



## **UWL REPOSITORY**

**repository.uwl.ac.uk**

Investigating a mobility-aware QoS model for multimedia streaming rate adaptation

Sardis, Fragkiskos, Mapp, Glenford, Loo, Jonathan ORCID logoORCID: <https://orcid.org/0000-0002-2197-8126>, Aiash, Mahdi and Vinel, Alexey (2015) Investigating a mobility-aware QoS model for multimedia streaming rate adaptation. Journal of Electrical and Computer Engineering, 2015. p. 548638. ISSN 2090-0147

<http://dx.doi.org/10.1155/2015/548638>

This is the Published Version of the final output.

UWL repository link: <https://repository.uwl.ac.uk/id/eprint/3522/>

**Alternative formats:** If you require this document in an alternative format, please contact: [open.research@uwl.ac.uk](mailto:open.research@uwl.ac.uk)

**Copyright:** Creative Commons: Attribution 3.0

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy:** If you believe that this document breaches copyright, please contact us at [open.research@uwl.ac.uk](mailto:open.research@uwl.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.

## Research Article

# Investigating a Mobility-Aware QoS Model for Multimedia Streaming Rate Adaptation

**Fragkiskos Sardis,<sup>1</sup> Glenford Mapp,<sup>2</sup> Jonathan Loo,<sup>2</sup> Mahdi Aiash,<sup>2</sup> and Alexey Vinel<sup>3</sup>**

<sup>1</sup>*Department of Informatics, King's College London, London WC2R 2LS, UK*

<sup>2</sup>*School of Science and Technology, Middlesex University, London NW44BT, UK*

<sup>3</sup>*Halmstad University, 301 18 Halmstad, Sweden*

Correspondence should be addressed to Fragkiskos Sardis; [fragkiskos.sardis@kcl.ac.uk](mailto:fragkiskos.sardis@kcl.ac.uk)

Received 20 May 2015; Accepted 26 July 2015

Academic Editor: Peppino Fazio

Copyright © 2015 Fragkiskos Sardis et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Supporting high quality multimedia streaming on wireless devices poses several challenges compared to wired networks due to the high variance in network performance encountered in the mobile environment. Although rate adaptation is commonly used in multimedia applications to compensate for fluctuations in network performance, it is a reactive mechanism which is not aware of the frequently changing connectivity that may occur on mobile devices. This paper proposed a performance evaluation model for multimedia streaming applications that is aware of user mobility and network performance. We presented an example of mathematical solution to the model and demonstrated the functionality using common mobility and connectivity examples that may be found in an urban environment. The proposed model is evaluated based on this functionality and how it may be used to enhance application performance.

## 1. Introduction

Multimedia content on the Internet has evolved from simple audio and images to high definition video and highly interactive video games. In recent years, this content has made its way to mobile devices such as smartphones and tablets. The capabilities of these devices have tremendously increased over the years in terms of processing power; however, in a mobile environment where network connectivity may be changing frequently, it is very difficult to achieve a high Quality of Experience (QoE) when it comes to streaming multimedia content. Although modern mobile devices are capable of streaming multimedia content at very high bit-rates, it is often the network's Quality of Service (QoS) that cannot deliver such content in a consistent and timely fashion. While this is also true for wired networks, wireless networks are more susceptible to performance degradation due to congestion at the access point or weak signaling, especially when dealing with mobile networks or other open-access wireless networks.

Multimedia streaming applications make up for varying QoS conditions by caching parts of the content opportunistically when possible or scaling the bit-rate and hence the quality of the audiovisual content in response to network conditions. These mechanisms provide a dampening effect on the varying QoS in wired networks; however, in wireless networks where the user may be moving and the device may be switching between heterogeneous networks, the QoS variance can be far greater than that found in wired networks even when not considering packet losses and jitter caused by handover events.

