

# Time Is Perception Is Money – Web Response Times in Mobile Networks with Application to Quality of Experience

Markus Fiedler, Patrik Arlos, Timothy Gonsalves, Anuraag Bhardwaj, Hans Nottehd

► **To cite this version:**

Markus Fiedler, Patrik Arlos, Timothy Gonsalves, Anuraag Bhardwaj, Hans Nottehd. Time Is Perception Is Money – Web Response Times in Mobile Networks with Application to Quality of Experience. Karin Anna Hummel; Helmut Hlavacs; Wilfried Gansterer. Performance Evaluation of Computer and Communication Systems (PERFORM), Oct 2010, Vienna, Austria. Springer, Lecture Notes in Computer Science, LNCS-6821, pp.179-190, 2011, Performance Evaluation of Computer and Communication Systems. Milestones and Future Challenges. <10.1007/978-3-642-25575-5\_15>. <hal-01586896>

**HAL Id: hal-01586896**

**<https://hal.inria.fr/hal-01586896>**

Submitted on 13 Sep 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



# Time is Perception is Money – Web Response Times in Mobile Networks With Application to Quality of Experience

Markus Fiedler<sup>1</sup>, Patrik Arlos<sup>1</sup>, Timothy A. Gonsalves<sup>2</sup>, Anuraag Bhardwaj<sup>1,3</sup>,  
and Hans Nottehd<sup>4</sup>

<sup>1</sup> Blekinge Institute of Technology, Karlskrona, Sweden, {mfi,pal}@bth.se

<sup>2</sup> Indian Institute of Technology Mandi, India, tag@iitmandi.ac.in

<sup>3</sup> Indian Institute of Technology Madras, Chennai, India, anuraagbady@gmail.com

<sup>4</sup> info24, Kista, Sweden, hans.nottehd@info24.se

**Abstract.** The number of mobile operators providing Internet access to end users is growing. However, irrespective of the access network, we observe a distinct sensitivity of user perception to response and download times, in particular for interactive services on the web. In order to facilitate the choice of the right network for a given task, this paper presents a systematic study of web download time and corresponding throughput as a function of the file size. Based on measurement data from three Swedish mobile operators and a particular strategy of choosing file sizes, we find surprisingly simple, yet sufficiently accurate approximations of download times. These approximations are based on simple-to-measure parameters and provide valuable quantitative insights into the acceleration of HTTP/TCP/IP-based data delivery. The paper discusses the emergence of these approximations and related errors. Furthermore, it correlates the findings with Quality of Experience, thus building bridges between performance, user perception and provisioning issues.

**Key words:** Download time; interactive service; web service; throughput; user perception; file size; measurements; Quality of Experience

## 1 Introduction

The success of interactive services depends on their responsiveness, as perceived by their users. Increasingly many such services are built from web services. These are using the HyperText Transfer Protocol (HTTP), running on top of the Transmission Control Protocol (TCP) and the Internet Protocol (IP), for exchanging data and configuration information in form of HyperText Markup Language (HTML) or eXtensible Markup Language (XML) files. This implies that the response time of a service is typically dominated by the *download time* of a corresponding HTML/XML document. The users find themselves at the end of service chains and have to wait until all the involved sub-services have executed

and delivered their results, e.g. via an InforMation eXchange (IMX) [1] that composes the information to be sent to the user.

We also observe that the number of service chains that include one or multiple mobile links is increasing. The simple exchangeability of mobile operators by replacing the SIM card in a USB dongle (or to switch between different dongles) has opened up new possibilities of choosing different operators if download times are perceived too long and/or throughput is felt to be insufficient. As waiting times increase, so does the risk that the user churns, i.e. leaves one operator for another that offers better performance, eventually entailing loss of revenue for the original provider. On the other hand, assigning more resources to a user increases cost for the provider. Thus, the latter needs to find a reasonable balance between user-perceived performance and cost.

This background motivates us to take a closer look at response times of HTTP-based downloads via different (in our case Swedish) mobile operators, and to relate them to user-perceived *Quality of Experience* (QoE) [2,3]. The typical procedure to benchmark (mobile) networks is as follows: Large files are downloaded, the corresponding *quasi-stationary throughput* is calculated as amount of downloaded data divided by the download time, and the result is reported together with other performance measures such as round-trip times or jitter values [4]. However, quasi-stationarity is preceded by a *transient phase*. TCPs startup behaviour [5,6] causes the transient throughput to approach the quasi-stationary throughput from below, which points at an acceleration of the download process. While being of minor interest for large files, the transient phase has a pronounced impact on the first few seconds of a file transfer [6], which is just in the order of magnitude of user patience [7] and thus deserves specific attention. The use of the quasi-stationary throughput to estimate download times during the transient phase will lead to over-optimistic and undesirable underestimations of the download time, and should thus be avoided.

