# Hierarchical Resource Allocation: Balancing Throughput and Energy Efficiency in Wireless Systems

Bho Matthiesen*, Eduard A. Jorswieck†, and Petar Popovski‡*

*University of Bremen, Department of Communications Engineering, Germany, email: bho.matthiesen@tu-dresden.de
†TU Braunschweig, Department of Information Theory and Communication Systemes, Germany, email: e.jorswieck@tu-bs.de
‡Aalborg University, Department of Electronic Systems, Denmark, email: petarp@es.aau.dk

*Abstract*—A main challenge of 5G and beyond wireless systems is to efficiently utilize the available spectrum and simultaneously reduce the energy consumption. From the radio resource allocation perspective, the solution to this problem is to maximize the energy efficiency instead of the throughput. This results in the optimal benefit-cost ratio between data rate and energy consumption. It also often leads to a considerable reduction in throughput and, hence, an underutilization of the available spectrum. Contemporary approaches to balance these metrics based on multi-objective programming theory often lack operational meaning and finding the correct operating point requires careful experimentation and calibration. Instead, we propose the novel concept of hierarchical resource allocation where conflicting objectives are ordered by their importance. This results in a resource allocation algorithm that strives to minimize the transmit power while keeping the data rate close the maximum achievable throughput. In a typical multi-cell scenario, this strategy is shown to reduces the transmit power consumption by 65% at the cost of a 5% decrease in throughput. Moreover, this strategy also saves energy in scenarios where global energy efficiency maximization fails to achieve any gain over throughput maximization.

*Index Terms*—multi-objective programming, global optimization, hierarchical optimization, mixed monotonic programming

## I. MOTIVATION AND PROBLEM STATEMENT

The goal of resource allocation in communication networks is to best utilize the available resources ensuring good Quality of Service (QoS) to all users. While the QoS constraints are mainly determined by the user's requirements or network slice configuration, the choice of a suitable utility function is entirely up to the operator or system designer [1]–[3]. Common choices are maximizing the throughput (TP) to best utilize the available spectrum [4], minimizing the total transmit power to save energy [5], or maximizing the energy efficiency (EE) to obtain a trade-off between these two [6], [7]. In general, these are conflicting metrics that can not be maximized simultaneously. Indeed, the multi-objective optimization problem (MOP)

$$\max_{\boldsymbol{p} \in \mathcal{P}} \begin{bmatrix} f_1(\boldsymbol{p}), & f_2(\boldsymbol{p}), & \dots \end{bmatrix} \qquad (1)$$

with network utility functions $f_1, f_2, \dots$ is known to posses an infinite number of noninferior solutions [8]. The MOP (1) is usually solved by transforming it into a scalar optimization problem, e.g., with the scalarization approach [9] where the weighted sum of the objectives is maximized, i.e.,

$$\max_{\boldsymbol{p} \in \mathcal{P}} \sum_i w_i f_i(\boldsymbol{p}),$$

or by the utility profile approach [9] where the intersection of a ray in the direction $\boldsymbol{w}$ and the outer boundary of the performance region is computed, i.e.,

$$\max_{t, \boldsymbol{p} \in \mathcal{P}} t \quad \text{s.t.} \quad \forall i : t w_i \leq f_i(\boldsymbol{p}).$$

Both methods obtain Pareto optimal points but share the weakness that the weights $\boldsymbol{w}$ often have no operational meaning and need to be chosen heuristically or by experimentation.

For example, consider balancing the TP with the total transmit power. This problem is formally stated as

$$\begin{cases} \max_{\boldsymbol{p}, \boldsymbol{r}} & \left[ \sum_i r_i, \quad -\sum_i p_i \right] \\ \text{s.t.} & \boldsymbol{r} \in \mathcal{R}(\boldsymbol{p}) \cap \mathcal{Q}, \quad \boldsymbol{0} \leq \boldsymbol{p} \leq \boldsymbol{P} \end{cases} \qquad (2)$$

where $\boldsymbol{P}$ is the maximum transmit power, $\mathcal{R}(\boldsymbol{p})$ the achievable rate region, and $\mathcal{Q}$ contains the QoS constraints. After scalarization, the problem becomes

$$\begin{cases} \max_{\boldsymbol{p}, \boldsymbol{r}} & w_1 \sum_i r_i - w_2 \sum_i p_i \\ \text{s.t.} & \boldsymbol{r} \in \mathcal{R}(\boldsymbol{p}) \cap \mathcal{Q}, \quad \boldsymbol{0} \leq \boldsymbol{p} \leq \boldsymbol{P} \end{cases} \qquad (3)$$

with nonnegative weights $w_1, w_2$. By varying these weights such that $w_1 + w_2 = 1$, the convex hull of the Pareto boundary is obtained. However, these weights do not have much operational meaning and there is no other guidance than experience or experimentation to choose them for a given system. Another approach to balance TP and transmit power is the notion of global energy efficiency (GEE), which is defined as the benefit-cost ratio of system throughput and total dissipated power, i.e., $\text{GEE} = \frac{\sum_i r_i}{\sum_i \mu_i p_i + P_c}$, where $\mu_i \geq 0$ and $P_c > 0$ are modeling constants reflecting the power amplifier inefficiency and static circuit power consumptions. Maximizing the GEE results in a Pareto optimal solution of (2) [7, p. 241] and has a well defined operational meaning. With
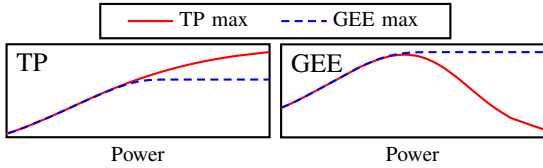
Fig. 1. Typical solution of TP and GEE maximization.

energy and spectrum being similarly scarce resources, the TP and GEE are considered to be the most important network utility functions in 5G and beyond networks.