In this paper we present a new approach at modeling the performance of multimedia applications based on the anticipated network connections caused by mobility and the rate of movement of the user. A mathematical solution to the multidimensional Markov model is presented along with examples given to demonstrate the functionality of the model. The novelty of this approach lies in the modeling of a multimedia application in a mobile environment from the perspective of the client rather than the perspective of each

individual network and therefore provides an overview of an applications performance over multiple networks.

The rest of the paper is outlined as follows: Section 2 includes background information in the field of scalable multimedia content, mobile service delivery, and mobile network technology. Section 3 presents the model under investigation and Section 4 provides an example solution for a two-dimensional model. Section 5 illustrates the functionality of the model with examples. Section 6 provides critical evaluation of the model and concludes this paper.

## 2. Background

In this section we present some background information in the fields of wireless network performance and multimedia content delivery mechanisms.

**2.1. Mobile Service Delivery.** In the context of mobile services, network performance is typically evaluated from the perspective of the network by means of queueing analysis [1–3]. Although the study of wireless network performance from the perspective of the network can help in optimizing their performance and as a result the delivery of Internet services to mobile clients, there is currently no focus on evaluating the performance of networks from the perspective of an application or a service. To this end, a mobile service delivery framework has been proposed [4, 5] in the context of 5th generation networks such as Y-Comm [6] where reliable and constant connectivity is achieved at all times by means of seamless horizontal and vertical handovers. The proposed framework uses network mechanisms to constantly probe networks adjacent to the user and select the best possible connection that satisfies the requirements of the user's applications. Thus far, only the traffic management aspects of the framework have been investigated through the use of Cloud technology for dynamic localization of services by means of Wide-Area-Network (WAN) migrations. In this paper, we explore the QoS aspects of the framework and therefore everything described in the following sections should be considered in the context of [4–6].

For the traffic management aspects of this framework, the network dwell time is considered based on user mobility patterns. However, in order to study the performance of an application over a series of networks, we need to express mobility as a rate of movement in terms of exiting the coverage of a network and entering another.

The cell outgoing rate  $\mu_{dwell}$  is defined as the average exit rate of uniformly distributed users equally likely to move in any direction with arbitrary distribution of moving speed [3]:

$$M_{cdwell} = \frac{E(v)L}{\pi A}, \quad (1)$$

where  $E(v)$  is the average velocity of the users,  $L$  is the length of the cell's perimeter, and  $A$  is the area of the cell.

Using this expression, we can determine the average exit rate for a given velocity and cell size and thus we can estimate the rate at which a user will be switching between cells. However, it should be noted that this is not an accurate

representation for each possible scenario of mobility and cell size but it is a good approximation of the rate without knowing in detail and with great accuracy where the user entered the cell, which direction they are going, and at what speed they are moving.

**2.2. Scalable Content Delivery.** Scalable multimedia content delivery encompasses mechanisms that dynamically adjust the quality of multimedia streams in such way that it adapts to varying network conditions. The main goal of the technology is to deliver content without interruptions at the cost of decreasing the audiovisual quality when necessary. Such mechanisms have been proposed for more than a decade [7] and are now widely used by online multimedia services such as YouTube [8].

As Rejaie et al. argue [7], multimedia streaming applications are subject to two conflicting requirements. The first requirement is that these applications are delay-sensitive and rate-based, and thus they require isochronous processing and end-to-end QoS guarantees. This stems from the fact that stored video has a predefined bit-rate which needs to be transmitted at a fixed rate and therefore requires constant bandwidth. On the other hand, the second requirement is that the Internet is a shared environment and therefore end systems are expected to cooperate by reacting to congestion properly and promptly by deploying congestion control mechanisms. Thus, the available bandwidth may vary in an unpredictable manner and more importantly large amounts of data streaming over the network may be the cause for triggering congestion control mechanisms. Rejaie et al. demonstrate that, by exploiting the flexibility of layered encoding, it is possible to maintain stable streaming by switching between the different encoding layers according to network performance. With layered encoding, each layer holds incrementally more fine details of the content. When sufficient bandwidth is available, the service streams the information from all the layers. When network QoS degrades and TCP congestion control is activated, the service drops some of the layers and maintains streaming of the more coarse layers which require less throughput.