This paper presents an original study of the download time and corresponding throughput as functions of the size of the downloaded file, as seen by the end user that does not have any insight into specific parameters and conditions of the mobile network used. The corresponding approximations are simple to parameterise from online measurements on application level and simple to deploy. They provide a useful basis for deciding which operator to choose for particular communication situations, here transfers of files upto some Megabytes). Vice versa, the formulae can be applied by operators and providers to adjust the trade-off between user perception, provisioning and cost, respectively.

The remainder of the paper is structured as follows. Section 2 discusses related work. Section 3 introduces notations and definitions to be used throughout the paper and the setup of the measurements. Section 4 presents measurement results for three different operators, followed by subsequent analysis and classification. Section 5 proposes and evaluates three- and two-parameter approximation formulae for download times as function of the file sizes. Section 6 presents and

discusses closed formulae relating QoE with file sizes through the presented approximation models. Section 7 concludes the paper and takes an outlook on future work.

## 2 Related Work

It has been recognised that users get the more distracted, the longer they have to wait [7,8]. Obviously, this implies a “negative impact on attitudes toward delay” [9]. In particular, it has been recognised that user perception and rating quickly decrease during the first seconds of waiting [3,10]. This effect can amongst others be modelled by differential equations [11,12], where the derivative of QoE with respect to the delay is negative. The resulting QoE–QoS relationships capture the effect of delays on QoE, typically described by negative-exponential [11–13] or logarithmic functions [10,12].

Much work has been done on the performance of TCP, the underlying protocol for web transfers. However, the steady-state loss and delay performance and throughput are targeted in the majority of contributions, see e.g. [3,14–17]. For example, [16] studies the impact of variable rate and variable delay on long-lived TCP performance in a Third Generation (3G) mobile network. The study was based on traces from a 3G 1 X network taken around 2002, and simulation results based on the network simulator ns-2. Also, [17] finds that the download time is proportional to the file size, which implies a constant throughput. The study employs a variable file size from some ten KB<sup>1</sup> on. Reference [18] defines web objects as between 1 and 100 KB, and [5] states that average and median of the flow sizes are smaller than 10 KB.

Motivated by the latter observation, [5] addresses explicitly the performance problems faced by short flows due to the start-up behaviour of TCP, and proposes an analytic model that explicitly covers such short flows. The model is validated by simulations. No measurements in real environments are taken into account. The latter provide the ground for reference [6], according to our knowledge the only reference to study measurements of the ramping-up of TCP in a real mobile network. For GPRS, the duration of the slow start phase was found to last at least six seconds before the congestion window is expanded sufficiently in order to allow for approaching the quasi-stationary throughput. No further mobile technologies (3G *etc.*) were analysed. We conclude that a study of the impact of the file size on the download time and thus on QoE, based on measurements on contemporary mobile networks and compiled into easy-to-deploy formulae, is missing. This paper closes the gap.

---

<sup>1</sup> 1 KB = 1000 B, 1 MB = 1000 KB; 1 KiB = 1024 B, 1 MiB = 1024 KiB; we will use the appropriate factor depending on the context

### 3 Notations, Definitions and Setup

We will now introduce the variables to be used throughout the paper; the *size of the file to be downloaded*  $X$ ; the *download time*  $T$ ; and the *perceived throughput on application level*  $R$ , with the corresponding binary logarithms interrelated through:

$$\text{lb}(R/\text{bps}) = \text{lb}(8 X/B) - \text{lb}(T/\text{s}) = 3 + \text{lb}(X/B) - \text{lb}(T/\text{s}). \quad (1)$$

In order to cover the whole range from small to large files, we increase the file sizes by factor two, starting from the smallest possible file size of one byte (just the end marker) in order to be able to see any particular behaviour for one-packet transmissions. We have chosen a maximal file size of 4 MiB, which comes closest to the file size of 3.7 MB used by the tool described in [4]. Table 1 contains the classification of file sizes that will be used throughout the paper. Furthermore, we denote the *quasi-stationary throughput* that can be observed for very large files as  $R_\infty$ , amongst others shown by the tool [4].

