

Radio resource management for the UMTS enhanced uplink in presence of QoS radio bearers

Andreas Mäder · Dirk Staehle

Published online: 13 September 2008
© Springer Science+Business Media, LLC 2008

Abstract The next step in the evolution of UMTS is the Enhanced Uplink or high speed uplink packet access (HSUPA), which is designed for the efficient transport of packet switched data. We propose an analytic modeling approach for the performance evaluation of the UMTS uplink with best-effort users over the enhanced uplink and QoS-users over dedicated channels. The model considers two different scheduling disciplines for the enhanced uplink: parallel scheduling and one-by-one scheduling. Resource Management in such a system has to consider the requirements of the dedicated channel users and the enhanced uplink users on the shared resource, i.e. the cell load. We evaluate the impact of two resource management strategies, one with preemption for dedicated channels and one without, on key QoS-indicators like blocking and dropping probabilities as well as user and cell throughput.

Keywords WCDMA · UMTS · Enhanced uplink · HSUPA · Radio resource management · Performance modeling · Radio network planning

1 Introduction and related work

The enhanced uplink (often also referred to as high speed uplink packet access—HSUPA) is the next step in the evolution process of the UMTS. Introduced with UMTS release 6 and specifically designed for the transport of packet switched data, it promises higher throughput, reduced packet delay and a more efficient radio resource utilization. The enhanced uplink introduces a new transport channel, the Enhanced-DCH (E-DCH) and three new signaling channels. The E-DCH can be seen as an “packet-optimized” version of the DCH. The major new features are: hybrid ARQ (automatic repeat request), implemented similarly as in the high speed downlink packet access (HSDPA), NodeB-controlled fast scheduling, and reduced transport time intervals (TTI) of 2 ms. In a UMTS network with enhanced uplink it

A. Mäder (✉) · D. Staehle
Dept. of Distributed Systems, University of Wuerzburg, Am Hubland, 97074 Würzburg, Germany
e-mail: maeder@informatik.uni-wuerzburg.de

D. Staehle
e-mail: dstaehle@informatik.uni-wuerzburg.de

is expected that QoS users will use the Rel. 99 dedicated channel (DCH) while best-effort users will use the new enhanced dedicated channel (E-DCH). A detailed overview can be found e.g. in 3GPP (2005a) or Parkvall et al. (2005).

The relocation of the scheduling control from the RNC to the NodeB introduces a new flexibility into the UMTS air interface, since it enables the vendor or operator to implement scheduling mechanisms which are between two fundamentally different scheduling paradigms: one-by-one scheduling and parallel scheduling. In Kumaran and Qian (2003) it is shown that the best scheduling strategy in terms of throughput is to schedule users which cannot utilize the total radio resource due to transmit power constraints in parallel, and the rest in a one-by-one manner. A similar conclusion has been found in Jäntti and Kim (2001) by means of dynamic programming. Both works (and this work as well) assume that an uplink synchronization mechanism exists which avoids that scheduled transmissions are interfering with each other. In Mäder and Staehle (2006), the authors investigate the impact of one-by-one and parallel scheduling on the performance of the enhanced uplink with an analytical queuing model which considers best-effort traffic as data flows.

The radio resource management (RRM) for the enhanced uplink benefits from the possibility to react faster on load variations in the cell. However, the question is how to perform RRM in a heterogeneous environment with enhanced uplink users which primarily generate best-effort traffic and QoS users which use dedicated channels (DCH) without rate control. Here, a vital part plays the admission control (AC): Should incoming QoS connections have precedence before already existing best-effort connections, or is it more beneficial to treat both connection types equal? The first policy is generally called *preemptive*, since a QoS connection may preempt resources before E-DCH connections. The second is a *preserving* policy since E-DCH connections are not dropped from the system but at most experience a slow-down (Mäder and Staehle 2007). Related work which can be found in the literature is e.g. Altman (2002), where a queueing analysis for the CDMA uplink with best-effort services is presented. A similar approach has been taken in Fodor and Telek (2005), which introduces a dynamic slow down approach for the best-effort users. In its successor (Fodor et al. 2006) the analysis is extended for multiple cells. Preemption for QoS-users is e.g. considered in Litjens and Boucherie (2002) for a GPRS/HSCSD system. In Boxma et al. (2006), the downlink of an integrated system with a preemptive admission control policy has been investigated.

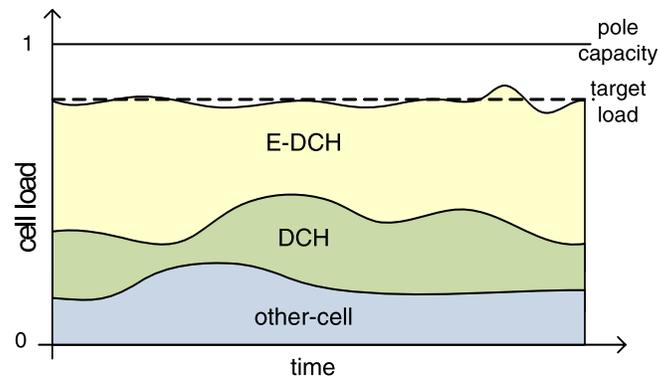
In this work we develop a capacity model for the enhanced uplink and investigate different options for scheduling and radio resource management. Our model is specific for the UMTS enhanced uplink such that it could be used in network planning or optimization. We look at the system on flow level, which means that we consider data traffic regardless of the protocol and content as a continuous flow of data, assuming that the E-DCH users use best effort applications which generate elastic traffic.

The rest of this paper is organized as follows: in Sect. 2 we define the basic radio resource management strategy which provides the frame for our calculations. This forms the base for the interference and cell load model in Sect. 3. In Sect. 4, we describe the E-DCH rate assignment and in Sect. 5 the admission control mechanism, which is then used for a queueing model approach in Sect. 6. In Sect. 7, we show some numerical examples and finally we conclude the paper with Sect. 8.

2 Radio resource management for the E-DCH best effort service

Radio resource management (RRM) for the E-DCH users is primarily done in the NodeBs, which control the maximum transmit power of the mobiles and therefore also the maximum

Fig. 1 Illustration of the RRM principles for the E-DCH best-effort service



user bit rate. The NodeBs send scheduling grants on the absolute or relative grant channel (AGCH and RGCH, resp.), which either set the transmit power to an absolute value or relative to the current value. The mobiles then choose the transport block size (TBS) which is most suitable to the current traffic situation and which does not exceed the maximum transmit power. The grants can be sent every TTI, i.e. every 2 ms, which enables a very fast reaction to changes of the traffic or radio conditions. Grants can be received from the serving NodeB and from non-serving NodeBs. However the latter may just send relative DOWN grants to reduce the other-cell interference in their cells. In our model, we consider grants from the serving NodeB only.