Raghuveer et al. [9] enhance the layered quality adaptation mechanisms by considering not only the status of the network but also the status of the buffer at the client. When near a buffer underflow, the proposed system increases the sending rate from the service. Conversely when near a buffer overflow, the system decreases the transmission rate even when there is no congestion on the network.

What becomes apparent from the above approaches is that, in order to provide scalable adaptation of content, there must be a strategy which defines how content quality is scaled according to network conditions. Typically, network condition is derived by monitoring the arrival rate of packets and comparing it to the consumption rate at the client which means that by definition, quality adaptation is reactive.

One attempt to probe the network in advance and subsequently at frequent intervals in order to determine the rate adaptation is presented by Li et al. in [10]. They propose a Probe and Adapt (PANDA) mechanism which compares the TCP throughput to the expected throughput

for uninterrupted streaming. By probing at frequent intervals they can dynamically adjust subsequent service requests and adapt the quality of the content so that the desired rate is met. This method of probing is quicker at responding to network changes and with sufficient buffering, helping to smooth out aggressive rate adaptation.

The main disadvantage of the state-of-art solutions is that they lack any form of QoS prediction in a mobile environment where even short-term probing may provide completely inaccurate information when network handover occurs. To this end, the optimal solution would be to monitor user mobility, predict future network configurations, probe each network in the user's path, and adjust the rate accordingly in a preemptive manner.

### 3. Mobile Service Performance Model

The structural complexity of the Internet along with factors such as distance and user demands leads to greatly varying levels of performance between different networks and locations. Latency, throughput, and response time are some of the determining factors to the performance of online multimedia applications and consequently to the perceived QoE. As discussed previously, there are various mechanisms on the network and on the service side that dynamically adapt the behavior of applications to mitigate this performance variance. However, these mechanisms often work independently of each other often resulting in duplicate functionality on the network and service side. For example, networks prioritize packets based on the application while at the same time multimedia applications attempt to adjust the quality of the content to adapt to network conditions.

In this section, we present a novel approach to modeling application performance based on user mobility, application requirements, and network service rates. The distinctive characteristic of the proposed model is that it evaluates the performance of a persistent connection from the perspective of the application rather than from the perspective of the network. As such, it offers insights as to how an application will behave over multiple networks across a user's path without considering the intrinsic details of network mechanisms. In other words, this model can be used to provide an overview of the performance of an application so that the application itself may decide how to best optimize the delivery of multimedia content according to the performance of anticipated network connections resulting from user mobility patterns.

Figure 1 shows a simple Markov chain representing a connection between a service and a client over the network. In this case a queue depth of two requests is depicted; however, this may be scaled to any buffer size desired. In order to introduce mobility and network handover as a factor of the performance, we will need to represent networks as individual chains and the user mobility as a rate of switching between those chains as illustrated in Figure 2.

The service request rate of the application is defined as  $\lambda$  and it is constant across all the networks if we assume a single application used constantly by the client. The service rate ( $\mu_n$ ) of each network is defined individually to represent the varying levels of QoS that each point of attachment

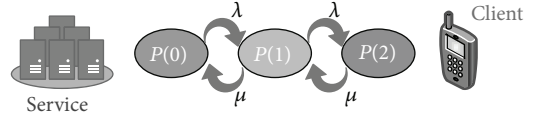


FIGURE 1: Markov chain with a buffer length of two requests.

