This is the authors' post-review version of paper:

# A distributed power-saving framework for LTE HetNets exploiting Almost Blank Subframes

Antonio Virdis, Giovanni Stea, Dario Sabella, Marco Caretti

*Abstract*—**Almost Blank Subframes (ABSs) have been defined in LTE as a means to coordinate transmissions in heterogeneous networks (HetNets), composed of macro and micro eNodeBs: the macro issues ABS periods, and refrains from transmitting during ABSs, thus creating interference-free subframes for the micros. Micros report their capacity demands to the macro via the X2 interface, and the latter provisions the ABS period accordingly. Existing algorithms for ABS provisioning usually share resources proportionally among HetNet nodes in a long-term perspective (e.g., based on traffic forecast). We argue instead that this mechanism can be exploited to *save power* in the HetNet: in fact, during ABSs, the macro consumes less power, since it only transmits pilot signals. Dually, the micros may inhibit data transmission themselves in some subframes, and optimally decide *when* to do this based on knowledge of the ABS period. This allows us to define a power saving framework that works *in the short term*, modifying the ABS pattern at the fastest possible pace, serving the HetNet traffic at reduced power cost. Our framework is designed using only standard signaling. Simulations show that the algorithm consumes less power than its competitors, especially at low loads, and improves the UE QoS.**

*Index Terms*—**Computer Networks, Algorithms, Computer Network Management, Computer Network Performance, Integer Programming, Scheduling**

## I. INTRODUCTION

Improving cell-edge performance in urban areas and increasing the energy efficiency of base stations have been among the main issues of LTE-Advanced research in recent years (e.g. [1],[2]). These requirements are met by using denser deployments, using a higher number of base stations, heterogeneous in terms of transmission power, coverage and power consumption: high-power nodes, called *macro*, cover large areas, whereas low-power ones, generally called *micro*, are used to boost capacity or extend coverage in specific zones. Such heterogeneous, dense deployments are called *HetNets*.

To allow efficient operations in HetNets, 3GPP introduced a set of techniques called *enhanced Inter-Cell Interference Coordination* (eICIC), since Rel-9. Notably, two mechanisms have been proposed. A fist one, *cell-range expansion* (CRE) can be used to favor the association of UEs to micro nodes, possibly offloading the macro. As a side effect, UEs at the edge of macro and micro cells may experience poor channel conditions due to interference. For this reason, a second mechanism has been devised, i.e. allowing the macro to inhibit data transmission during certain subframes (SFs), called *Almost Blank SFs*, to reduce interference to cell-edge UEs in micro cells. During ABSs, however, the macro still transmits pilot signals, hence ensuring correct operation of the UEs under its control. ABSs are arranged by the macro in *ABS periods* (APs) of 40 SFs, whose composition is transmitted to all the micros in the HetNet through the X2 interface. A slightly different embodiment of the ABS concept, discussed in [3], consists in having the macro transmit at a reduced power during ABSs. This allows the macro to serve *some* UEs (those nearer to it), possibly with a reduced SINR, and consuming less power. On the other hand, UEs attached to micros will perceive a mitigated interference from the macro, hence will be able to hear their serving micro.

Recent research [1] has shown that a node that only transmits pilot signals – such as a macro during ABSs – consumes *less power* than in a normal SF, since it can power down some circuitry. This creates a power saving opportunity: provisioning an AP so that the macro packs all its data traffic into as few SFs as possible, and then powers down some circuitry and transmits ABSs in the rest of the SFs, allows it to carry the same traffic at a reduced power cost. Moreover, micro nodes may adopt a similar behavior: when they receive an AP, they decide how many SFs they need to use, and of which type (whether ABSs or non-ABSs) to carry their traffic. Therefore, they too can pack their traffic into as few SFs as possible, and only transmit pilot signals in the other ones, thus saving some more power.

Now, both a macro and a micro need to estimate the load that they will carry in the next AP in order to decide when to power down. Since that load may vary in the short term, there

is a clear need for a *dynamic* ABS provisioning framework. The vast majority of the existing schemes (e.g. [4], [5]-[9]) consider a *semi-static* ABS partitioning, based on long-term capacity requirements. This way, they implicitly give up the opportunity of riding the peaks and valleys of the traffic demand at the nodes, which is instead the principal lever for power saving. There is, besides, no practical impediment to having a new, different AP at every AP boundary, if the required computations complete in time. Even the few dynamic schemes appeared so far, however, (e.g., [10]) do not consider power saving opportunities. Moreover, to the best of our knowledge, all schemes presented so far *either* assume that macro transmission is completely inhibited during ABSs, *or* assume low-power data transmission during them, and never consider that switching between *both* solutions, or employing them both simultaneously in the same AP, may improve performance and/or save power.

In this paper, we propose a practical framework for power saving that exploits *dynamic* ABS provisioning. Based on periodic reports from the micros and its own aggregate load and channel quality measures, the macro node selects how many "idle" ABSs (I-ABSs, those where the macro does not transmit any data), how many "low-power" ABSs (LP-ABSs, those where the macro does transmit data at a reduced power), and how many non-ABSs to provision in the next AP, so that all the HetNet load is carried, if at all possible, at the minimum power cost. When the network is overloaded (i.e., the capacity demands of the micros and the macro cannot be accommodated simultaneously), our framework degrades the service proportionally at all the coordinated nodes.

The points of strength of our framework are the following: first, to the best of our knowledge, ours is the first framework to leverage ABS provisioning for power saving. Moreover, it is the only one that takes into account *both* I-ABSs *and* LP-ABSs simultaneously in a unified framework. We show that slight variations in the HetNet scenarios - e.g., in the position or traffic of UEs - such as those that may occur dynamically at the same timescale at which our framework operates, make *different combinations of these mechanisms* preferable from a power saving standpoint. Using both I-ABSs and LP-ABSs allows more energy-efficient transmissions, hence greater power savings or higher throughputs, than using either of the two mechanisms in isolation. It is well known that power saving comes at a cost in terms of increased latency: however, our framework also *minimizes the latency* for a given level of power saving, thus offering the best possible trade-off. Moreover, since our our framework saves power by keeping nodes off as much as possible, it also reduces interference and increases the transmission efficiency, which further reduces latency for a given traffic load. Second, but not less important, our framework employs only *standardized signaling*, i.e. nodes are assumed to know *only* what the standard allows them to about the status of the HetNet. This makes our framework practically implementable. Third, it does not require complex computations: both the macro and the micro nodes

run only *simple algorithms*, that are independent of the number of users and scale well with the number of nodes. Fourth, it works under broad hypotheses, e.g., with arbitrary numbers of micros and network topologies, it accommodates a large class of power consumption models, and does not make hypotheses on the traffic or UE mobility model. Besides, it is orthogonal to algorithms running at both *slower* timescales, such as CRE bias selection for user association or network topology adaptation through selective node switching, and *faster* timescales, such as MAC-level scheduling. These will generate input data to our framework.

We evaluate our framework via simulation against a dynamic ABS provisioning scheme, showing that its power savings are remarkably higher, and that this does not come with any performance degradation: on the contrary, the cell throughput stays the same, and the user delay distribution improves considerably.

The rest of the paper is organized as follows: Section II reports background on LTE. Section III reviews the related work, and Section IV describes the system model. Our power-saving framework is explained in Section V, and Section VI reports performance results. Section VII concludes the paper.

## II. BACKGROUND ON LTE

In this section we describe those features of LTE that are more relevant to the problem at hand, i.e. downlink resource allocation at the MAC layer and the eICIC framework.

In an LTE network, an eNodeB (eNB) allocates resources to its UEs, composing transmission schedules (called *subframes*) periodically. A *cell* is an area where UEs are associated to an eNB, hence share resources taken from the same SF. Scheduling within a cell occurs on every Transmission Time Interval (TTI), whose duration is 1ms, and consists in the eNB allocating a vector of *Resource Blocks* (RBs) to UEs. One RB goes to one UE only (Multi-user Multiple-Input/Multiple-Output (MIMO) techniques are outside the scope of this paper). The number of bytes per RB is indirectly determined by the Channel Quality Indicator (CQI) reported by the UE to which that RB is allocated. The CQI depends on the measured Signal-to-Noise-and-Interference Ratio (SINR), and is reported periodically or on demand. Downlink transmissions are protected by a Hybrid ARQ (H-ARQ) scheme: the UE reports an ACK/NACK, and the eNB may reschedule the same transmission up to four times in a future TTI. H-ARQ errors often originate when the reported channel state (as per the UE CQI) is better than the one at the time of transmission.

The eNBs are often categorized according to their role and radiation power. Large-scale cellular coverage is provided essentially by *macro* eNBs, i.e. high-power, large-coverage eNBs. Lower-power eNBs are normally called *micro*, or *pico*, or *femto* eNBs (depending on their power and possibly other factors), and provide *additional localized* capacity, to increase the UE data rate where needed. For this reason, we will henceforth refer to all of them as *micro* eNBs for simplicity, any dif-
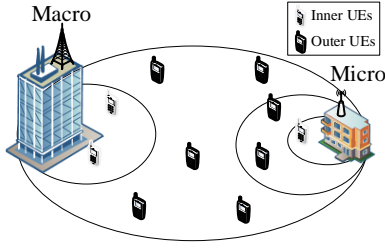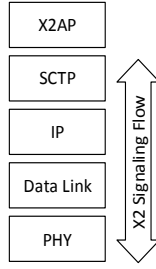
Fig. 1 – Example of HetNet deployment.
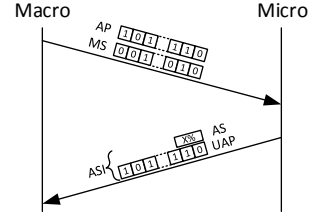


Fig. 2 - Description of the X2 stack.



Fig. 3 - Standard-based message exchange between macro and micro.

TABLE 1 - LIST OF RELATED WORKS AND RESPECTIVE TOPICS

| | [4] | [5] | [6] | [7] | [8] | [9] | [10] | [15] | [16] | [17] | [18] | [20] | [21] | [22] | [23] | [24] | [25] | [26] |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ABS ratio | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | | | ✓ | | | ✓ | | |
| UE association | ✓ | ✓ | | | ✓ | | | ✓ | | ✓ | ✓ | ✓ | ✓ | | | ✓ | | |
| Transmission power | | ✓ | | ✓ | | ✓ | | | | | | | | ✓ | ✓ | ✓ | ✓ | ✓ |

ference between them being immaterial for the purpose of this paper. Micro cells, thus, are embedded within macro cells, and UEs in micro cells suffer interference from the macros (see Fig. 1).

The eNBs communicate with each other – or, possibly, with other network entities – through the so-called X2 interface [11], which includes both a control and a data plane. The protocol stack for the control plane is shown in Fig. 2. Signaling information are generated by the X2 Application Protocol (X2AP) [12], which defines a large set of procedures and messages for supporting inter-eNB operations, including those related to eICIC, which we describe next. The X2 uses the Stream Control Transmission Protocol (SCTP) at layer 4, which establishes and maintains the association between two peering eNBs, and whatever link-level technology is available below the IP layer.

To limit inter-cell interference, the standard [13] allows a macro to provision ABSs, i.e. SFs where it either refrains from transmitting downlink data altogether, limiting to common reference and pilot signals [14], or transmits downlink data at a reduced power, so that micro UEs will experience a better SINR. We have already distinguished *idle ABSs* (I-ABSs), and *low-power ABSs,* (LP-ABSs). The term ABS will denote both when no distinction is needed.

ABSs are organized into an *ABS pattern* (AP), composed of 40 consecutive SFs in an FDD deployment, which is repeated periodically. The AP can be varied at the end of the period by the macro. An AP is sent to the micros within an X2 message, as a vector of binary values, stating whether a SF is/is not an ABS, as shown in Fig. 3. Moreover, the macro sends another bit vector to complement the first, called the *Measurement Subset (MS)*, which indicates to the micros *in which SFs their UEs should measure* interference. A macro may thus instruct a micro to evaluate CQIs in ABSs, to assess whether the interference in these is acceptable or not. On request of their macro node, micro nodes report periodic ABS Status IEs (ASIs), still using the X2 interface (see again Fig. 3). The latter include:

– ABS Status (AS): A value from 0 to 100 stating the *percentage of employed RBs* over the total available on the subset of ABSs defined in the next field. By reporting

AS=100, micros signal that they need more ABSs.
– Usable ABS Pattern Info (UAP): a string of binaries indicating the set of ABSs over which the previous percentage is computed. This must be a subset of the ABSs of the last AP, possibly their whole set, including ABSs where the measured interference is below a configurable value.

## III. RELATED WORK

The study of eICIC techniques has received considerable attention in the last few years, and many research papers are being published on the topic. Almost all works on eICIC pursue one or more of the following objectives: i) selecting an ABS ratio; ii) influencing or selecting the UEs association to macro and micro nodes, and iii) selecting the transmission power used by either or both the macro and the micros. Table 1 matches some recent works with the above topics.

