)\right)^{\frac{1}{\alpha_m+2}}. \qquad (32)$$

Specially, this bound is quite tight as $U_o/(W_m - w_m) \to 0$.

***Proposition 5*** The repulsive scheme is more benificial for the heavily loaded networks.

Eq. (32) indicates that the ratio of sleeping SCs under the repulsive scheme is inversely proportional to $\rho_s^{\frac{2}{\alpha_m+2}}$. As $\alpha_m \in (2,4]$, the performance of the repulsive scheme is less sensitive to the variation of the traffic load compared with the random scheme, which explains Proposition 5. In fact, the advantages of the repulsive scheme mainly comes from the shorter distance and smaller path loss of the offloaded UEs. However, this advantage degrades as $R_s$ increases. That is why the repulsive scheme is more advantageous for the heavily loaded networks where few SCs can be turned off and the sleeping radius $R_s$ is small.

*(2) With CB*

If CB is supported and $\alpha_m = \alpha_s$, the available bandwidth for the offloaded UEs $w_o = W_m + W_s - w_m - w_s$, where $w_m$ and $w_s$ are given by Eq. (29). In this case, the upper bound of the sleeping radius becomes

$$\tilde{R'}_s = \left(\frac{(\alpha_m+2)\eta_o P^m}{2\sigma^2(1+I)U_o\pi\lambda_s}(W_m+W_s-w_m-w_s)\right)^{\frac{1}{\alpha_m+2}}. \qquad (33)$$
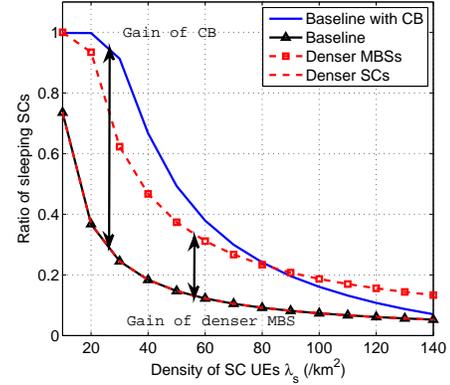
Notice that Proposition 3 and Proposition 4 also hold for the repulsive scheme.

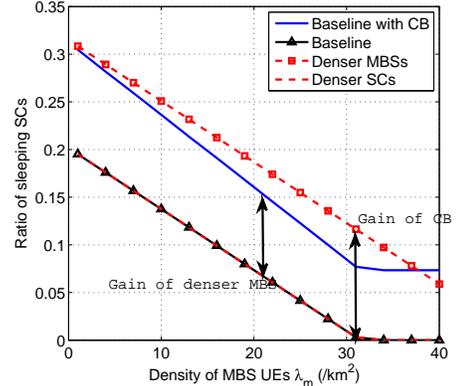### C. Numerical Results and Analysis

In this part, we analyze the relationship between the maximum sleeping ratio and system parameters (traffic load and BS density) through the optimal solutions of Problems (24) and (28) obtained by dichotomy. $\alpha_m$ is set to 3.5, and other parameters can be found in Table I.

TABLE II: Network density

| Parameter | Baseline | Denser SCs | Denser MBSs |
|---|---|---|---|
| $D$ | 500m | 500m | 400m |
| $\rho_s$ | 25/km$^2$ | 50/km$^2$ | 25/km$^2$ |



(a) Varying density of SC UEs $\lambda_s$ ($\lambda_m = 20$/km$^2$)



(b) Varying density of MBS UEs $\lambda_m$ ($\lambda_s = 100$/km$^2$)

Fig. 7: Maximum ratio of sleeping SC under the random scheme.

The maximum sleeping ratio versus traffic load $\lambda_s$ and $\lambda_m$ under random scheme is presented in Fig. 7, where the influences of BS density and CB are both considered. The settings of network density is listed in Table II. Besides, the results of the repulsive scheme are demonstrated in Fig. 8 Notice that the black solid lines with triangles are totally overlapped by the red dashed lines in Fig. 7 and Fig. 8.

