

Impacts of Buffering of Voice Calls in Integrated Voice and Data Services

Eser Gemikonakli, Orhan Gemikonakli, Enver Ever and Glenford Mapp
Middlesex University, London
e.gemikonakli, o.gemikonakli, e.ever, g.mapp@mdx.ac.uk

Abstract—In this study, we aim to analyse the relationship between various characteristics of a communication system with data and voice call requests. Queuing theory and Markov chain analysis are effectively used for this purpose. Such a study is useful for understanding how the proposed mathematical models behave which represents a system with integrated voice and data calls in homogenous wireless networks. We also propose to optimise the system characteristics in an attempt to provide better Quality of Service (QoS) for systems with integrated voice and data calls. The proposed models have two dimensions; one for voice calls and one for data calls. A channel is assigned for two input traffic call, namely, voice and data calls. The incoming voice and data calls are queued when the channel is busy. Since voice calls are delay-sensitive, priority is given to voice calls. Also, since there is only one channel, data calls are only serviced if there are no voice calls in the system. For such systems, it is important to analyse the impact of buffering the voice calls as well as data calls for various mean rates of call requests, and mean service times. The analytical models presented are generic which is applicable for various systems with similar characteristics. Numerical results are also provided. The results show that the proposed models can be used for optimisation of the performance of a given network.

Keywords-component; analytical modelling; performance analysis; integrated voice and data calls; QoS

I. INTRODUCTION

The increasing demand for unified multimedia communication on a single IP Network diverted network designers of both service providers and enterprises to use integrated voice and data systems. An Integrated system is a way to reduce the cost and take advantage of underused network capacity. Voice over Internet Protocol (VoIP) allows both voice and data communication to be run over a single network, which significantly reduces infrastructure cost.

The overall effect of service performance which decides the level of contentment of a user of the service can be defined as Quality of Service (QoS) [1]. In communication systems, service performance is closely related to network performance. Voice calls are delay sensitive. Hence, the impact of using various buffer sizes for voice calls, on overall system performance has significant importance. It is essential to specify service thresholds in this sense, in order to provide better QoS in terms of both data and voice communications.

Exact solutions are computationally very expensive for heavy load multi-dimensional systems because of the state space explosion problem. In the analysis of integrated voice and data systems a decomposition-approximated method for

the integrated voice/data has been studied extensively in the literature. The GPRS model for the analysis is characterized as a one dimensional decomposing Markov chain with fixed buffer size, and the impact of buffer assignment on the GPRS traffic is also investigated [2], [3], [4]. Although the decomposition method for infinite queuing capacity fails in the overload region, it provides a significant approximation solution technique for finite buffering performance analysis.

The effect of buffering on the QoS parameters for each traffic type is examined in [5]. A channel allocation scheme with dynamic threshold for two different types of traffic; voice calls and data calls is presented. This shows that, the proposed scheme deal with distinct types of traffic flows. Also a finite queuing capacity scheme is used for efficient resource utilization. Simulation results show that the proposed model gives better performance in resource management.

In [6], two types of traffic; real-time (voice/or video) and non-real-time (data) traffic flows are considered as a single multimedia stream. The authors highlight the necessity of considering relationship between two types of flows within the same multimedia stream. Arrivals are described by the Batch Marked Markov Arrival Process (BMMAP) where service time distribution is PH. The queue is finite and the corresponding queuing system's behaviour is explained as multi-dimensional continuous time skip-free to the left Markov chain. In most of the studies, in performance analysis of integrated systems, Poisson process is used for the inter-arrival time of voice and data calls and exponential distributions are used for the service times of the channels [1], [5], [6], [7], [8].

The rest of this paper is organised as follows: Section II presents a model description of the integrated system. In section III, exact product form approach for the model with finite queuing capacity is presented. Performance measures presented here are calculated using both the proposed approach and spectral expansion. Section IV shows the results obtained comparatively. In section V, and VI, conclusion and challenging tasks for further studies are presented respectively.