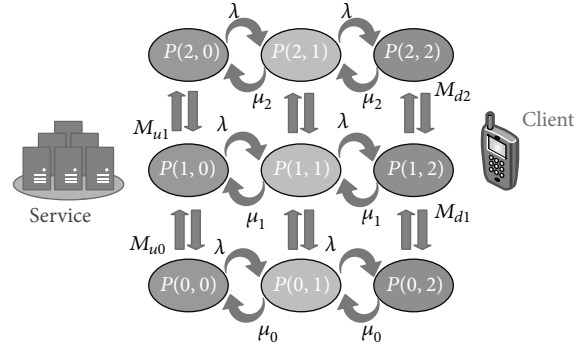
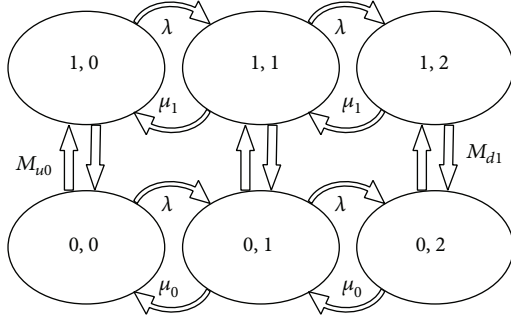


FIGURE 2: Multidimension model for evaluating application performance over a series of networks.

may deliver. The mobility rate of the user is defined as  $M_{un}$  and  $M_{dn}$  individually for each chain in order to accurately represent different coverage areas of networks and therefore different handover rates from one network to another. Mobility between chains is represented as upwards ( $M_u$ ) and downwards ( $M_d$ ) in order to address scenarios where a user may be switching between networks of the same provider or simply to address mobility cases where a user may be moving in such a pattern that causes an oscillating behavior between networks.

Being able to model the performance of an application in this manner presents new possibilities for the delivery of multimedia services over mobile networks. For example, using this model can help predetermine how an application will behave along a user's path provided that we know that the path and velocity are constant and which networks are present along that path. Consequently, we could instruct a multimedia service to preemptively adjust the quality of a stream by means of precaching some sections at lower quality so that they will be ready when needed. Alternatively, in the context of content-centric networks, we could instruct different parts of a multimedia file to be delivered by different locations that may provide better performance on networks that have insufficient service rate for a particular service instance. The same may be applied on service-centric networks where alternate component services may be used for networks where the performance is insufficient for the existing composite service.

To achieve the above, the model assumes that probe connections are used to gather performance metrics for each network along the path. Furthermore, mechanisms that identify where the performance degradation comes from are also needed in order to identify cases where congestion at the access point is causing the problem and therefore there is nothing that can be done on the network or the service side that will improve the performance. Therefore, in this

FIGURE 3:  $2 \times 3$  mobile network QoS model.

paper we explore this model at a theoretical level and assume that performance degradation occurs within the network's backbone infrastructure rather than the client access points.

The next section shows how this model can be solved mathematically followed by some examples of how this model can be used.

#### 4. Example Solution

In this section we present an example solution for a  $2 \times 3$  model which represents how service requests are being queued in two networks along a user's path and how the user's mobility pattern affects the overall service rate received at the client. The model is illustrated in Figure 3 where  $\lambda$  is the service rate requested by the client,  $\mu_n$  is the perceived service rate of each network at the client,  $M_{u0}$  is the mobility rate at which the client leaves the network represented by the bottom chain and enters the one above, and finally  $M_{d1}$  is the mobility rate at which the client leaves the network represented by the top chain and enters the one below.

To solve this model, we start by expressing each state as a function of its inbound and outgoing rates from and to other states. Thus, we have

$$\begin{aligned}
 (M_{u0} + \lambda) P_{0,0} &= \mu_0 P_{0,1} + M_{d1} P_{1,0}, \\
 (M_{u0} + \lambda + \mu_0) P_{0,1} &= M_{d1} P_{1,1} + \mu_0 P_{0,2} + \lambda P_{0,0}, \\
 (M_{u0} + \mu_0) P_{0,2} &= M_{d1} P_{1,2} + \lambda P_{0,1}, \\
 (M_{d1} + \lambda) P_{1,0} &= M_{u0} P_{0,0} + \mu_1 P_{1,1}, \\
 (M_{d1} + \lambda + \mu_1) P_{1,1} &= M_{u0} P_{0,1} + \mu_1 P_{1,2} + \lambda P_{1,0}, \\
 (M_{d1} + \mu_1) P_{1,2} &= M_{u0} P_{0,2} + \lambda P_{1,1}.
 \end{aligned} \tag{2}$$