A qualitative solution of TP and GEE maximization in wireless interference networks is displayed in Fig. 1. While leading to similar operating points in the low signal-to-noise ratio regime, it is characteristic for the GEE to saturate. The link and power budget in a wireless network often allow for an operating point far in this saturation region. In such a scenario, selecting the operating point by TP or GEE maximization either results in poor EE or in low spectral efficiency. Thus, it has been proposed in [10] to balance TP and GEE with multi-objective programming theory. While the obtained performance region provides valuable insights for system design, the weights still have little operational meaning. A more straightforward method is to maximize the GEE under QoS constraints which is expected to provide the best rate-energy trade-off while still providing satisfactory service to all users.

Taking the operator's perspective, saving energy is just a secondary concern, while generating revenue from their costly equipment and spectrum licenses is the primary goal. This requires good service quality to outperform competitors and thereby ensure customer loyalty. Satisfying QoS constraints and providing good connectivity is undoubtedly the foundation for good service but from there it's up to the operator to choose an operating point in the resource allocation design space. A viable strategy is to prioritize high service quality and minimize energy consumption as a secondary objective to reduce operational expenditures and further increase revenue. This could be achieved by solving (3) with $w_1 \gg w_2$. A more rigorous approach is to use lexicographic ordering [11, §4.2], a recursive multi-objective programming technique where objectives are strictly ordered by priority. In the context of this paper and (2), a lexicographic ordering approach is to maximize the TP first and then select the solution with lowest transmit power, i.e., $\boldsymbol{p}^\star = \min\{\sum_i p_i \mid \boldsymbol{p} \in \mathcal{T}^\star\}$ where $\mathcal{T}^\star$ is the set of throughput optimal power allocations, i.e., $\mathcal{T}^\star = \arg\max\{\sum r_i \mid \boldsymbol{r} \in \mathcal{R}(\boldsymbol{p}), \, \boldsymbol{0} \le \boldsymbol{p} \le \boldsymbol{P}\}$. When the solution to the TP maximization problem is (almost) unique, i.e., the volume of $\mathcal{T}^\star$ is close to zero, the possible power reduction due to this approach is negligible. However, significant gains are possible by slightly relaxing this strict ordering of the objectives. For example, the goal could be to achieve at least 95 % of the maximum TP instead of strictly maximizing it, i.e., $\boldsymbol{p}^\star = \min\{\sum_i p_i \mid \sum_i r_i \ge 0.95 \cdot r_\Sigma^\star, \, \boldsymbol{r} \in \mathcal{R}(\boldsymbol{p})\}$, where $r_\Sigma^\star$ is the optimal value of the TP maximization problem. Selecting a power allocation within this tight TP region leaves more

freedom than lexicographic ordering, while still ensuring high service quality.

Leaving economical considerations aside, there are plenty of other technical motivations to strictly prioritize a high TP over other metrics. One application arises from cross-layer optimization where the queue of a base station (BS) needs to be stabilized. Regardless of the underlying queuing model, the total storage capacity is essentially limited by the BS's installed memory. The TP determines the maximum departure rate of this joint queue and, hence, maximizing the TP ultimately enlarges the stability region. Please refer to [4] for further application examples.

The goal of this paper is to obtain a hierarchical Pareto optimal solution of (2) for wireless interference networks, and to evaluate the benefits of this approach over GEE maximization numerically. As the resulting optimization problem is NP-hard and numerically very challenging, this requires the careful design of a solution algorithm. We show that, by reducing the TP by just 5 %, almost 65 % of transmit power can be saved in a typical wireless network.

### A. System Model

We consider a Gaussian interference network with power allocation $\boldsymbol{p} = (p_1, p_2, \dots)$ and average power constraint $\boldsymbol{P}$. The receive signal to interference plus noise ratio (SINR) is $\frac{\alpha_i p_i}{\sum_{j \ne i} \beta_{ij} p_j + \sigma_i^2}$ and, under the assumption that interference is treated as noise, asymptotic error free communication is possible at all rates $\boldsymbol{r}$ satisfying

$$r_i \le B \log \left( 1 + \frac{\alpha_i p_i}{\sum_{j \ne i} \beta_{ij} p_j + \sigma_i^2} \right)$$

for all $i$, where $B$ is the communication bandwidth. In this setting, $\alpha_i$ is the effective channel gain of the direct channel from transmitter $i$ to receiver $i$, $\beta_{ij}$ are the effective channels from transmitter $j$ to receiver $i$, and $\sigma_i^2$ is the variance of circularly-symmetric complex Gaussian noise.

This adequately models the effective channel for multi-antenna transmission in 5G networks after precoder matrix selection [12, §11], for multi-cell networks with overlapping frequencies, and for dense low earth orbit (LEO) satellite constellations [13]. Other applications include, e.g., massive MIMO and relay-assisted CoMP networks [14].

## II. HIERARCHICAL OPTIMIZATION

Hierarchical optimization [11, §4.2.2], [15] is a solution method for the MOP (1) where the objectives are arranged a priori by their absolute importance. Without loss of generality, assume that $f_i$ is more important to the system designer than $f_{i+1}$. The optimization is carried out recursively by first maximizing $f_1$ and ignoring all other objectives $f_2, f_3, \dots$. Then, the next objective $f_2$ is maximized with additional constraint that the value of $f_1$ is close to the optimal value of the previous optimization. Mathematically, the $i$th optimization problem is

$$\max_{\boldsymbol{p} \in \mathcal{D}_i} f_i(\boldsymbol{x}) \quad \text{with} \quad \mathcal{D}_i = \{\boldsymbol{p} \in \mathcal{D}_{i-1} \mid f_{i-1}(\boldsymbol{p}) \ge \omega_{i-1} f_{i-1}^\star\}$$

for all $i > 1$ and some initial feasible set $\mathcal{D}_1$. Here, $f_i^\star$ denotes the optimal value of the $i$th problem and $\omega_1, \omega_2, \dots$

are so-called worsening factors. These are selected a priori by the system designer and have, contrary to the weights in the multi-objective programming solution approaches discussed in Section I, a clearly defined operational meaning in many engineering problems. Lexicographic ordering [11, §4.2] is a special case of this approach obtained by setting all worsening factors to one. For a MOP with two objectives, the second (and final) optimization step is equivalent to the $\varepsilon$-constraint method [11, §3.2] and its solution is a strictly Pareto optimal point if it is unique [11, Thm. 3.2.4].