Label	Size $X(i) = 2^i$ B	Range of $i$
S	1 B ... 1 KiB	0 ... 10
M	2 KiB ... 256 KiB	11 ... 18
L	512 KiB ... 4 MiB	19 ... 22

**Table 1.** Classification of file sizes.

The setup used for the measurements is as follows: Host P, a PC with Windows XP connected to the mobile operators using a Huawei E220 USB modem, requests a file via HTTP (not FTP) from the Linux-based host S that is running an Apache web server. Host S is not publicly known, so all HTTP traffic destined to S originates from P, and disturbances caused by competing traffic are avoided. For each file size as specified in Table 1, forty subsequent downloads are performed ( $j = 1 \dots 40$ ) without caching the file, before moving to the next file size ( $i = i + 1$ ). The download time for file size  $i$  and replication  $j$ , defined as  $T_j(i)$ , is obtained through time stamp  $T_j^s(i)$  just before retrieving the file using the `get` function found in the Perl module `LWP::Simple`. Completion of the transfer triggered a new time stamp  $T_j^e(i)$ . From this, we calculate  $T_j(i) = T_j^e(i) - T_j^s(i)$  and the throughput  $R_j(i) = X(i)/T_j(i)$  according to (1).

Three Swedish mobile operators are considered, denoted by A, B and C, each of them offering HSDPA (High-Speed Downlink Packet Access). Taking the role of an end user, no specific insights into configurations and load conditions of the mobile networks were available. The experiments were conducted during December 2009 during business days. All units were stationary during the tests, and no obvious network congestion was observed during the series of tests.

## 4 Results and Analysis

### 4.1 Median values of download time and throughput

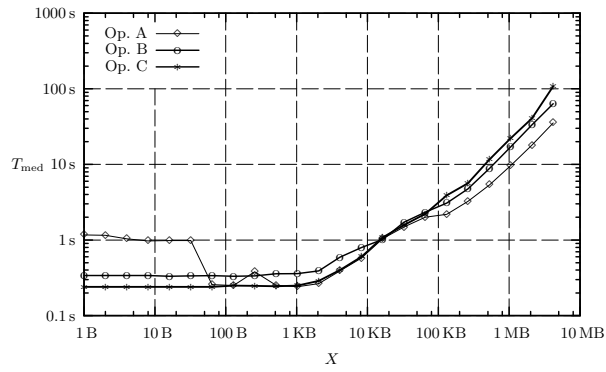
Typically, the spread of the distribution of  $R_j(i)$  over  $j$  is quite small. In 20 out of a total of 2640 measurements, we observed download times exceeding their neighboring values by up to ten seconds. This effect is rather common and discussed a.o. in [19]. As these extraordinary long download times have a significant impact on the average, we use the median (and skip that notion frequently for sake of brevity).

The quasi-stationary throughput per operator was approximated by the median of the throughput obtained for the largest investigated file size of 4 MiB, i.e.  $R_\infty = R_{\text{med}}(22)$ . The corresponding values are given in Table 2, accompanied by average  $m_R$ , standard deviation  $s_R$  and coefficient of variation  $c_V$ . Furthermore, the median of a set of round-trip time measurements  $\text{RTT}_{\text{med}}$  performed with the tool [4] is shown. Obviously, operator A provides the highest quasi-stationary throughput in combination with the smallest relative variation, while the reverse holds for operator C. Operator B, which is actually sharing the 3G network infrastructure with operator A, yields a considerably lower throughput than operator A with comparable relative variation, despite of quite similar RTT values.

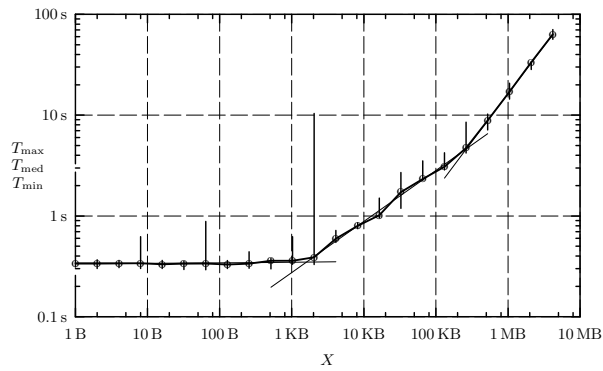
Op.	$R_\infty/\text{bps}$	$\text{lb}(R_\infty/\text{bps})$	$m_R(22)/\text{bps}$	$s_R(22)/\text{bps}$	$c_V(22)$	$\text{RTT}_{\text{med}}/\text{ms}$
A	949584	19.9	935374	38594	0.04	125
B	530114	19.0	530987	28409	0.05	130
C	311365	18.3	323128	43934	0.13	336