*(1) Traffic load*

Generally, the sleeping ratio decreases with traffic load, whereas the slopes of the curves are quite different. Under the random scheme, the sleeping ratio is inversely proportional to the density of SC UEs in Fig. 7(a), whereas it shows linear relation with the density of MBS UEs in Fig. 7(b). Under the repulsive scheme, the sleeping ratio decreases more slowly with $\lambda_s$ (Fig.8(a)), but there is a rapid decline when the traffic load of the MBS is high (Fig.8(b)). These results are consistent with the analytical results of Eqs. (25),(26),(32), and (33).
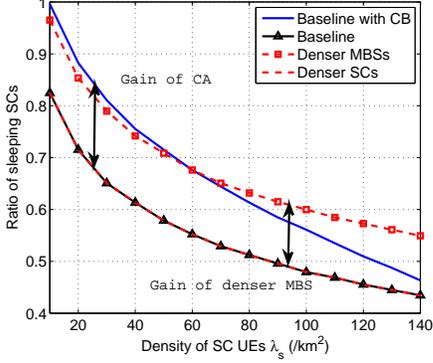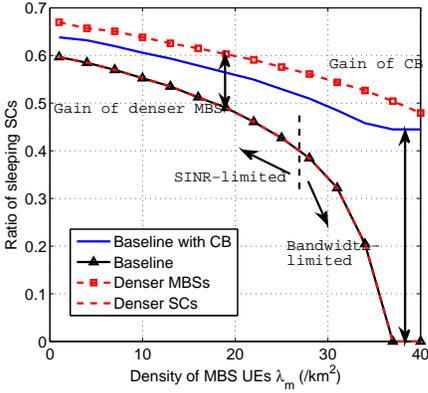
Notice that the sleeping ratio under the repulsive scheme

(a) Varying density of SC UEs $\lambda_s$ ($\lambda_m = 20$/km$^2$)



(b) Varying density of MBS UEs $\lambda_m$ ($\lambda_s = 100$/km$^2$)

Fig. 8: Maximum ratio of sleeping SC under the repulsive scheme.



Fig. 9: Comparison of two schemes ($\lambda_m = 200$/km$^2$).

baseline BS density are compared as shown in Fig. 9. Generally, the sleeping ratio decreases more slowly with $\lambda_s$ under the repulsive scheme, which is consistent with Proposition 5. When CB is not supported, repulsive scheme performs better than the random scheme. However, the random scheme is more advantageous than the repulsive scheme when $\lambda_m$ is smaller and CB is conducted. This can be explained from two aspects: (1) The offloaded UEs enjoy higher spectrum efficiency under the repulsive scheme due to the smaller path loss, especially for the heavily loaded networks (Proposition 5); (2) Due to the lower inter-cell interference level, the SC-layer can provide more residual bandwidth for the offloaded UEs through CB under the random scheme, especially when the network is lightly-loaded and more SCs are turned off.

Based on the above analysis and numerical results, we summarize our findings as follows.

1) Under the random scheme, the ratio of sleeping SCs is inversely proportional to $\lambda_s$ while decreases linearly with $\lambda_m$;
2) Without CB, the repulsive scheme performs better;
3) If CB is conducted, repulsive scheme performs better when the traffic load is high, otherwise the random scheme is better;
4) Deploying more MBSs can help to improve network energy efficiency.
5) CB brings higher performance gain for the networks with heavily-loaded MBS-layer and lightly-loaded SC-layer.

## VI. PERFORMANCE EVALUATION UNDER DAILY TRAFFIC PROFILES

In this part, we evaluate the performance of the two SC sleeping schemes under daily traffic profiles. Two typical traffic patterns are considered as shown in Fig. 10, where the x-axis denotes time and y-axis denotes the density of active UEs. We assume 80% UEs require high data rate service, while the others are served at low data rate. Traffic pattern 1 describes the daily traffic variations of places like bus stations, whose two peaks correspond to the rush hours.
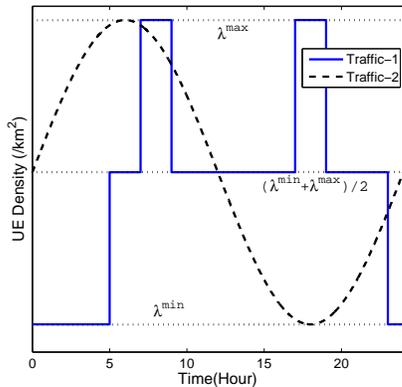
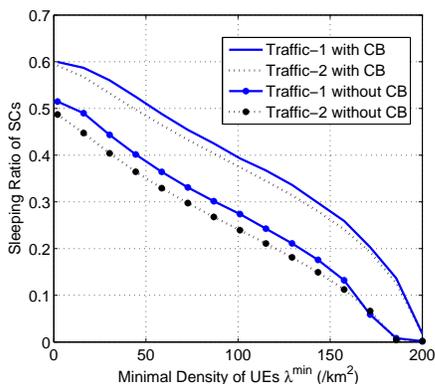can be divided into two cases: (1) high sleeping ratio with large sleeping radius, when the received SINR of the offloaded UEs is relatively low and the sleeping ratio increases slowly with available spectrum resource (SINR-limited); (2) low sleeping ratio with small sleeping radius, when the received SINR is high and the sleeping ratio mainly depends on the available spectrum resource (bandwidth-limited). When $\lambda_m$ is high, increasing bandwidth for the offloaded UEs can significantly improve the sleeping ratio as shown in Fig. 8(b).

### (2) BS density

As the black solid curves with triangles are completely overlapped by the red dashed curves in Figs. 7 and 8, deploying denser SCs does not improve network performance (consistent with Proposition 2).