Applying this approach to the MOP (2) and prioritizing the TP over the transmit power, we obtain two scalar optimization problems[1]

$$
\begin{cases}
\max_{\boldsymbol{p},\boldsymbol{r}} & \sum_i \log\left(1 + \dfrac{\alpha_i p_i}{\sum_{j\neq i}\beta_{ij}p_j + \sigma_i^2}\right) & \text{(4a)} \\[2mm]
\text{s.t.} & \forall i: \log\left(1 + \dfrac{\alpha_i p_i}{\sum_{j\neq i}\beta_{ij}p_j + \sigma_i^2}\right) \geq r_{i,\min} & \text{(4b)} \\[2mm]
& \boldsymbol{0} \leq \boldsymbol{p} \leq \boldsymbol{P} & \text{(4c)}
\end{cases}
$$

for minimum rate constraints $r_{i,\min} \geq 0$, and

$$
\begin{cases}
\min_{\boldsymbol{p},\boldsymbol{r}} & \sum_i p_i & \text{(5a)} \\[2mm]
\text{s.t.} & \sum_i \log\left(1 + \dfrac{\alpha_i p_i}{\sum_{j\neq i}\beta_{ij}p_j + \sigma_i^2}\right) \geq \omega r_\Sigma^\star & \text{(5b)} \\[2mm]
& \forall i: \log\left(1 + \dfrac{\alpha_i p_i}{\sum_{j\neq i}\beta_{ij}p_j + \sigma_i^2}\right) \geq r_{i,\min} & \text{(5c)} \\[2mm]
& \boldsymbol{0} \leq \boldsymbol{p} \leq \boldsymbol{P} & \text{(5d)}
\end{cases}
$$

where $r_\Sigma^\star$ is the optimal value of (4) and $\omega \in [0,1]$ is the worsening factor that determines the acceptable TP reduction. Clearly, it is necessary to solve (4) before (5).

Both problems (4) and (5) are challenging global optimization problems due to the nonconvexity of the objective in (4) and constraint (5b). In particular, (4) is known to be NP-hard [16], and, hence, (5) is also NP-hard due to constraint (5b). While (4) can be solved efficiently using the mixed monotonic programming (MMP) framework as discussed next, problem (5) needs a novel algorithm that is developed in Section III.

### A. Solution of Problem (4)

MMP is a global optimization framework that exploits partial monotonicity in the objective and constraints [17]. It is much more versatile than classical monotonic optimization [18] and shows tremendous performance gains over state-of-the-art algorithms for global optimal power allocation in interference networks and other scenarios [17, §IV].

The concept of mixed monotonic (MM) functions generalizes differences of increasing functions. Let $\mathcal{M}_0$ be a box in $\mathbb{R}^n$, i.e., $\mathcal{M}_0 = [\boldsymbol{r}^0, \boldsymbol{s}^0] = \{\boldsymbol{x} \in \mathbb{R}^n \,|\, \forall i: r_i^0 \leq x_i \leq s_i^0\}$. A continuous function $F : \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}$ is called MM function if it satisfies

$$
\begin{aligned}
F(\boldsymbol{x},\boldsymbol{y}) &\leq F(\boldsymbol{x}',\boldsymbol{y}) && \text{if } \boldsymbol{x} \leq \boldsymbol{x}', \\
F(\boldsymbol{x},\boldsymbol{y}) &\geq F(\boldsymbol{x},\boldsymbol{y}') && \text{if } \boldsymbol{y} \leq \boldsymbol{y}'.
\end{aligned}
$$

[1]The constant $B$ is inessential and moved into $r_{i,\min}$ for notational clarity.

for all $\boldsymbol{x}, \boldsymbol{x}', \boldsymbol{y}, \boldsymbol{y}' \in \mathcal{M}_0$ and a continuous optimization problem $\max_{\boldsymbol{x}\in\mathcal{D}} f(\boldsymbol{x})$ with compact feasible set $\mathcal{D} \subseteq \mathbb{R}^n$ is called MMP problem if there exists an MM function $F$ such that $F(\boldsymbol{x},\boldsymbol{x}) = f(\boldsymbol{x})$ for all $\boldsymbol{x} \in \mathcal{M}_0$, where $\mathcal{M}_0 \supseteq \mathcal{D}$ encloses $\mathcal{D}$. The MMP framework [17] solves such a problem very efficiently with global optimality using a branch and bound (BB) procedure.

Applying the MMP framework requires MM representations of the objective and constraint functions in (4). For the objective, such a function is $\boldsymbol{x}, \boldsymbol{y} \mapsto \sum_i R_i(\boldsymbol{x},\boldsymbol{y})$ with [17, §IV-A]

$$
R_i(\boldsymbol{x},\boldsymbol{y}) = \log\left(1 + \frac{\alpha_i x_i}{\sum_{j\neq i}\beta_{ij}y_j + \sigma_i^2}\right). \tag{6}
$$

Likewise, the QoS constraints have MM representation $\boldsymbol{x}, \boldsymbol{y} \mapsto r_{i,\min} - R_i(\boldsymbol{x},\boldsymbol{y})$. Theoretically, such MM constraints lead to an algorithm without guaranteed finite convergence. This is, because for general MM constraints and some boxes $\mathcal{M}$, it is impossible to determine whether $\mathcal{M} \cap \mathcal{D}$ contains feasible points or not [17, §III-A]. However, in practise this is seldom a problem for typical minimum rate constraints as in (4b).

