# An Adaptive Bandwidth Reservation Scheme in Multimedia Wireless Networks

Xiang Chen and Yuguang Fang
Department of Electrical and Computer Engineering
University of Florida, Gainesville, FL 32611

*Abstract*—**Next generation wireless networks target to provide quality of service (QoS) for multimedia applications. In this paper, the system supports two QoS criteria, i.e., the system should keep the handoff dropping probability always less than a predefined QoS bound, while maintaining the relative priorities of different traffic classes in terms of blocking probability. To achieve this goal, a dynamic multiple-threshold bandwidth reservation scheme is proposed, which is capable of granting differential priorities to different traffic class and to new ad handoff traffic for each class by dynamically adjusting bandwidth reservation thresholds. Moreover, in times of network congestion, a preventive measure by use of throttling new connection acceptance is taken. Another contribution of this paper is to generalize the concept of relative priority, hence giving the network operator more flexibility to adjust admission control policy by incorporating some dynamic factors such as offered load. The elaborate simulation is conducted to verify the performance of the scheme.**

## I. INTRODUCTION

With the increasing demands for mobile multimedia services such as audio, video and data, next generation wireless networks are expected to provide quality of service (QoS) for such multimedia applications to users on the move. Since the services have inherently different traffic characteristics, their QoS requirements may differ in terms of bandwidth, delay and connection dropping probabilities. It is the networks' responsibilities to fairly and efficiently allocate network resources among different users in order to satisfy such differentiated QoS requirements of each type of service independent of the others.

In order to guarantee QoS requirements of these services while accommodating rapidly growing population of mobile users, efficient connection admission control (CAC) schemes have to be used. In [1]-[3], the well-known Guard Channel scheme and some of its variations were proposed to give higher priority to handoff connections over new connections in voice traffic, whose performance depends largely on the choice of the number of guard channels, which is mainly based on *a priori* knowledge of the traffic patterns. Some dynamic bandwidth allocation schemes are investigated in [4] [5]. However, in [4], the priority among different traffic classes, either in handoff traffic or new traffic, was not

addressed while [5] did not take fairness into consideration. [6] uses predefined call blocking probability profile to maintain the relative priorities among different classes of traffic, whose underlying assumption, however, may not hold in fast changing networks.

In this paper, we propose a dynamic, multiple-threshold bandwidth reservation scheme (DMTBR) for multimedia mobile wireless systems. Three bandwidth reservation thresholds, $G_1$, $G_2$, and $G_3$ ($G_1 < G_2 < G_3$) are used to guarantee QoS requirements. When the network is under heavy traffic load, in order to guarantee QoS provisioning, further step by throttling new connection acceptance is taken. In summary, this scheme has the following features:

- It gives differential priorities to both new and handoff connections with different types of services.
- It maintains the relative priorities and fairness among traffic classes by taking into account both user QoS profile and real traffic condition, which generalizes the concept of relative priorities among traffic classes.
- It updates the reservation thresholds periodically, hence is able to respond to the changing network conditions quickly and effectively.

The rest of this paper is organized as follows. The next section describes the traffic model considered. Our scheme, including the target QoS criteria, is presented in Section III. Then, the calculation issues involved in the scheme are addressed in Section IV. In Section V, the scheme is verified through simulations combined with result analysis. Finally, this paper is concluded in Section VI.

## II. TRAFFIC MODEL

The system under consideration is a multimedia wireless network with a cellular infrastructure, comprising a number of cells. We assume the system uses fixed channel assignment (FCA), which means each cell has a fixed amount of capacity. No matter which multiple-access technology (FDMA, TDMA or CDMA) is used, we could interpret system capacity in terms of bandwidth. In this paper, we assume a single number, "effective bandwidth" [7], is adequate for guaranteeing desired QoS for any connection with certain traffic characteristics. Hereafter, whenever we refer to the bandwidth of a connection, we mean its effective bandwidth. We assume each cell has $C$ bandwidth units (BU). There are two classes of incoming traffic: 1) Class I—real-time traffic and 2) Class