### (3) Influence of CB

The gain brought by CB is marked in Figs. 7 and 8. Fig. 7(a) and Fig. 8(a) both reflect that CB brings higher performance gain when $\lambda_s$ is small, which validates Proposition 3. In addition, CB can greatly improve the sleeping ratio under the repulsive scheme when the performance is bandwidth-limited (high $\lambda_m$), as shown in Fig. 8(b).

Furthermore, the results of the two schemes under the

Fig. 10: Daily traffic profiles.



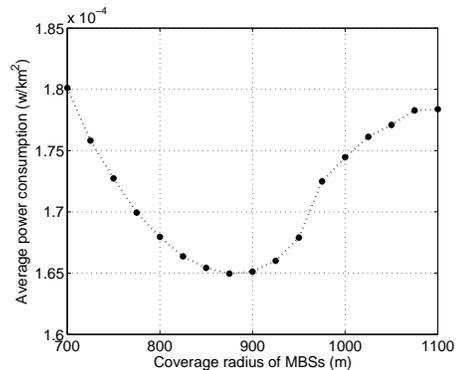Fig. 11: Average sleeping ratio of SCs.



Fig. 12: Network power consumption for different MBS density.

on-demand service is provided through flexible SC sleeping and traffic offloading.

Furthermore, we consider a network planning case when the maximal and minimal user densities are $\lambda^{\max} = 50/\text{km}^2$ and $\lambda^{\min} = 0.5/\text{km}^2$. To guarantee the basic coverage and service, the coverage radius of the MBSs should be no larger than 1100m. Fig. 12 demonstrates the average power consumption per unit area for different MBS coverage radius, under traffic pattern 2. It can be found that the energy-optimal coverage radius of the MBSs is 900m instead of 1100m, which indicates that deploying more MBSs can help to save energy. Although deploying denser MBSs causes higher energy consumption of the MBS-layer, more SCs can be turned off via vertical offloading, which helps to reduce energy consumption of the SC-layer. Thus there exists a tradeoff relation between the energy consumption of the two layers, and the energy-optimal MBS density may be higher than its minimal value required for basic coverage. The problem of energy-optimal network planning is left for our future work, due to the space limitations.

## VII. CONCLUSIONS

In this paper, the expected sleeping ratio of SCs with the time-varying traffic is obtained for two SC sleeping schemes (random and repulsive schemes), under which the UEs of the sleeping SCs are vertically offloaded to the MBS-layer for outage probability guarantee. Under the random scheme, the ratio of sleeping SCs is inversely proportional to the density of SC UEs whereas decreases linearly with the density of MBS UEs. Compared with the random scheme, the repulsive scheme performs better when the traffic load exceeds certain thresholds. In addition, the analytical results suggest that deploying more MBSs can help to save energy by offering more opportunities for SC sleeping. Furthermore, if the spectrum resource can be dynamically borrowed between layers, more SCs can be turned off especially when the MBS-layer is heavily-loaded. Numerical results show that half of the SCs can be turned off on average under two typical daily traffic profiles, and 10% more SCs can further go into sleep if channel borrowing is allowed. For future work, the energy-optimal network planning should be

Besides, traffic pattern 2 with sine function is often used to evaluate the effectiveness of energy saving algorithms [20] [22], [35] and [36].

For any given traffic load $\lambda(t)$, the maximum ratio of sleeping SCs under the two schemes can be obtained by solving Problems (24) and (28), and we dynamically choose the better one which turns off more SCs. The average ratio of sleeping SCs is shown in Fig. 11 with parameters of Table I. The peak traffic load of both traffic profiles is set as $\lambda^{\max} = 200/\text{km}^2$ (equivalent to the network capacity), whereas the minimal traffic load $\lambda^{\min}$ is set as a varying parameter reflecting the non-uniformity of the traffic load in time domain. Notice that different values of $\lambda^{\min}$ reflect different traffic non-uniformity.

The ratio of sleeping SCs decreases with $\lambda_{\min}$, which indicates traffic dynamics help to turn off more SCs. As the network power consumption increases linearly with the density of active SCs without power control, this also suggests that networks with more fluctuated traffic can save more energy through SC sleeping. For the real systems, the traffic load after mid-night usually goes to zero. In this case, about 50% SCs can be turned off on average without CB, and 10% more SCs can go into sleep if CB is conducted. Besides, more SCs are expected to be turned off if effective SC sleeping algorithms are designed. Hence, HCN is an energy-efficient network architecture, under which

analyzed based on the obtained sleeping ratio. In addition, detailed energy-efficient SC sleeping schemes should be designed, where the vertical and horizontal offloading schemes can be jointly optimized, and more realistic network scenarios (such as non-uniform traffic load) should be considered.