We then proceed to unzip the model by expressing each state as a function of its adjacent states. We have

$$P_{1,2} = a_{1,2} P_{0,2} + b_{1,2} P_{1,1}, \tag{3}$$

where

$$\begin{aligned}
 \frac{M_{u0}}{(M_{d1} + \mu_1)} &= a_{1,2}, \\
 \frac{\lambda}{(M_{d1} + \mu_1)} &= b_{1,2}.
 \end{aligned} \tag{4}$$

Now we can express  $P_{1,1}$  by substituting  $P_{1,2}$ :

$$P_{1,1} = a_{1,1} P_{0,1} + b_{1,1} P_{0,2} + c_{1,1} P_{1,0}, \tag{5}$$

where

$$\begin{aligned}
 a_{1,1} &= \frac{M_{u0}}{(M_{d1} + \lambda + \mu_1 - \mu_1 b_{1,2})}, \\
 b_{1,1} &= \frac{\mu_1 a_{1,2}}{(M_{d1} + \lambda + \mu_1 - \mu_1 b_{1,2})}, \\
 c_{1,1} &= \frac{\lambda}{(M_{d1} + \lambda + \mu_1 - \mu_1 b_{1,2})}.
 \end{aligned} \tag{6}$$

We proceed using the same methodology for all the states, at each step substituting the solved state in each equation as an expression of rates. In this example, we derive a final expression for each state as a function of  $P_{0,0}$  and from there we proceed by defining the sum of all the state probabilities to be equal to 1:

$$1 = P_{0,0} + P_{0,1} + P_{0,2} + P_{1,0} + P_{1,1} + P_{1,2}. \tag{7}$$

Thus, we can solve

$$\begin{aligned}
 1 &= P_{0,0} + f_{0,1}(P_{0,0}) + f_{0,2}(P_{0,0}) + f_{1,0}(P_{0,0}) \\
 &\quad + f_{1,1}(P_{0,0}) + f_{1,2}(P_{0,0}), \\
 P_{0,0} &= \frac{1}{f_{0,1} + f_{0,2} + f_{1,0} + f_{1,1} + f_{1,2}}.
 \end{aligned} \tag{8}$$

At this point, we have defined every state probability in the model as a function of  $P_{0,0}$  and  $P_{0,0}$  itself as a function of the rates. We can now input values for the different rates that we wish to solve for and examine how the model behaves under different performance and mobility scenarios.

It should be noted that the model may be solved for any  $N \times M$ ; however, as the number of chains increases, the number of variables also increases thus making a closed-form solution very difficult to derive. It would be easier to consider equal mobility rates between chains; however, it would also provide a less accurate model. The following section presents some common mobility scenarios that may be encountered in daily usage and how the model may be used to determine the overall performance of a persistent connection over multiple networks.

#### 5. Common Scenario Results

This section includes some examples of user mobility and network coverage that may be commonly encountered in the real world. The results of these examples can assist in understanding how the model works and what insights it may offer in the context of network performance and multimedia content quality adaptation. Furthermore, these results prove analytically that the model is functioning as expected.

**5.1. Fixed-Path Mobility with Overlapping Networks.** The first scenario to look at covers a mobile node moving on a fixed



TABLE 1: Fixed-path mobility with overlapping networks.