Generally, the WCDMA uplink is interference limited. Therefore, following Holma and Toskala (2001), we define the load in a cell as

$$\hat{\eta} = \frac{\hat{I}_D + \hat{I}_E + \hat{I}_{oc}}{\hat{I}_0 + W\hat{N}_0}, \quad (1)$$

with \hat{I}_D and \hat{I}_E as received powers from the DCH and E-DCH users¹ within the cell, \hat{I}_{oc} as other-cell interference from mobiles in adjacent cells, W as system chip rate, \hat{N}_0 as thermal noise power spectral density and $\hat{I}_0 = \hat{I}_D + \hat{I}_E + \hat{I}_{oc}$. It can be readily seen that this load definition allows the distinction of the cell load after its origin, hence we define

$$\hat{\eta} = \frac{\hat{I}_D}{\hat{I}_0 + W\hat{N}_0} + \frac{\hat{I}_E}{\hat{I}_0 + W\hat{N}_0} + \frac{\hat{I}_{oc}}{\hat{I}_0 + W\hat{N}_0} = \hat{\eta}_D + \hat{\eta}_E + \hat{\eta}_{oc} \quad (2)$$

subject to $\hat{\eta} < 1$. The goal of the RRM is now twofold: first, the cell load should be below a certain maximum load in order to prevent outage. Second, the RRM tries to maximize the resource utilization in the cell to provide high service qualities to the users. The second goal allows also the interpretation of the maximum load as a target load, which should be met as close as possible. Since the DCH-load and the other-cell load cannot be influenced in a satisfying way, the E-DCH load can be used as a means to reach the target cell load. The fast scheduling gives operators the means to use the E-DCH best-effort users for “water-filling” the cell² load at the NodeBs up to a desired target. This radio resource management strategy is illustrated in Fig. 1. The total cell load comprises the varying other-cell load, the load generated by DCH users and the E-DCH load. The received power for the E-DCH users

¹Variable \hat{x} is in linear and x is in dB scale.

²Corresponding to a sector in case of multiple sectors per NodeB.

is adapted such that the total cell load is close to the maximum load. However, due to the power control error and the other-cell interference there is always the possibility of a load “overshoot”. The probability for such an event should be kept low.

The cell load is a random variable due to fast fluctuation of the received E_b/N_0 values. We define that the goal of the RRM is to keep the probability of the total cell load below a maximum tolerable probability p_t :

$$P\{\hat{\eta} \geq \hat{\eta}^*\} \leq p_t. \quad (3)$$

This means that the received signal power (i.e. the E-DCH interference) of the E-DCH users depends on the amount of dedicated channel and other-cell interference. More precisely, the E-DCH users are slowed down if the DCH or the other-cell load is growing, or are speed up, if more radio resources are available for the E-DCH users. If we now assume that the buffers in the mobiles of the E-DCH users are always saturated, we can use this relation to calculate the grade-of-service the E-DCH users receive depending on the scheduling strategy.

3 Interference and load model

Let us consider a NodeB in a UMTS network corresponding to a single sector or cell, respectively. The cell serves a number of DCH users, each connected with a service class $s \in \mathcal{S}$. The service classes are defined by bitrate and target- E_b/N_0 -value. Additionally, n_E E-DCH users are in the system. The state vector \bar{n} comprises the users per DCH service class, n_s , and the E-DCH users n_E :

$$\bar{n} = (n_1, \dots, n_{|\mathcal{S}|}, n_E). \quad (4)$$

Each mobile power controlled by the NodeB perceives a bit-energy-to-noise ratio (E_b/N_0), which is given by

$$\hat{\varepsilon}_k = \frac{W}{R_k} \frac{\hat{S}_k}{W\hat{N}_0 + \hat{I}_0 - \hat{S}_k}. \quad (5)$$

In this equation, W is the chip rate of 3.84 Mcps, R_k is the radio bearer information bit rate, \hat{N}_0 is the thermal noise power density, \hat{S}_k is the received power of mobile k and \hat{I}_0 is the multiple-access interference (MAI) including the own- and other-cell interference.

We assume imperfect power control, so the received E_b/N_0 is a lognormally distributed r.v. with the target- E_b/N_0 -value ε_k^* as mean value (Viterbi and Viterbi 1993) and parameters $\mu = \varepsilon_k^* \cdot \frac{\ln(10)}{10}$ and $\sigma = \text{Std}[\varepsilon_k] \cdot \frac{\ln(10)}{10}$. The received power of each mobile is calculated from (5) as

$$\hat{S}_k = \hat{\omega}_k \cdot (W\hat{N}_0 + \hat{I}_0) \quad \text{with} \quad \hat{\omega}_k = \frac{\hat{\varepsilon}_k R_k}{W + \hat{\varepsilon}_k R_k}. \quad (6)$$

We define the r.v. ω_k as *service load factor* (SLF) depending on the bit rate and the E_b/N_0 -value. The sum of all concurrently received powers constitutes the received own-cell interference, i.e.

$$\hat{I}_D(\bar{n}) = \sum_{s \in \mathcal{S}} \sum_{k \in n_s} \hat{S}_k \quad \text{and} \quad \hat{I}_E(\bar{n}) = \sum_{j \in n_E^a} \hat{S}_j. \quad (7)$$

\hat{I}_D is the total received power of the DCH users and \hat{I}_E of the E-DCH users. Note that the set of currently active E-DCH users n_E^a depends on the scheduling discipline. For parallel

scheduling, $n_E^a = n_E$, since all users are concurrently active. For one-by-one scheduling, $|n_E^a| = 1$ since in this case only one E-DCH user is transmitting at the same time.