## APPENDIX A
## PROOF OF THEOREM 1

The inter-cell interference varies with the UE location, which makes the outage probability unsolvable. To derive closed-form expression, we average the inter-cell interference of all UE locations to approximate the uncertain value. Consider a MBS$_i$, the average inter-cell interference $\bar{\mathcal{I}}$ is

$$\bar{\mathcal{I}} = \int_{a \in \mathcal{A}_i} \sum_{j \in \mathcal{B}^{\mathrm{m}}, j \neq i} P_j^{\mathrm{m}} d_{aj}^{-\alpha_{\mathrm{m}}} f_{\mathcal{A}}(a) \mathrm{d}a, \tag{34}$$

where $\mathcal{A}_i$ is the coverage area of MBS$_i$, and $a$ is the locations of UEs. Denote by $I = \bar{\mathcal{I}}/\sigma^2$ the ratio of inter-cell interference to noise for simplicity.

Assume the number of residual UEs in the target MBS cell is given as $N_{\mathrm{m}}$, then the probability that the data rate requirement of UE$_u$ can be satisfied is given by

$$\mathbb{P}\left\{\gamma_u^{\mathrm{m}} \geq 2^{(N_{\mathrm{m}}+1)\frac{U_{\mathrm{m}}}{w_{\mathrm{m}}}} - 1\right\}$$

$$= \int_0^D \mathbb{P}\left\{h_u^{\mathrm{m}} \geq \frac{(I+1)\sigma^2}{P^m} d^{\alpha_{\mathrm{m}}}\left(2^{(N_{\mathrm{m}}+1)\frac{U_{\mathrm{m}}}{w_{\mathrm{m}}}} - 1\right)\right\} \frac{2d}{D^2}\mathrm{d}d$$

$$= \int_0^D \exp\left(-\frac{(I+1)\sigma^2}{P^m} d^{\alpha_{\mathrm{m}}}\left(2^{(N_{\mathrm{m}}+1)\frac{U_{\mathrm{m}}}{w_{\mathrm{m}}}} - 1\right)\right) \frac{2d}{D^2}\mathrm{d}d \tag{35a}$$

$$= \int_0^D \left(1 - \frac{(I+1)\sigma^2}{P^m} d^{\alpha_{\mathrm{m}}}\left(2^{(N_{\mathrm{m}}+1)\frac{U_{\mathrm{m}}}{w_{\mathrm{m}}}} - 1\right)\right) \frac{2d}{D^2}\mathrm{d}d \tag{35b}$$

$$= 1 - \frac{2D^{\alpha_{\mathrm{m}}}}{\alpha_{\mathrm{m}}+2} \frac{(I+1)\sigma^2}{P^m}\left(2^{(N_{\mathrm{m}}+1)\frac{U_{\mathrm{m}}}{w_{\mathrm{m}}}} - 1\right),$$

where MBS UEs are approximated to be uniformly distributed within a circle of radius $D$. Eq. (35a) holds as $h_u^{\mathrm{m}}$ follows exponential distribution, and (35b) is due to the assumption $\frac{\sigma^2}{P_{\mathrm{m}}} \to 0$. Although (35b) does not hold for the strong interference case (i.e., $I \to \infty$), (35b) applies to the MBS-layer, where most MBS UEs can receive high signal to interference ratio due to the large coverage radius.

Recall that the probability distribution of $N_{\mathrm{m}}$ follows Poisson distribution. By substituting Eq. (35) into Eq. (4), the outage probability of a typical MBS UE is:

$$1 - G_{\mathrm{m}} = \sum_{N=0}^{\infty} \mathbb{P}\left(\gamma_u^{\mathrm{m}} \geq 2^{(N+1)\frac{U_{\mathrm{m}}}{w_{\mathrm{m}}}} - 1\right) P_{N_{\mathrm{m}}}(N)$$

$$= \sum_{N=0}^{\infty} \mathbb{P}\left(\gamma_u^{\mathrm{m}} \geq 2^{(N+1)\frac{U_{\mathrm{m}}}{w_{\mathrm{m}}}} - 1\right) \frac{(\frac{3\sqrt{3}}{2}D^2\lambda_{\mathrm{m}})^N}{N!} e^{-\frac{3\sqrt{3}}{2}D^2\lambda_{\mathrm{m}}}$$

$$= 1 - \frac{2D^{\alpha_{\mathrm{m}}}(I+1)\sigma^2}{P^m(\alpha_{\mathrm{m}}+2)}\left(2^{\frac{U_{\mathrm{m}}}{w_{\mathrm{m}}}} \exp\left(\frac{3\sqrt{3}}{2}D^2\lambda_{\mathrm{m}}\left(2^{\frac{U_{\mathrm{m}}}{w_{\mathrm{m}}}} - 1\right)\right) - 1\right). \tag{36}$$

Hence, Theorem 1 is proved.