$\lambda$	$\mu_0$	$\mu_1$	$\mu_2$	$M_{u0}$	$M_{d1}$	$M_{u1}$	$M_{d2}$
60	80	40	20	0.0088	0.0001	0.0088	0.0001

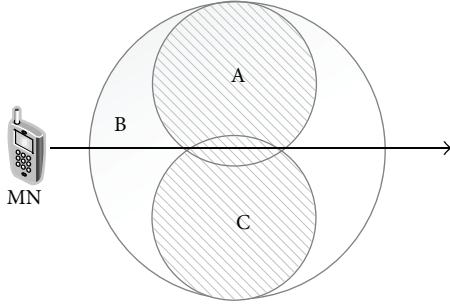


FIGURE 4: Linear mobility example with overlapping networks.

path while being connected to network B as shown in Figure 4. Along the path and within the coverage of network B, there are two smaller networks A and C. The node will enter a small area that is covered by all three networks and subsequently will exit the coverage of A and C and return to network A. This scenario expresses a case where a user may be connected to a large coverage area network such as LTE and reaches an area where smaller Wi-Fi networks are available. Assuming a constant velocity, the user will eventually leave the area of Wi-Fi coverage and fall back to LTE.

In this case, we shall consider a user moving at 5 km/h as a representation of average walking speed. The LTE radius is 500 meters and the Wi-Fi radius is 50 meters. Because in this scenario there are three networks involved, we consider a  $3 \times 3$  queueing model. The starting position of the user is in the middle chain of the model while networks A and C represent the top and bottom chains.

Table 1 shows the values considered in this scenario. The service rate of the LTE network ( $\mu_1$ ) is lower than the service rates of the other two networks. Additionally, due to the sizes of the coverage areas in the configuration presented in Figure 4, the  $M_{u0}$  and  $M_{d2}$  rates are equal since we are dealing with equally sized networks. Furthermore,  $M_{d1}$  and  $M_{u1}$  are also equal since they represent the exit rate from the LTE network towards two equally sized smaller networks. The service request rate and service rates are arbitrary but they could represent packets, frames, or any other metric significant to the application's performance.

As we see in Figure 5 most of the user's requests will be carried over the LTE network in the middle chain ( $P_{1,0}$ ,  $P_{1,1}$ , and  $P_{1,2}$ ) based on this mobility pattern and network configuration. Additionally, we see from the state probabilities in each network that all networks can provide sufficient service to support the application, while the Wi-Fi networks can provide better performance. Based on these results we can determine that the device may connect to either of the Wi-Fi networks temporarily to improve the performance of the application. If we consider a multimedia application such as video streaming, based on this model we can determine at which points we may enhance the quality

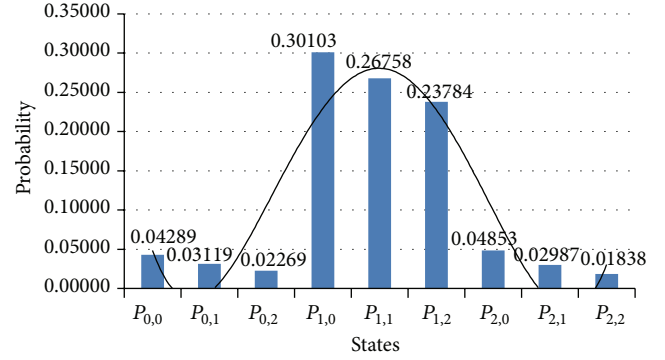


FIGURE 5: Fixed-path example with overlapping networks.

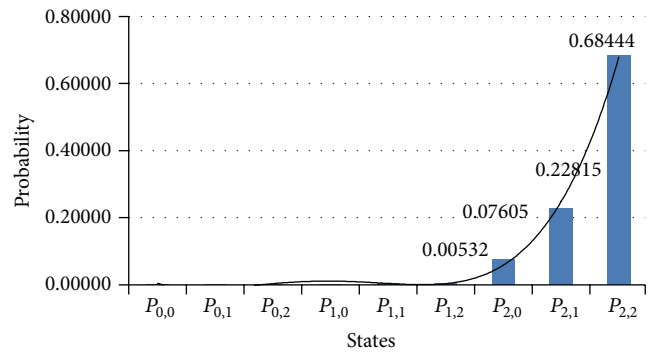


FIGURE 6: Fixed-path example with nonoverlapping networks.

of the video without having to wait for feedback from the device. This can help multimedia service providers cache the appropriate segments of the video at different bitrates or even preconfigure dynamically the sources of the video segments in a content-centric context.