The substitution of \hat{I}_D and \hat{I}_E in (2) with (7) gives us the load definitions depending on \bar{n} :

$$\hat{\eta}_D(\bar{n}) = \sum_{s \in \mathcal{S}} \sum_{k \in n_s} \hat{\omega}_k \quad \text{and} \quad \hat{\eta}_E(\bar{n}) = \sum_{j \in n_E^a} \hat{\omega}_j, \quad (8)$$

and the total load as

$$\hat{\eta}(\bar{n}) = \hat{\eta}_D(\bar{n}) + \hat{\eta}_E(\bar{n}) + \hat{\eta}_{oc}. \quad (9)$$

We assume the service load factors as lognormal r.v.'s which follows from the power control error of the received E_b/N_0 around the target- E_b/N_0 , see e.g. Viterbi and Viterbi (1993). The parameters μ , σ are then derived from the mean and variance of the E_b/N_0 distributions. These parameters depend on the service class of the users, but are equal for all users within one class. So we can write $E[\hat{\omega}_k] = E[\hat{\omega}_s]$ for all mobiles k with the same service class s . The other-cell load η_{oc} is modeled as a lognormal r.v. with constant mean and variance.

Since the total load $\hat{\eta}$ is a sum of independent lognormally distributed r.v.'s, we assume that $\hat{\eta}$ also follows a lognormal distribution (Fenton 1960). We get the distribution parameters from the first moment and variance of the cell load which can be calculated directly from the moments of the SLFs:

$$E[\hat{\eta}(\bar{n})] = \sum_{s \in \mathcal{S}} n_s \cdot E[\hat{\omega}_s] + n_E^a \cdot E[\hat{\omega}_E] + E[\hat{\eta}_{oc}]. \quad (10)$$

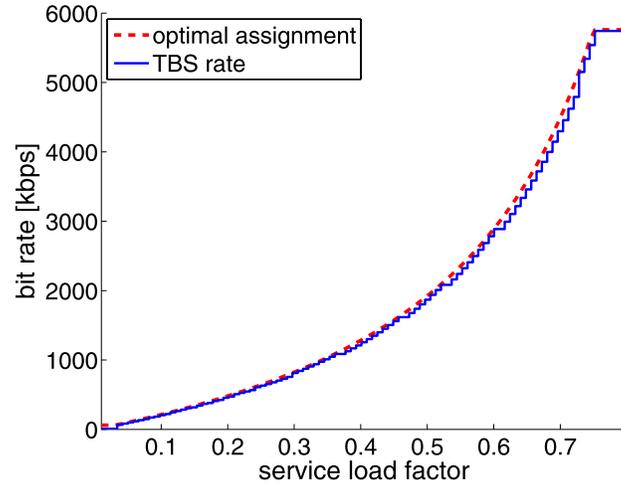
The variance is calculated analogously. The accuracy of this approach is validated e.g. in Staehle et al. (2003). Another novelty of the E-DCH is Hybrid ARQ (HARQ), which combines the automatic-repeat-request protocol with code combining techniques. The effect of HARQ can be modeled as a constant gain which is included in the target- E_b/N_0 of the E-DCH and with an additional overhead on the mean data volumes of the E-DCH.

4 Rate assignment

The available E-DCH load depends on the DCH and other-cell load. The task of the RRM is to assign each E-DCH mobile a service load factor ω such that the E-DCH load is completely utilized if possible. Due to the very flexible scheduling mechanism of the E-DCH, this can be reached in several ways. We consider two fundamentally different scheduling disciplines: the first one is parallel equal-rate scheduling, which means that every E-DCH user gets the same service load factor in every TTI. The second is equal-rate one-by-one-scheduling, where each E-DCH user gets the maximum possible service load factor in a round-robin-fashion. Note that we assume for the latter discipline that the E-DPDCHs (the enhanced dedicated physical data channels) are perfectly synchronized, hence do not generate any interference to each other. Although this is a strong assumption for current specifications of the Enhanced Uplink, the results can be seen as an upper bound for a more realistic system. We further assume that the network is dimensioned such that the transmit powers of the mobiles are sufficient to reach the maximum bitrate.

Generally, the user bit rate depends on the magnitude of the E-DCH cell load which may be generated without violating the RRM target in (3). The channel bit rate of the E-DCH is defined by the amount of information bits which can be transported within one TTI. This quantity is defined in 3GPP (2005b) by the set of transport block sizes \mathcal{TBS} . With

Fig. 2 Mapping of service load factors to bit rates



a TTI of 2 ms, the information bit rate per second follows as $R_{i,E}^I = TBS_i \cdot 500$, where $i = 1, \dots, |TBS|$ indicates the index of the TBS. We further define $R_{0,E}^I = 0$. With this interpretation we can map the E-DCH bit rate to a service load factor according to (6) as

$$\hat{\omega}_{i,E} = \frac{\hat{\varepsilon}_E R_{i,E}^I}{\hat{\varepsilon}_E R_{i,E}^I + W}, \tag{11}$$

where $\hat{\varepsilon}_E$ is the E_b/N_0 for the E-DCH RAB. Note that here we assume that the target- E_b/N_0 -values are equal for all rates. However, this restriction can be easily avoided by introducing individual target- E_b/N_0 -values for each TBS.

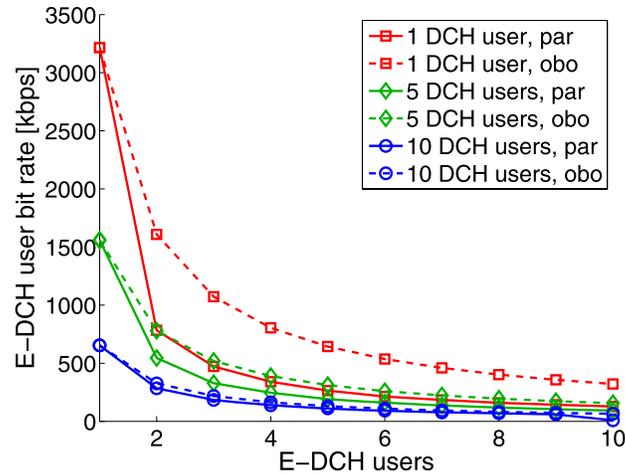
The next step is to select the information bit rate such that (3) is fulfilled:

$$R_E^I(\bar{n}) = \max\{R_{i,E}^I | P(\hat{\eta}_D(\bar{n}) + |n_E^a| \cdot \hat{\omega}_{i,E} + \eta_{oc} \geq \hat{\eta}^*) \leq p_t\}. \tag{12}$$

The actual user bit rates are now calculated according to the scheduling mechanism under the condition that the rate is higher than a certain minimum bit rate $R_{\min,E}$: In case of parallel scheduling, the user bit rate is simply the information bit rate. In case of one-by-one scheduling, the user bit rate is approximated by dividing the information bit rate by the number of E-DCH users:

$$R_E(\bar{n}) = \begin{cases} R_E^I(\bar{n}), & \text{if } R_E^I(\bar{n}) \geq R_{\min,E} \text{ and parallel scheduling,} \\ \frac{R_E^I(\bar{n})}{|n_E|}, & \text{if } \frac{R_E^I(\bar{n})}{|n_E|} \geq R_{\min,E} \text{ and one-by-one scheduling,} \\ 0 & \text{else.} \end{cases} \tag{13}$$