The MMP framework is also applicable to (5). However, the minimum sum rate constraint in (5b) is very tight and leads to a tiny feasible set compared to $\mathcal{M}^0 = [\boldsymbol{0}, \boldsymbol{P}]$. This results in impractically slow convergence of the MMP procedure. In the next section, we develop an algorithm with much faster and provably finite convergence.

### III. Successive Incumbent Transcending Scheme

The main challenge in solving (5) with the MMP framework is constraint (5b). An efficient solution to this problem is the successive incumbent transcending (SIT) scheme developed in [19]. The main idea is to solve a sequence of easily implementable feasibility problems. Specifically, given a real number $\gamma$, the core problem of the SIT algorithm is to check whether (5) has a feasible solution $\boldsymbol{p}$ satisfying $\sum_i p_i \leq \gamma$, or, else, establish that no such $\boldsymbol{p}$ exists. In this manner, a sequence of feasible points ("incumbents") with decreasing objective value is generated until no point with lesser objective value than the current best solution $\gamma$ exists.

Consider the optimization problem

$$
\min_{\boldsymbol{x}\in\mathcal{M}_0} f(\boldsymbol{x}) \quad \text{s.t.} \quad g(\boldsymbol{x}) \leq 0 \tag{7}
$$

which generalizes (5) and assume that $f$ is a nondecreasing function, $g$ has an MM representation, and $\mathcal{M}_0$ is a box. The outlined SIT scheme for this problem is given in Algorithm 1.

---

**Algorithm 1** SIT Scheme [20, Sect. 7.5.1]

---

**Step 0** Initialize $\bar{\boldsymbol{x}}$ with the best known feasible solution and set $\gamma = f(\bar{\boldsymbol{x}}) - \eta$; otherwise do not set $\bar{\boldsymbol{x}}$ and choose some $\gamma \leq f(\boldsymbol{x}) \,\forall \boldsymbol{x} \in \mathcal{M}_0 : g(\boldsymbol{x}) \leq 0$.

**Step 1** Check if (5) has a feasible solution $\boldsymbol{x}$ satisfying $f(\boldsymbol{x}) \geq \gamma$; otherwise, establish that no such feasible $\boldsymbol{x}$ exists and go to Step 3.

**Step 2** Update $\bar{\boldsymbol{x}} \leftarrow \boldsymbol{x}$ and $\gamma \leftarrow f(\bar{\boldsymbol{x}}) - \eta$. Go to Step 1.

**Step 3** Terminate: If $\bar{\boldsymbol{x}}$ is set, it is an $\eta$-optimal solution; else Problem (5) is infeasible.

---

Implementing the feasibility check in Step 1 of Algorithm 1 efficiently is crucial. Consider the optimization problem

$$\min_{\boldsymbol{x} \in \mathcal{M}_0} \quad g(\boldsymbol{x}) \quad \text{s.\,t.} \quad f(\boldsymbol{x}) \leq \gamma \tag{8}$$

which is dual to (7) in the sense that if the optimal value of (8) is greater than zero, the optimal value of (7) is greater than $\gamma$ [20, Prop. 7.13]. Thus, any point $\boldsymbol{x}'$ in the feasible set of (8) with objective value less than zero is also a feasible point in (7) with objective value less than $\gamma$. We can solve (7) sequentially by solving (8) with a BB method.

At first, this approach seems to increase the computational complexity significantly because if (7) is nonconvex, then so is (8). However, given that $f$ has favorable properties,[2] problem (8) might be considerably easier to solve than (7). Moreover, the SIT scheme can be combined with the BB procedure that solves (8). This eliminates the need to solve (8) multiple times.

Exploiting the properties of MM functions, we can obtain a lower bound on the objective value of (8) over a box $\mathcal{M} = [\boldsymbol{r}, \boldsymbol{s}]$ from its MM representation $G$ as

$$\min_{\boldsymbol{x} \in \mathcal{M}: f(\boldsymbol{x}) \leq \gamma} g(\boldsymbol{x}) \geq \min_{\boldsymbol{x} \in \mathcal{M}} G(\boldsymbol{x}, \boldsymbol{x}) \geq \min_{\boldsymbol{x}, \boldsymbol{y} \in \mathcal{M}} G(\boldsymbol{x}, \boldsymbol{y}) = G(\boldsymbol{r}, \boldsymbol{s}).$$

Together with an exhaustive rectangular subdivision [20], this bound leads to a convergent BB procedure that can be incorporated into the SIT scheme.

The complete algorithm is stated in Algorithm 2. It involves a parameter $\varepsilon$ that is related to the concept of $\varepsilon$-essential feasibility explained in [21]. Its primary roles are to exclude numerically instable points from the feasible set and ensure finite convergence of the algorithm. The latter is established in the theorem below. This is the first algorithm that combines the MMP approach with the SIT scheme.

*Theorem 1:* Algorithm 2 converges in finitely many steps to the $(\varepsilon, \eta)$-optimal solution of (7) or establishes that no such solution exists.