## APPENDIX B
## PROOF OF THEOREM 2

When no SCs go into sleep, the distribution of the received SINR of a typical SC UE$_u$ in the interference-limited networks is given by [32]:

$$\int_0^\infty \mathbb{P}\{\bar{\gamma}_u^{\mathrm{s}} \geq T\} f_{d_{\mathrm{s}}}(d)\mathrm{d}d$$

$$\approx \frac{1}{1 + T^{\frac{2}{\alpha_{\mathrm{s}}}} \int_{T^{-\frac{2}{\alpha_{\mathrm{s}}}}}^{\infty} \frac{1}{1+x^{\frac{\alpha_{\mathrm{s}}}{2}}}\mathrm{d}x}$$

$$\geq \frac{1}{1 + T^{\frac{2}{\alpha_{\mathrm{s}}}} \int_{T^{-\frac{2}{\alpha_{\mathrm{s}}}}}^{\infty} x^{-\alpha/2}\mathrm{d}x} \tag{37a}$$

$$= \frac{1}{1 + \frac{2}{\alpha_{\mathrm{s}}-2}T}.$$

and the equality of (37a) holds when $T \to 0$.

Therefore, the probability that the data rate requirement of UE$_u$ can be satisfied is given by

$$\mathbb{P}\left\{\gamma_u^{\mathrm{s}} \geq 2^{(N_{\mathrm{s}}+1)\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}}} - 1\right\}$$

$$\approx \left(1 + \frac{2}{\alpha_{\mathrm{s}}-2}\left(2^{(N_{\mathrm{s}}+1)\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}}} - 1\right)\right)^{-1} \tag{38a}$$

$$= \frac{\frac{\alpha_{\mathrm{s}}-2}{2} 2^{-(N_{\mathrm{s}}+1)\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}}}}{\left(\frac{\alpha_{\mathrm{s}}-2}{2} - 1\right) 2^{-(N_{\mathrm{s}}+1)\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}}} + 1}$$

$$= \frac{(\alpha_{\mathrm{s}}-2)}{2} 2^{-\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}}} \sum_{m=0}^{\infty} \left(\frac{4-\alpha_{\mathrm{s}}}{2}\right)^m 2^{-m} 2^{-m\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}}} 2^{-(m+1)N_{\mathrm{s}}\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}}} \tag{38b}$$

where Eq. (38a) holds for the assumption $\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}} \to 0$. As the path loss factor generally satisfies $\alpha_{\mathrm{s}} \in (2,4]$, $1 - \frac{\alpha_{\mathrm{s}}-2}{2} \in [0,1)$ and Eq. (38b) holds.

Furthermore, as $N_{\mathrm{s}}$ follows Poisson distribution with parameter $\lambda_{\mathrm{s}}A_{\mathrm{s}}$, we have

$$\sum_{N=0}^{\infty} \mathbb{P}\left\{\gamma_u^{\mathrm{s}} \geq 2^{(N+1)\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}}} - 1\right\} P_{N_{\mathrm{s}}}(N)$$

$$= \sum_{N=0}^{\infty} \frac{(\lambda_{\mathrm{s}}A_{\mathrm{s}})^N}{N!} \exp\{-\lambda_{\mathrm{s}}A_{\mathrm{s}}\} \mathbb{P}\left\{\gamma_u^{\mathrm{s}} \geq 2^{(N+1)\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}}} - 1\right\}$$

$$\approx \frac{(\alpha_{\mathrm{s}}-2)}{2} 2^{-\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}}} \sum_{m=0}^{\infty} \left(\frac{4-\alpha_{\mathrm{s}}}{2} 2^{-\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}}}\right)^m \exp\left\{\lambda_{\mathrm{s}}A_{\mathrm{s}}\left(2^{-(m+1)\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}}} - 1\right)\right\}. \tag{39}$$

Recall that $A_{\mathrm{s}}$ follows Gamma distribution with shape $K$ and scale $\frac{1}{K\rho_{\mathrm{s}}}$:

$$f_{A_{\mathrm{s}}}(A) = A^{K-1} \exp\{-K\rho_{\mathrm{s}}A\} \rho_{\mathrm{s}}^K \frac{K^K}{\Gamma(K)}, \tag{40}$$

then we have

$$\int_0^\infty \sum_{N=0}^{\infty} \mathbb{P}\left\{\gamma_u^{\mathrm{s}} \geq 2^{(N+1)\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}}} - 1\right\} P_{N_{\mathrm{s}}}(N) f_{A_{\mathrm{s}}}(A)\mathrm{d}A$$

$$= \frac{(\alpha_{\mathrm{s}}-2)}{2} 2^{-\frac{U_{\mathrm{T}}}{w_{\mathrm{T}}}} \sum_{m=0}^{\infty} \left(\frac{4-\alpha_{\mathrm{T}}}{2} 2^{-\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}}}\right)^m \left(\frac{\rho_{\mathrm{s}}K}{\rho_{\mathrm{s}}K + (1-2^{-(m+1)\frac{U_{\mathrm{s}}}{w_{\mathrm{s}}}})\lambda_{\mathrm{s}}}\right)^K. \tag{41}$$