More in detail, works that select the ABS ratios either operate at a network-wide scale, i.e. coordinate many macros together, or at the macro-cell scale, by coordinating *one* macro with its micros. Our work can be framed within this last stream. In [15] instead, a hybrid approach is defined: groups of macros coordinate with the same target micro, so that each micro UE is protected from the highest-interfering macro. Unlike ours, however, most of these works assume that the ABS ratios can only be chosen within a small, predefined set (e.g. 1/8, 1/4, etc.), and/or are configured statically or semi-statically, at a pace of tens of minutes [4]. The only works we found that tackle the selection of the ABS dynamically are [10] and [16]. Work [10] defines an algorithm to select the ABS ratio in an AP based on the number of resources requested by the macro UEs and micro *inner* and *outer* UEs, these being, respectively, those UEs that can/cannot hear the serving micro well when the macro transmits. The number of resources for each group $X$ of UEs is computed as $V_X = \sum_{i \in X} B_i / R_i$, where $B_i$ is the amount of data in the buffer for user $i$, and $R_i$ is the channel rate. This scheme thus requires the algorithm to be *omniscient*, i.e. to know both the buffer status and the CQIs for each UE anywhere in the HetNet. This requires a non-negligible communication overhead, which cannot be mapped on the standard ABS signaling. The

ABS ratio is then computed as $\theta = \sqrt{V_{\mu-outer}} \Big/ \left( \sqrt{V_{\mu-outer}} + \sqrt{V_{\mu-inner} + V_M} \right)$ ($\mu$ and $M$ denoting the micro and the macro, respectively). Work [16] advocates selecting *almost blank resource blocks (ABRB)* instead of ABSs, and shows how to do so in a near-optimal way, with the objective of maximizing a concave function of the overall physical-layer data rate, assuming that the central controller possesses a HetNet-wide topology graph, stating which eNB interferes with which UE. On one hand, ABRBs are not supported in LTE, and standard UE CQI reporting alone does not allow one to construct a topology graph. On the other, the above scheme requires that, on each AP, per-UE information is conveyed to a central HetNet controller, which solves several convex optimization problems, whose size depends on the number of eNBs and UEs in the HetNet. Convex optimization has superquadratic complexity in general, and no computational results are reported in [16] to investigate at which timescales and HetNet sizes or populations this is feasible.

Works that optimize the UE association usually aim at offloading macro traffic to micros by varying the CRE Bias [17], and are often associated to static ABS-selection algorithms. A small subset of these, instead, propose methods for directly *selecting* the serving node of a UE, based on cell load and/or channel quality information [15]. Work [18] considers a heterogeneous TDMA network with cellular and Wi-Fi base stations occupying non-interfering spectra, and selects to which infrastructure a UE should connect, so that the cellular operator minimizes the energy cost of its network, while still getting revenue from it. Energy cost reduction is done by switching off cellular eNBs, something that is infeasible at per-AP timescales (the switch-on/off of a cellular node takes several tens of seconds, and forces mass re-associations of UEs [19]). Work [20] shows that setting a load-aware CRE bias leads to higher spectral efficiency than associations based on the received power only. Work [21] tackles joint UE association and ABRB provisioning in a unified framework, assuming massive MIMO and Joint Transmission from clusters of eNBs, and computes an upper bound on network performance.

Works that focus on the transmission power generally aim at selecting the power of macros (during I- or LP-ABSs), and/or micros (during non-ABSs) [22]. The transmission power is generally set *semi-statically* based on the network topology. Work [5] proposes an analytical framework based on stochastic geometry to set the power level of LP-ABSs, the ABS ratio, the CRE bias and the thresholds for inner/outer UE classification. Work [23], instead, advocates a dynamic approach and varies the transmission power at a faster pace, using the UE positions as an input. Work [24] uses reinforcement learning to have a micro learn what RBs to use to serve its UEs, at what power, and using what CRE, implicitly positing ABRBs. Moreover, it leverages dynamic Carrier Aggregation to allow a UE to be served by both a macro and a micro, on non-overlapping carriers. It is worth noting that varying the power (hence the network coverage) and/or the CRE at a fast pace is often criticized (see e.g. [10]). Such approach, in fact, may

lead to unexpected ripples on the network at large, such as interference fluctuation, etc., which in turn may reduce the effectiveness of channel prediction mechanisms and lead to higher error rates [22],[25]. Finally, within this group, [26] proposes to use arrays of antennas on the micro side, to adapt the radiation pattern thus protecting macro edge users from interference.

Schemes for selecting UE association, node transmission power and UE classification can indeed be used to set the input values for our framework. As far as works targeting ABS ratios are concerned, ours differs from them in several respects. First of all, it is meant to work in the short term (i.e., per-AP timescales), whereas all the above ones (except [10] and [16]) consider long-term decisions. Long-term decisions can indeed leverage more complex algorithms, but they require *modeling* the aspects of the problem – specifically, forecasting channel qualities and traffic arrivals – and assuming that models hold over long timespans. Our work, instead, relies on *measures* rather than models, and employs fast algorithms that run in few milliseconds. Second, almost all the above works (including [16]) investigate how *capacity* or *spectral efficiency* are affected by ABS provisioning. To do so, they make assumptions (e.g., full-buffer traffic) that are incompatible with those under which a power-saving scheme makes sense at all. Third, our work takes into account the constraints of the protocol layers of the LTE standard. In a real LTE network, UE reporting is limited to periodic CQIs, and the X2AP clearly defines what information can be shared and transmitted among the coordinated eNBs for ABS management. Only few works from the above list ([25], [26], [16]) do consider communications among the eNBs. Some of them just mention the X2 interface, others define their own (non-standard) signaling architecture for ABS provisioning. Only two works use known simulators to test their results, namely [27] uses the Vienna Simulator [28], whereas [29] uses Lte-Sim [30]. All the other use ad-hoc system-level simulators, often purported to be "3GPP-compliant" as they rely on [31] for modeling physical-layer aspects. However, none of them provide any detail as to how, if at all, what is *above* the PHY (notably resource allocation at the MAC level or inter-eNB communication) is modeled. Finally, all the above works consider either I- or LP-ABSs, but not both simultaneously, and none discuss optimal SF placement.

## IV. SYSTEM MODEL

This section details the hypotheses and assumptions underlying our work. We focus on a HetNet composed of a macro node, acting as a master for one or more micro nodes, with which it performs dynamic eICIC using ABSs as a tool.

UEs may be associated to either the macro or micro nodes (but not to both). Such association – which changes at timescales much larger than those where our framework operates – is an input datum, and the means by which it is obtained (e.g., selection of a suitable CRE bias) are outside the scope of this
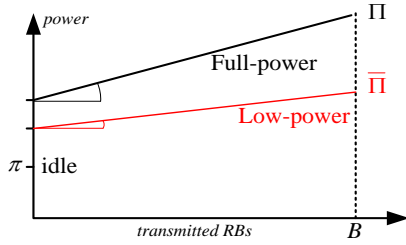
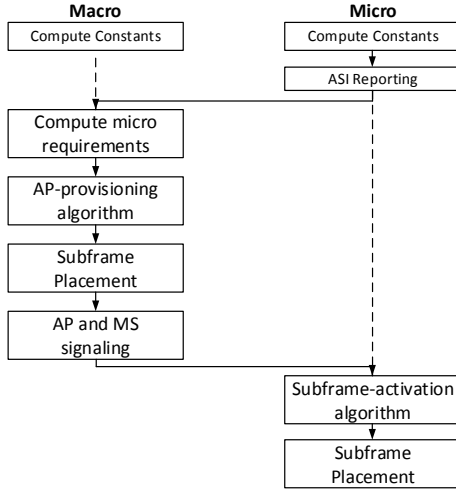Fig. 4 – Example of a power model for a node.



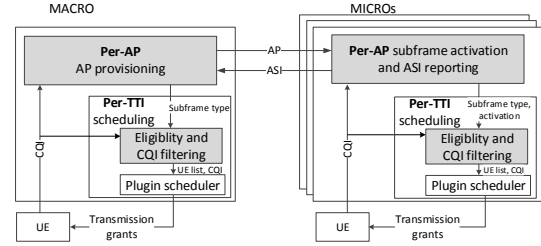Fig. 5 - Main operations of the proposed framework.



Fig. 6 - High-level representation of the proposed framework.
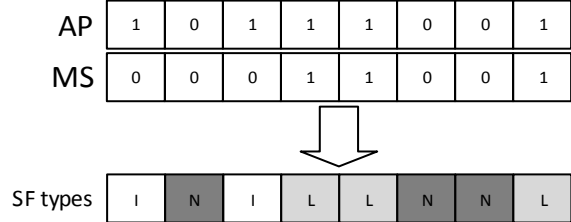


Fig. 7 - Example of SF types in an AP

TABLE 2 – TYPES OF SFS, ELIGIBLE UES AND POWER REGIMES

| Node | | Non-ABS | I-ABS | LP-ABS |
|------|------|---------|-------|--------|
| **Macro** | Eligible UEs | Outer/inner | - | Inner only |
| | Power regime | Full-power | idle | Low-power |
| **Micro (active)** | Eligible UEs | Inner only | Outer/inner | Outer/inner |
| | Power regime | Low-power | Full-power | Full-power |
| **Micro (idle)** | Eligible UEs | - | - | - |
| | Power regime | idle | idle | idle |

paper. We assume that the eNBs (both macro and micros) can classify their UEs into *inner* and *outer*: the former are those that can hear their eNB sufficiently well even when interference from the other HetNet nodes is considered, and the latter are those that cannot, hence must be protected. The criteria according to which this is done are outside the scope of this paper: for instance, a micro eNB may classify its UEs based on their attenuation or channel reporting (UEs reporting a high average CQI both during ABSs and non-ABSs will be considered as *inner*). Alternatively, outer UEs at a micro may be those that would not have associated to it, if it were not for the CRE bias. Dually, a macro eNB may classify as *inner* those UEs that have a smaller attenuation, or that report a high average CQI during both non-ABSs and LP-ABSs. Other schemes can be found in [5],[10]. Although UEs are mobile, we assume that their being *inner* or *outer* does not change within an AP. This is reasonable, since an AP is few milliseconds long (40 ms in an FDD deployment, according to the standard). Each node schedules its UEs on each TTI using its own scheduling algorithm. Since per-TTI scheduling works at a faster pace than our framework's, we deliberately abstain from placing any requirement on scheduling algorithms (nor we require that they be the same at all nodes in the HetNet), other than that they must be able to inhibit the scheduling of *outer* UEs when required. All nodes manage SFs of $B$ RBs.

We assume that a (macro or micro) node has *three* different regimes, in which it can do different things, consuming different amounts of power.

- *idle*, in which case it can *only* transmit pilot signals, and cannot transmit data traffic. Its power consumption $p$ in the *idle* state is a constant value $p = \pi$.

- *active*, operating at either *low power* or *full power*. An active node can transmit data traffic, and its power consumption is a function of the number of transmitted RBs $p = f_\tau(n)$, where $0 \le n \le B$ and $\tau \in \{low, full\}$.

Our assumptions on the power regimes are the following:

1) $\pi \le f_\tau(0)$, i.e., a node always consumes less when idle than when it is active and transmitting no RBs;

2) The power consumption of an active node (in both low-power and full-power regimes) is a *monotonically increasing affine* function of the number of RBs.

3) $f_{low}(n) \le f_{full}(n)$, $0 \le n \le B$, i.e. the low-power regime consumes less than the full-power one for the same RB number.

Assumption 1) implies that it is beneficial to have a node be *idle* when there is nothing to transmit. Besides the obvious point of monotonicity, 2) implies that – from a power-saving standpoint – once a node decides to be active in a SF, the optimal thing to do is to allocate as many RBs as required, given the traffic demand: in fact, the marginal power cost of an RB is a wide-sense decreasing function. This implies that a work-conserving per-TTI scheduler can be used. Finally, 3) is a self-evident consistency condition on the two power regimes of active nodes. Power models that verify all the above conditions are those used in [1], [32]-[33], also shown in Fig. 4. As we show later on, the only relevant data at each node, however, are the *maximum* power it consumes in active full-power and low-power SFs, call them $\Pi, \overline{\Pi}$, and its power consumption in the *idle* state $\pi$, all of which are easily measurable. The actual values of functions $f_\tau()$ will depend on the *node*, hence will differ for a macro and a micro, and may also be different for micros in the same HetNet.

We assume that nodes communicate through the X2 inter-

face, using *only* the standard signaling described in Section II. Aside from that, nodes only possess *local* knowledge, i.e. a macro node is *not* aware of its micros' position or power model, or the number, inner/outer classification and required load of micros' UEs, etc.

## V. DISTRIBUTED POWER-SAVING FRAMEWORK

This section details our power saving framework (PSF). We first provide a high-level view of PSF and its main operation, sketched respectively in Fig. 6 and Fig. 5, and then present each block in detail. PSF operates at both the macro and the micros, and on two different timescales: the *AP timescale*, and the *TTI timescale*. At the AP timescale, the macro collects the ASI reports from the micros, computes a new AP and transmits it to the micros. The micros, in turn, obtain the new AP. Based on the information reported therein and on a forecast of their own traffic, they select *which SFs to be active in*, and which to be idle in instead. Finally, at the end of the AP, micros prepare the ASI for the macro. At the TTI timescale, both the macro and the micro select *which subset* of UEs to schedule, according to the settings of the ongoing AP.