*Proof sketch:* By virtue of [20, Prop. 7.14] a BB procedure for solving (8) with pruning criterion $G(\boldsymbol{r}, \boldsymbol{s}) > -\varepsilon$ and stopping criterion $g(\boldsymbol{r}) < 0$ or $\mathcal{R}_k = \emptyset$ implements Step 1 in Algorithm 1. Thus, start with the MMP algorithm in [17, Alg. 1] for (8) and modify it according to the previous sentence. Establishing finite convergence is a minor modification of [17, Thm. 1]. Next, integrate the SIT scheme in Algorithm 1 into this procedure: move the termination criterion $g(\boldsymbol{r}) < 0$ into the incumbent update in Step 3 and update $\gamma_k$ if a box satisfies this criterion. It remains to show that continuing the procedure after updating $\gamma_k$ preserves convergence. This part of the proof follows along the lines of the proof of [21, Thm. 1]. ∎

The purpose of the reduction in Step 2 is to speed up the convergence. This is achieved by replacing the box under consideration by a smaller one that still contains all candidate solutions and, thereby, improves the quality of the computed bounds. One approach to determine this procedure for Algorithm 2 is to replace $\mathcal{M}$ by $\mathcal{M}' = [\boldsymbol{r}', \boldsymbol{s}']$ with

$$r_i' = \min_{\boldsymbol{x} \in \mathcal{M}: f(\boldsymbol{x}) \leq \gamma_k} x_i, \qquad s_i' = \max_{\boldsymbol{x} \in \mathcal{M}: f(\boldsymbol{x}) \leq \gamma_k} x_i \tag{10}$$

[2]Such favorable properties could be, e.g., linearity, convexity, or being increasing.

---

**Algorithm 2** SIT Algorithm for (7)

**Step 0 (Initialization)** Set $\varepsilon, \eta > 0$, Let $k = 1$ and $\mathcal{R}_0 = \{\mathcal{M}_0\}$. If available, initialize $\bar{\boldsymbol{x}}^0$ with the best known feasible solution and set $\gamma_k = f(\bar{\boldsymbol{x}}) - \eta$. Otherwise, do not set $\bar{\boldsymbol{x}}^0$ and choose $\gamma \geq f(\boldsymbol{x})$ for all feasible $\boldsymbol{x}$.

**Step 1 (Branching)** Let $\mathcal{M}_k = [\boldsymbol{r}^k, \boldsymbol{s}^k]$ be the oldest box in $\mathcal{R}_{k-1}$. Bisect $\mathcal{M}_k$ via $(\boldsymbol{v}^k, j_k)$ with $j_k \in \arg\max_j s_j^k - r_j^k$ and $\boldsymbol{v}^k = \frac{1}{2}(\boldsymbol{s}^k + \boldsymbol{r}^k)$, i.e., compute

$$\mathcal{M}^- = \{\boldsymbol{x} \mid r_j^k \leq x_j \leq v_j^k, \ r_i^k \leq x_i \leq s_i^k \ (i \neq j)\}$$
$$\mathcal{M}^+ = \{\boldsymbol{x} \mid v_j^k \leq x_j \leq s_j^k, \ r_i^k \leq x_i \leq s_i^k \ (i \neq j)\},$$

and set $\mathcal{P}_k = \{\mathcal{M}_-^k, \mathcal{M}_+^k\}$.

**Step 2 (Reduction)** Replace each box in $\mathcal{M} \in \mathcal{P}_k$ with some $\mathcal{M}'$ such that $\mathcal{M}' \subseteq \mathcal{M}$ and

$$\min\{g(\boldsymbol{x}) \mid f(\boldsymbol{x}) \leq \gamma_k, \ \boldsymbol{x} \in \mathcal{M}\}$$
$$= \min\{g(\boldsymbol{x}) \mid f(\boldsymbol{x}) \leq \gamma_k, \ \boldsymbol{x} \in \mathcal{M}'\} \tag{9}$$

**Step 3 (Incumbent)** Let $\mathscr{I} = \{\boldsymbol{r} \mid [\boldsymbol{r}, \boldsymbol{s}] \in \mathcal{P}_k, \ g(\boldsymbol{r}) \leq 0\}$. If not empty, set $\boldsymbol{r}^k = \arg\min_{\boldsymbol{r} \in \mathscr{I}} f(\boldsymbol{r})$. If $\bar{\boldsymbol{x}}^{k-1}$ is not set or $f(\boldsymbol{r}^k) < \gamma_{k-1} + \eta$, set $\bar{\boldsymbol{x}} = \boldsymbol{r}^k$ and $\gamma_k = f(\boldsymbol{r}^k) - \eta$. In all other cases, set $\bar{\boldsymbol{x}}^k = \bar{\boldsymbol{x}}^{k-1}$ and $\gamma_k = \gamma_{k-1}$.

**Step 4 (Pruning)** Delete every $[\boldsymbol{r}, \boldsymbol{s}] \in \mathcal{P}_k$ with $f(\boldsymbol{r}) \geq \gamma_k$ or $G(\boldsymbol{r}, \boldsymbol{s}) > -\varepsilon$. Let $\mathcal{P}_k'$ be the collection of remaining sets and set $\mathcal{R}_k = \mathcal{P}_k' \cup (\mathcal{R}_{k-1} \setminus \{\mathcal{M}_k\})$.

**Step 5 (Termination)** Terminate if $\mathcal{R}^k = \emptyset$: If $\bar{\boldsymbol{x}}^k$ is not set, then (7) is $\varepsilon$-essential infeasible; else $\bar{\boldsymbol{x}}^k$ is an essential $(\varepsilon, \eta)$-optimal solution of (7). Otherwise, update $k \leftarrow k + 1$ and return to Step 1.