**5.2. Fixed-Path Mobility with Nonoverlapping Networks.** For this scenario, we consider a fast-moving user passing through a series of networks that do not have overlapping coverage areas. Such a scenario may be envisioned by considering a car driving by an area with multiple LTE networks. To further demonstrate how this model behaves in different scenarios, we study a case where the LTE networks have different service rates. Once again, we are using a  $3 \times 3$  model for this example. Table 2 shows the values used for a user speed of 50 km/h and uniform coverage areas with a radius of 1 km.

As we see from Figure 6, below, in this scenario the user is moving very rapidly across the three networks and therefore very few of the service requests are covered by each network. Because this model only represents three networks, in this example we see that the requests converge on the third chain. Since the network represented by the third chain does not have adequate service rate for the user's application, the requests are being queued in  $P_{2,2}$  which tells us that as the user is moving quickly, the final network along their path is the one that will have to service any pending requests from other networks and hence the performance of the final network is an important factor to the overall QoS.

TABLE 2: Fixed-path mobility with nonoverlapping networks.

$\lambda$	$\mu_0$	$\mu_1$	$\mu_2$	$M_{u0}$	$M_{d1}$	$M_{u1}$	$M_{d2}$
40	55	45	65	0.015	0.0018	0.0018	0.015

TABLE 3: Random mobility example.

$\lambda$	$\mu_0$	$\mu_1$	$M_{u0}$	$M_{d1}$
60	30	60	0.0177	0.0177

**5.3. Random Mobility Examples.** The final scenario to look at represents a case of random mobility such as when a user does not exhibit any kind of linear or predictable movement. For this scenario we set a walking speed of 5 km/h and a cell radius of 50 m. We are using a  $2 \times 3$  model in this case. Table 3 summarizes the values used.

The model shows (Figure 7) that, based on the user mobility, the two networks will have equal probability of receiving service requests as the sum of probabilities for each chain is equal to 0.5. However, we see that the first chain does not have adequate service rate to support the application and hence there is a higher probability of queueing requests for that network. Based on this example it would be difficult to determine exactly the level of QoS delivered to the client at each point and therefore it would be difficult to adjust the quality of multimedia applications accordingly. Nevertheless, this result may also be used as an indication of which networks should be avoided in a particular area so that the user experience will not degrade as the user is moving.

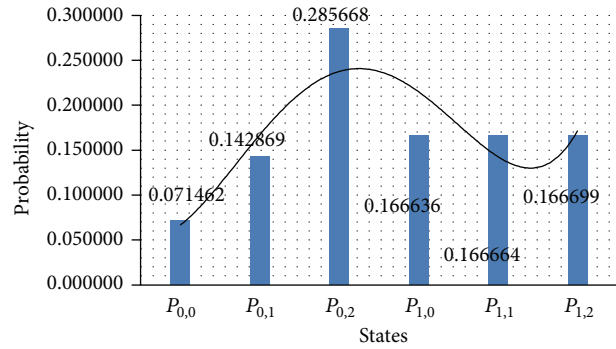
## 6. Evaluation and Conclusion

The proposed model offers a new approach at evaluating the performance of streaming applications in different mobility and network coverage scenarios; however, there are some limitations that must be highlighted in the interest of further improving the model and understanding its applications.

The mobility rate equation is not accurate for every scenario as it only represents an average approximation. For greater accuracy, we can consider the exact coordinates for the user and cell location and derive the outgoing rate based on the user's direction and speed.

The model relies on knowing in advance and with high confidence the exact route that the user will take and therefore the exact sequence of network handover that will occur. This may be impossible to achieve in real life scenarios but at a theoretical level, it can help analyze the performance of an application under certain network conditions and mobility patterns.