Assume that an AP consists of $T$ consecutive SFs, and that $AP[j]$, $MS[j]$ are the bits in the AP and MS vectors related to SF *j*. The macro – and, specifically, its *AP provisioning* (APP) algorithm – periodically selects SFs to be I-ABS, LP-ABS or non-ABS, and communicates that decision to the micros. The three possible SF types are identified by different combinations of the bits in the AP and MS vectors: if $AP[j]=1$, then $MS[j]=0$ will identify that SF as an I-ABS, and $MS[j]=1$ will identify it as an LP-ABS. If, on the other hand, $AP[j]=0$, then the SF will be a non-ABS one. An example of the above combinations with the resulting SF types is shown in Fig. 7.

The macro will transmit at its full-power regime to both *inner* and *outer* UEs during non-ABSs. It will instead transmit data at a *low-power* regime in LP-ABSs, to its *inner* UEs only. Finally, it will be *idle* in I-ABSs.

On the other hand, the micros will serve both inner and outer UEs during ABSs (both I- and LP-ABSs), using their full-power regime, and serve only inner UEs during non-ABSs, using their low-power regime to reduce the interference made on the macro's outer UEs. Moreover, a micro may decide to be *idle* in some SFs (of any type), if it does not need all the capacity, so as to reduce its own power consumption. The different combinations are summarized in Table 2.

The APP at the macro selects the type of each SF in the AP so as to minimize the macro's power consumption, provided that the load of *both* the macro *and* the micros is carried, if this is possible at all. Following Fig. 4 and Table 2, the macro should declare *as few non-ABSs as possible*, e.g., just enough to carry the load of its own *outer* UEs. The rest of the AP should then consist of – possibly – some LP-ABSs, if there are still *inner* UEs that could not be served during non-ABSs, and then as many I-ABSs as possible, to reap the energy benefits

of going idle. By doing this, in turn, the macro creates an ideal environment for the micros, which in turn see an AP with plenty of ABSs, at least when the network load is low, and can exploit them to serve their own UEs at reduced interference. This paves the way to further power saving opportunities *at the micros* as well: micros may in fact decide to go idle themselves at some SFs, *both* ABSs and non-ABSs, if they do not need the available downlink capacity. Mirroring what happens at the macro, *mutatis mutandis*, a micro should stay active as little as necessary to carry its load, and select the optimal combination of SFs where it is active, based on its inner/outer UE load and distribution, and on the interference suffered by the macro.

UE scheduling at both the macro and the micros will be done by a per-TTI scheduler. The latter is a plug-in for our framework, to which an *eligibility* module will feed the list of eligible UEs, depending on the type of SF as per Table 2. The eligibility module should also present the scheduler with the most recent CQIs reported by each UE *in the same type of SF as the current one*, since the interference measured by the UEs – hence their CQIs – will be different depending on the type of SFs. For instance, CQIs reported by micro UEs would typically be much lower in non-ABSs than in I-ABSs.

Hereafter, we first describe the algorithms run at the micro, and then move to describing those at the macro.

### A. Algorithms run at the micro

All the micros run the same algorithms. When the micro receives an AP, it decides which SFs to be active or idle in, according to its *SF activation* (SA) algorithm. This decision is made so as to *carry the micro's load* (if this is possible at all) trying to *minimize its power consumption*.

In order to do this, we need to estimate the micro's *expected load*, i.e. the number of bytes that we expect to transmit in the next AP. We do this by measuring what happens in the past AP, and use the measure thus obtained as a forecast. Since we ultimately need to decide whether or not we need to use the RBs of the various SFs in the AP, it is convenient to normalize the expected load to the average per-RB capacity. The latter, however, depends on the *SF type*, which determines the power level of the micro, the interference suffered from the macro, and which UEs are targeted. Let $C_N, C_I, C_L$ be the average per-RB capacity in non-ABS, I-ABS, and LP-ABS, respectively, measured in the past AP by logging the number of transmitted bytes and allocated RBs, and let $\mathbf{X} \equiv \{N,I,L\}$ be the set of SF types. Define coefficients $\alpha_N = C_N/C_I$, $\alpha_L = C_L/C_I$, i.e., the capacity per-RB *normalized to I-ABSs'*. For symmetry, one may also define $\alpha_I = C_I/C_I = 1$. Since $C_N, C_I, C_L$ will change over time, $\alpha_N$ and $\alpha_L$ will change as well. Thus, all three coefficients are initialized to a default of 1 when the system starts, and they are updated at the end of each AP. If no data is available to update them (e.g., because inner UEs have all been served during ABSs in the last AP), the *last* computed value is carried over to the next AP. $\alpha_N$ and $\alpha_L$ can be expected to be smaller than one. However, this is not

necessarily the case, due to e.g. fading peaks that change CQIs considerably in an AP, or the fact that a significant fraction of UEs may appear/disappear within an AP, so that the set over which the measurements are taken changes significantly from one SF to the other. In any case, our algorithms work regardless of the values of $\alpha_N$ and $\alpha_L$.

The *expected load* of the micro is represented by two values $K, K_{inner}$, representing the *overall* expected load and the load of *inner-UEs* only, both measured in multiples of $C_I$. The algorithm that estimates $K$ for the *next* AP is shown in Fig. 8, and it computes the number of RBs that would be needed to clear all the micro's backlog in the *ongoing* AP. $K$ is initialized to zero at the start of an AP, and is increased by the number of allocated RBs (including those for retransmissions), times $\alpha_x$, in every type-*x* SF. At the end of the AP, $K$ is increased by the number of RBs required to clear the residual backlog at the micro. $K_{inner}$ is computed the same way as $K$, however, only on *inner* UEs.

```
1.  @AP start:
2.       let K=0
3.  @type-x SF:
4.     let b=no. allocated RBs in the SF
5.     let K=K+b* α_x
6.  @AP end:
7.     let K=K+total backlog/ C_I
```

Fig. 8 – Algorithm to compute the micro's expected load in the next AP.

*1) Subframe activation algorithm*

Once the micro's expected load has been computed, the SA algorithm is quite straightforward. Let $T_x$, $x \in \mathbf{X}$, be the number of type-*x* SFs in the next AP (provisioned by the macro and identified by the micro via straightforward bitwise operations on the AP and MS vectors received in the X2 message). The SA computes $t_x$, i.e., how many type-*x* SFs to be *active* in, so as to minimize the micro power consumption. This is done by solving at optimality the following integer-linear problem (ILP):

$$\min\left\{\sum_{x\in\mathbf{X}} W_x \cdot t_x\right\}$$

*s.t.*

$$t_x \leq T_x \qquad\qquad\qquad x \in \mathbf{X} \quad (i)$$
$$t_I + t_L \geq \min\{T_{on}, T_I + T_L\} \qquad\qquad (ii)$$
$$t_N \leq \left\lceil K_{inner}/(B\cdot\alpha_N) \right\rceil \qquad\qquad (iii)$$
$$B\cdot t_N \cdot\alpha_N \leq K_{inner} + (1-b)\cdot\Omega \qquad\qquad (iv)$$
$$B\cdot t_N \cdot\alpha_N + b\cdot\Omega \geq K_{inner} \qquad\qquad (v)$$
$$B\cdot\sum_{x\in\mathbf{X}} t_x \cdot\alpha_x + (1-b)\cdot\Omega \geq \min\left\{K, B\cdot\sum_{x\in\mathbf{X}} T_x \cdot\alpha_x\right\} \quad (vi)$$
$$B\cdot(t_I \cdot\alpha_I + t_L \cdot\alpha_L) + b\cdot\Omega \geq \min\left\{K - K_{inner}, B\cdot(T_I \cdot\alpha_I + T_L \cdot\alpha_L)\right\} \quad (vii)$$
$$t_x \in \square^+ \qquad\qquad\qquad x \in \mathbf{X} \quad (viii)$$
$$b \in \{0,1\} \qquad\qquad\qquad (ix)$$

(1)

The objective function to be minimized is the variable part of the power consumed by the micro. A fixed baseline of $\pi\cdot T$ will be consumed in any case, even if the micro is always idle. For each type-*x* SF when the micro is active, an *additional* power consumption of $W_x$ is required, with $W_N = \bar{\Pi} - \pi$,

$W_L = W_I = \Pi - \pi$, (see again Table 2) to have the micro transmit all the $B$ RBs in that SF. Note that the decision on *how many* RBs are actually allocated in a SF is *not* taken by the SA, nor could it reasonably be: rather, it is taken on a per-TTI basis by the scheduler, which is the only one that knows the *actual* traffic demand at that TTI. Thus, the value of the objective function, plus the constant term $\pi\cdot T$, is an *upper bound* on the overall power that will be consumed by the micro in the next AP, and as tight as an upper bound can be unless clairvoyance of future loads and CQIs is assumed. Constraint (*i*) models the upper bound on the number of active SFs of each type given by the AP provisioning. Constraint (*ii*) guarantees that the micro is active for a minimum number of SFs $T_{on} \geq 0$ to serve both inner and outer UEs, if this is compatible with the AP decisions made at the macro (hence the *min* on the right-hand side). This may help a network engineer to strike a better trade-off between energy consumption and latency. In fact, an AP is typically long (i.e., 40ms), and if a micro has no traffic during one AP, then its expected load would be null for the next AP as well, hence the SA would set the micro to be idle for the whole *next* AP. Thus, any downlink traffic arriving in that AP would be neglected until the successive AP (when it would be counted in $K$, and possibly $K_{inner}$, as remaining backlog). By setting $T_{on}$ to a non-null value, one guarantees that enough active SFs always occur in an AP, hence the delay is reduced at the expenses of a higher power consumption. Constraint (*iii*) states that, since you can only serve inner UEs in non-ABSs, these must not exceed the number necessary to carry expected load $K_{inner}$. Constraints (*iv-vii*) state that the overall normalized capacity must be sufficient to carry the expected load $K$, if it can be carried at all. The easiest way to explain these constraints is to understand that (*iii*) allows that $B\cdot t_N \cdot\alpha_N > K_{inner}$ (as per the ceiling operator), but this does not mean that the leftover capacity $B\cdot t_N \cdot\alpha_N - K_{inner}$ can be made available to *outer* UEs, which are ineligible in non-ABSs. Therefore, we must discriminate whether $B\cdot t_N \cdot\alpha_N \leq K_{inner}$ or $B\cdot t_N \cdot\alpha_N \geq K_{inner}$, and we do this by employing a helper binary variable $b$, such that $b=1 \Rightarrow B\cdot t_N \cdot\alpha_N \leq K_{inner}$, and $b=0 \Rightarrow B\cdot t_N \cdot\alpha_N \geq K_{inner}$. The above two implications are written in constraints (*iv-v*), by multiplying alternatively $b$ or $1-b$ by a large positive constant $\Omega$, so that either constraint is inactive depending on the value of $b$. Now, if $b=1$ (i.e., the capacity in non-ABSs does not exceed the requests of inner UEs), then the total capacity in active SFs is sufficient if $B\cdot\sum_{x\in\mathbf{X}} t_x \cdot\alpha_x \geq K$. However, it may be that $K$ exceeds the allocable capacity $B\cdot\sum_{x\in\mathbf{X}} T_x \cdot\alpha_x$, hence that constraint must be reformulated as (*vi*). On the other hand, if $b=0$ (i.e., the capacity in non-ABSs may exceed the request of inner UEs), then we only need to check whether there are enough active LP-ABSs and I-ABSs to carry the expected load of *outer* UEs, if that load can be carried at all. This is written in constraint (*vii*). Note again that either (*vi*) or (*vii*) is inactive, depending on the value of $b$.

Problem (1) is an ILP with four variables and nine constraints, which is always feasible. We prove in [35] that it has

no more than $2 \cdot (T/3+1)^3$ feasible solutions. That number is independent of the number of UEs, their traffic, and the system bandwidth (i.e., constant $B$). Since the number of *constraints* is independent of all the above as well, an upper bound on the complexity of solving (1) at optimality is $O(T^3)$. In fact, a brute-force algorithm could just find the optimum by testing all constraints on all the feasible solutions. Since $T=40$, there are no more than 5890 solutions, hence even this would only take milliseconds on a modern CPU. However, state-of-the-art solution algorithms, such as those implemented in commercial solvers [36], converge rapidly towards the optimum, instead of testing solutions exhaustively.

Once problem (1) is solved, a micro knows *how many* SFs of each type to be active in, but it must still decide *which* SFs to activate. This apparently simple problem is in fact non-trivial, hence we postpone addressing it until Section V.C, i.e., after we present the algorithms for the macro, to avoid disrupting the flow of the discussion.

*2) Eligibility*

The *eligibility* module at the micro runs at every TTI and presents the scheduler with the list of eligible UEs for the current SF, according to Table 2. Moreover, it stores the CQIs of the UEs *together with the SF type* they are related to. When presenting the scheduler with the list of eligible UEs for a type-*x* SF, it also communicates to the scheduler the most recent CQI of that UE *for type-x SFs*, as shown in Fig. 9. CQIs can be expected to differ for the same UE in different SF types (this is, in fact, the very reason for employing ABSs). Hence, using a CQI measured in a different type of SF would lead to either H-ARQ errors or resource underutilization, depending on whether you overestimate or underestimate it.