II. MODEL DESCRIPTION

The proposed models consider modelling a single cell in a wireless network for different traffic behaviours and service requirements to attain the best trade-off. A homogenous wireless system is considered. The proposed models consider systems with finite as well as infinite queuing capacities. A single channel system is considered for both models in order to analyse the behaviour of the channel with priority given to voice calls. The channel is

assigned two different types of traffic; real-time (voice) and non-real-time (data) traffic flows.

When the channel is busy, a buffer is used for the incoming traffic flows. Voice traffic has priority over data traffic since data calls are usually less delay sensitive. It is essential to adapt the delay requirements for such systems in order to have better QoS. The queuing capacity is limited for voice calls arriving at a mean rate of σ_{vc} , and given as B_{vc} as shown in Fig. 1 and Fig. 2. The maximum number of voice calls allowed in the system is equal to one voice call assigned to the channel ($C=1$) in the system plus the queuing capacity B_{vc} . The maximum number of voice calls in the system is given by L_{vc} , where $L_{vc} = C + B_{vc}$. On the other hand, data calls arriving at a mean rate of σ_{dt} , cannot be lost but can tolerate some delay or jitter [4]. In Fig. 1, the number of data calls accepted in the system is equal to one data call being serviced plus queuing capacity B_{dt} . The maximum number of data calls in the system is given by L_{dt} , where $L_{dt} = C + B_{dt}$. In order to prevent data loss for the proposed model, queuing capacity for data calls has been assumed to be large. The model considered has Markov processes for arrival and departure of voice as well as data calls, where voice and data call requests share a common queue.

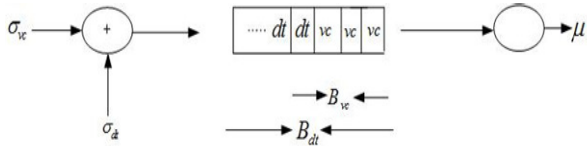


Figure 1. The model of the system with finite buffering

In Fig. 2, a model with an infinite buffering capacity is shown. These systems are also considered in order to analyse the effects of queuing delay. In other words L_{dt} goes to infinity.

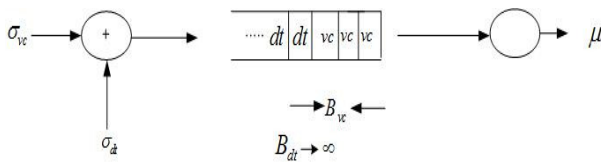


Figure 2. The model of the system with infinite buffering

The systems considered in Fig. 1 and Fig. 2 consists of a single channel with a common queue. The call requests arrive at the system in a Poisson stream and join the queue when no server is available. The voice call requests as well as data call requests are stored in the queue as they arrive, however since voice requests have priority over data requests, it is possible to assume that voice calls are always in front of the data calls as illustrated in Fig. 1 and Fig. 2.

For wireless communication systems, exponentially distributed service times, and Poisson arrivals are commonly used [8], [9], [10], [11]. For both models, the arrivals of voice and data calls are assumed to be independent and follow Poisson processes with mean arrival

rates σ_{vc} and σ_{dt} respectively. The service time of voice and data calls, are assumed to follow exponential distribution with means $1/\mu_{vc}$ for voice calls and $1/\mu_{dt}$ for data calls respectively.

The calls in the system receive service in a pre-emptive fashion. If a data call is pre-empted due to the arrival of a voice call, its service restarts whenever the channel is available. In this case the service policy is repeated with re-sampling.

The algorithms used for the admission of voice and data calls are illustrated in Fig. 3 and Fig. 4 for systems with finite and infinite queuing capacities respectively. When there is a new request; either voice or data call, and a channel is available the request is assigned to the channel. If the channel is not available, voice calls are queued ahead of data calls as long as the amount of the occupied buffer in the wireless cells is less than or equal to the total capacity, B_{vc} , as shown in Fig. 1 and Fig. 2. If there are more than one voice calls, they are serviced according to the first come, first served (FCFS) queuing discipline. When the voice calls reach to limited buffer size, the incoming call arrivals are rejected. In Fig. 3, since finite queue is assumed for both traffic types, data calls can also be rejected similarly. Since the proposed model has a single channel, data calls are only serviced when there is no voice call in the system. In the absence of voice calls, if there are more than one data calls in the system they also receive service according to the FCFS queuing discipline.