II—non-real-time traffic. Typically, class I traffic includes voice and video service while class II traffic is comprised of data services like email, file transfer and web browsing. The arrivals are Poisson processes, with respective arrival rates $\lambda_{rt}$ and $\lambda_{nrt}$. The connection duration times of both connections follow exponential distribution, with means $1/\mu_{rt}$ and $1/\mu_{nrt}$. Furthermore, we assume that the cell residence time distributions of these two kinds of connections are also exponentially distributed, with means $1/\gamma_{rt}$ and $1/\gamma_{nrt}$. The number of BUs required by each real-time and non-real-time connection is $BW_{rt}$ and $BW_{nrt}$, respectively. These assumptions are appropriate and commonly used in literature. Since time spent by a connection in a cell is the minimum of connection duration time and cell resident time, following the assumptions, we can easily obtain that, for these two type of traffic, the distribution of the time spent by a connection in a cell is also exponentially distributed, with the mean of $d_{rt} = 1/(\mu_{rt} + \gamma_{rt})$ and $d_{nrt} = 1/(\mu_{nrt} + \gamma_{nrt})$.

## III. PROPOSED SCHEME

### A. QoS Criteria

The first criterion is about the handoff connection dropping probability (CDP), i.e., the probability of a handoff connection being dropped when a handoff has to be made due to user roaming from one cell to another. In wireless networks, we usually set an upper bound for this probability, like the following:

$$P_{d,rt} \leq QoS_{rt}$$
$$P_{d,nrt} \leq QoS_{nrt} \tag{1}$$

where $P_{d,rt}$, $QoS_{rt}$ and $P_{d,nrt}$, $QoS_{nrt}$ are CDP and the allowable maximum dropping probability for real-time and non-real-time traffic, respectively.

The second criterion is to maintain the relative priority among different type of traffic in terms of new connection blocking probability (CBP). Obviously, if there is no such criterion, it is unfair to the traffic classes that require larger bandwidth. To address this problem, we assume in each traffic class's profile, there exists a parameter named traffic priority weight, $W$, indicating which priority level the traffic class will have. This parameter is set during the negotiation between the user and the network operator, taking traffic characteristics into account. A smaller weight means a higher priority. To achieve a better fairness in CBPs among all traffic classes, the network may keep the CBPs satisfying the following equation:

$$\frac{P_{b,rt}}{P_{b,nrt}} = \frac{W_{rt}}{W_{nrt}} \tag{2}$$

where $P_{b,rt}$, $W_{rt}$ and $P_{b,nrt}$, $W_{nrt}$ are CBP and the predefined traffic priority weight for the two traffic classes, respectively.

Generally speaking, many factors play a role in determining the CBP. Each traffic class's actual CBP depends on system capacity, offered traffic load of the traffic class, the priority of the traffic class, the admission policy adopted to fulfill the QoS criteria related to handoff traffic, the action the network may take in times of congestion, and so on. While some factors like system capacity or pre-assigned traffic

priority could be static, offered load, actions taken to deal with network congestion are dynamic. In this sense, the criterion in Eq. (2) is static and not fair enough since it fails to reflect the real time network situation. Therefore, we generalize the concept of relative priority and propose a more general way to maintain the relative priority among different traffic class, using the following formula:

$$\frac{P_{b,rt}}{P_{b,nrt}} = \alpha \frac{W_{rt}}{W_{nrt}} \tag{3}$$

Compared with Eq. (2), we add one factor, $\alpha$, on the right-hand side of Eq. (3). $\alpha$ can be thought of as a function of some of the dynamic factors described above, representing network's real traffic conditions or some procedures responding to traffic changes or QoS status.

