

Bounds on QoS-Constrained Energy Savings in Cellular Access Networks with Sleep Modes

Balaji Rengarajan*, Gianluca Rizzo* and Marco Ajmone Marsan*[†]

{balaji.rengarajan, gianluca.rizzo, marco.ajmone}@imdea.org

*Institute IMDEA Networks, Madrid, Spain

[†]Politecnico di Torino, Italy

Abstract—Sleep modes are emerging as a promising technique for energy-efficient networking: by adequately putting to sleep and waking up network resources according to traffic demands, a proportionality between energy consumption and network utilization can be approached, with important reductions in energy consumption. Previous studies have investigated and evaluated sleep modes for wireless access networks, computing variable percentages of energy savings. In this paper we characterize the *maximum* energy saving that can be achieved in a cellular wireless access network under a given performance constraint. In particular, our approach allows the derivation of realistic estimates of the energy-optimal density of base stations corresponding to a given user density, under a fixed performance constraint. Our results allow different proposals to be measured against the maximum theoretically achievable improvement. We show, through numerical evaluation and simulation, the possible energy savings in today’s networks, and we further demonstrate that even with the development of highly energy-efficient hardware, a holistic approach incorporating system level techniques is essential to achieving maximum energy efficiency.

I. INTRODUCTION

The ethical imperative to reduce their carbon footprint, combined with the financial realities of increasing energy costs, and the difficulties of network deployment in developing countries with unreliable power grids, has telecommunication network operators keenly interested in energy saving approaches. In cellular networks, reducing the power consumed by base stations is, by far, the most effective means to streamline energy consumption. As an example, in the case of UMTS, one typical Node-B consumes around 1500 W, and the multitude of these devices accounts for between 60 and 80% of the network’s energy consumption [1], [2], often representing the main component of an operator’s operational expenditures.

The bulk of the research on energy savings in wireless network was initially focused on the case of ad-hoc and sensor networks [3], and in the context of hand-held, battery operated devices [4], [5]. Not much attention was paid until recently to reducing energy consumption of base stations, since these were assumed to rely on access to a reliable supply of energy with acceptable cost. Both assumptions are challenged in the networking context of today. While equipment manufacturers are working to produce more energy-efficient hardware [6], as we show, system-level approaches are called for, to obtain networks with the lowest possible energy consumption.

Base stations are deployed according to dimensioning strategies that ensure acceptable user performance at peak (worst-

case) traffic loads. However, traffic loads fluctuate throughout the day. For example, we expect diurnal patterns in the rate of user requests that mirror human patterns. Additionally, as the users of the network move during the day, they cause fluctuations in the spatial traffic load seen by base stations serving different locations. In [1] and [7], the possibility of reducing power consumption in cellular networks by reducing the number of active cells in periods of low traffic was considered, but the degradation in performance experienced by users in such a scenario, due to active base stations having to serve larger numbers of users that are located farther away from their serving base station was not explicitly taken into account. However, an important requirement for any energy saving measure, such as the introduction of sleep modes for base stations, is that they must be (almost) transparent to users. This means that the user-perceived performance must be above the target threshold at peak hours, when the load on the network is the highest, and all base stations are active, as well as in non-peak periods, when the load is lower, but the network is operating with reduced resources. In other words, the performance sacrifices that are implied by the introduction of energy-saving measures must be compatible with the target design objectives.

Recently, heuristics have been proposed to turn off base stations to conserve energy [8]. Approaches to vary cell sizes through changing the base station transmit power and in the limit turning off base stations have also been proposed [9], [10]. However, to the best of our knowledge, the maximal energy savings that can be achieved under some predefined performance constraint are not known.

In this paper, our objective is to obtain a realistic characterization of the potential energy savings that can be achieved by sleep mode schemes under fixed user performance constraints, and study the impact of base station topology, power consumption model, and user density on the energy-optimal configuration of the access network. The metric we use to capture performance is the *per-bit delay* [11] (whose inverse approximates the throughput) perceived by a typical best-effort user. The network is constrained to maintain, at all times, the average per-bit delay across users below a predetermined threshold. This focus on user-perceived performance is one of the key contributions of this paper.

Our contributions are as follows:

- For a given base station topology, we develop a method

for estimating the density of base stations that minimizes energy consumption and which is sufficient to serve a given set of active users, with fixed performance guarantees.

- For base stations whose power consumption is independent of load (not unlike current hardware), we derive a topology-independent lower bound on the density of base stations required to support a particular user density and thus an upper bound on energy savings.
- Through numerical evaluation and simulations, we compute bounds on the maximum energy saving and illustrate the impact of various system parameters. We demonstrate that even with highly energy efficient hardware, system level techniques are crucial to minimizing energy consumption. We find that the variability in performance across users is sufficiently low, validating the choice of the mean of the per-bit delay as a suitable metric for capturing user performance.

Our results are *bounds* with respect to what can be achieved in real networks, since we assume that *any base station density is achievable*, although this is clearly not possible in practice, since in real networks base stations can be turned off, but their locations cannot be rearranged according to traffic variations. The relevance of our bounds lies in that they indicate what are the theoretical minimum base station densities and energy consumption, allowing the effectiveness of different proposals to be measured against the maximum theoretically achievable improvement.

The paper is organized as follows. In Sec. II, we present our model for the distribution of users and of base stations, and we state the main assumptions underlying our approach. In Sec. III, we derive the average and the variance of the per-bit delay. In Sec. IV, we use the results of the previous sections to compute the energy-optimal base stations density for a given user density, and to estimate the achievable energy savings. Sec. V presents lower bounds on the base station densities required to satisfy the performance constraints. In Sec. VI, we present numerical and simulation results, and we conclude the paper in Sec. VII.