---

for all $i$. For $f$ nondecreasing, the solution to the first problem is always $r_i$ unless it is infeasible. For the upper bound in (10), recall that $\boldsymbol{r}$ minimizes $f(\boldsymbol{x})$ over $\mathcal{M}$. Thus, the optimal solution to this optimization problem is to set $x_j = r_j$ for all $j \neq i$. Then, the optimal $x_i = \min\{\tilde{x}_i, s_i\}$ where $\tilde{x}_i$ satisfies

$$f(\boldsymbol{r} + (\tilde{x}_i - r_i)\boldsymbol{e}_i) = \gamma_k. \tag{11}$$

*Remark 1 (Branch selection):* Most BB procedures select the box with the largest bound for further partitioning. The rationale is that this choice leads to fastest convergence. In practice, when the number of boxes in $\mathcal{R}_k$ grows very large, this selection rule might become the performance and memory bottleneck of the algorithm. First, it tends to store suboptimal boxes longer than necessary and therefore increases memory consumption. Second, inserting new boxes into $\mathcal{R}_k$ has complexity $O(\log |\mathcal{R}_k|)$. Instead, with the *oldest-first* rule employed in Algorithm 2 inserting new boxes has constant complexity. Also, every box is visited after a fixed amount of time and, thus, likely to be pruned much earlier than with the best-first rule [17]. Since Algorithm 2 is essentially memory limited, the oldest-first rule performs much better than the standard best-first rule.

*Remark 2 (Other SIT applications):* Despite its tremendous numerical advantages, the SIT approach is currently not widely used. Besides the applications to DC and monotonic optimization problems in [19], [20], it is only employed in [21] where it is applied to resource allocation problems with fractional objectives and partial convexity. The implementation most closely related to Algorithm 2 is the monotonic optimization

variant in [20, §11.3]. The key advantage of Algorithm 2 over this procedure is that cumbersome transformations and an auxiliary variable are required to bring (5) into a suitable form for [20, §11.3]. This leads to much slower convergence due to the extra variable and much looser bounds on the constraints.

### A. Solution of Problem (5)

Identify $\mathcal{M}_0 = [\mathbf{0}, \boldsymbol{P}]$ and $f(\boldsymbol{p}) = \sum_i p_i$. Note that $f(\boldsymbol{p})$ is an increasing function. MM representations of (5b) and (5c) are $\boldsymbol{x}, \boldsymbol{y} \mapsto \omega r_\Sigma^\star - \sum_i R_i(\boldsymbol{y}, \boldsymbol{x})$ and $\forall i : \boldsymbol{x}, \boldsymbol{y} \mapsto r_{i,\min} - R_i(\boldsymbol{y}, \boldsymbol{x})$, respectively, with $R_i(\boldsymbol{x}, \boldsymbol{y})$ as in (6). They can be merged into a single inequality constraint $\max_i g_i(\boldsymbol{x}) \le 0$ with MM representation

$$G(\boldsymbol{x}, \boldsymbol{y}) = \max\left\{\omega r_\Sigma^\star - \sum_i R_i(\boldsymbol{y}, \boldsymbol{x}),\ \max_i\left\{r_{i,\min} - R_i(\boldsymbol{y}, \boldsymbol{x})\right\}\right\}$$

due to [17, Eq. (9)]. In the reduction step, the solution to (11) is $\tilde{x}_i = \gamma_k - \sum_{j \ne i} r_j$. Thus, every box $\mathcal{M} = [\boldsymbol{r}, \boldsymbol{s}]$ in Step 2 can be replaced by $[\boldsymbol{r}, \boldsymbol{s}']$ with $s_i' = \min\{s_i, \gamma_k - \sum_{j \ne i} r_j\}$. With these choices, Algorithm 2 solves (5) in a finite number of iterations.

### IV. NUMERICAL EVALUATION

We consider uplink transmission in a single-input single-output multi-cell system. User equipments (UEs) are placed randomly in a rectangular area with edge length 1 km. This area is divided into four equal sized cells with BSs located at the center of their cell. Path-loss is modeled according to the Hata-COST231 [22], [23] urban scenario with carrier frequency 1.9 GHz, 30 m BS height and 8 dB log-normal shadow fading. Small scale effects are modeled as Rayleigh fading. Each UE is associated to the BS with the best channel. Scenarios where more than one UE is associated to a BS are dropped. The receivers have noise spectral density $N_0 = -174$ dBm and noise figure $F = 3$ dB. The communication bandwidth is $B = 180$ kHz and the noise power is calculated as $\sigma_i^2 = N_0 F B$. The UEs RF chains have a static power consumption $P_c = 400$ mW and power amplifiers with an efficiency of 25 %. No cooperation between BSs is assumed, i.e., interference from other cells is treated as noise.

TP and GEE are maximized using the MMP framework [17]. Algorithm 2 is used to solve (5) for $\omega = 0.95$, i.e., the obtained resource allocation uses the minimum total transmit power under the constraint that the system TP is not less than 95 % of the maximum achievable system TP. We call this resource allocation high throughput energy efficiency (HTEE) for reasons that will become apparent below. All algorithms obtain the global optimal solution within an absolute tolerance of $\eta = 0.01$. In Algorithm 2, we set $\varepsilon = 10^{-5}$. All results are averaged over 1000 independent and identically distributed channel realizations.