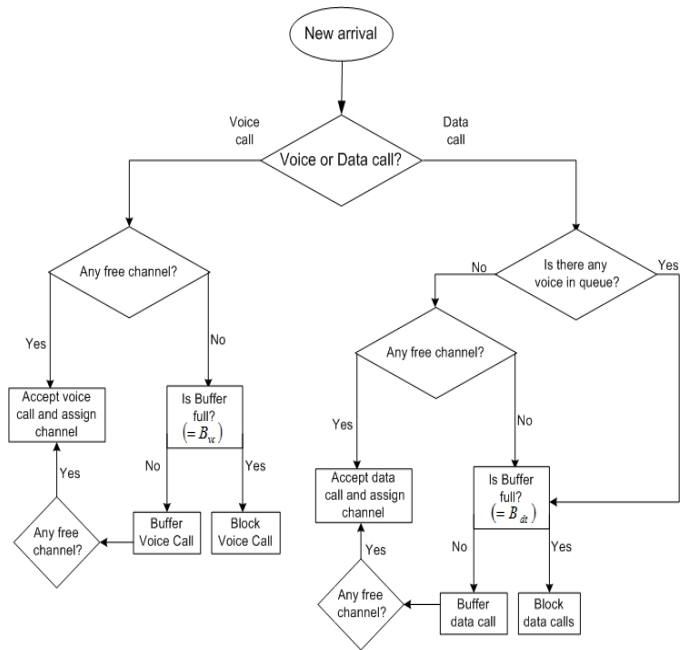


Figure 3. Flowchart of the traffic flows with finite buffering

In Fig. 4, data calls have infinite buffering capacity. As B_{dt} goes to infinity due to the infinite queuing capacity, L_{dt} goes to infinity as well. Therefore, data calls are not blocked.

$$C_{0,j} = \frac{1}{\mu_{dt}} \left((\sigma_{vc} + \sigma_{dt} + \mu_{dt}) C_{0,j-1} - \sigma_{dt} C_{0,j-2} - \mu_{vc} C_{1,j-1} \right) \quad i=0, 1 < j < L_{dt} - L_{vc} \quad (11)$$

$$C_{0,L_{dt}-L_{vc}} = \frac{1}{\mu_{dt}} \left((\sigma_{vc} + \sigma_{dt} + \mu_{dt}) C_{0,L_{dt}-L_{vc}-1} - \sigma_{dt} C_{0,L_{dt}-L_{vc}-2} - \mu_{vc} C_{1,L_{dt}-L_{vc}-1} \right) \quad i=0, j=L_{dt}-L_{vc} \quad (12)$$

$$C_{0,j} = \frac{1}{\mu_{dt}} \left((\sigma_{vc} + \sigma_{dt} + \mu_{dt}) C_{0,j-1} - \sigma_{dt} C_{0,j-2} - \mu_{vc} C_{1,j-1} \right) \quad i=0, L_{dt} - L_{vc} < j \leq L_{dt} \quad (13)$$

$$d_{i,j} = \frac{\sigma_{vc}}{(\sigma_{vc} + \sigma_{dt} + \mu_{vc}(1.0 - d_{i+1,j}))} \quad 0 < i < L_{vc}, 0 \leq j < L_{dt} - L_{vc} \quad (14)$$

$$d_{i,j} = \frac{\sigma_{vc}}{(\mu_{vc} + \sigma_{dt})} \quad i=L_{vc}, 0 \leq j < L_{dt} - L_{vc} \quad (15)$$

$$d_{i,j} = \frac{\sigma_{vc}}{\mu_{vc}} \quad i=L_{dt}-j, L_{dt}-L_{vc} \leq j < L_{dt} \quad (16)$$