Furthermore, the model relies on knowing in advance and with high confidence the service rate of each network that will carry application traffic. This may be achieved with network mechanisms that report the achievable QoS between a content source and an access point but once more, in a real life environment, it would add to the infrastructure and application complexity. The model also assumes a constant service request rate by the streaming application which may also be unrealistic in real-world scenarios; however, the average

FIGURE 7: Random mobility example using a  $2 \times 3$  model.

rate may be considered for the purposes of the model in order to provide an approximation of the performance.

The current version of the model has the disadvantage of not taking into account the performance cost of a handover between heterogeneous networks. However, this disadvantage is eliminated in the context of seamless handover technologies such as proposed by Y-Comm.

Despite these restrictions, the model presents a theoretical approach at evaluating the overall performance of a connection in a mobile environment with handover rate awareness and therefore provides a fine-grained analysis of the impact of mobility and network performance on multimedia streaming applications. As this model is still at its early stages conceptually, the authors will appreciate feedback and are open to collaboration.

## Conflict of Interests

The authors declare that there is no conflict of interests regarding the publication of this paper.

## References

- [1] X. Chen and M. R. Lyu, "Message queueing analysis in wireless networks with mobile station failures and handoffs," in *Proceedings of the IEEE Aerospace Conference Proceedings*, vol. 2, pp. 1296–1303, IEEE, March 2004.
- [2] A. Rico-Páez, C. B. Rodríguez-Estrello, F. A. Cruz-Pérez, and G. Hernández-Valdez, "Queueing analysis of mobile cellular networks considering wireless channel unreliability and resource insufficiency," in *Managing Traffic Performance in Converged Networks*, pp. 938–949, Springer, Berlin, Germany, 2007.
- [3] J. Wang, Q.-A. Zeng, and D. P. Agrawal, "Performance analysis of integrated wireless mobile networks with queueing handoff scheme," in *Proceedings of the IEEE Radio and Wireless Conference (RAWCON '01)*, pp. 69–72, IEEE, Waltham, Mass, USA, August 2001.
- [4] F. Sardis, G. Mapp, J. Loo, M. Aiash, and A. Vinel, "On the investigation of cloud-based mobile media environments with service-populating and QoS-aware mechanisms," *IEEE Transactions on Multimedia*, vol. 15, no. 4, pp. 769–777, 2013.
- [5] F. Sardis, G. Mapp, J. Loo, and M. Aiash, "Dynamic edge-caching for mobile users: minimising inter-as traffic by moving cloud services and VMs," in *Proceedings of the 28th IEEE International Conference on Advanced Information Networking and Applications Workshops (WAINA '14)*, pp. 144–149, May 2014.

- [6] G. Mapp, D. N. Cottingham, F. Shaikh et al., "An architectural framework for heterogeneous networking," in *Proceedings of the International Conference on Wireless Information Networks and Systems (WINSYS '06)*, pp. 5–12, August 2006.
- [7] R. Rejaie, M. Handley, and D. Estrin, "Layered quality adaptation for Internet video streaming," *IEEE Journal on Selected Areas in Communications*, vol. 18, no. 12, pp. 2530–2543, 2000.
- [8] Computerphile, "How YouTube Works," 2014, <https://www.youtube.com/watch?v=OqQk7kLuaK4>.
- [9] A. Raghuv eer, E. Kusmierek, and D. H. C. Du, "Network-aware rate adaptation for video streaming," in *Proceedings of the IEEE International Conference on Multimedia and Expo (ICME '04)*, vol. 2, pp. 1039–1042, IEEE, June 2004.
- [10] Z. Li, X. Zhu, J. Gahm et al., "Probe and adapt: rate adaptation for HTTP video streaming at scale," *IEEE Journal on Selected Areas in Communications*, vol. 32, no. 4, pp. 719–733, 2014.