II. MODEL AND ASSUMPTIONS

We consider the downlink information transfer in a cellular access network. Users form a homogeneous planar Poisson point process, Π_u , with intensity λ_u users per square km, while base stations form a planar point process, Π_b , with density λ_b base stations per square km. While the methodology introduced in this paper is quite general, and can be extended to many different base station configurations, we restrict ourselves to the following models for base station distribution across the service area:

- Manhattan layout: base stations lie on the vertices of a square grid, where the side of each square is $l_b = \frac{1}{\sqrt{\lambda_b}}$ km.
- Hexagonal layout: base stations lie at the centers of a hexagonal tessellation of side $l_H = \left(\frac{2}{3\sqrt{3}\lambda_b}\right)^{\frac{1}{2}}$ km.

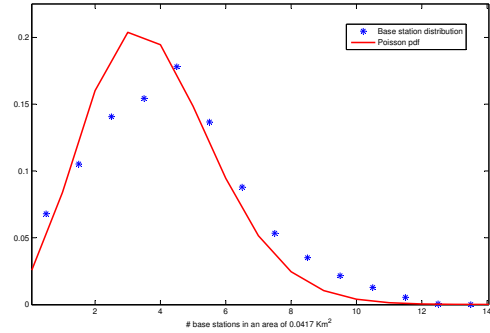


Fig. 1. Empirical distribution of the number of base stations in a rectangular area of downtown Sydney (AU), and Poisson distribution with equal average.

- Poisson layout: base stations are distributed over the service area according to a two-dimensional homogeneous Poisson point process with density λ_b .

The first two distributions above reflect regular topologies used for the analysis and design of cellular networks, while the third reflects the result of real life constraints on the base station locations. For example, we examined the distribution of the base stations operated by an important international operator in the bay area of Sydney, Australia [12]. The area we chose is densely populated, with an average base station density of 81.64 base stations per square km, and is a good candidate for reducing the density of active base stations in periods of low load. Fig. 1 displays the empirically determined distribution of the number of base stations within a randomly centered $124 \text{ m} \times 336 \text{ m}$ rectangle, along with a Poisson pdf with an expected value matching the average number of base stations found within the rectangle. While the Poisson pdf is not an exact fit, it reasonably approximates the variability introduced by practical constraints on base station location.

We assume that all base station densities are feasible. In the case of the Manhattan and hexagonal layout of base stations, since only a subset of existing base stations can be turned off, only a discrete subset of densities corresponding to those that maintain the structure of the topology can be achieved. However, in the homogeneous Poisson process layout of base stations, if each base station independently makes a decision to either turn off, or stay on, according to some probability, the resulting point process of base stations is a thinned homogeneous Poisson process, and all base station densities are indeed achievable.

The end user performance metric that we use is per-bit delay of best effort data transfers.

Definition 2.1 (Per-bit delay): The per-bit delay, τ , that a user perceives is defined as the inverse of the user throughput, i.e., the actual rate at which the user is served, taking into account both the capacity to the user, as well as the sharing of the base station time across all associated users.

Both the average and the variance of the per-bit delay will be computed and used as performance metrics in this paper. The performance constraint that is enforced is as follows: if the per-bit delay experienced by a *typical* user, $\bar{\tau}$, is less than a predefined threshold $\bar{\tau}^0$ seconds, then users are said to perceive

satisfactory performance, and the corresponding base station distribution is feasible. Here, the interpretation of a typical user is that provided by Palm theory [13], and $\bar{\tau}$ is computed as the expectation of τ with respect to the Palm distribution P^0 associated with Π_u . Intuitively, the Palm distribution is the conditional distribution given that there is a point belonging to Π_u at the origin. The variance of the per-bit delay allows the characterization of the spread of the performance perceived by different end users at a given time instant. It should be however observed that user mobility makes the performance of each individual user vary over time, reducing variance across users in the long run. For this reason, we just use the average as a performance constraint, but we also observe the variance, in order to verify that performance differences across users remain acceptable.

A. Channel and Service Model

In this paper, we do not consider the effect of shadowing and only take into account distance-dependent path loss. We assume that users are served by the base station that is closest to them, i.e., by the base station that corresponds to the strongest received signal, as it normally happens in reality. Denote by $S(x)$, the location of the base station that is closest to a user located at x , and by $D(x)$ the distance between the user and the closest base station. The number of active users associated with base station $S(x)$ is denoted $N(S(x))$. We denote the capacity to a user located at a distance r from the base station by $C(r)$ bit/s per Hertz. The capacity can be modeled, for example, using Shannon's capacity law or other models such as a quantized set of achievable rates. In this paper, we focus on the case where the network only serves best-effort traffic. The analysis can be extended to the case of a mixture of best-effort and delay-sensitive traffic, however we do not report the results here due to space limitations. We assume that base stations use a processor sharing mechanism to divide capacity among all the connected best-effort users. By doing so, a notion of fairness is imposed, since all best effort users associated with a particular base station are served for an identical fraction of time.

B. Energy Consumption Model

We assume that base stations always transmit at a fixed transmit power. When the base station density is higher than that required to achieve the threshold expected per-bit delay $\bar{\tau}^0$, we assume that base stations only serve users for the fraction of time required to satisfy the performance constraint, and remain idle (i.e., not transmitting to any user) for the rest. We denote with U the utilization of base stations, i.e., U is the average fraction of time in which the base station is transmitting.