Since offered load is one of the commonly used measures of network traffic load, as one way to make $\alpha$ concrete, we let $\alpha$ be a function of offered load per cell for each traffic class. Offered load can be defined as the product of each traffic class's traffic arrival rate, call holding time and normal bandwidth, i.e.,

$$OL_{(n)rt} = \frac{\lambda_{(n)rt} BW_{(n)rt}}{\mu_{(n)rt}} \tag{4}$$

Replacing $\alpha$ with the ratio of offered load of real-time and non-real-time traffic, we obtain the following:

$$\frac{P_{b,rt}}{P_{b,nrt}} = \frac{OL_{rt}}{OL_{nrt}} \times \frac{W_{rt}}{W_{nrt}} \tag{5}$$

In this way, we take into account the offered load of each traffic class. In other words, we maintain the relative priority by keeping the ratio of the CBP of each traffic class equal to the product of their corresponding ratio of traffic load and the weight pre-defined. The advantage of this approach will be shown in Section V. Hereafter, we term the scheme satisfying QoS Eq. (1) and (2) DMTBR_A, and term the scheme satisfying QoS Eq. (1) and (5) DMTBR_G.

### B. Connection Admission Policy

Based on the thresholds $G_1$, $G_2$, and $G_3$, the admission policy, including the adoption of throttling new connections acceptance in case of heavy congestion, is shown in Fig. 1. Note that for new connections, the thresholds $G_2$ and $G_3$ are not fixed for either of the two traffic classes; instead, their roles switch depending on the network's instantaneous situation. Parameter $prob_{rt}$ (or $prob_{nrt}$) denotes the probability of throttling and function $rand()$ generates a random number belonging to $[0, 1)$. $switch$ can be considered a Boolean sign, which indicates the roles of $G_2$ and $G_3$.

### C. Cooperations Among Cells

In cellular networks, traffic in different cells has correlation. Hence, it is necessary and more efficient to deal with network congestion in a cooperative manner to prevent this from happening through admission control. We adopt this idea in this scheme to cope with the situation where the network is undergoing heavy traffic load.

For a period of time, each cell measures CBP and CDP, i.e., $P_{b,rt}$ and $P_{d,rt}$ (or $P_{b,nrt}$ and $P_{d,nrt}$). We will count the times of increasing the reservation thresholds for handoff traffic. Once one reservation threshold is consecutively increased for a certain number of times, say three times, the cell is deemed experiencing heavy handoff traffic. In this case, to reduce the potential incoming handoff traffic hence keeping CDP below the upper bound, the cell will inform all of its neighbors to further throttle the acceptance of new connections of the same traffic class as the handoff traffic class in the current cell. One method to achieve this is to admit the new connection request with a certain probability generated online, which is called the probability of throttling new connections. Details on how to generate the probability are given in Section IV.

## IV.  CALCULATION ISSUES

### A.  Calculation of $G_1$, $G_2$, and $G_3$

The scheme requires accurately adjusting the values of these three thresholds, $G_1$, $G_2$, and $G_3$, for every period of time. We assume $G_1$ are the number of BUs that needs to be reserved to deal with handoff real-time connections that will arrive in a period of $d$ from now to the future, where $d$ is the corresponding expectation of channel holding time for real-time traffic. If during period $d$, there are $m$ real-time connections in cell $i$ which will leave cell $i$ due to completion or handoff, and $n$ handoff real-time connections which will enter cell $i$. Therefore, a total of $m + n$ events will happen in cell $i$ during period $d$. Let $s$ be a sequence of these $m + n$ events and $S(m, n)$ be the set of all possible sequences which may take place in $d$. Let $Y(s)$ denote the maximum net change in the number of BUs allocated to real-time connections in cell $i$ in $d$ corresponding to each specific $s$. We set $G_1$ equal to the expected value of $Y(s)$, which can be obtained as shown in [4].

$$G_1 = \sum_{s \in S(m,n)} \frac{Y(s)}{|S(m,n)|} \qquad (6)$$

where $|S(m, n)|$ is the cardinality of $S(m, n)$. According to the traffic model described in Section II, we can easily obtain the parameters in Eq. (6). Because of space limit, we do not give the calculation details.

In a similar way, we can calculate $G_2'$, which is the number of BUs that needs to be reserved to deal with handoff non-real-time connections that will arrive in a period of $d'$, the corresponding expectation of channel holding time for non-real-time traffic.