$$d_{i,j} = \frac{\sigma_{vc}}{(\sigma_{vc} + \sigma_{dt} + \mu_{vc}(1.0 - d_{i+1,j}))} \quad 0 < i < L_{dt}-j, L_{dt}-L_{vc} \leq j < L_{dt} \quad (17)$$

$$a_{L_{vc},j} = \frac{\sigma_{dt}}{(\mu_{vc} + \sigma_{dt})} C_{L_{vc},j-1} \quad i=L_{vc}, 1 \leq j < L_{dt} - L_{vc} \quad (18)$$

$$a_{i,j} = \frac{(\mu_{vc} a_{i+1,j} + \sigma_{dt} C_{i,j-1})}{(\sigma_{vc} + \sigma_{dt} + \mu_{vc}(1.0 - d_{i+1,j}))} \quad 1 \leq i < L_{vc}, 1 \leq j < L_{dt} - L_{vc} \quad (19)$$

$$a_{i,j} = \frac{\sigma_{dt}}{\mu_{vc}} C_{i,j-1} \quad i=L_{vc}, 1 \leq j < L_{dt} - L_{vc} \quad (20)$$

$$a_{i,j} = \frac{(\mu_{vc} a_{i+1,j} + \sigma_{dt} C_{i,j-1})}{(\sigma_{vc} + \sigma_{dt} + \mu_{vc}(1.0 - d_{i+1,j}))} \quad 1 \leq i < L_{dt}-j, L_{dt}-L_{vc} \leq j < L_{dt} \quad (21)$$

$$P_{00} + \sum_{i=0}^{L_{vc}} P_{i0} + \sum_{j=0}^{L_{dt}} P_{0j} + \sum_{i=0}^{L_{vc}} \sum_{j=0}^{L_{dt}-i} P_{ij} = 1 \quad (22)$$

$$P_{00} = \frac{1}{1 + \sum_{i=0}^{L_{vc}} C_{i0} + \sum_{j=0}^{L_{dt}} C_{0j} + \sum_{i=0}^{L_{vc}} \sum_{j=0}^{L_{dt}-i} C_{ij}}$$

The solution is given for a finite lattice. Performance measures such as, Mean Queue Length (MQL), Blocking Probability, Throughput, and Response Time can be calculated using state probabilities. In (23) – (28), the notations *vc* and *dt* are used for voice and data calls respectively.

$$MQL_{vc} = \sum_{i=0}^{L_{vc}} \sum_{j=0}^{L_{dt}} iP_{i,j} \quad (23)$$

$$MQL_{dt} = \sum_{i=0}^{L_{vc}} \sum_{j=0}^{L_{dt}} jP_{i,j} \quad (24)$$

$$Throughput_{vc} = \sum_{i=1}^{L_{vc}} \sum_{j=0}^{L_{dt}} \mu_{vc} P_{i,j} \quad (25)$$