Figure 2 shows the mapping of the service load factors to information bit rates in case of a target- E_b/N_0 of 3 dB. The optimal case indicated by the dashed line is calculated from the definition of the service load factors as $R_{\text{opt}} = \frac{\hat{\omega} \cdot W}{\hat{\varepsilon}^* (1 - \hat{\omega})}$. The solid line shows the corresponding rate calculated from the TBS. Both curves are very close to each other, and we see that for high SLFs, a small change means a large change on the bit rate. The non-linear dependency between bit rate and SLF is the basis for the argument that a slow-down (in terms of bit rate) of the users leads to an increased system capacity in terms of admissible sessions if an admission control based on the cell load is used (Altman 2002; Fodor and Telek 2005). However, if we define capacity as the cumulated bit rate per cell, the capacity

Fig. 3 Rate selection for varying DCH and E-DCH loads

shrinks with the number of parallel transmitting users due to the increased interference, as we can see in Fig. 3 and in Sect. 7.

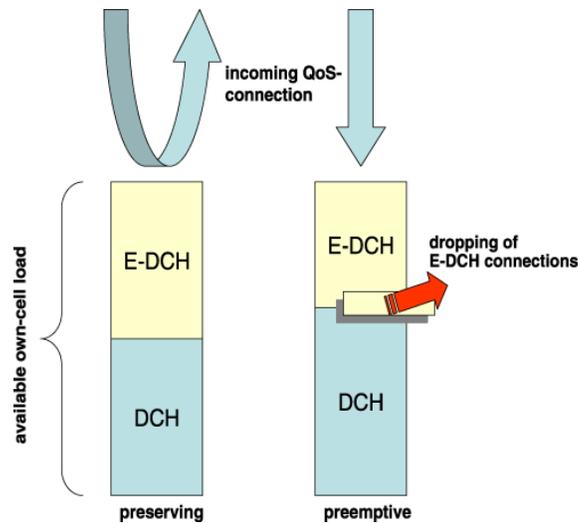
Figure 3 shows the E-DCH bit rate per user for different numbers of DCH and E-DCH users. The solid lines show the parallel scheduling case and the dashed ones the one-by-one scheduling. Different colors indicate different numbers of concurrently active DCH users. We see that for only one E-DCH user, parallel and one-by-one scheduling have naturally the same throughput. However, with more E-DCH users the gain of the one-by-one scheduling over the parallel scheduling increases since the users do not interfere with each other and thus are able to utilize the radio resources more efficiently, i.e. get high SLFs and correspondingly also high bit rates. This gain depends on the number of DCH users in the system: With more DCH users, the gain shrinks such that with 10 DCH users, there is nearly no gain for the one-by-one scheduling. In this case the available resources for the E-DCH are already quite low and thus only SLFs with low transmission rates are possible.

5 Admission control

The admission control (AC) is responsible for keeping the cell load below the maximum load. Generally, we model the AC on basis of the RRM target condition. We distinguish between two RRM policies for incoming QoS users: The first, which we call *preserving* treats E-DCH and QoS equally, which means that an incoming connection of either class is blocked if there are not enough resources available. The second, which we call *preemptive*, gives priority to QoS users, which means that eventually active best effort connections may be dropped from the system in order to make room for the incoming QoS user. In both policies existing E-DCH connections are slowed-down if the number of QoS-connections increases. However, with the preserving strategy incoming QoS-calls are blocked if the RRM cannot slow-down the E-DCH connections any more. With the preemptive strategy, one or more E-DCH connections are dropped from the system in this case, meaning that blocking for the QoS users occurs only if nearly all resources are occupied by QoS connections, cf. Fig. 4.

If a new connection is to be established to the network, the AC is done in two steps: At first, the amount of resources ω which the incoming connection will occupy is identified. In case of a QoS-connection, this is simply ω_s . In case of an E-DCH connection, incoming connections are admitted if a minimum bit rate $R_{\min,E}$ can be guaranteed. The corresponding SLF is denoted with $\hat{\omega}_{\min,E}$. Let us further denote with \bar{n}^+ the state vector \bar{n} plus the

Fig. 4 Principle of the preserving and preemptive policy



incoming connection with service class s or with an additional E-DCH connection. The second step is then to estimate the probability for exceeding the maximum load with the new connection included. This step depends on the implemented policy:

Preserving policy: In the preserving case, we calculate the parameters for the distribution of the expected cell load η_{AC} as in (10), but with $\omega_{\min,E}$ for the E-DCH users:

$$\eta_{AC}(\bar{n}^+) = \eta_D(\bar{n}^+) + n_E^+ \cdot \hat{\omega}_{\min,E} + \eta_{oc}, \quad (14)$$

where n_E^+ is the number of E-DCH mobiles with the incoming mobile included, if any. So, if the probability $P(\eta_{AC} \geq \eta^*)$ is higher than the target probability p_t , the connection is rejected, otherwise the connection is admitted.

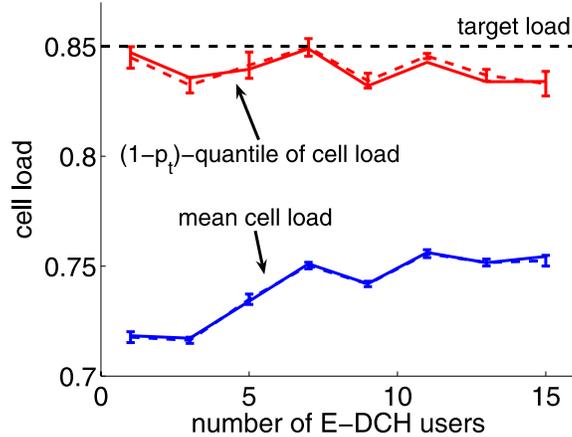
Preemptive policy: With preemption, the incoming call is admitted if enough resources are available such that $P(\eta_{AC} \geq \eta^*) \leq p_t$, as in the preserving case. However, if the resources are insufficient, we distinguish two cases: If the incoming call belongs to an E-DCH user, the call is blocked. If the incoming call belongs a QoS user, the RRM calculates from the service requirement ω_s the number of E-DCH connections with minimum rate $R_{E,\min}$ which must be dropped from the system such that the incoming call can be admitted. The number of E-DCH connections $n_d(\bar{n}, s)$ which must be dropped depends on the current state and on the SLF of the incoming QoS-connection. It is given by the following rule:

$$n_d(\bar{n}, s) = \min\{n | P(\eta_D(\bar{n}^+) + (n_E - n) \cdot \omega_{\min,E} + \eta_{oc} \geq \eta^*) \leq p_t\}. \quad (15)$$

Note that $0 \leq n_d \leq \lceil \frac{\omega_s}{\omega_{\min,E}} \rceil$. Blocking for QoS-users occurs if the number of E-DCH connections is too low to meet the requirements of the service class, i.e. if $n_d(\bar{n}, s) > n_E$. Blocking for E-DCH users occurs if the existing connections cannot be slowed down any further, due to the constraint on the minimum bit rate.