Finally, due to the well-known exponential limit,

$$\lim_{\frac{U_s}{w_s}\to 0}\left\{\frac{\rho_s K}{\rho_s K+\left(1-2^{-(m+1)\frac{U_s}{w_s}}\right)\lambda_s}\right\}^K$$

$$=\lim_{\frac{U_s}{w_s}\to 0}\left\{1+\frac{1}{\frac{1}{\left(1-2^{-(m+1)\frac{U_s}{w_s}}\right)\lambda_s}\rho_s K}\right\}^{-K}$$

$$=\lim_{\frac{U_s}{w_s}\to 0}\exp\left\{-\frac{\left(1-2^{-(m+1)\frac{U_s}{w_s}}\right)\lambda_s}{\rho_s}\right\} \quad (42)$$

$$=\exp\left\{-(m+1)\frac{U_s}{w_s}\frac{\lambda_s}{\rho_s}\log 2\right\}=2^{-(m+1)\frac{U_s}{w_s}\frac{\lambda_s}{\rho_s}}.$$

By substituting (42) into (41), we have

$$\int_0^\infty\sum_{N=0}^\infty\mathbb{P}\left\{\gamma_u^s\ge 2^{(N+1)\frac{U_s}{w_s}}-1\right\}P_{Ns}(N)f_{A_s}(A)\mathrm{d}A$$

$$\approx\frac{(\alpha_s-2)}{2}2^{-\frac{U_s}{w_s}(1+\frac{\lambda_s}{\rho_s})}\sum_{m=0}^\infty\left\{\left(\frac{\alpha_s-2}{2}-1\right)2^{-\frac{U_s}{w_s}(1+\frac{\lambda_s}{\rho_s})}\right\}^m \quad (43)$$

$$=\frac{\frac{\alpha_s-2}{2}2^{-\frac{U_s}{w_s}(1+\frac{\lambda_s}{\rho_s})}}{1-\frac{4-\alpha_s}{2}2^{-\frac{U_s}{w_s}(1+\frac{\lambda_s}{\rho_s})}}.$$

Hence, Theorem 2 is proved.

## REFERENCES

[1] J. Andrews, et al., "What will 5G be?", *IEEE J. Sel. Areas Commun.*, vol. 32, no. 6, pp. 1065-1082, May, 2013.

[2] X. Ge, H. Cheng, M. Guizani, T. Han, "5G wireless backhaul networks: challenges and research advances," *IEEE Network*, vol. 28, no. 6, pp. 6-11, Nov., 2014.

[3] N. Zhang, N. Cheng, A. Gamage, K. Zhang, J. Mark, and X. Shen, "Cloud assisted HetNets toward 5G wireless networks", *IEEE Commun. Magazine*, to appear.

[4] L. Suarez, L. Nuaymi, and J. Bonnin, "An overview and classification of research approaches in green wireless networks," *EURASIP J. Wireless Commun. and Netw., Special Issue: Green Radio*, Apr., 2012.

[5] X. Ge, X. Huang, Y. Wang, M. Chen, Q. Li, T. Han and C. Wang, "Energy efficiency optimization for MIMO-OFDM mobile multimedia communication systems with QoS constraints," *IEEE Transactions on Vehicular Technology*, vol. 63, no. 5, pp. 2127-2138, June 2014.

[6] Z. Niu, "TANGO: traffic-aware network planning and green operation", *IEEE Wireless Commun.*, vol. 18, pp. 22-29, Oct. 2011.

[7] Z. Niu, S. Zhou, S. Zhou, X. Zhong, and J. Wang, "Energy efficiency and resource optimized hyper-cellular mobile comminication system architecture and its technical challenges," *Scientia Sinica(Informationis)*, vol. 42, no. 10, pp. 1191-1203, 2012.

[8] A. Capone, A. F. dos Santos, I. Filippini, and B. Gloss, "Looking beyond green cellular networks," in *IEEE WONS'12*, Courmayeur, Italy, Jan. 2012.

[9] H. Ishii, Y. Kishiyama, and H. Takahashi, "A novel architecture for LTE-B C-plane/U-plane split and phantom cell concept," in *IEEE GLOBECOM'12*, Anaheim, CA, USA, pp. 624-630, 2012.

[10] 3GPP TR 36.842 V0.2.0, "Study on small cell enhancements for E-UTRA and E-UTRAN - Higher layer aspects," 2013.

[11] X. Xu, G. He, S. Zhang, Y. Chen, and S. Xu, "On functionality separation for green mobile networks: concept study over LTE," *IEEE Commun. Mag.*, vol. 51, no. 5, pp. 82-90, 2013.