$$Throughput_{dt} = \sum_{j=1}^{L_{dt}} \mu_{dt} P_{0,j} \quad (26)$$

$$Response\ Time_{vc} = \frac{MQL_{vc}}{Throughput_{vc}} \quad (27)$$

$$Response\ Time_{dt} = \frac{MQL_{dt}}{Throughput_{dt}} \quad (28)$$

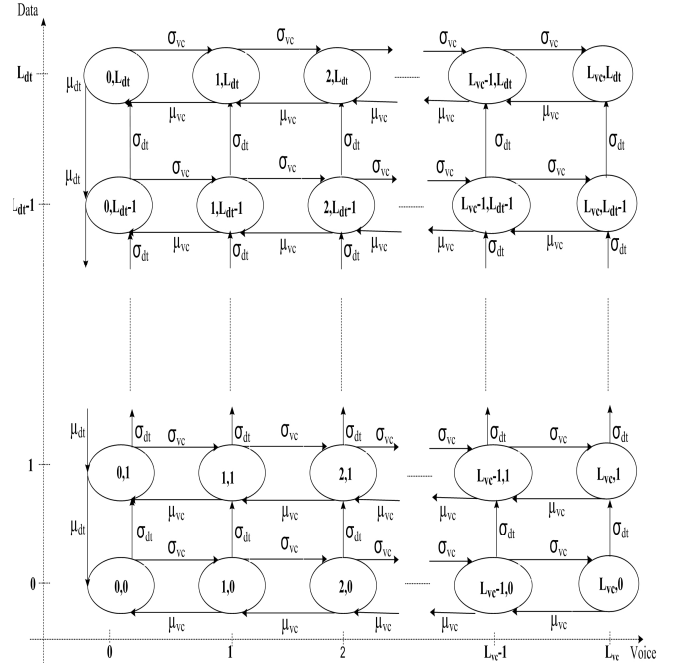


Figure 6. The state transition diagram for the performance model of the system considered with infinite buffering

IV. NUMERICAL RESULTS AND DISCUSSION

In this section, numerical results are presented. The results show the behaviour of the system with integrated voice and data calls. Also in order to validate the product form solution presented, results are presented comparatively with the spectral expansion method. In order to show the accuracy, the discrepancy between the results from the Spectral Expansion method (Fig. 6), and results obtained by using the proposed solution technique are presented.

The parameters used in the mathematical calculations are as follows: σ_{vc} range from 0.01 job/ms to 0.08 job/ms while L_{vc} changes from 1 to 80. The other parameters are taken as $S=1$, σ_{dt} from 0.01 job/ms to 0.06 job/ms, $\mu_{vc}=0.08$ job/ms and $\mu_{dt}=0.06$ job/ms. Since the unbounded systems are considered, very large L_{dt} values are used for the product form approach and unbounded queuing capacity is used for Spectral Expansion method (semi infinite lattice in Fig. 6). The results show that when large L_{dt} values are considered for the product form solution used in this study, the discrepancy between the Spectral Expansion method and the proposed approach is negligibly small (See table 1).

Fig. 7 shows the blocking probabilities and response time of voice calls as a function of L_{vc} for different σ_{vc} values. The results show that as the queuing capacity of voice calls increases, the blocking probability of voice calls decreases significantly in acceptable ranges where response time is less than 100 ms. Since we give priority to voice calls, the blocking probability decreases when the queuing capacity of the voice calls is increased. The applications with voice call requests can tolerate blocking probabilities which are less than 0.01 [15].

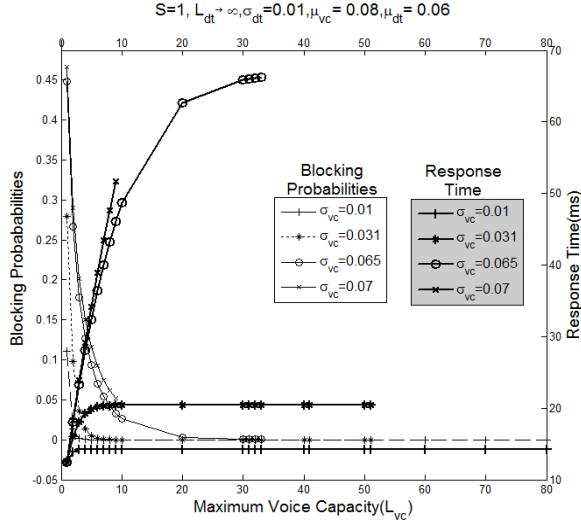


Figure 7. Blocking Probability and Response Time as function of L_{vc}

Fig. 8 shows the effects of L_{vc} on the MQL of voice and data calls for $\sigma_{vc}=0.031$. The results show a sharp increase in MQL as L_{vc} increases for $L_{vc}<8$. However, if the response time is acceptable, increasing the queuing capacity of the voice calls further, does not affect the system. Also, having various σ_{dt} values does not affect the MQL of the voice calls, since the priority is given to the voice call requests.