*3) ASI reporting*

At the end of an AP, the micro reports two values, namely the AS and UAP. According to the standard, the UAP should report the number of ABSs where the interference is "small enough" (no quantitative definition is provided). We embody this concept by measuring the average per-RB capacity: this can be expected to be smaller in LP-ABSs than in I-ABSs, due to interference from the macro, hence all it takes is to compare the per-RB capacity of an LP-ABS against a threshold $\sigma \cdot C_I$ to determine whether that SF will contribute to the UAP or not. We justify later on that $\sigma$ plays a very limited role in our framework, and we set $\sigma = 0.7$. Thus, ASI reporting boils down to the self-explanatory algorithm of Fig. 10.

Note that the per-ABS part of the algorithm is run also when the micro is idle, in which case no RBs are allocated.

The ASI reporting algorithm is coherent with the fact that the macro advertises LP-ABSs to the micros by setting the relevant bit in the MS vector: this means instructing the micros to measure the interference in those SFs.

*B. Algorithms at the macro*

This section details the algorithms being run at the macro. Since many aspects of these are similar, and specular, to those of the micros, we will reuse the same notation whenever possible, there being no ambiguity as to the fact that we are dealing with quantities related to the macro instead.

In order to provision an AP, the macro needs to gather the ASI reports from its micros, based on which it can infer their requirements, and to compute its own expected load as well. Within a micro's ASI, two data are relevant:

− whether AS=100 or not: in fact, AS=100 denotes a *possible* overload situation, whereas AS<100 indicates that there are unused RBs in active SFs. We say "possible" since it may be that the micro node is actually able to carry its load (e.g., by serving inner UEs during non-ABS frames). However, the standard X2-based signaling does not allow to discriminate the two cases, hence we must assume a worst-case scenario.

− whether or not UAP is equal to the number of ABSs in the previous period $T_L + T_I$. If it is not, then the macro should prefer I-ABSs to LP-ABSs for that micro, given a choice, since the interference in the latter is presumably too high.

The cross-product of the above yields four different combinations, described below, together with the inferences that the macro should reasonably make on the micro's requirements as a consequence − assuming that that micro's load in the next period stays the same:

a) (UAP == $T_L + T_I$) and (AS<100):
 • tolerable interference in LP-ABSs;
 • the available ABSs (both I- and LP- ) are enough to carry the micro's load.

b) (UAP == $T_L + T_I$) and (AS==100):
 • tolerable interference in LP-ABSs;
 • not enough ABSs to carry the current load in the micro. This micro will fare better if the number of either LP-ABSs or I-ABSs is increased in the next AP.

c) (UAP< $T_L + T_I$) and (AS<100):
 • intolerable interference in LP-ABSs;
 • the available I-ABSs are enough for the current load in that micro, hence that micro does not need LP-ABSs.

d) (UAP< $T_L + T_I$) and (AS==100):

UEs attached to the micros and related CQIs

| UE | category | CQI | | | | AP | | | | |
|----|----------|-----|-----|---------|--|-----|---------|---------|---------|--|
| | | I-ABS | LP-ABS | Non-ABS | | ... | I-ABS | Non-ABS | LP-ABS | ... |
| A | inner | 15 | 13 | 5 | | | A, 15 | A, 5 | A, 13 | |
| B | inner | 14 | 11 | 6 | Eligible UEs and related CQIs | | B, 14 | B, 6 | B, 11 | |
| C | outer | 9 | 6 | - | | | C, 9 | | C, 6 | |
| D | outer | 8 | 4 | - | | | D, 8 | | D, 4 | |

Fig. 9 – Eligible UEs per SF type.

```
1.   @AP start:
2.      set UAP=T_L+T_I ; set R=0
3.   @ABS SF:
4.      measure per-RB capacity C
5.      let b=number of allocated RBs in SF
6.      if (SF is LP-ABS and C<σ·C_I)
7.         decrease UAP
8.      else
9.         let R=R+b
10.  @AP end:
11.     report UAP
12.     if (UAP)>0) report AS=R/UAP
13.     else report AS=0
```

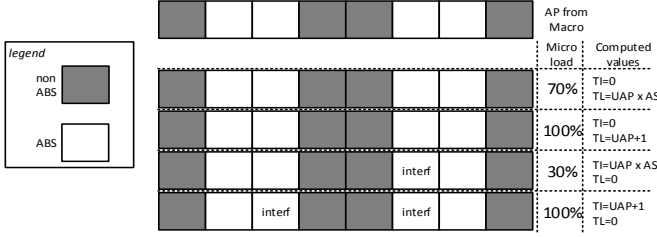Fig. 10 – ASI reporting at the micros.

```
1.   for each micro node i
2.      let req=ceiling(UAP[i]*AS[i]/B)
3.      if (AS[i]=100) increase req
4.      if (UAP[i]<T_L+T_I)
5.         let τ_I[i]=req
6.         let τ_L[i]=0
7.      else   //(UAP[i]==T_L+T_I)
8.         let τ_I[i]=0
9.         let τ_L[i]=req
10.  let τ_I=max_i{τ_I[i]}
11.  let τ_L=min{max_i{τ_L[i]}-τ_I,0}
```

Fig. 11 – Computing the minimum requirements $\tau_I, \tau_L$ at the macro.



Fig. 12 - Micro feedback to the macro node

- intolerable interference in LP-ABSs;
- There is not enough capacity in ABSs to carry that micro's load. This micro will fare better *only* if the number of *I-ABSs* is increased in the next AP.

Based on the above inferences, the macro computes two values $\tau_I, \tau_L$, corresponding to the *minimum* number of I-ABS and LP-ABSs necessary to satisfy the requirements of all the micros simultaneously, using the algorithm in Fig. 11.

For each micro *i,* the algorithm computes a required number of ABSs (line 2) based only on those where the interference is tolerable (thus, erring on the safe side). That number is increased by one if the micro signals overload (line 3). That number is then assigned to either $\tau_L[i]$ or $\tau_I[i]$ based on whether the interference in LP-ABSs is tolerable or not (lines 4-9). Finally (line 10), the minimum number of I-ABSs $\tau_I$ is computed as the maximum of $\tau_I[i]$ among all the micros. If a micro *j* requested a higher number of ABSs than $\tau_I$, and tolerates that they are LP-ABSs, then it is $\tau_L[j] > \tau_I$, hence the maximum difference $\tau_L[j] - \tau_I$, if positive, should make up $\tau_L$ (line 11). This way, some of micro *j*'s ABS will be I-ABSs in any case, which can only improve its performance. Fig. 12 shows an example of computation for $\tau_I$ and $\tau_L$, obtained from four ASI reports sent by as many micros. Note that all the above four cases are covered.

The macro also needs to compute its own expected load, in the same way as micros do, by measuring its own $K, K_{inner}$ using a similar algorithm as the one in Fig. 8. The macro computes its expected load as a multiple of $C_N$, its average per-RB capacity in non-ABSs, and RBs allocated in LP-ABSs are rescaled by coefficient $\alpha_L = C_L / C_N$. We leave to the alert reader the straightforward adaptation of the algorithm in Fig. 8 to compute $K, K_{inner}$ at the macro.

*4) AP provisioning*

The AP provisioning at the macro consists in selecting how many SFs of each type will occur in the next AP, i.e., selecting

$T_I, T_L, T_N$ such that $\sum_{x \in \mathbf{X}} T_x = T$. The following inequalities constrain the choice:

- *micro requirements*: in order to meet them, it must be $T_I \geq \tau_I$ and $T_I + T_L \geq \tau_I + \tau_L$. This last constraint allows LP-ABSs to be upgraded to I-ABSs, if the macro has no need for them (e.g., at low loads).
- *maximum latency constraints*: $T_N \geq T_{on}$. Similarly to what we do at the micros, we need to guarantee that the macro is always able to dedicate some capacity to its *outer* UEs, even if it was completely unloaded in the previous AP. Hence, at least $T_{on}$ non-ABSs should be set in an AP.
- *macro capacity requirements*: the overall capacity at the macro must be sufficient to carry its own expected load. This requires that $T_L \leq \lceil K_{inner} / (B \cdot \alpha_L) \rceil$, and that $T_N \cdot B + \min\{T_L \cdot B \cdot \alpha_L, K_{inner}\} \geq K$. The "min" at the left-hand side of the last inequality is necessary because $T_L$ may be larger than strictly necessary to serve all the *inner* UEs, as per the ceiling in the first inequality, but this does not mean that any leftover capacity can be used for *outer* UEs in LP-ABSs (like with problem (1) at the micro).

This said, the power-optimal AP provisioning at the macro is the optimum of the following ILP:

$$\min\left(W_L \cdot T_L + W_N \cdot T_N\right)$$
$$s.t.$$
$$\sum_{x \in \mathbf{X}} T_x = T \qquad\qquad (i)$$
$$T_I \geq \tau_I \qquad\qquad (ii)$$
$$T_I + T_L \geq \tau_I + \tau_L \qquad\qquad (iii)$$
$$T_N \geq T_{on} \qquad\qquad (iv)$$
$$T_L \leq \lceil K_{inner} / (B \cdot \alpha_L) \rceil \qquad\qquad (v)$$
$$B \cdot T_L \cdot \alpha_L \leq K_{inner} + (1-b) \cdot \Omega \qquad\qquad (vi)$$
$$B \cdot T_L \cdot \alpha_L + b \cdot \Omega \geq K_{inner} \qquad\qquad (vii)$$
$$B \cdot (T_N + T_L \cdot \alpha_L) + (1-b) \cdot \Omega \geq K \qquad\qquad (viii) \qquad (2)$$
$$B \cdot T_N + b \cdot \Omega \geq K - K_{inner} \qquad\qquad (ix)$$
$$T_x \in \square^+ \qquad\qquad x \in \mathbf{X} \quad (x)$$
$$b \in \{0,1\} \qquad\qquad (xi)$$

Problem (2) is indeed similar to (1). Its objective function includes only the extra power consumed in LP-ABSs and non-ABSs (the macro is in fact idle in I-ABSs), using weights $W_x$ defined as those in the objective of (1) and obtained from the macro's power model. We remark again that the value of the objective (plus a constant term $\pi \cdot T$) is an upper bound on the

power that the macro will consume in the next AP. Constraint (*i*) bounds the length of the AP. Constraints (*ii-iii*) are the micro requirements[1], (*iv*) is the maximum latency requirement, and (*v-ix*) are the macro requirements, rewritten as linear constraints with the help of binary helper variable $b$ as in (1). The problem is an ILP with four variables (one of which is redundant, since *(i)* is an equality) and nine constraints. We prove in [35] that it admits no more than $(T+1)\cdot(T+2)$ feasible solutions. Since the number of constraints is constant, solving (2) at optimality is again $O(T^2)$. Moreover, since $T = 40$, there are fewer than 1722 solutions. The same considerations as for (1) regarding solving times apply here *a fortiori*. However, unlike (1), problem (2) may actually be *infeasible*. This may occur for three reasons:

1) $T_{on} > T - (\tau_I + \tau_L)$, i.e. the micro requirements leave too few non-ABSs available for the maximum latency constraint to be met. This follows from (*i*), (*iii*) and (*iv*).

2) $K - K_{inner} > \left[ T - (\tau_I + \tau_L) \right] \cdot B$, i.e., the maximum number of non-ABSs (i.e., the term between square brackets) is not sufficient to carry $K_{out} = K - K_{inner}$, i.e. the expected load of *outer* UEs at the macro.

3) $K > \left[ T - (\tau_L + \tau_I) \right] \cdot B + \alpha_L \cdot \tau_L \cdot B$, i.e., the expected load at the macro is larger than its total exploitable capacity, in both non-ABSs (i.e. the first term at the right-hand side) and LP-ABSs (second term).