We model the power in watts consumed by a base station as $k_1 + k_2 U$, where k_1 is the power consumed by keeping a base station turned on with no traffic, and k_2 is the rate at which the power consumed by the base station increases with the utilization. The first energy model that we study reflects the current base station design, and assumes that the bulk of the

energy consumption at the base stations is accounted for by just staying on, while the contribution to energy consumption due to base station utilization is negligible (i.e., $k_2 = 0$). We also study energy consumption models with k_1 and k_2 chosen to reflect a more energy-proportional scenario i.e., $k_1 \ll k_2$.

III. MODELING USER PERCEIVED PERFORMANCE

We characterize the per-bit delay perceived by a typical best-effort user who is just beginning service, as a function of the density of users and base stations under the different base station topologies.

Theorem 3.1: The average per-bit delay, $\bar{\tau}$, perceived by a typical best-effort user joining the system when the density of base stations is λ_b and the density of users is λ_u , is given by:

- Hexagonal layout:

$$\bar{\tau}_H = 6\lambda_u \int_0^{\left(\frac{1}{2\sqrt{3}\lambda_b}\right)^{\frac{1}{2}}} \int_{-\frac{y}{\sqrt{3}}}^{\frac{y}{\sqrt{3}}} \frac{1}{C(\sqrt{x^2 + y^2})} dx dy \quad (1)$$

- Manhattan layout:

$$\bar{\tau}_M = \lambda_u \int_{-\frac{1}{2\sqrt{\lambda_b}}}^{\frac{1}{2\sqrt{\lambda_b}}} \int_{-\frac{1}{2\sqrt{\lambda_b}}}^{\frac{1}{2\sqrt{\lambda_b}}} \frac{1}{C(\sqrt{x^2 + y^2})} dx dy \quad (2)$$

- Poisson layout:

$$\bar{\tau}_P = \frac{\int_0^\infty \left(\int_0^\infty \int_0^{2\pi} e^{-\lambda_b A(r,x,\theta)} \lambda_u x d\theta dx \right) e^{-\lambda_b \pi r^2} \lambda_b 2\pi r dr}{C(r)} \quad (3)$$

where $A(r, x, \theta)$ is the area of the circle centered at (x, θ) with radius x that is not overlapped by the circle centered at $(0, -r)$ with radius r .

Proof Sketch: We leverage Slivnyak's theorem [13], and derive a formula for the mean per-bit delay experienced by adding a point at the origin to Π_u . The mean per-bit delay depends on the capacity at which the user at the origin can be served, which in turn depends on the distance between the user and the serving base station (the one that is closest to the origin). Further, the per-bit delay perceived by any user is affected by the number of users that share the serving base station. The mean per-bit delay experienced by the user at the origin can be computed as:

$$E^0[\tau] = E^0 \left[\left(\frac{C(D(0))}{N(S(0))} \right)^{-1} \right] = E^0 \left[\frac{N(S(0))}{C(D(0))} \right]. \quad (4)$$

Here, E^0 denotes the expectation with respect to the Palm distribution associated with Π_u . A detailed proof including the formula to compute $A(r, x, \theta)$ is in Appendix A. ■

Further, we characterize the variance in the user-perceived per-bit delay through the following theorem.

Theorem 3.2: The variance of the per-bit delay, $(\sigma)^2$, perceived by a typical best-effort user joining the system when the density of base stations is λ_b and the density of users is λ_u , is given by:

- Hexagonal layout:

$$(\sigma_H)^2 = -(\bar{\tau}_H)^2 + (6\lambda_u + 9\sqrt{3}l_H^2(\lambda_u)^2) \int_0^{\frac{\sqrt{3}}{2}l_H} \int_{-\frac{y}{\sqrt{3}}}^{\frac{y}{\sqrt{3}}} \frac{1}{(C(\sqrt{x^2+y^2}))^2} dx dy \quad (5)$$

- Manhattan layout:

$$(\sigma_M)^2 = -(\bar{\tau}_M)^2 + \left(\lambda_u + \frac{(\lambda_u)^2}{\lambda_b} \right) \int_{-\frac{1}{2\sqrt{\lambda_b}}}^{\frac{1}{2\sqrt{\lambda_b}}} \int_{-\frac{1}{2\sqrt{\lambda_b}}}^{\frac{1}{2\sqrt{\lambda_b}}} \frac{1}{(C(\sqrt{x^2+y^2}))^2} dx dy \quad (6)$$

- Poisson layout:

$$(\sigma_P)^2 = \int_0^\infty \left[\left(\int_0^\infty \int_0^{2\pi} e^{-\lambda_b A(r,x,\theta)} \lambda_u d\theta dx \right)^2 + \int_0^\infty \int_0^{2\pi} e^{-\lambda_b A(r,x,\theta)} \lambda_u d\theta dx \right] \frac{e^{-\lambda_b \pi r^2} \lambda_b 2\pi r}{C(r)^2} dr - (\bar{\tau}_P)^2. \quad (7)$$

Proof Sketch: The proof is similar to that of Theorem 3.1, and additionally makes use of the fact that the users form a Poisson point process, thus

$$\text{Var}^0[N(S(0))] = E^0[N(S(0))].$$

We skip the detailed proof due to space limitations. ■

IV. OPTIMIZING BASE STATION ENERGY CONSUMPTION

In the case of the energy model with $k_2 = 0$, energy consumption is minimized by using the lowest base station density that can achieve the desired user performance. Given λ_u and λ_b , the per-bit delay perceived by a typical user can be evaluated using the results from Sec. III. $E^0[\bar{\tau}]$ is decreasing in λ_b . Thus, we can set the expressions equal to the target per-bit delay, $\bar{\tau}^0$, to determine the minimum required base station density λ_b^* . For this case, that approximates current base station power consumption trends, we determine lower bounds for the required base station density and thus energy consumption, irrespective of base station distribution, in the following section.