Figures 2 and 3 display TP and GEE, respectively, with very typical behavior. With increasing transmit power budget, the maximum TP increases. Instead, the GEE saturates at some point and the transmit power stays constant in the GEE optimal resource allocation. Increasing the transmit power beyond the GEE saturation point, as is done in the TP optimal
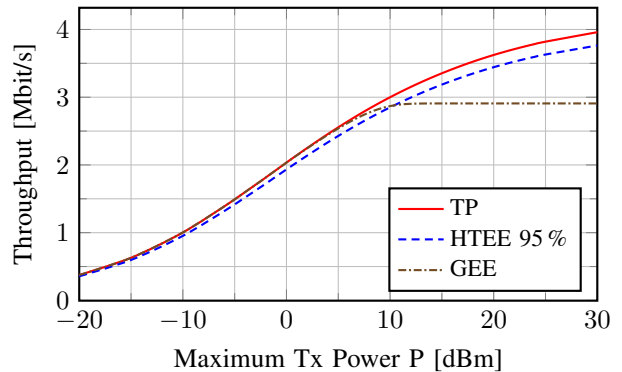


Fig. 2. Achievable throughput with different resource allocation approaches.
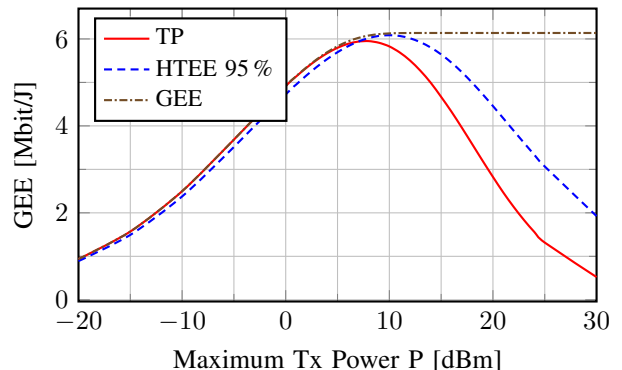


Fig. 3. Global energy efficiency of the discussed resource allocation strategies.

allocation, decreases the GEE. For a maximum transmit power of 23 dBm, which corresponds to the typical UE power budget [12], [24], the GEE optimal allocation achieves 22.4 % or 0.84 Mbit/s less TP than possible. Instead, the HTEE resource allocation is within 95 % of the maximum achievable TP and achieves a 97 % higher GEE than the maximum TP allocation at 23 dBm. This corresponds to a gain of 1.8 Mbit/J at the cost of 0.19 Mbit/s.

However, the GEE is not the optimal metric to evaluate transmit power savings. Consider a second operating point at $-10$ dBm, the median transmit power of 4G UEs in urban scenarios [24]. The TP and GEE optimal strategies both achieve almost the same TP and GEE. Figure 4 displays the power consumption relative to the TP optimal resource allocation. It can be observed that the GEE strategy consumes almost as much transmit power as the TP strategy, and, thus, is unable to exploit the "rate reduction budget" of the system designer. Instead, the HTEE strategy uses almost 40 % less transmit power at a TP cost of 50 kbit/s, which is less than the data rate of classical digital telephone line modem. Nevertheless, its GEE is worse than that of the other strategies, despite the tremendous transmit power reduction.

Returning to our previous scenario with 23 dBm maximum transmit power, it can be seen from Fig. 4 that the HTEE strategy consumes only 35.7 % of the transmit power necessary to achieve the maximum TP. Of course the GEE strategy saves even more transmit power but at a much higher cost to the throughput. This trade-off is illustrated in Fig. 5 where the
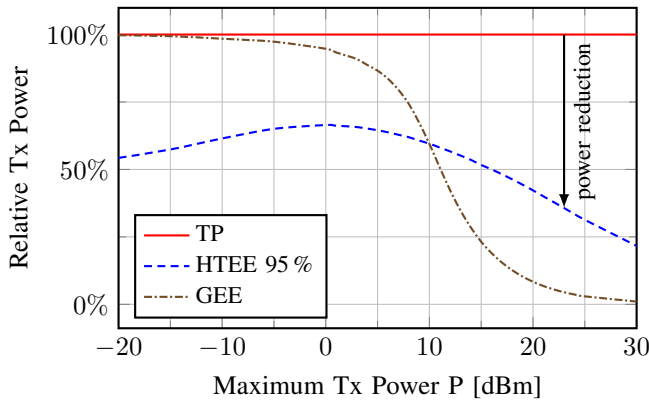
Fig. 4. Total power consumption relative to the throughput optimal strategy.
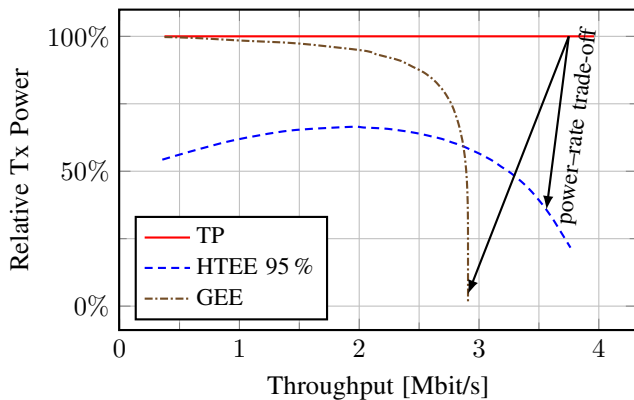


Fig. 5. Relative power consumption as a function of the achievable throughput.

relative transmit power is plotted over the achievable data rate. It can be observed that a major advantage of the HTEE strategy over the GEE optimal power allocation is that the TP does not saturate and any data rate is achievable given a sufficient transmit power budget. Thus, it results in an energy-efficient resource allocation while still ensuring high TP.

Finally, to support the statement at the end of Section II-A that (5) is hard to solve with a traditional BB method, we have also employed the MMP framework to solve (5). Out of 1000 problem instances that ran on an Intel Xeon E5-2680 v3 CPU with a memory usage limit of 21 GB, 483 problem instances ran out of memory and 517 problem instances did not complete within 24 hours, i.e., not a single problem instance of (5) could be solved by a traditional BB algorithm with reasonable usage of computational resources. In contrast, the same problem instances could be solved with Algorithm 1 using a maximum of 50 MB memory and not taking longer than 752 ms to complete. The median computation time among all problem instances was 1.75 ms.