After admission control, the RRM executes the rate assignment as in (12) to adjust the bit rate of the E-DCH users to the new situation. Figure 5 illustrates the principle of admission control and rate selection. It shows the mean and the $(1 - p_t)$ -quantile (here $p_t = 5\%$) of the cell load distribution for 5 DCH users and an increasing number of E-DCH users. The target load is $\hat{\eta}^* = 0.85$. Note that the results from a Monte-Carlo-simulation which uses random E_b/N_0 -values, denoted by dashed lines, are very close to the analytical results, which shows

Fig. 5 Mean cell load and 95%-quantiles



the accuracy of the lognormal approximation. Due to the discretization of the available rates, the $(1 - p_t)$ -quantile does not exactly meet the target-load, but stays just below. Since the coefficient of variation of the cell load is decreasing with the number of users in the system, the mean load comes closer to the target load with an increasing number of E-DCH users.

6 Performance evaluation

Now we assume that all calls arrive with exponentially distributed interarrival times with mean $\frac{1}{\lambda}$. The users choose a DCH service class or the E-DCH with probability p_s and p_E , hence the arrival rates per class are $\lambda_s = p_s \cdot \lambda$ and $\lambda_E = p_E \cdot \lambda$. The holding times for the DCH calls are also exponentially distributed with mean $\frac{1}{\mu_s}$. For the E-DCH users we assume a volume based user traffic model (Hossfeld et al. 2006). With exponentially distributed data volumes, the state-dependent departure rates of the E-DCH users are then given by

$$\mu_E(\bar{n}) = n_E \cdot \frac{R_E(\bar{n})}{E[V_E]}, \tag{16}$$

where $E[V_E]$ is the mean traffic volume of the E-DCH users.

The resulting system is a multi-service $M/M/n - 0$ loss system with state dependent departure rates for the E-DCH users. We are now interested in calculating the steady-state distribution of the number of users in the system. Since the joint Markov process is not time-reversible which can be instantly verified with Kolmogorov’s reversibility criterion, no product form solution exists. The steady state probabilities follow by solving

$$Q \cdot \bar{\pi} = 0 \quad \text{s.t.} \quad \sum \pi = 1 \tag{17}$$

for $\bar{\pi}$, where Q is the transition rate matrix. The rate matrix Q is defined with help of the bijective index function $\phi(\bar{n}) : \Omega \rightarrow N$, which maps the state vector \bar{n} to a single index number. The transition rate $q(\phi(\bar{n}), \phi(\bar{n} \pm \bar{1}))$ in the rate matrix between states \bar{n} and $\bar{n} \pm \bar{1}$ is then

$$\begin{aligned} q(\phi(\bar{n}), \phi(\bar{n} + \bar{1}_s)) &= \lambda_s, \\ q(\phi(\bar{n}), \phi(\bar{n} + \bar{1}_E)) &= \lambda_E, \\ q(\phi(\bar{n}), \phi(\bar{n} - \bar{1}_s)) &= n_s \cdot \mu_s, \\ q(\phi(\bar{n}), \phi(\bar{n} - \bar{1}_E)) &= \mu_E(\bar{n}) \end{aligned} \tag{18}$$

for all valid states in the state space Ω and $q(\phi(\bar{n}), \phi(\bar{n} \pm \bar{1})) = 0$ otherwise. The sets of $\Omega_{ps,b}^+$ states where blocking occurs in the preserving case are defined by the condition $P(\eta(\bar{n}^+) \geq \eta^*) > p_t$, i.e. they form the 'edges' of the state space. With preemption, an E-DCH connection is dropped if $P(\eta(\bar{n}^{+s}) \geq \eta^*) > p_t$ and $n_d(\bar{n}, s) \geq \lceil \frac{\omega_s}{\omega_{\min,E}} \rceil$, i.e. in case of an incoming QoS connection. We define this set as $\Omega_{pe,d}^{+s}$. Blocking occurs then in the set $\Omega_{pe,b}^{+s} = \Omega_{ps,b}^{+s} \setminus \Omega_{pe,d}^{+s}$. The set of blocking states for E-DCH connections is the same for both policies. For the preemptive policy, an additional entry in the transition rate matrix is generated for states where preemption may occur:

$$q(\phi(\bar{n}), \phi(\bar{n} + \bar{1}_s - \bar{n}_d(\bar{n}, s))) = \lambda_s. \tag{19}$$

As performance measures we choose the service-dependent call blocking probabilities P_s , the call dropping probability P_d which applies only in the case of the preemptive strategy and the mean user bit rate $E[R_U]$ achieved by the E-DCH users. The call blocking probabilities are easily calculated as the sum of all states probabilities in which blocking may occur:

$$P_s = \sum_{\bar{n} | \bar{n} \in \Omega_b^{+s}} \pi(\bar{n}). \tag{20}$$

Note that we omit the qualifier for the admission control policy. We define the call dropping probability in our analysis as the probability that an E-DCH connection is dropped if a QoS-call is arriving in the system. This probability is given by

$$P_d = \sum_{\bar{n} | \Omega_{pe,d}^{+s}} \frac{\pi(\bar{n}) \cdot \sum_{s' \in \mathcal{S}} \frac{\lambda_{s'}}{\sum_{s'' \in \mathcal{S}} \lambda_{s''}} \cdot \frac{n_d(\bar{n}, s')}{n_E}}{\sum_{\bar{n}' | n_E > 0} \pi(\bar{n}')} . \tag{21}$$

We define further the mean throughput per user at a random time instance as

$$E[R_U] = \sum_{R_E > 0} R_E \cdot \frac{\sum_{\bar{n} | R_E(\bar{n}) = R_E} n_E \cdot \pi(\bar{n})}{\sum_{\bar{n}' | n_E > 0} n'_E \cdot \pi(\bar{n}')} , \tag{22}$$