Since SF selection affects the performance of the whole network, such *overload* situations must be managed by rescaling the input constants $K, K_{inner}, \tau_I, \tau_L$, so that all the constraints become feasible. A proportional scaling is as follows:

1. If $T_{on} > T - (\tau_I + \tau_L)$ then compute the new $\hat\tau_I$ and $\hat\tau_L$ as

$$\hat{y} = y\cdot(T-T_{on})/(\tau_I + \tau_L), \text{ with } y\in\{\tau_I, \tau_L\};$$

2. If $K - K_{inner} > \left[ T - (\tau_I + \tau_L) \right]\cdot B$ then compute the new $\hat\tau_I$, $\hat\tau_L$ and $\hat{K}$ as:

$$\hat{K} = \frac{T\cdot\min(T, K_{out}/B)\cdot B}{\min(T, K_{out}/B)+\tau_I+\tau_L} + K_{inner},$$

$$\hat{y} = \frac{T\cdot y}{\min(T, K_{out}/B)+\tau_I+\tau_L}, \text{ with } y\in\{\tau_I, \tau_L\};$$

3. If $K > \left[ T - (\tau_L + \tau_I) \right]\cdot B + \alpha_L\cdot\tau_L\cdot B$ then compute the new $\hat{K}, \hat{K}_{inner}, \hat\tau_I, \hat\tau_L$ as:

$$\hat{K}_{inner} = \frac{(T - K_{out}/B)\cdot\min(T, K_{inner}/B)\cdot B}{\min(T, K_{inner}/B)+\tau_I+\tau_L(1-\alpha_L)},$$

---

[1] Note that the way constraints *(ii-iii)* are written implies that the ASI reporting threshold $\sigma$ plays a very limited role. In fact, it only discriminates between cases b) and d) above, when the micro has a high load (ASI=100). A higher $\sigma$ will lean towards d), limiting the possibility of using LP-ABSs. At low loads, i.e., when $\tau_L + \tau_I + \lceil K/B \rceil$ is well below $T$, the macro will select $T_I \geq \tau_L + \tau_I$ regardless of the value of $\tau_L$, hence of the threshold $\sigma$, since this is preferable power- and capacity-wise. At low loads, the macro may end up using LP-ABSs because these are more power efficient than non-ABSs to serve its own *inner* UEs, and not because of the micro requirements.

$$\hat{K} = (K - K_{inner}) + \hat{K}_{inner},$$

$$y = \frac{(T - K_{out}/B)\cdot y}{\min(T, K_{inner}/B)+\tau_I+\tau_L(1-\alpha_L)}, \text{with } y\in\{\tau_I, \tau_L\}.$$

We terminate this section by mentioning that the eligibility module at the macro works the same as the one for the micros, *mutatis mutandis* according to Table 2.

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| N | N | N | L | L | I | I | I | I | I |

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| N | I | L | N | I | I | N | I | L | I |

Fig. 13 – Two different SF placements in an AP.

### C. *The problem of optimal subframe placement*

The macro must *place* all types of SFs - whose numbers it has computed by solving (2) - within the AP. Similarly, each micro must select which SFs to be active in, given the AP provisioned by the macro and the numbers obtained by solving (1). The two problems share the same objective, i.e. to enable the relevant node to serve its users *as often as possible*, so as to minimize the maximum latency: with reference to Fig. 13, where a repeating AP of 10 TTIs is assumed for ease of reading, the topmost placing is such that an *outer* macro UE will have to wait up to seven TTIs before being served (i.e., if its traffic arrives at time 3, it will have to wait until time 10 for a non-ABS), whereas the maximum latency for an *inner* UE will be five TTIs (i.e., for traffic arriving at time 5). However, arranging SFs as in the bottom AP reduces these latencies to three and two, respectively. To the best of our knowledge, no paper dealing with ABS describes *how to place* SFs in the AP. The algorithm for optimal placing at the macro is intuitively straightforward, but slightly tricky to formalize. Assume that the SFs in the AP are numbered from 0 to $T-1$. First, all the *non-ABSs* are placed so as to minimize the *maximum distance* between consecutive ones (keeping into account wrap-around at the end of the AP). This guarantees that *outer* UEs can be served as frequently as possible, and is achieved by placing non-ABSs at positions $\lfloor i\cdot T/T_N \rfloor$, with $i$ going from 0 to $T_N - 1$. This is what happens at the bottom of Fig. 13, where $T_N = 3$ is assumed, and leaves the AP with $T_N$ disjoint intervals (i.e., 1-2, 4-5, 7-9). Some of these intervals ($T \bmod T_N$, in fact) have a length $\lceil T/T_N \rceil - 1$, whereas the remaining ones include $\lfloor T/T_N \rfloor - 1$ SFs, i.e. one less than the former, unless $T$ is an integer multiple of $T_N$. In our example, there is one "longer" interval (7-9), and two "shorter" ones (1-2 and 4-5).

Then, all the *LP-ABSs* are placed in these intervals so as to minimize the maximum distance between consecutive *active* SFs (whether non-ABSs or LP-ABSs) at the macro, thus increasing the frequency of transmission opportunities for *inner* UEs. This is done by assigning either $\lfloor T_L/T_N \rfloor$ or $\lfloor T_L/T_N \rfloor + 1$ LP-ABSs to the intervals, taking care to assign *more* LPs to larger intervals first, and splitting each interval evenly. In the example of Fig. 13, we need to place two LP-ABSs, hence we will have either zero or one LP-ABS per interval. The only

"larger" interval (7-9) will get one LP-ABS, and the remaining one will go to the first of the two smallest intervals, i.e., 1-2.

The pseudocode for the placement algorithm is reported in Fig. 14. Lines 2-5 compute all the involved constants (i.e., how many long and short intervals, and how many of these will get one more LP-ABS than the other intervals of the same length). The outer *for* cycle assigns non-ABSs, and the inner *for* cycle assign LP-ABSs within an interval. The alert reader will notice that the placement algorithm at the macro is $O(T)$.

At the micro, the added complication is that the $t_I$, $t_L$ and $t_N$ activations selected by the SA algorithm can occur only in SFs of matching type, whose positions are *fixed*. Finding the "most regular spacing" of activations in these settings is challenging – to the point that even formally defining what a "most regular spacing" should be is all but trivial. Moreover, it would involve solving at optimality an ILP with $O(T)$ variables, which would clearly be impractical. Since what really matters is to avoid *clustering* activations in the same region of the AP (thus unduly leaving the micro idle for too many consecutive SFs), we use a simple heuristic that leverages the even spacing already done by the macro. For each SF *type* $x$, we activate the $\lfloor i \cdot T_x / t_x \rfloor$-th SF *of that type*, with $i$ ranging from 0 to $t_x - 1$. To reduce the chance of activations of adjacent SFs of different types, the AP vector is scanned circularly starting from a *different* offset for each SF type, e.g., $r, \lfloor r + \lfloor T/3 \rfloor \rfloor_{\mathrm{mod}\,T}, \lfloor r + \lfloor 2 \cdot T/3 \rfloor \rfloor_{\mathrm{mod}\,T}$, if all three $t_I$, $t_L$ and $t_N$ are non null, where $r$ is a random offset taken uniformly in $[0, T[$ at each AP. Value $r$ is used to desynchronize activations among micros, which further reduces the interference at low loads. The placement algorithm at the micro is also $O(T)$.

### D. Communication overhead and reporting periods

As shown in the previous subsections, computing a new AP involves exchanging messages and running algorithms at both the macro and the micros. The whole timing is shown in Fig. 15. Algorithms at the macro and the micros cannot run in parallel, since the latter take the output of the former as an input, and – simple and fast as they may be – cannot be supposed to run in zero time. Assume that AP $j+1$ starts at $(j+1) \cdot T$, and call $x$ the *reporting offset*, i.e. the time it takes for our framework to prepare a new AP. The computations for AP $j+1$ must start by $(j+1) \cdot T - x$, hence must make reference to the most recent AP completed by then (i.e., the $j-1$-th in the figure, since the $j$-th is ongoing).

However, the communication and computation overhead of our framework is shorter than one AP, as we show later on. Thus, one may ask if having fresher reports from the micros may be beneficial, and how this can be accomplished. All it takes is to interpret the standard liberally, and assume that an AP means a *generic* period of $T$ consecutive SFs with an arbitrary offset. Define a *reporting period* $[j \cdot T - x, (j+1) \cdot T - x]$ (consisting of $T$ SFs), and assume that the micros make their report based on *reporting periods* instead of APs. This would allow a new AP to be computed using the most recent possible traffic, CQI and interference estimates. This setting generaliz-

es the standard, which can be seen as the particular case $x = T$. For this to happen, the macro and its micros would then need to agree on the value of $x$, using extra (non-standard) X2 signaling. Note that the standard does not specify the *period* at which APs can be modified. The latter must be a multiple of the AP, and we call it *reconfiguration period* henceforth.

```
1.  init: all SFs are I-ABS
2.  let q=⌊T_L/T_N⌋ , r=T_L mod T_N
3.  let NLongInts=T mod T_N , NShortInts=T_N −(T mod T_N)
4.  let ExtraLPLong=min(r,NLongInts)
5.  let ExtraLPShort=min(r-ExtraLPLong,NShortInts)
6.  for i=0 to T_N −1
7.      mark SF ⌊i·T/T_N⌋ as non-ABS
8.      let NLPs=q
9.      let I =⌊(i+1)·T/T_N⌋−⌊i·T/T_N⌋−1 // int. lgth
10. let O =⌊i·T/T_N⌋+1              // int. offset
11. if (I == ⌈T/T_N⌉−1) and (ExtraLPLong>0)
12.     let NLPs=q+1
13.     let ExtraLPLong=ExtraLPLong-1
14. if (I == ⌊T/T_N⌋−1) and (ExtraLPShort>0)
15.     let NLPs=q+1
16.     let ExtraLPShort=ExtraLPShort-1
17. for j = 1 to NLPs
18.     mark SF O+⌊j·I/(NLPs+1)⌋ as LP-ABS
```
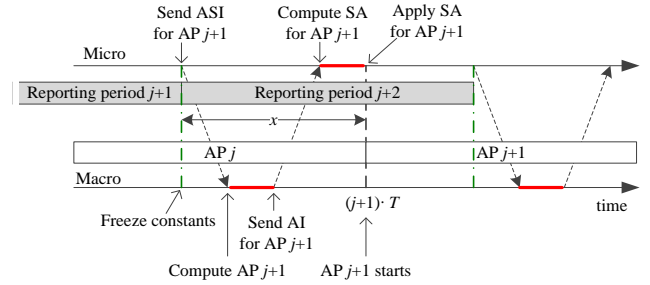
Fig. 14 – SF placement at the macro.



Fig. 15 - Message exchange for the preparation of a new AP.

## VI. PERFORMANCE EVALUATION

In this section, we evaluate the performance of PSF. Our evaluation is carried out using SimuLTE [37]-[38], a system-level simulator developed for the OMNeT++ simulation framework [39]. SimuLTE simulates the data plane of the LTE/LTE-A radio access network. The SimuLTE protocol stack includes all the LTE protocol layers, i.e., a Packet Data Convergence Protocol – Radio Resource Control (PDCP-RRC), Radio Link Control (RLC), MAC and PHY. It also includes models for *macro-*, *micro-*, *pico-*eNBs, with different radiation profiles (both isotropic and anisotropic), functions for MAC-level scheduling in downlink and uplink directions and X2-based inter-eNB communication. Note that the benefits of having an OMNeT++-based simulator include leveraging the INET library [40], which already includes validated models of well-known Internet protocols (therein including SCTP for the X2, IP, TCP, UDP, etc.).

We extended SimuLTE to support ABSs and ABS signaling

TABLE 3 – MAIN SIMULATION PARAMETERS

| Parameter | Value |
|---|---|
| Bandwidth | 10 MHz – 50 RBS |
| Carrier Frequency | 2 GHz |
| Path Loss Model | ITU Urban Macro [31] |
| Fading Model | Jakes |
| Simulation Time | 30 s |
| Warm-up time | 3 s |
| # Indep. Replicas | 5 |
| Packet size (default) | 50 bytes |
| Inter-packet time | 20ms |
| Active period duration | Weibull, mean 0.8s |
| Silence period duration | Weibull, mean 1.0s |

TABLE 4 - POWER MODEL

| Scenario | Macro | | Micro | |
|---|---|---|---|---|
| | 1 Micro | 2 Micros | 1 Micro | 2 Micros |
| Tx Power (Low) | 30 dBm | 21 dBm | 14 dBm | 14 dBm |
| Tx Power (High) | 40 dBm | 40 dBm | 30 dBm | 21 dBm |
| Antenna pattern | Anisotropic | | Anisotropic | |
| $\pi$ | 174 | | 7 | |
| $\Pi$ | 445.5 W | | 342.5 W | |
| $\Pi$ | 38.5 W | | 21 W | |
| $f(0)$ | 320 W | | 16 W | |

over X2, and we included an implementation of the PSF framework. Both the ILPs are built dynamically during the simulation, and solved at runtime using the CPLEX solver [36]. The latter uses *presolving* techniques [41] to reduce the size of a problem and improve the tightness of its formulation, and runs *branch-and-cut* [42] afterwards. An optimality gap of 0.5% is set, hence solutions are at least 99.5% optimal.

UE classification is performed at the start of the simulation, based on the attenuation to the serving node, with a threshold of 80dB. This allows you to have both inner and outer UEs at the macro and the micro. Only downlink traffic is simulated, and we assume that the traffic is generated by applications running on a server and forwarded by a router to the serving cell of the receiver. Applications alternate between active and inactive states, whose duration is extracted from a Weibull distribution. When active, the application generates fixed-length packets every 20 ms, unless otherwise specified. In our evaluations, we often vary the offered load by varying the packet size, whose default value is 50 bytes. Scheduling at the MAC layer is done according to a MaxCQI policy, i.e., eligible UEs are sorted by decreasing CQI and each one is allocated enough RBs to empty its queue, until either all RBs are occupied or no eligible UEs remain. Propagation and transmis-

sions delays of the X2 are assumed to be negligible, given the small inter-node distance and size of the ABS X2 messages.

We measure the end-to-end application delay, the HetNet throughput at the MAC layer, and the HetNet power consumption, computed as the sum of the consumption for all the eNBs in the system. For the latter we only display the part that can be actually affected by the ABS decision process, i.e. the one exceeding the baseline $\pi$. Confidence intervals at the 95% level are shown only when visible. The main simulation parameters are reported in Table 3, whereas the values used in the power models are described in Table 4, and are taken from [33]. With these values, a macro (micro) in its low-power regime consumes the same power as a micro (pico) at its full-power regime. A more extensive performance evaluation, including more scenarios, is reported in [35].