## V. Conclusions

We have introduced the novel concept of hierarchical resource allocation and applied it to minimize energy consumption while still ensuring high spectrum utilization. The numerical results show a transmit power reduction of 65 % in a multi-cell communication scenario at the cost of a 5 % drop

in TP. Instead, state-of-the-art GEE maximization results in a TP reduction of almost 25 %. Moreover, this strategy also saves energy in scenarios where GEE optimization fails to provide a gain over TP maximization. The developed algorithms solve the involved optimization problems with global optimality and, therefore, rigorously demonstrate the gains of hierarchical resource allocation and high-throughput energy efficiency maximization over state-of-the-art approaches.

## References

[1] Z. Han and K. J. R. Liu, *Resource Allocation for Wireless Networks: Basics, Techniques, and Applications*. Cambridge Univ. Press, 2008.

[2] H. Zhang, N. Liu, X. Chu, K. Long, A. Aghvami, and V. C. M. Leung, "Network slicing based 5G and future mobile networks: Mobility, resource management, and challenges," *IEEE Commun. Mag.*, vol. 55, no. 8, pp. 138–145, Aug. 2017.

[3] P. Popovski, K. F. Trillingsgaard, O. Simeone, and G. Durisi, "5G wireless network slicing for eMBB, URLLC, and mMTC: A communication-theoretic view," *IEEE Access*, vol. 6, pp. 55 765–55 779, Sep. 2018.

[4] P. C. Weeraddana, M. Codreanu, M. Latva-aho, A. Ephremides, and C. Fischione, *Weighted Sum-Rate Maximization in Wireless Networks: A Review*, ser. FnT Netw. Now, 2012, vol. 6, no. 1-2.

[5] S. A. Grandhi, R. Vijayan, D. J. Goodman, and J. Zander, "Centralized power control in cellular radio systems," *IEEE Trans. Veh. Technol.*, vol. 42, no. 4, pp. 466–468, Nov. 1993.

[6] C. Isheden, Z. Chong, E. Jorswieck, and G. Fettweis, "Framework for link-level energy efficiency optimization with informed transmitter," *IEEE Trans. Wireless Commun.*, pp. 1–12, 2012.

[7] A. Zappone and E. Jorswieck, *Energy Efficiency in Wireless Networks via Fractional Programming Theory*, ser. FNT in Communications and Information Theory. Now Publishers, 2015, vol. 11, no. 3-4.

[8] L. A. Zadeh, "Optimality and non-scalar-valued performance criteria," *IEEE Trans. Autom. Control*, vol. 8, no. 1, pp. 59–60, Jan. 1963.

[9] R. Zhang and S. Cui, "Cooperative interference management with MISO beamforming," *IEEE Trans. Signal Process.*, vol. 58, no. 10, pp. 5450–5458, Oct. 2010.

[10] O. Aydin, E. A. Jorswieck, D. Aziz, and A. Zappone, "Energy-spectral efficiency tradeoffs in 5G multi-operator networks with heterogeneous constraints," *IEEE Trans. Wireless Commun.*, vol. 16, no. 9, pp. 5869–5881, Sep. 2017.

[11] K. Miettinen, *Nonlinear Multiobjective Optimization*. Springer, 1999.

[12] E. Dahlman, S. Parkvall, and J. Sköld, *5G NR: The Next Generation Wireless Access Technology*, 1st ed. Academic Press, 2018.

[13] I. Leyva-Mayorga *et al.*, "LEO small-satellite constellations for 5G and beyond-5G communications," *IEEE Access*, vol. 8, Oct. 2020.

[14] A. Zappone, L. Sanguinetti, G. Bacci, E. Jorswieck, and M. Debbah, "Energy-efficient power control: A look at 5G wireless technologies," *IEEE Trans. Signal Process.*, vol. 64, no. 7, pp. 1668–1683, Apr. 2016.

[15] D. Bestle and P. Eberhard, "Dynamic system design via multicriteria optimization," in *Multiple Criteria Decision Making: Proceedings of the Twelth International Conference*, ser. Lecture Notes in Economics and Mathematical Systems, G. Fandel and T. Gal, Eds. Springer, 1997.

[16] Z.-Q. Luo and S. Zhang, "Dynamic spectrum management: Complexity and duality," *IEEE J. Sel. Areas Commun.*, vol. 2, no. 1, Feb. 2008.

[17] B. Matthiesen, C. Hellings, E. A. Jorswieck, and W. Utschick, "Mixed monotonic programming for fast global optimization," *IEEE Trans. Signal Process.*, vol. 68, pp. 2529–2544, Mar. 2020.

[18] H. Tuy, "Monotonic optimization: Problems and solution approaches," *SIAM J. Optimization*, vol. 11, no. 2, pp. 464–494, Feb. 2000.

[19] ——, "$\mathcal{D}(\mathcal{C})$-optimization and robust global optimization," *J. Global Optim.*, vol. 47, no. 3, pp. 485–501, Oct. 2009.

[20] ——, *Convex Analysis and Global Optimization*. Springer, 2016.

[21] B. Matthiesen and E. A. Jorswieck, "Efficient global optimal resource allocation in non-orthogonal interference networks," *IEEE Trans. Signal Process.*, vol. 67, no. 21, pp. 5612–5627, Nov. 2019.

[22] 3GPP, "Digital cellular telecommunications systems (phase 2+); radio network planning aspects," Tech. Rep. TR 43.030 V9.0.0 R9, Feb. 2010.

[23] T. S. Rappaport, *Wirless Communications*, 2nd ed. Prentice-Hall, 2002.

[24] P. Joshi, D. Colombi, B. Thors, L.-E. Larsson, and C. Törnevik, "Output power levels of 4G user equipment and implications on realistic RF EMF exposure assessments," *IEEE Access*, vol. 5, pp. 4545–4550, Mar. 2017.