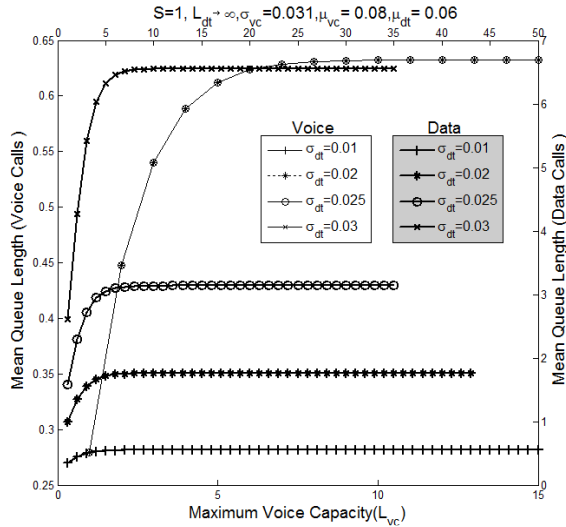


Figure 8. MQL as a function of L_{vc}

Fig. 9 shows the Blocking Probability and Response Time results for voice calls as a function of σ_{vc} . The results show that the arrival rate of the incoming voice calls, affect the systems with low queuing capacities quite significantly. For $L_{vc}=1$, the blocking probability increases rapidly and reaches up to 50%. Similar to the trends of figures considered above, the effects of queuing capacity L_{vc} is less evident for $L_{vc}>20$ for both response time and blocking probability measures. Also please note that for $L_{vc} =1$, the response time does not change since most of the incoming requests are blocked.

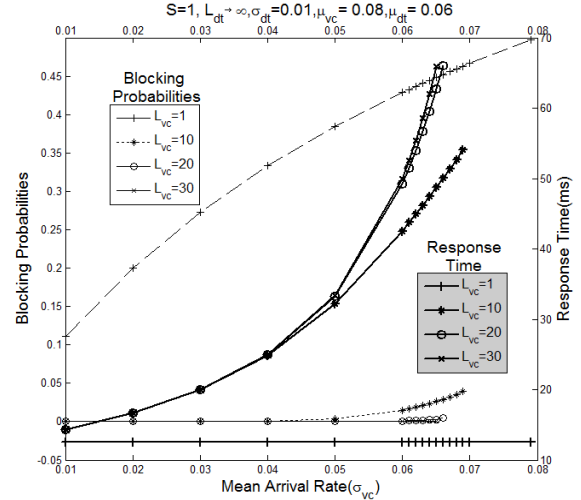


Figure 9. Blocking Probability and Response Time of voice calls as a function of σ_{vc}

Similar to Fig. 9, Fig.10 shows the response time results this time for data calls as a function of σ_{vc} . The results show that the arrival rate of the incoming voice calls, affect the data calls quite significantly. This is mainly because the data calls do not receive service if there are voice calls in the system. Because of the same reason, the system can provide better service to incoming data calls when L_{vc} value is lower. Since the amount of buffered voice requests is smaller, the response time of data calls decreases as L_{vc} decreases.

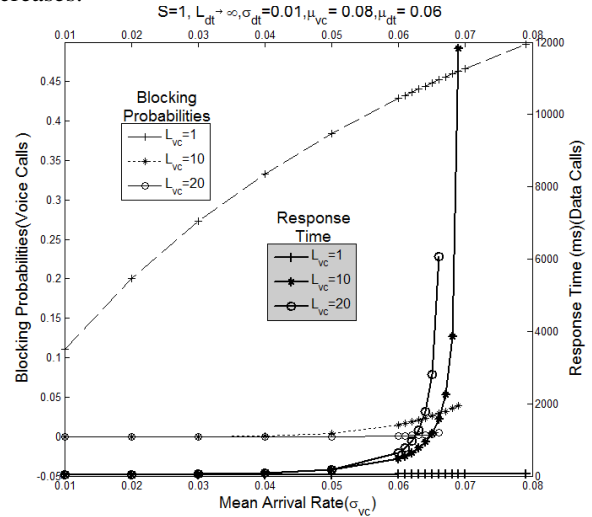


Figure 10. Blocking Probability and Response Time of data calls as a function of σ_{vc}

Further results are considered in order to validate the accuracy of the exact solution approach. Lattice strip provided in Fig. 6 is employed for the Spectral Expansion method. Table 1 clearly shows that when the accuracy of the method is considered, the discrepancy between the product form results and the results obtained from the Spectral Expansion is less than 0.001% for both MQL_{vc} and MQL_{dt} results.