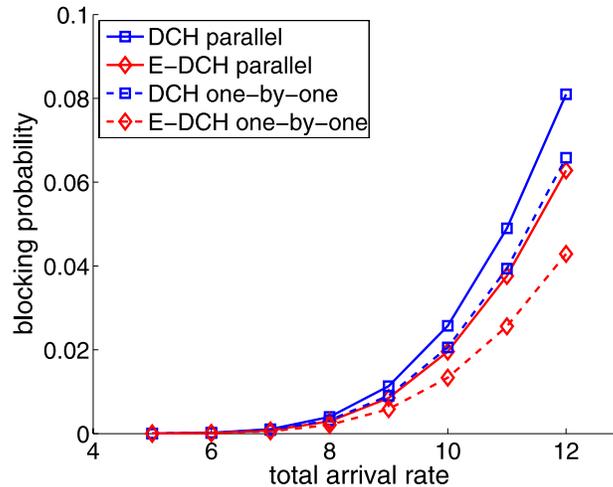
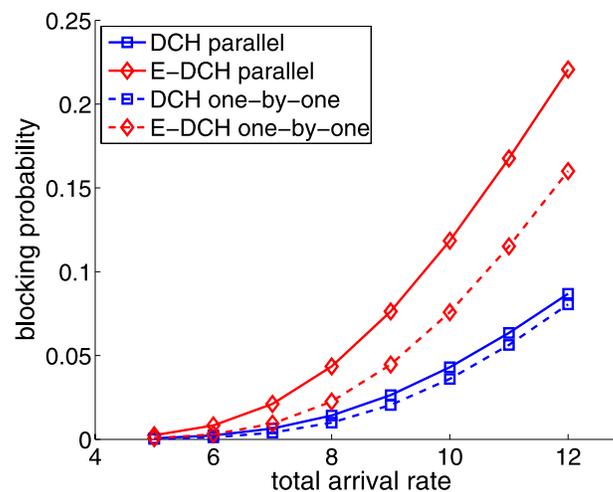
which is conditioned with the probability that at least one E-DCH user is in the system.

7 Numerical results

In this section we give some numerical examples for our model. Our scenarios, if not stated otherwise, consist of two service classes: 64 kbps QoS-users (i.e. DCH users) with a target- E_b/N_0 of 4 dB and the E-DCH best effort users with a target- E_b/N_0 of 3 dB. The service probabilities are $p_1 = 0.4$ and $p_E = 0.6$.

7.1 Comparison of parallel and one-by-one scheduling

In the first scenario we compare the blocking probabilities between parallel scheduling and one-by-one scheduling for preserving AC. In Fig. 6, $R_{\min,E}$ is 60 kbps. $E[V_E]$ is 72 kbit. Red lines indicate the blocking probabilities for the E-DCH users, blue lines for the DCH users. We see that the blocking probabilities for the DCH users are higher than for the E-DCH users. Due to the low minimum E-DCH bit rate and the resulting low minimal

Fig. 6 Blocking probabilities for 60 kbps min. E-DCH bit rate**Fig. 7** Blocking probabilities for 200 kbps min. E-DCH bit rate

service load factor, E-DCH users may still connect to the system whereas DCH users are already blocked. The comparison of the parallel (solid lines) with the one-by-one scheduling case (dashed lines) shows that the throughput gain of the one-by-one users leads to lower blocking probabilities, and also to a higher difference between DCH and E-DCH users.

In Fig. 7, the scenario is equal to the previous one with the exception that $R_{\min,E}$ is 200 kbps. In this case, the minimal service load factors for the E-DCH user is higher than the load requirements of the DCH users. Consequently, the E-DCH blocking probabilities are now higher than the DCH blocking probabilities. We see further that the DCH blocking probabilities for both scheduling disciplines are now very close to each other.

In Fig. 8, the total mean E-DCH throughput is shown, which is given by $E[R_T] = \lambda_E \cdot (1 - P_{b,E}) \cdot E[V_E]$. The throughput for the one-by-one scheduling is as expected always higher than for the parallel scheduling. The gain gets higher with increasing load due to the lower blocking probabilities. This also explains why the total throughput with $R_{\min,E} = 200$ kbps is lower than with $R_{\min,E} = 60$ kbps, although the per-user-bit rate is higher, as in Fig. 9, where at a low offered load, the bit rates for $R_{\min,E} = 60$ kbps and $R_{\min,E} = 200$ kbps are close to each other. This changes with increasing loads, since the AC prevents the decline for the 200 kbps case stronger than for the 60 kbps case. We see further that the gain of the one-by-one scheduling over the parallel scheduling shrinks with a higher load. This is due to the effect of the concurrently increased DCH load, as already seen in Fig. 3.

Fig. 8 Mean total E-DCH throughput

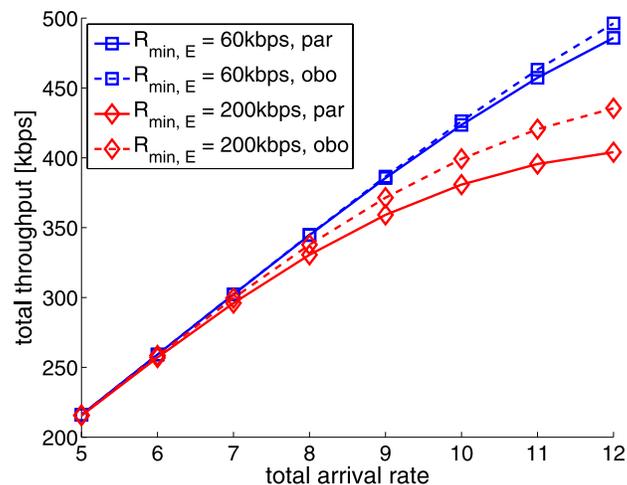
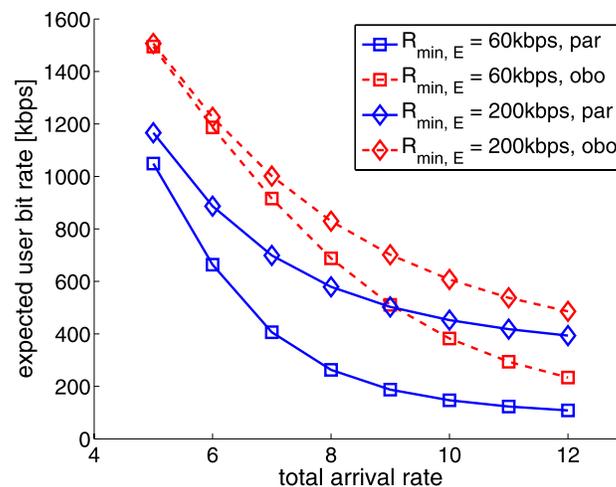


Fig. 9 Mean E-DCH bit rates per user



7.2 Comparison of preserving and preemptive admission control

Now we compare preemptive and preserving admission control. We restrict our scenario to the case of parallel scheduling, but the conclusions do also hold for the case of one-by-one scheduling. The guaranteed bandwidth for E-DCH users is $R_{\min, E} = 60$ kbps.