When $k_1 \ll k_2$, the utilization of the base stations in the network plays a key role in determining the energy consumed. Again, $\bar{\tau}$ can be evaluated given λ_u and λ_b using the results from Sec. III. In this case, it is easy to see that the desired user performance can be achieved by the base stations only actively serving best-effort users for a time fraction $\frac{\bar{\tau}}{\bar{\tau}^0}$ of the time originally used, provided that $\bar{\tau} < \bar{\tau}^0$. If, instead, $\bar{\tau} > \bar{\tau}^0$, the base station density λ_b cannot meet the performance constraint. Thus, the base station serving the typical user will be serving actively for a time fraction $\frac{\bar{\tau}}{\bar{\tau}^0}$. From this, we can calculate the energy consumed in order to satisfy the performance constraint at any feasible base station density. By inspection, we can then determine the base station density that minimizes energy consumption.

V. A LOWER BOUND ON BS DENSITY

Clearly, the density of base stations required to support a particular population of users depends on the geometry of the base station layout. In this section, we determine a lower bound on the base station density required to achieve the target average per-bit delay across all base station distributions. This lower bound corresponds to the base station density that minimizes energy consumption in the case of the energy model with $k_2 = 0$.

Theorem 5.1: A lower bound on the minimum density of base stations sufficient to serve a population of users with density λ_u with an average per-bit delay $\bar{\tau}^0$ is given by λ_b^* that satisfies

$$\bar{\tau}^0 = 2\pi\lambda_u \int_0^{\frac{1}{\sqrt{\lambda_b^* \pi}}} \frac{1}{C(r)} r dr \quad (8)$$

Also, there exists a configuration with base station density less than $1.173\lambda_b^*$ that is feasible.

Proof: see Appendix B.

VI. NUMERICAL EVALUATION

In this section we estimate numerically, in some simple scenarios, the potential energy savings that can be obtained by turning off base stations in periods of low load, while still guaranteeing quality of service. Base station transmit power p is assumed to be 30W. Base stations work at a frequency of 1 GHz, and use a bandwidth of 10 MHz. We use a log distance path loss model, with path loss at a reference distance of one meter calculated using Friis equation, and with a path loss exponent $\alpha = 3.5$. We assume that the rate perceived by users is given by Shannon's capacity law. Thus, the capacity to a user located at a distance r from the base station is given by $C(r) = 10^7 \log_2 \left(1 + \frac{pr^{-\alpha}}{N_0} \right)$ bit/s, where $N_0 = -174$ dBm/Hz is the power spectral density of the additive white Gaussian noise. However, the maximum rate at which the base station can transmit data is limited to 55 Mbps.

We considered different choices for the parameters of the base stations energy model while always keeping the total power consumed by a base station with utilization 100% at 1500W. In one setting, the total energy consumption does not vary with the base station utilization. In this setting, we choose $k_1 = 1500$ W and $k_2 = 0$ W, in accordance with typical values found in the literature. We refer to this setting as the *on-off* setting. This choice of parameters approximately models the behavior of base stations currently deployed, in which the dependency of the energy consumed on load is negligible. Moreover, as current trends in base stations design aim at tying power consumption to base station utilization, we considered a few settings in which the energy consumed by a base station depends on the utilization of the base station. These *energy proportional* (EP) settings allow us to examine how strategies for turning off base stations could evolve in the future. We distinguish them by the ratio $\frac{k_2}{k_1+k_2}$ that we use as a metric for energy proportionality. For instance, a setting with $k_1 = 500$ W and $k_2 = 1000$ W is denoted EP 66.6% and one with $k_1 = 100$ W and $k_2 = 1400$ W is denoted EP 93.4%.

In Fig. 2, we plot the optimal base stations density (i.e. the one that minimizes the average power consumption per Km^2 due to base stations, as described in Section IV) versus user density, for various base stations layouts and energy settings. We also plot the lower bound on base station density obtained as described in Section V.

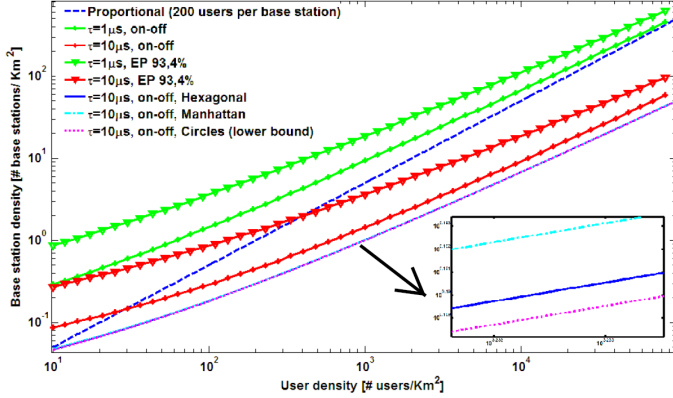


Fig. 2. Energy-optimal base stations density versus user density, for Poisson base stations layout (unless otherwise indicated).