Once we get $G_1$ and $G_2'$, $G_2$ is obtained as followes:

$$G_2 = G_1 + G_2' \qquad (7)$$

Note that we grant higher priority to handoff real-time traffic over non-real-time traffic by use of the calculation order of $G_1$ and $G_2$ as described above.

Before calculating $G_3$, we calculate $G_3'$, which could be thought of as the reservation threshold that could be used either to reserve bandwidth for new real-time traffic against new non-real-time traffic, or to reserve bandwidth for new non-real-time traffic against new real-time traffic, depending on the instantaneous relative priority status for the traffic

```
if (the incoming handoff connection is real-time)
        if (available BUs >= BWₙ)          accept;
        else          reject;
else // the incoming handoff connection is non-real-time
        if (available BUs >= Gₜ + BWₙᵣₜ)  accept;
        else          reject;

if (switch == true)
        if (the incoming new connection is real-time)
          if (available BUs >= G₂ + BWₙ && rand( ) <= probₙ)
            accept;
          else     reject;
        else // the incoming new connection is non-real-time
          if (available BUs >= Gₜ + BWₙᵣₜ && rand( ) <= probₙᵣₜ)
            accept;
          else     reject;
else // (switch == false)
        if (the incoming new connection is real-time)
          if (available BUs >= G₃ + BWₙ && rand( ) <= probₙ)
            accept;
          else     reject;
        else // the incoming new connection is non-real-time
          if (available BUs >= G₃ + BWₙᵣₜ && rand( ) <= probₙᵣₜ)
            accept;
        else     reject;
```

Figure 1. Admission policy

classes. The initial value for $G_3'$ could be set as $BW_{rt}$ or $BW_{nrt}$. Therefore, $G_3$ can be estimated like the following:

$$G_3' = \begin{cases} BW_{rt}, & if \quad for \quad RT \\ BW_{nrt}, & if \quad for \quad NRT \end{cases} \qquad (8)$$

$$G_3 = G_2 + G_3'$$

### B.  Adaptation of Bandwidth Reservation Thresholds

The techniques we used to estimate $G_1$, $G_2'$ and $G_3'$ only serves to provide a good initial value. To meet the QoS criteria in a dynamically changing network environment, further adaptation of these thresholds is needed.

In Fig. 2, $up\_th_1$, $down\_th_1$ ($0 < down\_th_1 < up\_th_1 < 1$) are the threshold factors indicating when the measured CDP is above $up\_th_1 * QoS_{rt}$ or below $down\_th_1 * QoS_{rt}$, the threshold will increase or decrease. Once the threshold is consecutively increased for a certain number of times, denoted by $time\_th$, the cell will inform all of its neighbors to do throttling as we described before. $Pow(up_1, v\_index)$ refers to the $v\_index$ power of $up_1$ ($> 1$), in which $v\_index$ is an integer. We also notice that when the measured CDP exceeds $up\_th_1 * QoS_{rt}$, we immediately boost $v\_index$ to zero if it was negative in previous step. In this way, this scheme is always able to be responsive enough to fulfill the QoS bound criterion. The portion of how to adapt $G_2'$ is omitted since it is similar to that of adapting $G_1$.