In Fig. 10, blocking and dropping probabilities for both admission policies are shown. The curves with a circle marker indicate the blocking probabilities for the 64 kbps QoS users, while the curves with a square marker show the blocking probabilities for the E-DCH users. The dashed line with diamond markers shows the dropping probabilities in case of preemption. Although a system with such high blocking probabilities would be considered as heavily overloaded, we show these results for a better understanding of the effect of preemption. It can be stated that preemption leads to an enormous performance gain for the QoS users, which is caused by the substantially smaller sets of states where blocking can occur at all. The blocking probabilities for the E-DCH users, however, are nearly identical and only begin to differ from each other under very high load. The dropping probabilities do not exceed approx. 10% because in high load regions the system is nearly fully occupied by QoS users.

The impact of preemption on the user and cell bit rates (defined as the cumulated bit rates of all users at any time) is shown in Fig. 11. The user throughputs have solid lines,

Fig. 10 Blocking and dropping probabilities for QoS and E-DCH users

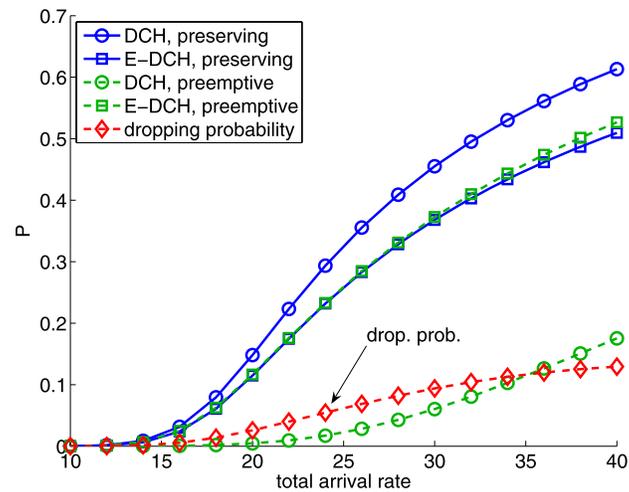
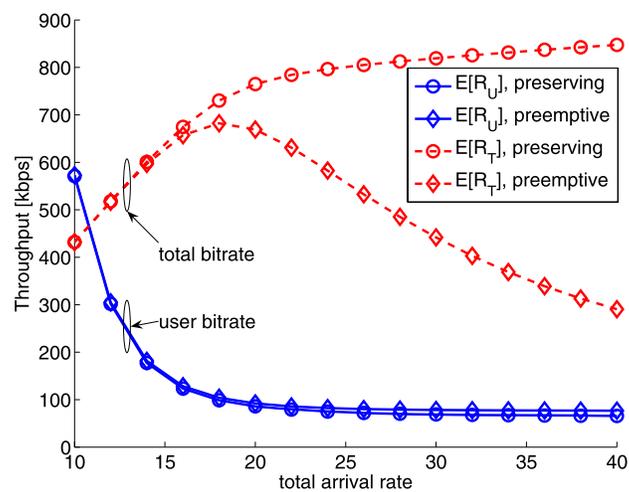


Fig. 11 Mean user and cell bit rates



while the cell throughputs have dashed lines. The expected user throughputs in both cases $E[R_U]$ are nearly identical with a slight advantage for the preemptive case. However, due to the dropping of users the total cell throughput $E[R_T]$ in the preemptive case is significantly lower than in the preserving case. Since the cell throughputs also consider the case if no E-DCH user at all is in the system, the curves are first increasing and then decreasing. In the next scenario we fix the total arrival rate to 15 and vary the ratio between DCH and E-DCH arrivals from 10%/90% to 90%/10%. The results are shown in Fig. 12. They show that in situations with a high fraction of best-effort traffic preemption leads to a substantial decrease of the blocking probabilities for the QoS users with dropping probabilities smaller than 1%. However, if the ratio is shifted to the QoS side, the decreasing load available to the E-DCH users leads to increased dropping probabilities.

Figure 13 shows the sensitivity of the system to different volume size distributions for the E-DCH users. The results are calculated with an event-based simulation which was also used for the validation of the analytical results. Three cases are presented: Constant volume size, exponentially and Pareto distributed volume sizes (with parameters $k = 1.5$ and $x_m = 2.4 \cdot 10^4$), all with the same mean. As expected (see e.g. Litjens and Boucherie 2003), a higher variance leads to lower dropping probabilities, although in this very low load regions the differences are quite small, which may lead to the conclusion that the exponential assumption may be a sufficient approximation in these cases.