**Table 2.** Quasi-stationary throughput estimations and round-trip measurements for different operators.

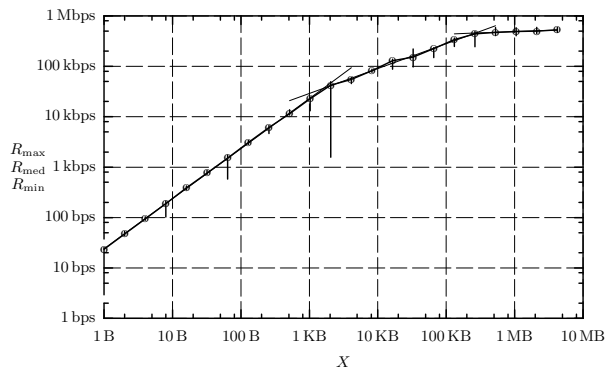
Figure 1 provides a first impression on the median of the download time obtained from the three different operators. For small files, the download time is quite independent of the file size, with fastest delivery provided by operator C, followed by B and A. A discontinuity is observed for operator A: For very small sizes ( $X \leq 32 \text{ B}$ ), the download time is found between 1 s and 1.2 s, while it then drops to values between 0.25 s and 0.38 s for  $64 \text{ B} \leq X \leq 1 \text{ KB}$ , which means that small packets may take more time than large packets to get delivered, an observation also reported and discussed in [19]. For medium-sized files, the download time is growing with the file size for all operators, however with different gradients. Operator B initially displays the longest download times. From around 20 KB on, operator C takes most time to deliver the file. Operator A provides fastest delivery for files larger than 30 KB. For large files, the order of delivery is operator A, followed by operators B and C.



**Fig. 1.** Median of the download time  $T_{\text{med}}$  via different operators versus file size  $X$ .



**Fig. 2.** Median of the download time  $T_{\text{med}}$  via operator B versus file size  $X$ , with error bars and potential regressions.



**Fig. 3.** Median of the perceived throughput  $R_{\text{med}}$  via operator B versus file size  $X$ , with error bars and potential regressions.

As particular examples, we are taking a closer look at the medians of download times  $T_{\text{med}}$  and throughputs  $R_{\text{med}}$  provided by operator B, accompanied by error bars indicating the asymmetrical span between maximum and minimum of the 40 contributing values. From figures 2 and 3, we can see three areas in which  $T_{\text{med}}$  and  $R_{\text{med}}$  are almost linear in log-log representation. This is illustrated piecewise by linear regressions and points at *power-type relationships* between download time, throughput and file size, respectively. Quite similar behaviours are obtained for operator A and C, cf. Figure 1. For small files, cf. Table 1, the download time is practically constant, and consequently, the throughput grows proportionally with the file size. For large files, we observe saturation of  $R_{\text{med}}$ . Its logarithm shows a quite weak, yet linear dependence on  $\text{lb}(X/B)$ , while  $\text{lb}(T_{\text{med}})$  grows in almost the same manner as  $\text{lb}(X/B)$ , with a gradient around 1. For medium-sized files, both  $\text{lb}(T_{\text{med}})$  and  $\text{lb}(R_{\text{med}})$  raise in a similar fashion, with a gradient around 1/2. Obviously, we are facing a transition area in which both download time and throughput grow as the file size grows, which indicates acceleration while downloading due to TCP start-up. The throughput has not reached its quasi-stationary value yet, which implies that the real download time is higher than any estimation that builds upon the quasi-stationary value  $R_{\infty}$ .

## 4.2 Regressions

We now construct linear least-square regressions on the logarithms of  $X$ ,  $T$  and  $R$  for the areas S, M and L for all operators of the types

$$\text{lb}(\hat{T}/s) = a_T \text{lb}(X/B) + b_T; \quad (2)$$

$$\text{lb}(\hat{R}/s) = a_R \text{lb}(X/B) + b_R; \quad (3)$$

and evaluate the coefficient of determination  $\mathcal{R}^2$  [3], which mostly signals good matches ( $\mathcal{R}^2 \rightarrow 1$ ) with exception of some weak trends ( $a \rightarrow 0$ ).