To guarantee the second QoS criterion, $G_2$ and $G_3$ are used to make Eq. (2) or (5) hold. There are three parameters, namely, switch, percentage and adj_index. switch is defined as before. Percentage refers to the deviation error the scheme may tolerate and the second criterion is still considered being met. For instance, if percentage is set to be 0.1, this means, as long as the ratio of the right-hand side and the left-hand side of Eq.

```
//Assuming the initial values of G₁, G₂', G₃' are already obtained.
time₁ = 0, time₂ = 0;
v_index = 0, d_index = 0;
switch = true;
if (P_{d, rt} >= up_th₁ *QoS_{rt}) {
            if (v_index < 0)        v_index = 0;
            else         v_index++;
            G₁ = G₁*pow(up₁, v_index);
            time₁++;
            if ((time₁ % time_th) = 0)
                {asking neighboring cells to throttle; time₁ = 0;}
    }
else if (P_{d, rt} < down_th₁ *QoS_{rt}) {
            v_index--;
            G₁ = G₁*pow(up₁, v_index);
            time₂++;
            if ((time₂ % time_th) = 0)
                {ask neighbors to de-throttle; time₂ = 0; }
    }

G₂ = G₁ + G₂'; // Adaptation of G₂' is omitted.

if (switch = = true){
        if (P_{b, rt} / P_{d, nrt} >= W_{rt}/W_{nrt} * [OL_{rt}/OL_{nrt}]* (1+percentage))
            adj_index++;
        else if (P_{b, rt} / P_{d, rt} <= W_{rt}/W_{nrt} * [OL_{rt}/OL_{nrt}]* * (1-percentage))
            adj_index--;
        }
else{
        if (P_{b, rt} / P_{d, nrt} >= W_{rt}/W_{nrt} * [OL_{rt}/OL_{nrt}]* (1+percentage))
            adj_index--;
        else if (P_{b, rt} / P_{d, nrt} <= W_{rt}/W_{nrt} * [OL_{rt}/OL_{nrt}]* (1-percentage))
            adj_index++;
        }
G₃' = G₃'*pow(up₃, adj_index);
if (adj_index < adj_index_th) {
            adj_index = 0;
            switch = ! switch;
        }
G₃ = G₂ + G₃';
```

Figure 2. Reservation thresholds adaptation

(2) or (5) is within the range [0.9, 1.1], the equations hold and the QoS criterion is met. The role of *adj_index* is very similar to *v_index*. However, in the adaptation here, we change $up_3$ ( $>$ 1), according the value of *adj_index* in a way that, the larger the absolute value of *adj_index*, the faster the adaptation speed. This ensures the adaptation of $G_3$ can promptly respond to the change of the incoming traffic and (or) QoS status. Finally, when *adj_index* is less than a threshold, *adj_index_th*, which means $G_3'$ is nearly zero, the scheme will reverse the parameter *switch*, letting $G_3$ reserved for the other traffic class instead of the current traffic class it is for.

### C. Probability of Throttling

Each cell keeps a $J \times K$ non-negative integer array $A$ for each traffic class for its neighbors. $J$ is the number of traffic classes and $K$ is the number of neighboring cells. For real-time and non-real-time traffic considered in this paper, $J$ is equal to 2. If the cell's $i$th ($i = 0, 1, \ldots K-1$) neighbor sends a message to the cell to throttle or de-throttle a real-time traffic class, then $A[0][i]$ is incremented or decremented by 1. It is similar for non-real-time traffic. When making admission decision for an incoming new connection request, the cell will use the

following equation to generate the probability of throttling new connections:

$$prob_{rt} = b^{\max(A[0][i])}$$
$$prob_{nrt} = b^{\max(A[1][i])} \quad (9)$$

where $b$ is a real number less than and close to 1, say 0.9.

## V. SIMULATION RESULTS AND ANALYSIS

In this section, we present the performance of our proposed scheme through simulation carried out with OPNET Modeler 8.0. The simulation model is a wrap-around model [8]. The total number of BUs in each cell is 50. The number of BUs each real-time or non-real-time connection needs is $BW_{rt}$ = 1 or $BW_{nrt}$ = 4. The real-time connection may be voice calls and the non-real-time connection may represent file transfer or web browsing. For real-time traffic, the mean duration $1/\mu_{rt}$ = 300 seconds and the mean cell residence time $1/\gamma_{rt}$ = 150 seconds. For non-real-time traffic, $1/\mu_{nrt}$ =1500 seconds and $1/\gamma_{nrt}$ = 750 seconds. On average, each connection will handoff once during its lifetime. A handoff request will randomly choose a destination from the six neighboring cells. We assume that 25% of traffic is real-time traffic, and 75% of the traffic is non-real-time traffic. New connections arrive according to a Poisson process. According to the assumption, 86.96% of the new connection arrivals are real-time, and the rest are non-real-time. For both DMTBR_A and DMTBR_G, $QoS_{rt}$ = 0.01 and $QoS_{nrt}$ = 0.05. The ratio $W_{rt}/W_{nrt}$ is equal to 1, with the deviation error 10%.