We focus first on the curves that represent the on-off setting. Note that for this setting, energy consumption is directly proportional to base station density. We see that regular layouts (namely, the hexagonal and Manhattan layouts) are the most energy efficient, and they are only slightly worse than the lower bound derived from (8). The Poisson layout consumes more energy due to the variability in cell sizes. As we would expect, decreasing the target average per-bit delay results in layouts with increased base station densities. Fig. 2 also exhibits the base station density corresponding to the case where the number of users per base station is held constant, i.e., a case where base station density is directly proportional to user density. We can see that decreasing base stations density proportionally to user density results in a highly optimistic estimate of energy savings. When user performance constraints are taken into account, actual energy savings are much less.

Under the energy proportional model, the minimum base station density that achieves the target performance is not necessarily the one that minimizes energy consumption. As illustrated in the figure, the base station density that minimizes energy consumption is higher in this case than under the on-off model. This indicates that as hardware becomes increasingly energy proportional, cellular layouts would tend towards higher densities of smaller cells. The effect on energy consumption is discussed later.

We also observe that the gap in the energy-optimal base station density between the on-off energy model and the more energy proportional model decreases with increasing user density. To understand the reason behind this, we refer to Fig. 3. This figure shows that, at the energy-optimal base station density, base station utilization increases with user density. This increase is due to the non-linearly increasing inefficiency in serving users farther and farther away from

the base station. Thus, at higher user densities, base stations tend to operate closer to peak capacity and thus the difference between the two energy models diminishes. Note that the base

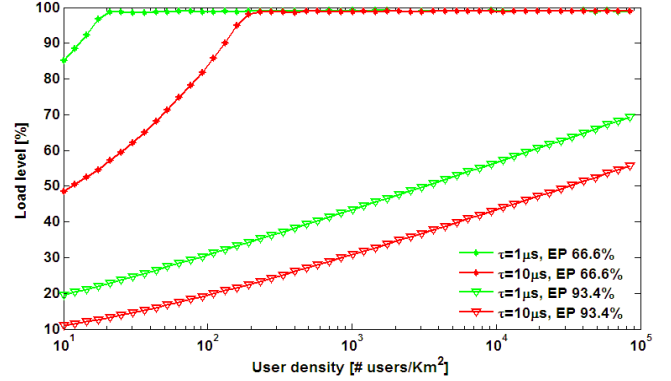


Fig. 3. Average utilization level of base stations at the optimal base stations density versus user density for a Poisson base stations layout.

station utilization under the on-off energy model (not shown) in the energy-optimal base station density is always 100%. For a given user density, this utilization decreases as base stations become increasingly energy proportional, indicating that base station densities increase and cells become smaller.

The amount of energy savings achievable with sleep modes is shown in Fig. 4. For a given energy model and a target average per-bit delay, we consider a network that is optimally planned for a peak user density of 10^5 users per Km^2 , and evaluate the amount of energy that can be saved by switching off base stations in periods of lower user density. We see that, when user density reduces from 10^5 to 10^3 , we can achieve energy savings of up to 95% by reducing accordingly the number of active base stations. Moreover, a reduction of user density by a factor of 10 is already sufficient to save more than 85% on the power consumed at peak load. We can also observe that energy savings exhibit little dependence on either the specific target average per-bit delay, or on the base station energy model.

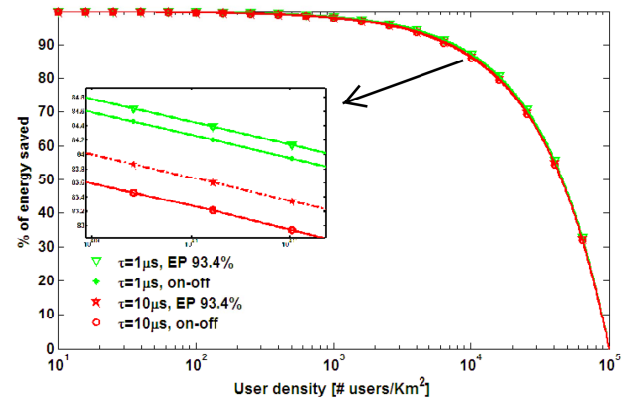


Fig. 4. Percentage of energy saved with sleep modes in a Poisson layout, with respect to the energy consumed at a peak user density of 10^5 users/ Km^2 .

The importance of sleep modes and system level techniques

is evident from Fig. 5, where we plot the average power consumed per square kilometer for the Poisson layout in two cases: i) when sleep modes are used to adapt the base station density to load, and ii) when the network is always provisioned for the peak load, so that power savings are only due to the energy proportionality of the base station power consumption.

We observe that in case i), when sleep modes are used, energy proportional base stations result in a slightly more energy efficient behavior at low user densities, as expected. However, we clearly see that much of the reduction in energy consumption is obtained through the intelligent use of sleep modes to adapt the active base station density to the user population, even in the absence of improved hardware.

On the contrary, in case ii), when sleep modes are not used, and the base station density remains at the level required to support the peak user density, energy proportional base stations do provide large energy savings with respect to current base stations whose power consumption is almost independent of utilization. However, the power consumption at low user densities is up to two orders of magnitude higher in this case with respect to case i), even under highly optimistic (and probably unrealistic) assumptions on energy proportionality. This highlights the need to tackle the problem of energy consumption in cellular access networks through both improved hardware and system level techniques. It also shows clearly that, even under futuristic assumptions on the energy efficiency of hardware, the intelligent use of sleep modes and other dynamic provisioning techniques can be crucial to achieving maximum energy efficiency.