When the data queuing capacity is assumed to be very large, the system considered shows the behaviour of a queuing system having an infinite capacity for data. Results presented in Table 1, where $\sigma_{vc}=0.05$ job/ms, $\sigma_{dt}=0.01$ job/ms, $\mu_{vc}=0.08$ job/ms and $\mu_{dt}=0.06$ job/ms, show that the proposed product form solution technique is useful to solve such models having either finite or infinite queuing capacity.

TABLE 1. COMPARATIVE RESULTS

L_{vc}	Discrepancy					
	Product Form		Spectral		$Mql_{vc}(\%)$	$Mql_{dt}(\%)$
	Mql_{vc}	Mql_{dt}	Mql_{vc}	Mql_{dt}		
1	0.3846	0.4373	0.3846	0.4373	0	0
2	0.6976	0.6964	0.6976	0.6964	0	0
3	0.9464	0.9426	0.9464	0.9426	0	0
4	1.1395	1.1556	1.1395	1.1556	8.78E-05	8.65E-05
5	1.2863	1.3280	1.2863	1.3280	7.78E-05	0.000377
6	1.3958	1.4612	1.3958	1.4612	0.000287	0.000137
7	1.4759	1.5607	1.4759	1.5607	0.000136	0
8	1.5337	1.6331	1.5337	1.6331	0.000326	0.000123
9	1.5748	1.6848	1.5748	1.6848	0.000127	0.000237
10	1.6037	1.7212	1.6037	1.7212	6.24E-05	0.000116
11	1.6238	1.7465	1.6238	1.7465	0.000123	0.000172
12	1.6377	1.7639	1.6377	1.76392	0.000183	0.000170
13	1.6472	1.7757	1.6472	1.7757	0	0
14	1.6536	1.7838	1.6536	1.7838	0.000302	0
15	1.6579	1.7892	1.6579	1.7892	0.000121	0.000224
16	1.660	1.7928	1.6609	1.7928	0.000301	5.58E-05
17	1.6628	1.7952	1.6628	1.7952	0	0
18	1.6641	1.7969	1.6641	1.7969	0	0

V. CONCLUSION

This paper presents a model and its solution technique for finite and infinite queuing capacities for studying the impact of buffering of voice calls as well as data calls having different arrival and service rates. Since the proposed models have only one channel and voice has priority over data, data calls are only served when there are no voice calls in the system. When the channel is not available for the incoming calls, both calls are queued provided that there is room for queuing. The product form solution approach is used to find the threshold of the proposed model in order to prevent voice being delayed and data being lost, providing better QoS. The results are validated using the Spectral Expansion solution. Moreover, product form solution approach has an advantage over spectral expansion; it can deal with larger L_{vc} .

In conclusion, this work shows that unlike current practice, it is possible to queue voice calls providing that a small portion of buffering is made available. This reduces the blocking probability and hence allows for better operation of both traffic services.

VI. FURTHERWORK

With the increasing popularity of integrated voice/data systems, customers are expecting high level QoS to deal with incoming jobs in a more cost-effective way. Traditional

pure performance model ignores failures and recovery of system components. To obtain realistic composite performance with failure and recovery behaviour, availability has to be taken into account. It is also important to evaluate the performance of multi-channel systems. Work is in progress. Using performability analysis we will investigate the limitations of integrated voice and data calls, develop and evaluate different algorithms and propose solutions for better QoS.