Fig. 3 and 4 show CBP and CDP for both traffic classes, as a function of average new connection arrival rate for both DMTBR_A and DMTBR_G. Through calculation, we know that arrival rate 0.1 connection/sec corresponds to about 110 Erlangs, which are 220% of the full load. In fig. 3, as expected, we can see that, for DMTBR_A, the CBP for the two classes are almost equal to each other. For DMTBR_G, since the ratio of offered load of each traffic class is taken into consideration, which is equal to 1:3, the ratio of CBP for the two traffic classes is also about 1:3. This is consistent with Eq. (2) or (5). Through direct calculation, we find out that as an average, the CBP for real-time traffic is reduced 58.23% in DMTBR_G compared to that in DMTBR_A, while the CBP for non-real-time traffic is only increased 9.89% compared to that in DMTBR_A. In fig. 4, both schemes successfully keep the CDP of both traffic classes under the predefined QoS bounds as expected, even when the network is experiencing heavy traffic. Also, there is no big difference in these two schemes in terms of CDP.

Next, the performance is investigated in terms of throughput. The system throughput is defined as followes:

$$TP = \frac{\sum_i BW_i * \text{time spent by each conn. } i \text{ in a cell}}{C * CELL\_NUM * ST} \quad (10)$$

where $C$, as mentioned before, is the total number of BUs available in each cell, $CELL\_NUM$ is the total number of cells in the entire network and $ST$ is the total simulation time. It can be observed that both schemes successfully achieve a stable system throughput even under heavy traffic situation, as

shown in fig. 5. The network throughput keeps increasing as the offered load increases, showing very little difference from each other. Combining the observation in fig 3 and 4, we see clearly the advantage of DMTBR_G over DMTBR_A, i.e., the benefit gained by generalizing the concept of relative priority. The network does not lose anything (in terms of network throughput); however, the user satisfaction for real-time new traffic (in terms of CBP) is significantly increased while the user satisfaction for non-real-time new traffic is only slightly affected. Meanwhile, the user satisfaction for handoff traffic (in terms of CDP) is well maintained for both schemes.

Finally, we consider the detailed throttling operations in each cell. Fig. 6 shows the throttling probabilities for both types of traffic, starting from the beginning of a simulation run

for arrival rate = 0.1 in cell 0 and cell 36. Cell 0 is in the center and cell 36 is located in the edge in the simulation model. As time passes, the throttling probabilities for real-time traffic are almost 1, which means the neighboring cells of cell 0 or 36 rarely throttle the new connection acceptance. This is consistent with Fig. 4, where the CDP for real-time traffic is well kept below the predefined QoS bounds, indicating there is no need to reduce the new connection admission for fulfilling the first QoS criterion. For non-real-time traffic, as time passes, the throttling probabilities first drop, then fluctuate around a certain value after the network enters into a steady state. Thus, we know that the cells keep the first QoS criterion for non-real-time traffic with the help of cooperative neighbors, which reduce the admission probability for new connections due to non-real-time traffic when necessary.

## VI. CONCLUSIONS

In this paper, a dynamic multiple-threshold bandwidth reservation (DMTBR) scheme is proposed to guarantee QoS provisioning in wireless multimedia networks. By dynamically updating the bandwidth thresholds according to network traffic situation and QoS criteria, this scheme works well to provide QoS guarantee and efficiently use network resource, as shown in the simulation. Some of the benefits acquired by generalizing the concept of relative priority are also shown. Since the proposed scheme involves slight changes to the architecture of the current wireless network, it can be easily adopted by 3G and beyond wireless systems.
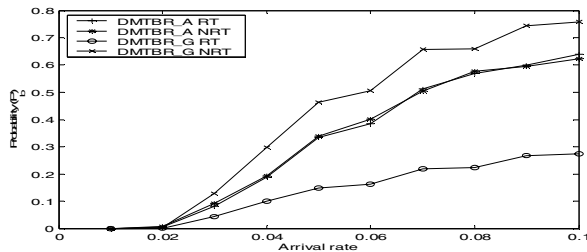