Fig. 12 Impact of preemption depends on ratio between DCH and E-DCH users

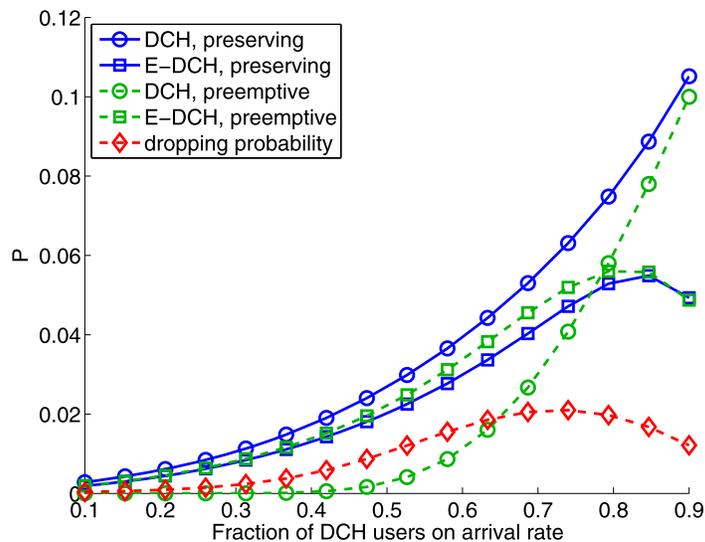
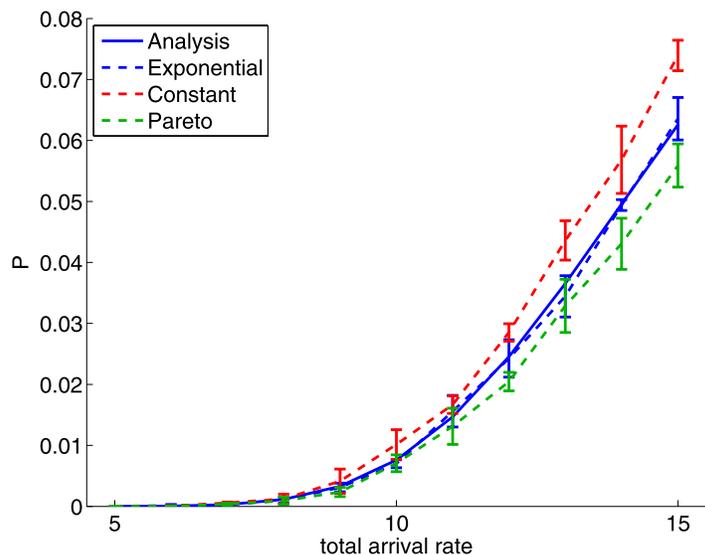


Fig. 13 Sensitivity of the dropping probabilities against volume size distribution



8 Conclusion

We presented an analytical model for a UMTS system with QoS traffic over DCH radio bearers and best-effort traffic over the Enhanced Uplink. The model considers two fundamental different scheduling mechanisms: Parallel scheduling, i.e. classical CDMA, and one-by-one scheduling, i.e. a mixture between TDMA and CDMA. The model considers the effects of imperfect power control and varying other-cell interferences. Because of these effects the bit rate selection for the E-DCH users and the admission control is based on a probabilistic metric, which states that the cell load should not exceed the allowed maximum cell load with a certain probability. We further evaluated two admission control policies, *preserving* and *preemptive*, where in the latter case incoming QoS connections may lead to the dropping of best-effort connections, if necessary.

The numerical results showed that one-by-one scheduling has the largest performance gain over the parallel scheduling if the number of DCH users is low, i.e. if the own-cell interference level in the cell is low. We further saw the effect of a minimum allowed bit rate

for the E-DCH users on the blocking probabilities and on the user bit rates. The evaluation of the two admission control policies showed that preemption can lead to a substantial decrease in blocking probabilities for the QoS users, but it should be generally carefully used since it can also lead to high dropping probabilities in scenarios with low quantities of best-effort traffic. A possible solution would be to reserve a certain amount of load to best-effort users only.

References

- 3GPP (2005a). *TS 25.309 V6.4.0 FDD enhanced uplink; overall description* (Technical report). 3GPP.
- 3GPP (2005b). *TS 25.321 V6.6.0 Medium Access Control (MAC) protocol specification* (Technical report). 3GPP.
- Altman, E. (2002). Capacity of multi-service cellular networks with transmission-rate control: a queueing analysis. In *Proc. of ACM MobiCom '02* (pp. 205–214). Atlanta, Georgia, USA.
- Boxma, O. J., Gabor, A. F., Núñez-Queija, R., & Tany, H.-P. (2006). Performance analysis of admission control for integrated services with minimum rate guarantees. In *Proc. of NGI 2006*, València, Spain.
- Fenton, L. F. (1960). The sum of log-normal probability distributions in scatter transmission systems. *IRE Transactions on Communication Systems*, 8(1), 57–67.
- Fodor, G., & Telek, M. (2005). Performance analysis of the uplink of a CDMA cell supporting elastic services. In *Proc of NETWORKING 2005* (pp. 205–216). Waterloo, Canada.
- Fodor, G., Telek, M., & Badia, L. (2006). On the tradeoff between blocking and dropping probabilities in CDMA networks supporting elastic services. In *Proc. of networking 2006* (pp. 954–965). Coimbra, Portugal.
- Holma, H., & Toskala, A. (Eds.). (2001). *WCDMA for UMTS*. New York: Wiley.
- Hossfeld, T., Mäder, A., & Staehle, D. (Eds.). (2006). When do we need rate control for dedicated channels in UMTS? In *Proc. of IEEE VTC spring '06* (pp. 425–429). Melbourne, Australia.
- Jäntti, R., & Kim, S.-L. (2001). Transmission rate scheduling for the non-real-time data in a cellular CDMA system. *IEEE Communication Letters*, 5(5), 200–202.
- Kumaran, K., & Qian, L. (2003). Uplink scheduling in CDMA packet-data systems. In *Proc. of IEEE INFOCOM '03* (pp. 292–300). San Francisco, CA, USA.
- Litjens, R., & Boucherie, R. J. (2002). Performance analysis of fair channel sharing policies in an integrated cellular voice/data network. *Telecommunication Systems*, 19(2), 147–186.
- Litjens, R., & Boucherie, R. J. (2003). Elastic calls in an integrated services network: the greater the call size variability the better the QoS. *Performance Evaluation*, 52(4), 193–220.
- Mäder, A., & Staehle, D. (2006). An analytical model for best-effort traffic over the UMTS enhanced uplink. In *Proc. of IEEE VTC fall '06* (pp. 1–5). Montréal, QC, Canada.
- Mäder, A., & Staehle, D. (2007). Comparison of preemptive and preserving admission control for the UMTS enhanced uplink. In *Proc. of 15th KiVS* (pp. 201–212). Bern, Switzerland.
- Parkvall, S., Peisa, J., Torsner, J., Sågfors, M., & Malm, P. (2005). WCDMA enhanced uplink—principles and basic operation. In *Proc. of IEEE VTC spring '05* (pp. 1411–1415). Stockholm, Sweden.
- Staehle, D., Leibnitz, K., Heck, K., Tran-Gia, P., Schröder, B., & Weller, A. (2003). Analytic approximation of the effective bandwidth for best-effort services in UMTS networks. In *Proc. of IEEE VTC spring '03* (pp. 1153–1157). Jeju, South Korea.
- Viterbi, A. M., & Viterbi, A. J. (1993). Erlang capacity of a power controlled CDMA system. *IEEE Journal on Selected Areas in Communications*, 11(6), 892–900.