[12] S. Chen, J. Zhao, "The requirements, challenges, and technologies for 5G of terrestrial mobile telecommunication," *IEEE Commun. Mag.*, vol. 52, no. 5, pp. 36-43, 2014.

[13] B. Bangerter, S. Talwar, R. Arefi, and K. Stewart, "Networks and devices for the 5G era," *IEEE Commun. Mag.*, vol. 52, no. 5, pp. 90-96, 2014.

[14] G. Auer, et al., "D2.3: energy efficiency analysis of the reference systems, areas of improvements and target breakdown," INFSO-ICT-247733 EARTH, Tech. Rep., Nov., 2010. [Online]. Available: www.ict-earth.eu/publications/ deliverables/deliverables.html.

[15] Z. Niu, Y. Wu, J. Gong, and Z. Yang, "Cell zooming for cost-efficient green cellular networks," *IEEE Commun. Mag.*, vol. 48, no. 11, pp. 74-79, Nov. 2010.

[16] S. Zhou, J. Gong, Z. Yang, Z. Niu and P. Yang, "Green mobile access network with dynamic base station energy saving," in *ACM MobiCom'09*, Beijing, China, Sept. 2009.

[17] M. Marsan, L. Chiaraviglio1, D. Ciullo1, and M. Meo, "Optimal energy savings in cellular access networks," in *IEEE ICC'09*, Dresden, Germany, June, 2009.

[18] J. Gong, S. Zhou, Z. Niu, and Y. Peng, "Traffic-aware base station sleeping in dense cellular networks," in *IEEE WiOpt'10*, Avignon, France, June, 2010.

[19] L. Chiaraviglio, D. Ciullo, M. Meo, and M. Marsan, "Energy-aware UMTS access networks," in *WPMC'08*, Lapland, Finland, Sept. 2008.

[20] Y. Wu, and Z. Niu, "Energy efficient base station deployment in green cellular networks with traffic variations," in *IEEE ICCC'12*, Beijing, China, Aug. 2012.

[21] Y. Soh, T. Quek, M. Kountouris, and H. Shin, "Energy efficient heterogeneous cellular networks," *IEEE J. Sel. Areas Commun.*, vol. 31, no. 5, pp. 840-850, May, 2013.

[22] S. Zhang, Y. Wu, S. Zhou, and Z. Niu, "Traffic-aware network planning and green operation with BS sleeping and cell zooming," in *IEICE Trans. Commun.* vol. E97-B, no. 11, pp. 2337-2346, Nov. 2014.

[23] D. Cao, S. Zhou, Z. Niu, "Optimal combination of base station densities for energy-efficient two-tier heterogeneous cellular networks," *IEEE Trans. Wireless Commun.*, vol. 12, no. 9, pp. 4350-4362, Sept. 2013.

[24] S. Cho, and W. Choi, "Energy-efficient repulsive cell activation for heterogeneous cellular networks," *IEEE J. Sel. Areas Commun.*, vol. 31, no. 5, pp. 870-882, May, 2013.

[25] S. Morosi, P. Piunti, and E. Re, "Improving cellular network energy efficiency by joint management of sleep mode and transmission power," in *IEEE TIWDC'13*, Genoa, Italy, Sept. 2013.

[26] S. Morosi, P. Piunti, and E. Re, "Sleep mode management in cellular networks: a traffic based technique enabling energy saving, "*Trans. Emerging Tel. Tech.*, vol. 24, issue 3, pp. 331-341, 2013. [Online], Available: http://onlinelibrary.wiley.com/doi/10.1002/ett.2621/abstract

[27] G. Wang, C. Guo, S. Wang, and C. Feng, "A traffic prediction based sleeping mechanism with low complexity in femtocell networks,", in *IEEE ICC'13*, Budapest, Hungary, June, 2013.

[28] L. Saker, S. Elayoubi, R. Combes, T. Chahed, "Optimal control of wake up mechanisms of femtocells in heterogeneous networks," *IEEE J. Sel. Areas Commun.*, vol. 30, no. 2, pp. 664-672, Apr., 2012.

[29] S. Zhang, J. Wu, J. Gong, S. Zhou, and Z. Niu, "Energy-optimal probabilistic base station sleeping under a separation network architecture," in *IEEE GLOBECOM'14*, Austin, USA, Dec., 2014.

[30] Z. Wang, and W. Zhang, "A separation architecture for achieving energy-efficient cellular networking," *IEEE Trans. Wireless Commun.*, vol. 13, no. 6, pp. 3113-3123, June, 2014.

[31] S. Zhang, S. Zhou, and Z. Niu, "Traffic aware offloading for BS sleeping in heterogeneous networks," in *Asilomar Conference on Signals, Systems, and Computers'15*, California, USA, Nov., 2014.