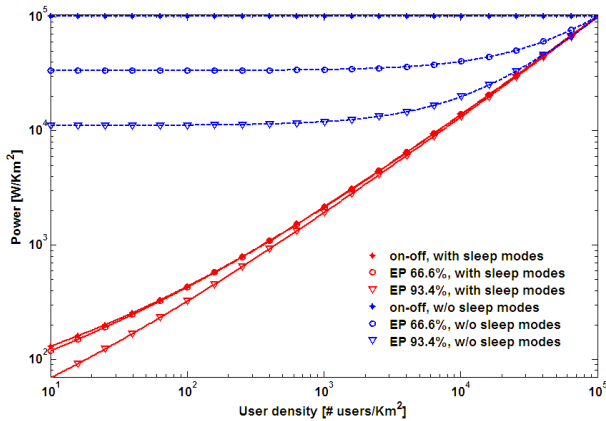


Fig. 5. Minimum power consumed by base stations per Km^2 , as a function of user density. Base stations layout is Poisson, and $\tau = 10\mu s$.

In Fig. 6 on the right y axis, we plot the minimum amount of power consumed per user, and on the left y axis, the optimal number of users per cell, both as a function of user density, for Poisson base station layouts. We observe how the per-user consumed power decreases with increasing user density. At high user densities, cells are small and base stations serve users that are relatively close. Therefore, as path losses are inferior on average, this represent a more energy efficient configuration. Moreover, as user density grows, the number

of users per cell in the energy-optimal configuration increases while the size of the cells decreases. We also note that the slope of these curves is higher at low user densities. This is again due to the inefficiency of serving users farther away from base stations, which increases non-linearly with the size of the cells. The inefficiency of serving low user densities suggests that operators could gain substantially by cooperating and sharing infrastructure in periods of low demand as suggested in [14].

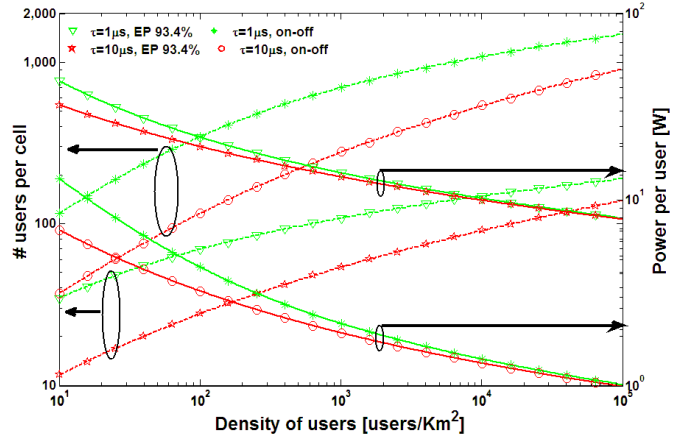


Fig. 6. Right y axis: Minimum amount of power consumed per user. Left y axis: Optimal number of users per cell as a function of user density. Base stations are distributed according to the Poisson layout.

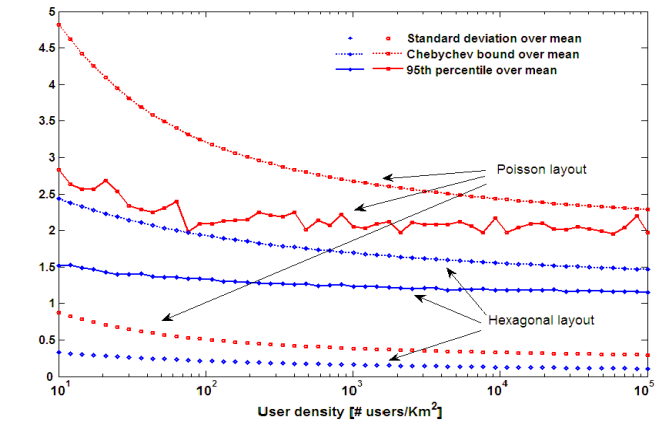


Fig. 7. Standard deviation, 95% Chebyshev bound and 95th percentile of the per-bit delay. All quantities are normalized over an average per-bit delay of $1\mu s$.

In order to validate our analytical results, we have run a large number of simulations, whose results (obviously) are in accordance with the numerical results derived from the formulas presented in this paper. In Fig. 7, we plot the ratio of the standard deviation of the per-bit delay (as derived in Theorem 3.2) to the average, and compare it to the 95th percentile of the per-bit delay derived from simulations, for the on-off energy model. We also plot the bound on the 95th percentile obtained using the Chebyshev bound normalized by the mean per-bit delay. As we can see, in the Poisson layout the

95th percentile is never larger than three times the average, and it does not vary significantly with user density. Also, the ratio of standard deviation and percentiles to the mean is very flat over the range of user densities. The curves for the hexagonal layout show that regular base station layouts translate into less variability in the per-bit delay across users. As these results on variability do not take into account the averaging effect on the user perceived per-bit delay induced by user mobility, we would expect variability in a more realistic situation with user mobility to be lower. Overall, these results suggest that the mean per-bit delay (possibly with a safety margin) is a reasonable design metric for sleep mode algorithms.

VII. CONCLUSIONS

In this paper, we presented a novel approach for estimating both the energy savings that can be achieved in cellular access networks by using sleep modes in periods of low traffic loads as well as the energy-optimal base station densities as a function of user density. By taking into account the quality of service perceived by end users, our approach allows the derivation of more realistic estimates that can be used to evaluate the efficacy of schemes utilizing sleep modes to save energy. The proposed approach can be applied to many base station configurations, and to many energy models for base stations. We demonstrated with numerical and simulation results that substantial energy savings are possible through schemes that adapt the density of base stations to the fluctuations in user density. We also showed that such system level schemes are essential even if base stations themselves become more energy proportional in the future.