*A.  Benchmarking*

We first benchmark the performance of PSF in various configurations, to explore trade-offs between power consumption and UE QoS. We consider a scenario with one macro and one micro node where UEs are dropped uniformly within a circle whose radius is 250 m, and are associated to a serving eNB on a highest-received-power basis, as we show in Fig. 16. The default CRE bias is zero. In this deployment, the number of inner UEs is around 5% of the total at the macro, and 25% at the micro. In Fig. 17 we show the distribution of the user ap-
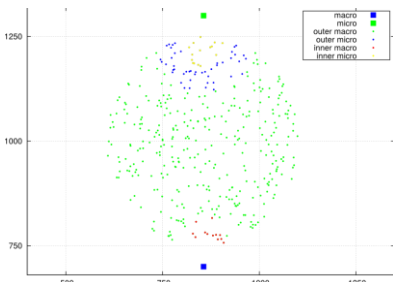


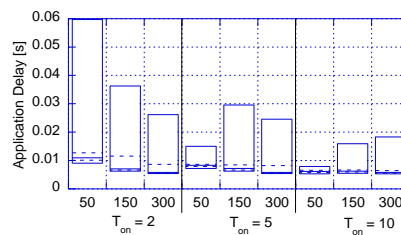Fig. 16 – Node deployment and UE association.



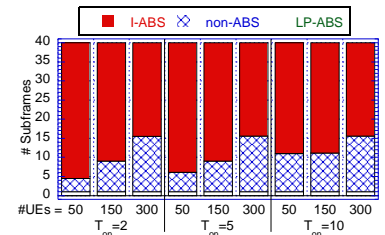Fig. 17 – Application delay with increasing values of $T_{on}$.



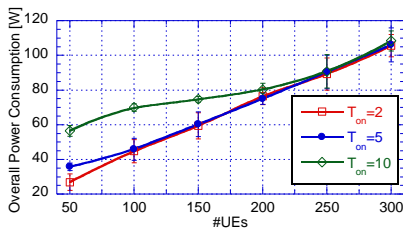Fig. 18 – Number of SFs of each type with increasing values of $T_{on}$.



Fig. 19 – HetNet power consumption (variable part) with increasing values of $T_{on}$.
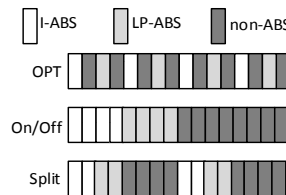


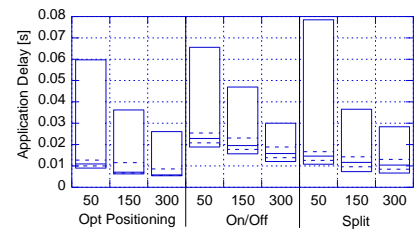Fig. 20 – Positioning of I- and LP-ABSs with three different algorithms, assuming an AP of 16 SF for convenience



Fig. 21 – Application delay for various positioning algorithms, $T$x
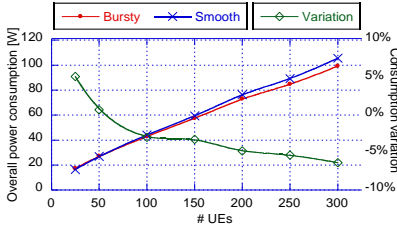
Fig. 22 – HetNet power consumption with smooth and bursty traffic (left *y* axis), and relative difference (right *y* axis), with an increasing number of UEs.
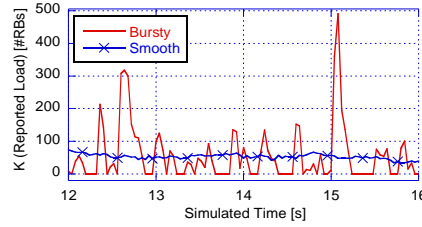


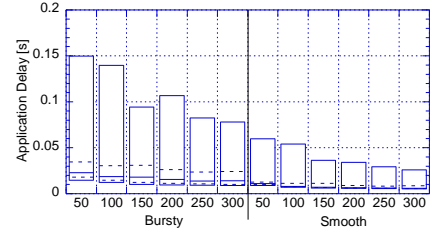Fig. 23 – Reported *K* at the macro with smooth and bursty traffic, in a scenario with 25 UEs.



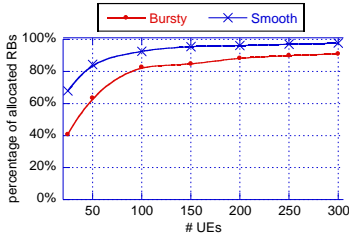Fig. 24 – Application delay with smooth and bursty traffic.



Fig. 25 – Percentage of RBs allocated by macro eNB during non-ABS, with smooth and bursty traffic.
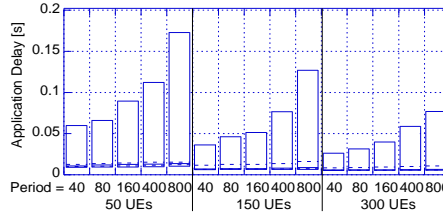


Fig. 26 – Application delay with an increasing value of the reconfiguration period, $T_{on}$=2.



Fig. 27 - HetNet power consumption against the number of UEs with several reconfiguration periods and $T_{on}$=2.



Fig. 28 - Number of SFs of each type with an increasing value of the reconfiguration period, $T_{on}$=2.



Fig. 29 - Application delay with an increasing value of *offset*, $T_{on}$=2, reconfiguration period 40 ms.

plication delay for three values of $T_{on}$ when the number of UEs ranges from 50 to 300. Each bar represents the interval between the 90[th] (upper edge) and 10[th] (lower edge) percentiles. The solid line marks the median, and the two dashed lines mark the 25[th] and 75[th] percentiles. From such graphs we can see that, when $T_{on}$ is low, all the percentiles decrease with the load. Although this might seem counterintuitive, it can be easily justified analyzing the number of SF per type, which is shown in Fig. 18: the higher the number of UEs, the higher the number of SF where the eNBs will be active, hence the faster it reacts to UEs traffic requests. If we increase the value of $T_{on}$, delay percentiles decrease globally at low and medium loads, and the number of non-ABSs increases as well. At higher loads, instead, increasing $T_{on}$ is less effective, as eNBs are already active most of the time, but it still reduces the 90[th] percentile. On the other hand, using larger values of $T_{on}$ increases the power consumption, especially at low loads, as demonstrated by Fig. 19. Increasing $T_{on}$, in fact, increases both the number of non-ABSs and the number of ABSs where the micro is active. By tuning $T_{on}$, a network operator can trade the

power consumption of the HetNet for the UE QoS.

We then show that placing ABSs smartly within an AP brings benefits. We compare our optimal placing against two baselines: one (*On/Off*) where SFs of the same type are clustered, and a second one (*split*) where clustering occurs in the two halves of the AP. An exemplary placement of I-ABSs and LP-ABSs using the three algorithms is shown in Fig. 20, whereas in Fig. 21 we show the distribution of the application delay for the three placing algorithms with an increasing number of UEs. As we can see, the optimal placing yields lower delays in general.

We now show how PSF performs with bursty traffic. We modify traffic generation so that an application packs into two bursts the same number of packets it would send during an active period. The average number of packet per burst is 20, and the mean rate of applications stays the same. As we can see in Fig. 22, the power consumption is similar to the case of smooth traffic, the relative difference (shown in green and referenced on the right *y* axis) being less than 10% either way depending on the load. With bursty traffic, PSF consumes

slightly *more* power at very low loads (e.g., 25 UEs), and *less* power as the load increases, having a break-even point around 50 UEs. This is due to the interplay of two different phenomena: on one hand, the fact that traffic is bursty allows the MAC scheduler to pack several RLC SDUS (i.e., IP packets with PCDP and RLC headers) into the same transport block, thus reducing the overhead of MAC and RLC headers on MAC-layer transmissions: this yields a slightly smaller MAC-level load for the same application-level load. This reduces power consumption, and this effect dominates at higher loads. On the other hand, bursty traffic affects the computation of $K$. Successive values of $K$ at the macro are shown in Fig. 23. The fact that a non-zero $K$ is sometimes followed by a null $K$ implies that the PSF uses the former in an AP where no traffic arrives, hence overestimates the expected load for that AP. Fig. 25 confirms this, showing that fewer RBs are allocated during non-ABSs at the macro when bursty traffic is transmitted, and such difference is more evident at low loads. Note that the knee in the RB utilization at low loads (regardless of whether traffic is smooth or bursty) is due to the presence of $T_{on}>0$, which forces more non-ABSs than strictly necessary. Finally, delays are higher with bursty traffic, as shown in Fig. 24, due to the fact that an arriving packet is also delayed by the part of burst arrived just ahead of it and already sitting in the queue.

The last two parameters whose effects we test are the *reporting offset* and *reconfiguration period*, both defined in Section V. In Fig. 26 we show the distribution of the application delay for values of the reconfiguration period ranging from 40 to 800 ms. As we can see, using smaller periods allows faster reactions to traffic variations, hence reduces the application delay. This is confirmed by Fig. 28 and Fig. 27, which show that the number of non-ABSs decreases slightly as the period increases, leading to small differences in terms of power consumption. We then evaluate the impact of reporting offsets. Fig. 29 reports the application delay for various offsets, demonstrating a negligible effect for up to 20 ms, a time largely sufficient to solve the two provisioning problems *and* to complete the handshake via X2, as we discuss below. From now on, we assume a reconfiguration period of 40ms, i.e., APs are changed every time, and a reporting offset of 10ms.

Finally, we quantify the algorithmic work required to solve our ILPs. CPLEX's presolver solves 80-90% of the instances.
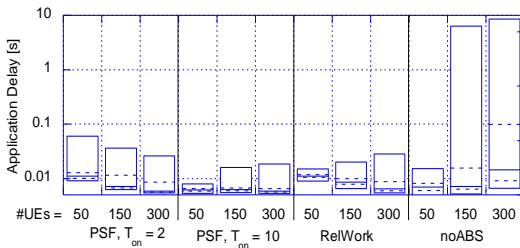
The remaining ones are solved with branch-and-cut, employing a total number of simplex iterations seldom exceeding ten [35]. The solving times, measured on an off-the-shelf PC equipped with an 8-core Intel i7-4790 CPU @ 3.60GHz and 16 GB RAM, were always between 0.5 and 5 ms.

### B. Comparison

We now compare PSF against a baseline with no ABS support and a work that advocates dynamic computation of ABSs, i.e., [10], and we will refer to them respectively as *noABS* and *RelWork*. To make the comparison more challenging, we enhanced [10] with our ABS-placing algorithm, and we allow a macro node to go to sleep during ABSs (something which [10] did not posit, and that does reduce that scheme's power consumption). We also remark that [10] assumes HetNet-wide *omniscience* at the macro, while our framework does not.

### 1) Single-Micro scenario

We consider a single-micro deployment like the one of the previous sub-section. We start by comparing the application delays for an increasing number of UEs, using two values of $T_{on}$. As Fig. 30 shows, when using $T_{on}=10$ our framework yields lower delays than all the other configurations, with a significant improvement at lower loads. With $T_{on}=2$ instead, delays for percentiles up to the 75th are close to the ones of the related work. Fig. 31 shows that PSF is the most energy-efficient, with a power saving of up to 72% with $T_{on}=2$, and up to 41% with $T_{on}=10$. The latter configuration in fact, achieves lower delays at a cost of an increased power consumption, i.e. allocates more non-ABSs. However, the number of I-ABSs is still significantly higher than the related work, which allocates SFs proportionally to the *ratio* of the macro and micros loads, regardless of their *absolute* offered load. This explains why the related work curve in Fig. 31 is approximately flat, and implies that, when the absolute load is low, an excess of non-ABSs is allocated. The fact that, especially at higher loads, using a power-saving algorithm yields lower powers *and* smaller delays (see Fig. 30 and Fig. 31) appears to be counter-intuitive: in fact, power saving is achieved by switching off nodes, and this creates interference-free environments for other nodes, which in turn employ fewer resources to serve their UEs with higher efficiency, hence serving them faster. This is why a solution without ABSs performs the worst in both re-



Fig. 30 - Application delays with an increasing number of UEs, packet size 50 bytes.
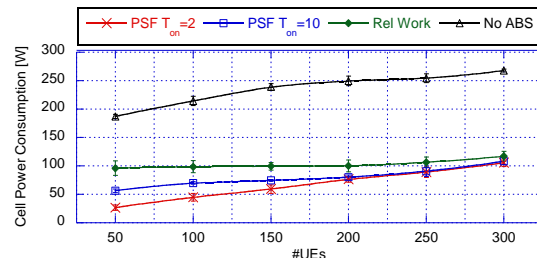


Fig. 31 - HetNet power consumption (variable part) with an increasing number of UEs, packet size 50 bytes
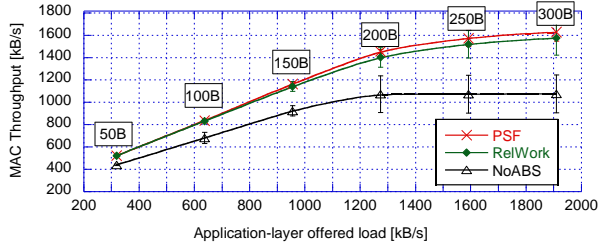
Fig. 32 - MAC-layer HetNet throughput against the application-layer offered load. For each offered load, the corresponding application-layer packet size is reported in the labels.