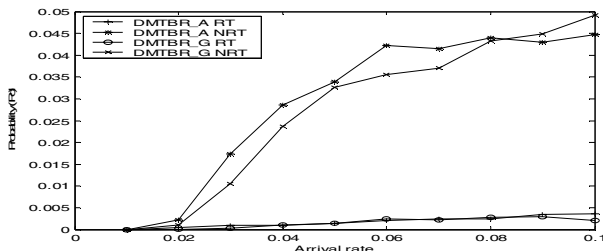Figure 3. CBP vs. Arrival rate



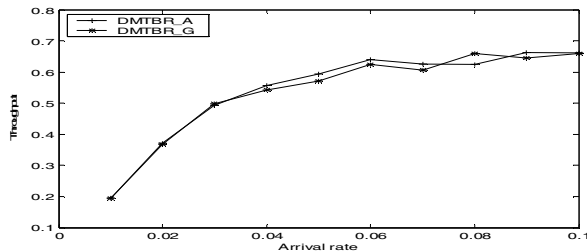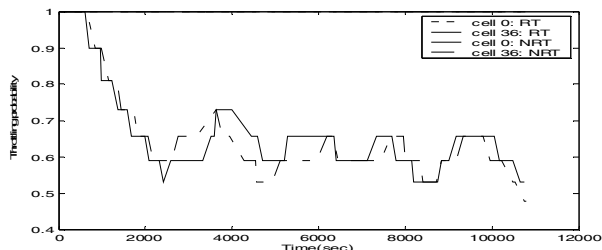Figure 4. CDP vs. Arrival rate



Figure 5. System throughput vs. Arrival rate



Figure 6. Throttling probability

REFERENCES

[1] D. Hong and S. S. Rappaport, "Traffic model and performance analysis for cellular mobile radio telephone systems with prioritized and nonprioritized handoff procedures," IEEE *Trans. Veh. Technol.*, vol. VT-35, no.5, pp. 77-92, Aug. 1986

[2] R. Guerin, "Queueing-blocking system with two arrival streams and guard channels," IEEE *Trans. Commun.*, vol. 36, pp. 153-163, Feb. 1988

[3] C.-J. Chang, T.-T. Su and Y.-Y. Chiang, "Analysis of a cutoff priority cellular radio system with finite queueing and reneging/dropping," IEEE/ACM *Trans. Networking*, vol. 2, no. 2, pp. 166-175, Apr. 1994

[4] P. Ramanathan, K. M. Sivalingam, P. Agrawal and S. Kishore, "Dynamic resource allocation schemes during handoff for mobile multimedia wireless networks," IEEE *J. Select. Areas Commun.*, vol. 17, no. 7, pp. 1270-1283, Jul. 1999

[5] C. Oliver, J. B. Kim and T. Suda, "An adaptive bandwidth reservation scheme for high-speed multimedia wireless networks," IEEE *J. Select. Areas Commun.*, vol. 16, no. 6, pp. 858-874, Aug. 1998

[6] B. M. Epstein and M. Schwartz, "Predictive QoS-based admission control for multiclass traffic in cellular wireless networks," IEEE *J. Select. Areas Commun.*, vol. 18, no. 3, pp. 523-534, Mar. 2000

[7] A. I. Elwalid and D. Mitra, "Effective bandwidth of general Markovian traffic sources and admission control of high-speed networks," IEEE/ACM *Trans. Networking*, vol. 2, no. 2, pp. 166-175, Apr. 1994

[8] I. Chlamtac, Y. Fang and H. Zeng, "Call blocking analysis for PCS networks under general cell residence time," in Proc. IEEE WCNC, Sept. 1999