We are currently working on extending this approach to mixed traffic scenarios, where voice and video traffic have higher priority than best effort traffic, and to clustered user populations. We also aim at incorporating mobility in our analysis, and investigating the impact that methods such as power control and opportunistic scheduling have on energy consumption.

REFERENCES

- [1] J. T. Louhi, "Energy efficiency of modern cellular base stations," in *29th International Telecommunications Energy Conference (INTELEC)*, Rome, Italy, october 2007, pp. 475–476.
- [2] (2008) Node b datasheets. [Online]. Available: <http://www.motorola.com/>
- [3] R. Zheng and R. Kravets, "On-demand power management for ad hoc networks," in *INFOCOM*, March 2003, pp. 481–491.
- [4] E. Shih, P. Bahl, and M. J. Sinclair, "Wake on wireless: An event driven energy saving strategy for battery operated devices," in *MobiCom*, 2002, pp. 160–171.
- [5] G. Anastasi, M. Conti, E. Gregori, and A. Passarella, "Performance comparison of power saving strategies for mobile web access," *Performance Evaluation*, vol. 53, pp. 273–294, 2003.
- [6] M. Hodes, "Energy and power conversion: A telecommunication hardware vendors perspective," Power Electronics Industry Group, 2007.
- [7] M. A. Marsan, L. Chiaraviglio, D. Ciullo, and M. Meo, "Optimal energy savings in cellular access networks," in *First International Workshop on Green Communications (GreenComm'09)*, June 2009.
- [8] K. Dufkova, M. Bjelica, B. Moon, L. Kencl, and J.-Y. Le Boudec, "Energy Savings for Cellular Network with Evaluation of Impact on Data Traffic Performance," in *European Wireless*, 2010.

- [9] S. Bhaumik, G. Narlikar, S. Chattopadhyay, and S. Kanugovi, "Breathe to stay cool: adjusting cell sizes to reduce energy consumption," in *Proceedings of the first ACM SIGCOMM workshop on Green networking*, 2010, pp. 41–46.
- [10] Z. Niu, Y. Wu, J. Gong, and Z. Yang, "Cell zooming for cost-efficient green cellular networks," *Communications Magazine, IEEE*, vol. 48, no. 11, pp. 74–79, november 2010.
- [11] T. Bonald, "Insensitive traffic models for communication networks," *Discrete Event Dynamic Systems*, vol. 17, pp. 405–421, 2007.
- [12] Australian geographical radiofrequency map. [Online]. Available: <http://www.spench.net>
- [13] D. Stoyan, W. S. Kendall, and J. Mecke, *Stochastic geometry and its applications*. Wiley, 1987.
- [14] M. A. Marsan and M. Meo, "Energy efficient management of two cellular access networks," *SIGMETRICS Perform. Eval. Rev.*, vol. 37, pp. 69–73, March 2010.

APPENDIX

A. Proof of Theorem 3.1

Proof:

1) Hexagonal layout:

$$\begin{aligned}
 \bar{\tau}_H &= \mathbb{E}^0 \left[\frac{N(S(0))}{C(D(0))} \right] = \mathbb{E}^0 [N(S(0))] \mathbb{E}^0 \left[\frac{1}{C(D(0))} \right] \\
 &= \frac{3\sqrt{3}}{2} l_H^2 \lambda_u \int_0^{\frac{\sqrt{3}}{2} l_H} \int_{-\frac{y}{\sqrt{3}}}^{\frac{y}{\sqrt{3}}} \frac{1}{C(\sqrt{x^2 + y^2})} \frac{4}{\sqrt{3} l_H^2} dx dy \\
 &= 6\lambda_u \int_0^{\frac{\sqrt{3} l_H}{2}} \int_{-\frac{y}{\sqrt{3}}}^{\frac{y}{\sqrt{3}}} \frac{1}{C(\sqrt{x^2 + y^2})} dx dy
 \end{aligned}$$

The first step above follows because the size of the hexagons in the tessellation is fixed, and the number of users served by the base station closest to the origin is independent of the distance to the origin, and only depends on the area of a hexagonal cell. The proof for the Manhattan layout follows closely the above methodology.

2) *Poisson layout:* The case where base stations are distributed as a homogeneous Poisson point process is more involved, since the size of the cell that the typical user belongs to is correlated with the distance between the user and the base station. For example, if the closest base station to a user is far away, that base station is likely to be serving a large cell with many users, and vice versa.

In the following, $B(c, r)$ denotes a ball of radius r centered at c .

$$\begin{aligned}
 \bar{\tau}_P &= \mathbb{E}^0 \left[\frac{N(S(0))}{C(D(0))} \right] \\
 &= \int_0^\infty \mathbb{E}^0 \left[\frac{N(S(0))}{C(D(0))} \middle| r \leq D(0) \leq r + dr \right] \\
 &\quad \mathbb{P}(r \leq D(0) \leq r + dr) \\
 &= \int_0^\infty \frac{\mathbb{E}^0 [N(S(0)) | r \leq D(0) \leq r + dr]}{C(r)} \\
 &\quad \mathbb{P}(B(0, r) = \phi) \lambda_b 2\pi r dr \\
 &= \int_0^\infty \frac{\mathbb{E}^0 [N(S(0)) | r \leq D(0) \leq r + dr]}{C(r)} \\
 &\quad e^{-\lambda_b \pi r^2} \lambda_b 2\pi r dr. \tag{9}
 \end{aligned}$$

where $P(B(0, r) = \emptyset)$ is the probability that a ball of radius r centered at the origin is empty.