Size	Op.	$a_T$	$b_T$	$\mathcal{R}^2$	$a_R$	$b_R$	$\mathcal{R}^2$
S	A	-0.276	0.566	0.785	1.276	2.434	0.987
	B	0.006	-1.582	0.211	0.994	4.582	0.999
	C	0.006	-2.072	0.602	0.994	5.072	1.000
M	A	0.520	-7.486	0.976	0.480	10.486	0.971
	B	0.506	-6.907	0.995	0.494	9.907	0.953
	C	0.627	-8.769	0.997	0.373	11.768	0.991
L	A	0.906	-14.83	0.999	0.094	17.83	0.933
	B	0.949	-14.89	0.999	0.051	17.89	0.953
	C	1.054	-16.57	0.987	-0.054	19.57	0.164

**Table 3.** Regressions for  $\hat{T}$  and  $\hat{R}$ .

Table 3 shows the obtained regressions. As expected from (1), we observe for each operator and file size class  $a_T + a_R = 1$  and  $b_T + b_R \simeq 3$ . A closer look



at the results provides the following insights. For small files, the download time regression for operator A has a coefficient  $a_T < 0$ , which is due to the above-described discontinuity. For the other operators B and C, there is hardly any dependence of the download time on the size for small files ( $a_T \simeq 0$ ). For medium-sized files, the coefficients  $a_T$  and  $a_R$  are found close to 0.5, which has shown to be the key for the approximation presented in the next section. For large files, operator C shows an unexpected negative trend in the throughput, seen from  $a_R < 0$ . The other two operators show slightly increasing trends ( $0 < a_R < 0.1$ ), i.e. the throughput still rises slightly as the file size increases.

## 5 Download Time Approximations

The regressions shown in Figure 2 suggest the use of maximum of the three approximations

$$\hat{T} = \max\{\hat{T}^S, \hat{T}^M, \hat{T}^L\} \quad (4)$$

to estimate the download time. The components of (4) belong to the different regimes of file sizes as defined in Table 1 and are given as follows:

- For small files

$$\hat{T}^S = \text{const.} \quad (5)$$

The time  $\hat{T}^S$  needed to send a one-packet file depends on the operator.

- For medium-sized files

$$\hat{T}^M/s = 8 \left( \frac{X/B}{R_\infty/\text{bps}} \right)^{a_T}. \quad (6)$$

This part builds upon the observation that  $b_R \simeq a_T \text{lb}(R_\infty/\text{bps})$ , cf. tables 2 and 3. For  $a_T = 0.5$ , the general power-type relationship (6) reduces to a simple square-root formula, which reminds of an accelerated movement with constant acceleration  $a$  in which the time to reach a distance  $s$  is given by  $t = \sqrt{2s/a}$ . Translated to our case ( $s = X$ ), the acceleration of the file transfer would amount to  $a = R_\infty/32$  bps.

- For large files

$$\hat{T}^L/s = 8 \frac{X/B}{R_\infty/\text{bps}}. \quad (7)$$

The motivation for this part of the approximation is found in the observation  $a_T \rightarrow 1$  and  $a_R \rightarrow 0$  in Table 3, i.e. the throughput is almost constant and can thus be approximated by  $R_\infty$ .

Obviously, the approximation depends on *three parameters*:

1. the *one-packet download time*  $\hat{T}^S$  to be measured from downloading files of a typical size (e.g. 1 KB);

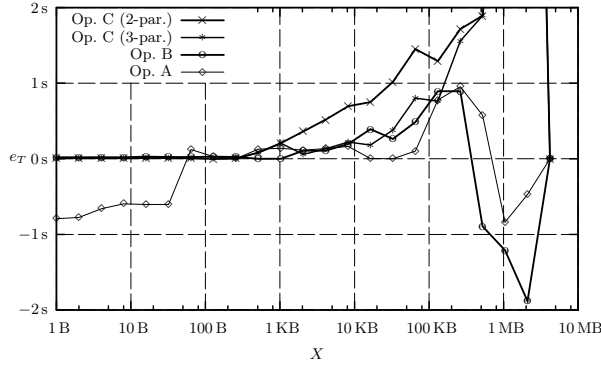
2. the quasi-stationary throughput  $R_\infty$  obtained from downloading large files of several MB and taking the median of the measured values; and
3. the acceleration factor  $a_T$  obtained from downloading a sequence of medium-sized files, followed by the regression (2).

The three-parameter approximation can be simplified by letting  $a_T = 0.5$ , which fits operators A and B almost exactly and operator C approximately ( $a_T \simeq 0.625$ ). This leaves us with *two parameters*  $\hat{T}^S$  and  $R_\infty$ .