[32] J. Andrews, F. Baccelli, and R. Ganti, "A tractable approach to coverage and rate in cellular networks," *IEEE Trans. Wireless Commun.*, vol. 51, no. 11, pp. 3122-3134, Nov. 2011.

[33] S. N. Chiu, D. Stoyan, W. S. Kendall, and J. Mecke, *Stochastic Geometry and Its Applications*, 3rd Edition, UK: John Wiley & Sons, 2013, pp. 48-51.

[34] M. Heanggi, J. Andrews, F. Baccelli, O. Dousse, and M. Franceschetti, "Stochastic geometry and random graphs for the analysis and design of wireless networks," invited paper, *IEEE J. Sel. Areas Commun.*, vol. 27, no. 7, pp. 1029-1046, Sept., 2009.

[35] M. Marsan, L. Chiaraviglio, D. Ciullo, and M. Meo, "Energy-efficient managment of UMTS access networks," in *IEEE ITC'09*, Paris, France, Sept., 2009.

[36] L. Chiaraviglio, D. Ciullo, G. Koutitas, M. Meo, and L. Tassiulas, "Energy-efficient planning and management of cellular networks," in *IEEE WONS'12*, Courmayeur, Italy, Jan., 2012.
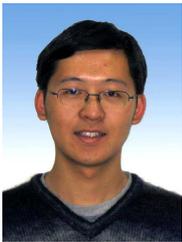
**Shan Zhang** received B.S. degree in Electronic Engineering from Beijing Institute Technology, Beijing, China, in 2011 and is currently a Ph.D. candidate in Department of Electronic Engineering, Tsinghua University. She received the Best Paper Award from the 19th Asia-Pacific Conference on Communication (APCC) in 2013. Her research interests include network planning, resource and traffic management for green communications.

**Jie Gong** received the B.S. and Ph.D. degrees from Tsinghua University, Beijing, China, in 2008 and 2013, respectively. He is currently a Postdoctoral Scholar with the Department of Electronic Engineering, Tsinghua University. From July 2012 to January 2013, he visited the Institute for Digital Communications, The University of Edinburgh, Edinburgh, U.K. His research interests include base station cooperation in cellular networks, energy harvesting, and green wireless communications. Dr. Gong was a co-recipient of the Best Paper Award from the IEEE Communications Society Asia Pacific Board in 2013.

**Sheng Zhou** (S'06, M'12) received his B.S. and Ph.D. degrees in Electronic Engineering from Tsinghua University, China, in 2005 and 2011, respectively. He is now a postdoctoral scholar in Electronic Engineering Department at Tsinghua University, Beijing, China. From January to June 2010, he was a visiting student at Wireless System Lab, Electrical Engineering Department, Stanford University, CA, USA. He is a co-recipient of the Best Paper Award from the 15th Asia-Pacific Conference on Communication (APCC) in 2009, and the 23th IEEE International Conference on Communication Technology (ICCT) in 2011. His research interests include cross-layer design for multiple antenna systems, cooperative transmission in cellular systems, and green wireless cellular communications.

**Zhisheng Niu** (M'98-SM'99-F'12) graduated from Northern Jiaotong University (currently Beijing Jiaotong University), Beijing, China, in 1985, and got his M.E. and D.E. degrees from Toyohashi University of Technology, Toyohashi, Japan, in 1989 and 1992, respectively. After spending two years at Fujitsu Laboratories Ltd., Kawasaki, Japan, he joined with Tsinghua University, Beijing, China, in 1994, where he is now a professor at the Department of Electronic Engineering, deputy dean of the School of Information Science and Technology, and director of Tsinghua-Hitachi Joint Lab on Environmental Harmonious ICT. He is also a guest chair professor of Shandong University. His major research interests include queueing theory, traffic engineering, mobile Internet, radio resource management of wireless networks, and green communication and networks. Dr. Niu has been an active volunteer for various academic societies, including Director for Conference Publications (2010-11) and Director for Asia-Pacific Board (2008-09) of IEEE Communication Society, Membership Development Coordinator (2009-10) of IEEE Region 10, Councilor of IEICE-Japan (2009-11), and council member of Chinese Institute of Electronics (2006-11). He is now a distinguished lecturer (2012-13) of IEEE Communication Society, standing committee member of both Communication Science and Technology Committee under the Ministry of Industry and Information Technology of China and Chinese Institute of Communications (CIC), vice chair of the Information and Communication Network Committee of CIC, editor of IEEE Wireless Communication Magazine. Dr. Niu received the Outstanding Young Researcher Award from Natural Science Foundation of China in 2009 and the Best Paper Awards (with his students) from the 13th and 15th Asia-Pacific Conference on Communication (APCC) in 2007 and 2009, respectively. He is now a fellow of both IEEE and IEICE.