Now, we turn to deriving the conditional expectation above. The expected number of users attached to the base station serving the user at the origin can be evaluated as follows:

$$\begin{aligned} & \mathbb{E}^0[N(S(0)) | r \leq D(0) \leq r + dr] \\ &= \mathbb{E}^0 \left[\int_0^\infty \int_0^{2\pi} \mathbf{1}_{(S(x, \theta) = S(0) | r \leq D(0) \leq r + dr)} \lambda_u d\theta dx \right] \\ &= \int_0^\infty \int_0^{2\pi} P(S(x, \theta) = S(0) | r \leq D(0) \leq r + dr) \lambda_u d\theta dx, \end{aligned}$$

where $\mathbf{1}_{(S(x, \theta) = S(0))}$ is the indicator function of the event that a user at location (x, θ) is served by same base station that serves the user at the origin.

For the purpose of computing the conditional probability, we assume without loss of generality that the base station closest to the origin is located at $(0, r)$. To evaluate the probability that a user at a given location is served by the same base station that serves a user at the origin, we use a simple change of coordinates, that moves the base station to the origin. In this shifted coordinate system, the typical user placed at the origin is now located at $(0, -r)$. A user at location (x, θ) will also be served by the base station at the origin, if there is no other base station that is closer, i.e., if there is no base station in a circle of radius x centered at (x, θ) . The probability that this is the case, given that there are no base stations in a circle of radius r centered at $(0, -r)$, is given by $\exp(-\lambda_b A(r, x, \theta))$, where $A(r, x, \theta)$ is the area of the circle centered at (x, θ) with radius x that is not overlapped by the circle centered at $(0, -r)$ with radius r . This non-overlapped area can be computed using standard trigonometric identities. Denoting the distance between the centres of the two circles by $d(r, x, \theta) = \sqrt{x^2 + r^2 + 2xr \sin(\theta)}$, we have:

$$\begin{aligned} A(r, x, \theta) &= \pi x^2 - \left[r^2 \arccos \left(\frac{2r^2 + 2xr \sin(\theta)}{2rd(r, x, \theta)} \right) + \right. \\ & x^2 \arccos \left(\frac{2x^2 + 2xr \sin(\theta)}{2xd(r, x, \theta)} \right) \\ & \left. - \frac{1}{2} (-d(r, x, \theta) + r + x)^{\frac{1}{2}} (d(r, x, \theta) + r - x)^{\frac{1}{2}} \right. \\ & \left. (d(r, x, \theta) - r + x)^{\frac{1}{2}} (d(r, x, \theta) + r + x)^{\frac{1}{2}} \right]. \end{aligned}$$

Using the above expression, we obtain

$$\begin{aligned} & \mathbb{E}^0[N(S(0)) | r \leq D(0) \leq r + dr] = \\ & \int_0^\infty \int_0^{2\pi} e^{-\lambda_b A(r, x, \theta)} \lambda_u d\theta dx. \end{aligned} \quad (10)$$

Finally, we obtain the mean per-bit delay experienced by a typical user by substituting expression (10) into (9). Note that this methodology can be applied to other base station layouts as well. ■

B. Proof of Theorem 5.1

First, we examine the case of a single base station and determine the shape of the cell that maximizes the area (users) covered while still satisfying the performance requirements.

Lemma A.1: When capacity to a user is a decreasing function of distance, a base station maximizes the area (number of users) covered while satisfying the performance constraint on per-bit delay by serving an area that is a circle with the base station at the center.

Proof: Consider a maximal service area that satisfies the per-bit delay constraint and is not a circle. There must exist a region at a distance d_1 from the base station that is not included in the service area while another at a distance $d_2 > d_1$ is. Let the average per-bit delay achieved by the maximal service area be $\bar{\tau}^m$. Consider swapping an area of measure ϵ at distance d_2 with an area of the same measure at distance d_1 . The expected per-bit delay for the new service area, $\bar{\tau}^n$ can be calculated as:

$$\bar{\tau}^n = \bar{\tau}^m - \frac{\lambda_u \epsilon}{C(d_2)} + \frac{\lambda_u \epsilon}{C(d_1)}$$

Since $C(d_1) > C(d_2)$, $\bar{\tau}^n < \bar{\tau}^m$. Thus, the new service area satisfies the per-bit delay constraint as well. We can continue this procedure until a region at a distance d' from the base station is included only if all regions at distance $d < d'$ are included. ■

Proof of Theorem 5.1: To determine a lower bound on the density of base stations, we determine r_c^* , the radius of the largest circular service area (users therein) that a single base station can serve while meeting the per-bit delay constraint. The area of this circle corresponds to the maximum area of a cell that satisfies the performance constraint. The density of base stations corresponding to cells of this size provides the lower bound. The expected user-perceived per-bit delay in a circular service area of radius r_c^* can be computed similar to the case of the hexagonal layout as:

$$\bar{\tau}_C = 2\pi \lambda_u \int_0^{r_c^*} \frac{1}{C(r)} r dr, \quad (11)$$

providing the lower bound when $\lambda_b^* = \frac{1}{\pi(r_c^*)^2}$.

Now, consider a hexagonal layout of base stations. If a base station can support users within the circle that superscribes a hexagon, then the base station can clearly support the users in the hexagon. Thus, an upper bound for the density of base stations required in a hexagonal layout, and thus an upper bound on the minimal density of base stations can be computed using the packing density of a hexagonal layout to be: $\lambda_b^U = \left(\frac{3\sqrt{3}(r_c^*)^2}{2} \right)^{-1}$, which proves the tightness result. ■