REFERENCES

- [1]. X. Massip-Bruin, M.Yannuzzi,J.Domingo-Pascual, A. Fonte, M.Curado, E.Monterio,F.Kuipers, P.Van Miegham, S.Avollene, G.Ventre,P.Aranda-Gutierrez, M.Hollick, R.Steintmetz, L.Iannone, K.Salamatian, "Research challenges in QoS routing" Computer Communication, 29, 2006, pp. 563-581.
- [2]. L. Hung-Huan; J.-L.C. Wu, H.Wan-Chih, "Delay analysis of integrated voice and data service for GPRS," *IEEE Communications Letters*, vol.6, no.8, 2002, pp. 319- 321.
- [3]. S. Ghani, M. Schwartz, "A decomposition approximation for the analysis of voice/data integration," *IEEE Transactions on Communications*, vol.42, no.7, 1994, pp.2441-2452.
- [4]. I. Candan, M. Salamah, "Analytical Modeling of a Time-Threshold Based Multi-guard Bandwidth Allocation Scheme for Cellular Networks," *Fifth Advanced International Conference on Telecommunications, 2009. AICT '09.*, vol., 24-28, pp.33-38.
- [5]. O.A Ojesanmi,; A. A. Ojesanmi, S. O. Ojesanmi, O. Makinde, "Enhanced Channel Allocation Scheme for Integrated Voice/Data Calls in Cellular Network" *International Journal of Intelligent Information Technology Application*,2009,2(2): pp.80-83
- [6]. K. Al-Begain, A. Dudin, V. Mushko," Novel Queuing Model for Multimedia over Downlink in 3.5G Wireless Networks. In Al-Begain, Bolch, Telek (eds.) "ASMTA05:12th International Conference on Analytical and Stochastic Techniques and Applications", ISBN 1-84233-116-7, Riga, Latvia, 1-4 June 2005. pp. 111-117.
- [7]. S. Wu, K. Y. M. Wong, L. Bo, "A new distributed and dynamic call admission policy for mobile wireless networks with QOS guarantee," *The Ninth IEEE International Symposium on Personal, Indoor and Mobile Radio Communications*,., vol.1, pp.260-264, 8-11 Sep 1998
- [8]. W. Shah, S. Soomro, F. Z. Khan, and G. D. Menghwar, "Performance Evaluation of Multistage Service System Using Matrix Geometric Method", *Fourth International Conference on Systems and Networks Communications*, 2009, pp: 265-269.
- [9]. G. Song, L. Cuthbert, and J. Schormans, "Modelling Cellular/Wireless LAN Integrated Systems with Multi-Rate Traffic Using Queueing Network". *Wireless Communications, Networking and Mobile Computing*, 2008. pp: 1-4,
- [10]. K. Trivedi, S. Dharmaraja, and X. Ma, "Analytic Modelling of Wireless Communication Systems", *Information Sciences*, vol. 148, 2002, pp: 155-166.
- [11]. W. Xia, L. Shen, "Modeling and Analysis of Hybrid Cellular/WLAN Systems with Integrated Service-Based Vertical Handoff Schemes". *The Institute of Electronics Information and Communication Engineers, IEICE Transactions*, vol. 92, 2009, pp: 2032-2043.
- [12]. R. Chakka, "Spectral Expansion Solution for Some Finite Capacity Queues", *Annals of Operations Research*, Vol.79, 1998, pp: 27-44.
- [13]. I. Mitrani, and R. Chakka, "Spectral expansion for a class of Markov models: Application and Comparison with the Matrix Geometric Method". *Performance Evaluation Journal*, 1995, 23(3), pp.241-260.
- [14]. I. Mitrani, "Approximate solutions for heavily loaded Markov-modulated queues". *Performance Evaluation*, 2005, 62(1-4): 117-131.
- [15]. A. Zahedi and K. Pahlavan (2000) "Capacity of a Wireless LAN with Voice and Data Services" *IEEE Transaction on Communications*, Vol.48, No.7, July 2000.