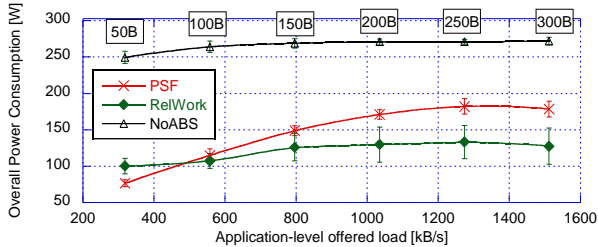Fig. 33 - Application delay with an increasing offered load.



Fig. 34 - HetNet power consumption (variable part) against application-layer offered load. For each offered load, the corresponding application-layer packet size is reported in the labels.
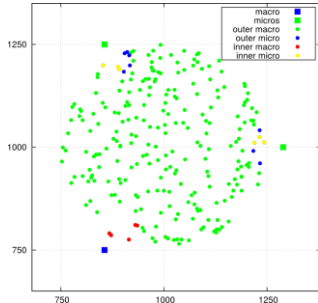


Fig. 35 - SF-type allocation with an increasing application-layer offered load, 200 uniformly distributed UEs.



Fig. 36 - Node deployment and UE association in a scenario with two micros.



Fig. 37 - Application delay with an increasing number of UEs.

spects. At low loads, instead, switching off nodes results in an increased latency. These two phenomena can be observed consistently in all the scenarios analyzed in the rest of the paper. In the following experiments, we use $T_{on}=2$.

As a next step, we perform a saturation analysis of the system. To that purpose, we place 200 UEs, which gives us enough spatial diversity, and we increase the size of the packets, until the system reaches saturation. Fig. 32 represents the overall MAC-level throughput, which shows that using ABSs increases the overall system performance, and that our solution achieves a slightly higher saturation throughput than the related work's. Moreover, we can see from Fig. 33 that the related work scheme cannot keep the *delays* within an acceptable range, already at moderate loads. The fact that its power consumption (shown in Fig. 34) is lower, which is due to the overabundance of I-ABSs (see Fig. 35), shows that the trade-off point between power saving and QoS is largely suboptimal for the related work. Our framework, instead, consumes less power at low loads, and allocates more power to preserve QoS
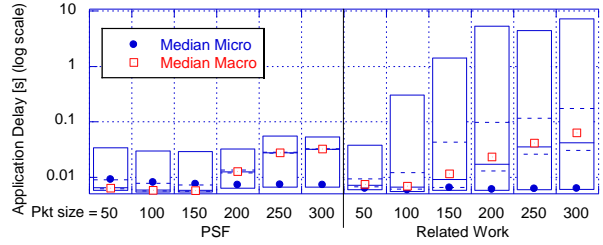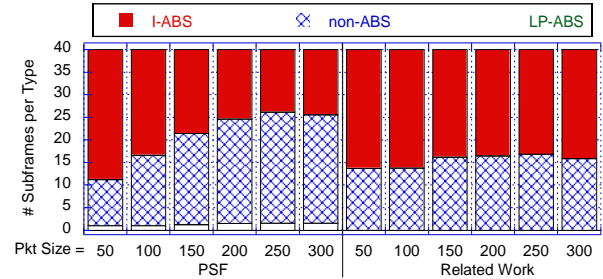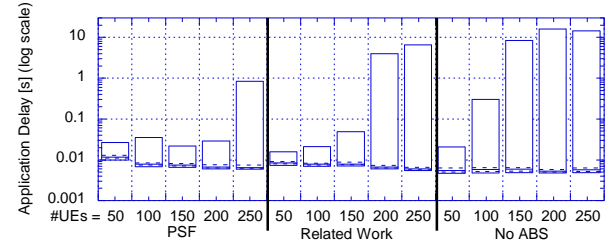
as the load increases.

*2) Two Micros*

We now compare the performance of our framework in a scenario with one macro and two micros, with UEs uniformly dropped, as shown in Fig. 36. First, we evaluate the QoS by considering an increasing number of UEs generating traffic with fixed bitrate (packets 50 bytes long). Fig. 37 shows the distribution of application delays in the three cases. As we can see, our solution outperforms both the baseline and the related work. However, we observe significantly different results with respect to the single micro scenario, and the main reason behind this phenomenon is the interference *between micros*. As explained in Section V.C, our SA algorithm forces micros to remain active for the minimum required time, *and* de-synchronizes micro activations by picking a random initial offset at each AP. Besides the (comparatively minor) power saving due to the switch-off of the single micros, this mechanism favors a multiplexing over time of micro activations, hence *reduces inter-micro interference*. In fact, we verified

that our solution achieves higher CQIs and has a lower rate of decoding errors at the MAC layer, as the rate of simultaneous micro transmissions – which also causes interference – is kept lower than the related work's (see Fig. 38). The fact that at high loads micros transmit *simultaneously* even less often when no ABSs are provisioned is due to the fact that micros cannot serve their outer UEs at all in that case, hence transmit less often altogether. Fig. 39 shows that the combined effect of micro switchoff and inter-micro interference reduction allows further power saving opportunities. The results are qualitatively similar if we change the number or position of micros.

*3) Two Micros: hot spot*

All the results we analysed so far showed an allocation of ABS that was exclusively composed of I-ABS and non-ABS, with only one LP-ABS allocated for CQI measurement. To demonstrate the usefulness of LP-ABS, we now assess the performance of our framework in a hot-spot scenario like the one shown in Figure 40. We used two micros as in the previous section, deploying UEs in two sets: a first one composed
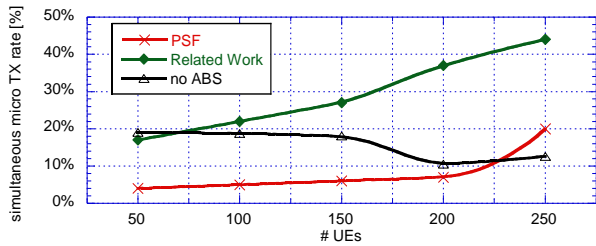
of 50 UEs uniformly distributed in the cell, and a second one of 150 UEs uniformly dropped in a hotspot close to the macro node. The traffic volume is varied by increasing the packet size only for UEs within the hotspot, until the system reaches saturation. Fig. 41 shows that the performance in terms of MAC cell throughput are similar in the three cases, as most of the traffic is generated in proximity of the macro node, thus benefiting less from the usage of ABSs. However, the power consumed by the three nodes, which is shown in Fig. 42, is significantly lower with our framework. The main reason for this is displayed in Fig. 43: in this case in fact, our framework adapts its allocation to the unbalanced deployment of UEs serving most of them at low power, thus using LP-ABS and reducing power consumption with no impact on QoS.

*4) Evaluation with multiple macros*

We now evaluate the scenario of Fig. 44, where a *central* cell like the one in is surrounded by three *external* macro nodes at 1km from its center, radiating towards it, to create a high interference. The central nodes always run the PSF



Fig. 38 – Percentage of TTIs where micros transmit simultaneously against the number of UEs.



Fig. 39 – HetNet power consumption (variable part) with an increasing number of UEs.

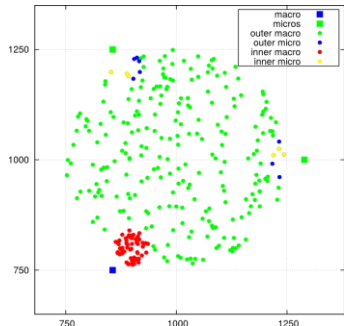

Figure 40 - Node deployment and UE association in a scenario with two micros and a hot spot close to the macro.



Fig. 41 - MAC-layer HetNet throughput against the application-layer offered load. The offered load is increased for hot-spot users only. For each offered load, the corresponding application-layer packet size is reported in the labels.
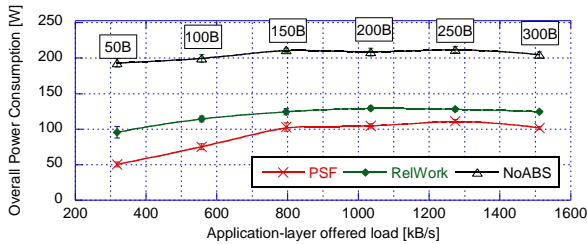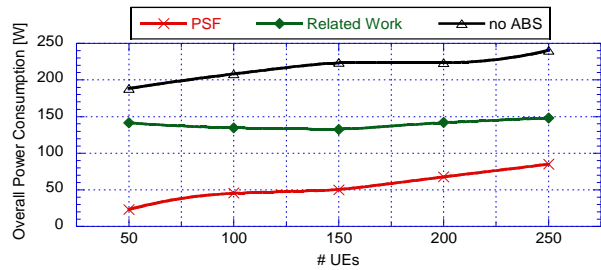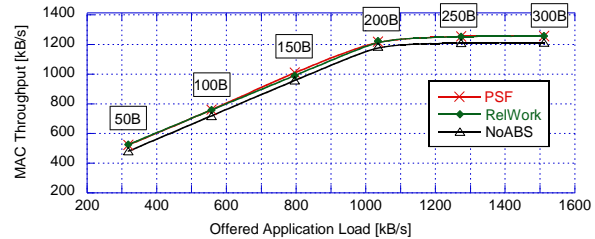


Fig. 42 - HetNet power consumption against application-layer offered load. The offered load is increased for hot-spot users only. For each offered load, the corresponding application-layer packet size is reported in the labels.
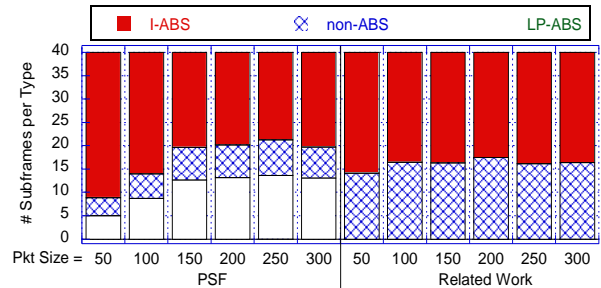


Fig. 43 - SF-type allocation with an increasing application-layer offered load

framework, whereas the *external* nodes may: a) not use the ABS mechanism; b) run *autonomous* instances of PSF, i.e., compute their own ABS pattern; c) run a *coordinated* PSF with the central cell. In the last case, each macro node computes its own values of $K$ and reports this and its micro load to a central entity. The latter takes the maximum of all the reported values, runs the ABS provisioning algorithm described in Section V.B, and returns the result to all the macros, which then apply it and send it to their micros. This makes all the coordinated entities agree on the same ABS pattern. We consider two load levels for the *external* nodes: low load, where the RB demand is 10% of the total, and high load, where it is 50%. We measure the performance in the central cell only.

As we can see in Fig. 45, external nodes increase the power consumption. This is because a higher load implies higher interference, hence lower CQIs, hence more RBs required to serve the same traffic. However, when external cells run *autonomous* PSF instances, the power consumption increases *less*. The justification can be found in Fig. 46. When no ABSs are used at the *external* nodes in fact, these generate a strong interference potentially in *all* the SFs, causing low CQIs of UEs served by the micro. On the other hand, running *coordinated* PSF, instead, achieves maximum protection for micro UEs, significantly increasing their channel quality. However, the same mechanism forces all the macro nodes to transmit during the same SFs (i.e. non ABS), thus increasing inter-macro interference (see left part of Fig. 46) and reducing the CQIs of the macro UEs, which are more numerous and served by a node with a higher power consumption. This justifies the higher power consumption of coordinated PSF. Running *autonomous* PSF instances yields a favorable trade-off, increasing the chance that both micro and macro users are scheduled in SFs where external interferers are inactive.

## VII. Conclusions

In this paper, we presented a framework for ABS provisioning that aims at reducing the power consumption in HetNets. Our framework exploits two types of ABSs, idle and low-power ABS, where a macro node respectively refrains from

transmitting data or still does so at a lower transmission power. The provisioning process is carried out dynamically at the fastest possible pace, selecting which ABS type to use depending on information on the network load, such as its volume and/or spatial distribution. The framework also accounts for macro-micro communication, which is realized via X2-interface and using standard-compliant signaling only, and its computational burden is independent of the number of UEs, their traffic, and the system bandwidth.

Finally, we evaluated our scheme by means of system-level simulations, comparing it against a legacy system with no ABS support, and a dynamic scheme taken from the literature. We showed that our framework consumes significantly less power at low loads, and preserves QoS as the system approaches saturation by dynamically adapting to different user deployments. Moreover, running autonomous instances of our scheme at macro nodes yields benefits in terms of consumed power and UE channel quality (both macro and micro).

## References

[1] Dario Sabella *et al.*, "Energy Management in Mobile Networks Towards 5G", in M.Z. Shakir et al. (eds.), Energy Management in Wireless Cellular and Ad-hoc Networks, Studies in Systems, Decision and Control, Springer, 2016

[2] Towards Real Energy-efficient Network Design (TREND) project. http://www.fp7-trend.eu/

[3] 3GPP, R1‑113482, Ericsson, ST-Ericsson System performance evaluations on FeICIC, Oct. 2011.

[4] S. Deb, P. Monogioudis, J. Miernik and J. P. Seymour, "Algorithms for Enhanced Inter-Cell Interference Coordination (eICIC) in LTE HetNets," *IEEE/ACM Trans. on Net.*, vol. 22, no. 1, pp. 137-150, Feb. 2014.

[5] A.Merwaday, S.Mukherjee, I.Güvenç, "Capacity analysis of LTE-Advanced HetNets with reduced power subframes and range
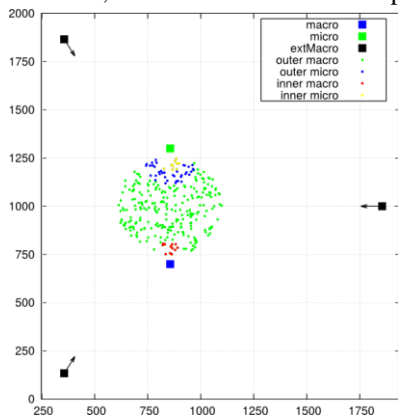


Fig. 44 – Node deployment and association in a scenario with interference coming from external cells.
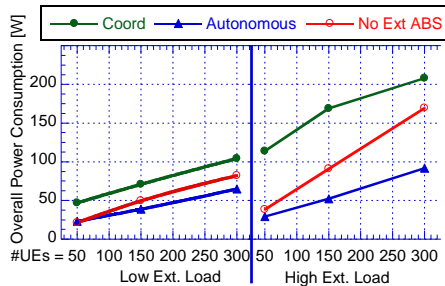


Fig. 45 – HetNet power consumption in a scenario with interference from external macros in two configurations having respectively low (left) and high (right) load in the external cells.
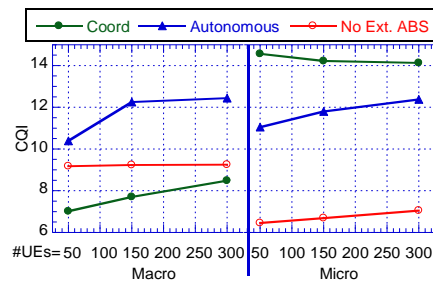


Fig. 46 – CQI of Macro (left) and Micro (right) UEs in a scenario with interference from external macros at high load.

expansion", *EURASIP J. Wireless Com. Network*, 2014:189, pp. 1-19, Dec. 2014

[6]   E. Visotsky *et al.*, "Joint Scheduling for CoMP and eICIC in Heterogeneous Network Deployments," in Proc. IEEE 77[th] VTC Spring, Dresden, DE, 2013

[7]   M. Wang, H. Xia and C. Feng, "Joint eICIC and dynamic point blanking for energy-efficiency in heterogeneous network," in Proc. WCSP 2015, Nanjing, CN, 2015.

[8]   F. Alfarhan, R. Lerbour and Y. Le Helloco, "An Optimization Framework for LTE eICIC and Reduced Power eICIC," in Proc. IEEE GLOBECOM, San Diego, US, 2015.

[9]   Yu Chen, Xuming Fang and Bo Huang, "Joint ABS power and resource allocations for eICIC in heterogeneous networks," in Proc. IWSDA, Tokyo, JP, 2013.

[10]  S. Vasudevan, R. N. Pupala and K. Sivanesan, "Dynamic eICIC — A Proactive Strategy for Improving Spectral Efficiencies of Heterogeneous LTE Cellular Networks by Leveraging User Mobility and Traffic Dynamics," *IEEE Trans. on Wireless Comm.*, vol. 12, no. 10, pp. 4956-4969, Oct. 2013.

[11]  *Evolved Universal Terrestrial Radio Access Network (E-UTRAN); X2 general aspects and principles (Release 13)*, 3GPP TS 36.420 v13.0.0, Dec. 2015.

[12]  *Evolved Universal Terrestrial Radio Access Network (E-UTRAN); X2 Application Protocol (X2AP) (Release 13)*, 3GPP TS 36.423 v13.2.0, Dec. 2015.

[13]  *Evolved Universal Terrestrial Radio Access (E-UTRA) and Evolved Universal Terrestrial Radio Access Network (E-UTRAN); Overall description*; Stage 2, 3GPP - TS 36.300, Sep. 2016

[14]  B. Soret, H. Wang, K. I. Pedersen and C. Rosa, "Multicell cooperation for LTE-advanced heterogeneous network scenarios," *IEEE Wireless Comm.*, vol. 20, no. 1, pp. 27-34, Feb. 2013.

[15]  K. I. Pedersen, *et al.* "Dynamic Enhanced Intercell Interference Coordination for Realistic Networks," *IEEE Trans. on Vehicular Tech.,* vol. 65, no. 7, pp. 5551-5562, July 2016.

[16]  A. Liu, V. Lau, L. Ruan, J. Chen, and D. Xiao, "Hierarchical radio resource optimization for heterogeneous networks with enhanced intercell interference coordination (eicic)," *IEEE Trans. Signal Processing*, vol. 62, no. 7, pp. 1684–1693, Apr. 2014.

[17]  H. Zhou, Y. Ji, X. Wang and B. Zhao, "ADMM based algorithm for eICIC configuration in heterogeneous cellular networks," in Proc. IEEE INFOCOM, Kowloon, HK, 2015, pp. 343-351.

[18]  S. Kim, S. Choi, and B. G. Lee, "A joint algorithm for base station operation and user association in heterogeneous networks," *IEEE Comm. Letters*, vol. 17, no. 8, pp. 1552–1555, August 2013.

[19]  M. Ajmone Marsan, L. Chiaraviglio, D. Ciullo, M. Meo, "Switch-off transients in cellular access networks with sleep modes", in Proc. ICC 2011, Kjoto, JP, 5-9 June 2011.

[20]  U. Siddique, H. Tabassum, E. Hossain, D. In Kim, "Channel-Access-Aware User Association With Interference Coordination in Two-Tier Downlink Cellular Networks," *IEEE Trans. on Vehicular Tech.*, vol. 65, no. 7, pp. 5579-5594, July 2016.

[21]  Q. Ye, O. Y. Bursalioglu, H. C. Papadopoulos, C. Caramanis, and J. G. Andrews, "User association and interference management in massive MIMO hetnets," *IEEE Trans. Commun.*, vol. 64, no. 5, pp. 2049–2065, May 2016.

[22]  X. Wang, *et al.*, "Reduced power centralized eICIC for LTE-advanced heterogeneous networks," in Proc. ICCC 2014, Shanghai, CN, 2014.

[23]  M. Amara and A. Feki, "Optimized ABS in LTE-Advanced Heterogeneous Networks with Adaptive Macro Cell Transmission," in Proc. 2015 IEEE GLOBECOM Workshops, San Diego, US, 2015.

[24]  M. Simsek, M. Bennis, İ. Güvenç, "Learning Based Frequency- and Time-Domain Inter-Cell Interference Coordination in HetNets", IEEE Trans. on Vehicular Techn., vol. 64, no. 10, pp. 4589-4602, Oct. 2015.

[25]  A. Morimoto, N. Miki and Y. Okumura, "Performance Evaluation of Reduced Power Inter-cell Interference Coordination for Downlink in

LTE-Advanced Heterogeneous Networks", in Proc. 19[th] European Wireless Conference, Guildford, UK, 2013.

[26]  A. M. Sadekar and R. H. Hafez, "LTE-A enhanced Inter-cell Interference Coordination (eICIC) with Pico cell adaptive antenna," in Proc. NOF 2015, Montreal, CA, 2015.

[27]  M.Yassin, *et al.*, "Survey of ICIC techniques in LTE networks under various mobile environment parameters", *Wireless Networks* vol. 23, no.2, pp. 403-418, Feb. 2017

[28]  C. Mehlführer *et al.* "The Vienna LTE simulators—enabling reproducibility in wireless communications research". *EURASIP J. Adv. Signal Process.* (2011) 2011: 29.

[29]  M. I. Kamel and K. M. F. Elsayed, "Performance evaluation of a coordinated time-domain eICIC framework based on ABSF in heterogeneous LTE-Advanced networks," in Proc. IEEE GLOBECOM 2012, Anaheim, US, 2012, pp. 5326-5331.

[30]  G. Piro *et al.*, "Simulating LTE Cellular Systems: An Open-Source Framework," *in IEEE Trans. on Vehicular Tech.*, vol. 60, no. 2, pp. 498-513, Feb. 2011..

[31]  Further advancements for E-UTRA physical layer aspects, Release 9, 3GPP TR 36.814, V.9.0.0, March 2010.

[32]  A. Prasad, A. Maeder, C. Ng "Energy Efficient Small Cell Activation Mechanism for Heterogeneous Networks", in Proc. IEEE GLOBECOM 2013, Atlanta, US, Dec. 2013, pp. 754-759.

[33]  K. I. Pedersen, Y. Wang, S. Strzyz and F. Frederiksen, "Enhanced inter-cell interference coordination in co-channel multi-layer LTE-advanced networks," in IEEE Wireless Com., vol. 20, no. 3, pp. 120-127, June 2013.

[34]  Energy Aware Radio and Nework Technologies (EARTH) project. Deliverable D2.3 - Energy efficiency analysis of the reference systems, areas of improvements and target breakdown. http://www.ict-earth.eu/

[35]  A. Virdis, G. Stea, D. Sabella, M. Caretti, "A distributed power-saving framework for LTE Het-Nets exploiting Almost Blank Subframes", Technical Report, University of Pisa, 2017. [Online]. Available: http://eprints.adm.unipi.it/id/eprint/2370

[36]  ILOG CPLEX Software, http://www.ilog.com

[37]  A. Virdis, G. Stea, G. Nardini, "SimuLTE: A Modular System-level Simulator for LTE/LTE-A Networks based on OMNeT++", in Proc. SimulTech 2014, Vienna, AT, August 28-30, 2014

[38]  SimuLTE webpage. http://www.simulte.com

[39]  OMNeT++, http://www.omnetpp.org

[40]  INET framework for OMNeT++: http://inet.omnetpp.org/

[41]  M. W. P. Savelbergh, "Preprocessing and probing techniques for mixed integer programming problems", *ORSA Journal of Computing*, Vol. 6, No. 4, pp. 445-454, Fall 1994

[42]  L. Wolsey, *Integer Programming*, John Wiley and Sons, Inc, 1998

**Antonio Virdis** is a post-doc researcher at the University of Pisa, where he obtained his MSc degree in Computer System Engineering in 2011, and his PhD in Information Engineering in 2015. His research interests include Quality of Service, scheduling and resource allocation in wireless networks, network simulation and performance evaluation. He has been and is currently involved in national, EU-funded and industry-funded research projects. He coauthored six patents and sixteen peer-reviewed papers in the field of cellular network modeling and algorithms.

**Giovanni Stea** is an Associate Professor at the Department of Information Engineering of the University of Pisa, Italy, where he also got his PhD in 2003. His current research interests include quality of service and resource allocation in wireline and wireless networks, performance evaluation through simulation and analytical techniques, traffic engineering. In

these fields he has coauthored more than 80 peer-reviewed papers and 15 patents. He has been involved in national and European research projects, and he has led joint research projects with industrial partners. He has served as a member of the technical and/or organization committees for several international conferences, including SIGCOMM, WoWMoM, and VALUETOOLS, and he is serving on the editorial board of the Springer Wireless Networks journal.

**Dario Sabella** received his MSc degree in electronic engineering from Politecnico di Torino (Italy) in 2002, and a post-degree specialization in Telecommunications in 2004. From 2002 to 2017 he was with TIM (Telecom Italia Group), first as an embedded application programmer and then as a member of the Wireless Innovation department, where he worked on wireless radio access technologies and system engineering for mobile networks toward 5G. In 2017 he joined Intel as Senior Standards and Research Engineer. He has been involved in European projects (FP7, EIT Digital, Horizon 2020), often also with leadership roles. He coauthored more than 20 patents and 40 publications. He has been a delegate within ETSI Energy Efficiency TC, and, more recently, he served as vice chair of ETSI MEC IEG. He is currently serving within ETSI

MEC as ISG Secretary and Lead of Industry Group Relationships.

**Marco Caretti** received a Laurea degree in Electronic Engineering from University of L'Aquila. Since 2000, he has been working with Telecom Italia Lab, the Research center of Telecom Italia group, where he has been involved in several internal R&D projects on performance study of wireless access communication systems (EDGE, UMTS, WiMAX and LTE). His research interests include design and performance study of radio access technologies (in particular LTE and LTE Advanced) and their evolution towards 5G systems. He coauthored several international patents and papers. Since 2005 he has participated in several standardization groups (WiMAX, 3GPP RAN1 and RAN4). He is currently in charge of an internal project focused on virtualized radio access solutions and 5G radio technologies.