We will now investigate the approximation errors, given as the difference of approximated and measured median of the download time

$$e_T(X(i)) = \hat{T}(X(i)) - T_{\text{med}}(X(i)) \quad (8)$$

and seen from Figure 4. Positive values of  $e_T(X(i))$  mean that the approximation overestimates the download time for the given operator and file size  $X(i)$ .



**Fig. 4.** Estimation errors for different operators.

For small files, the only significant deviation is seen for operator A due to the discontinuity mentioned above. For medium-sized files, the approximation estimates on the safe side. The deviations are less than one second for operators A and B, while the two-parameter approximation for operator C displays a steady drift towards two seconds. The performance of the three-parameter approximation for operator C is comparable to that of A and B until a file size of 128 KiB thanks to the better capture of the gradient in the log scale by the third parameter  $a_T = 0.625$ . Upon reaching the area of large file sizes, a trend towards underestimations of no more than two seconds for operator A and one second for operator B is observed. The approximation for operator C is very conservative, overestimating the median by up to 13s for a file size of 2 MiB, which indicates the need to use a refined value of  $R_\infty$  to better account for this domain. All errors vanish for a file size of 4 MiB, as this was the anchor point for  $R_\infty$ .

## 6 Links to Quality of Experience

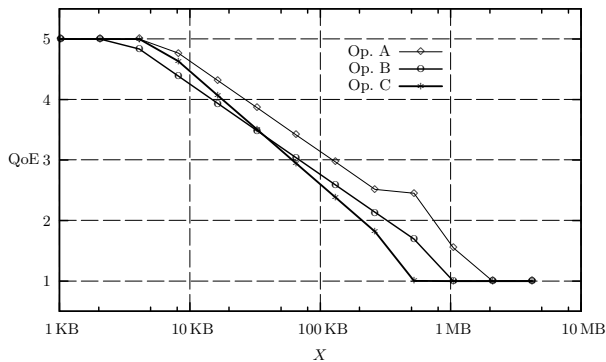
In the following, we assume that we can quantify the QoE by employing a numerical linear Mean Opinion Score (MOS)-type scale from 5 (= excellent) to 1 (= unacceptable). We employ a logarithmic relationship between QoE (user rating) and QoS (response time) found in standard [10] and recently re-confirmed *a.o.* by [3,20]. In our example, the numerical values stem from the case of a time scale of reference of 6 s and a mix of two user groups. They can (be) change(d) according to context and user expectations, as illustrated in [10]. Starting from

$$\text{QoE} = \max\{\min\{4.38 - 0.9\ln(T/s), 5\}, 1\} \quad (9)$$

and inserting the above results (5)–(7), we arrive at

$$\text{QoE} = \begin{cases} \min\{4.38 - 0.9\ln(\hat{T}^S/s), 5\} & \text{for S} \\ \max\{\min\{1.68 + 0.9a_T(\ln(R_\infty/\text{bps}) - \ln(X/B)), 5\}, 1\} & \text{for M} \\ \max\{\min\{1.68 + 0.9(\ln(R_\infty/\text{bps}) - \ln(X/B)), 5\}, 1\} & \text{for L} \end{cases} \quad (10)$$

As expected, the QoE decreases as file size and download time grow, and increases as throughput increases. However, the sensitivity of QoE to the parameters depend on the file size itself: Due to the acceleration behaviour, the sensitivity is smaller for medium-sized files than for large files. Figure 5 illustrates the QoE estimations from (10), based on the three-parameter approximation and Table 2. The discontinuity for operator A is due to the change of domain between medium-sized and large files. Best QoE is reached for operator A, followed by operator C for files up to 30 KB. From then on, operator B provides the second-best perception.



**Fig. 5.** QoE estimations versus file size  $X$  for different operators.

Figure 5 also shows that, in order not to challenge user patience beyond feasibility (i.e. keep the QoE above 3), downloaded files should not be larger than approximately 120 KB for operator A, 70 KB for operator B and 60 KB for operator C, respectively.

## 7 Conclusions and Outlook

This paper presented a measurement-based study of download performance and approximations of download times and corresponding throughputs as functions of the size of the downloaded file. Through a specific strategy of choosing file sizes, employing a factor of two, we were able to see linear relationships between the logarithms of download time, throughput and file size for small, medium-sized and large files. From a subsequent analysis of the parameters of the approximations, we found surprisingly simple approximation formulae with reasonably tight error bounds for the download times. For the small- and medium-sized files that represent the majority of the web downloads and that yield download times in the order of user patience, we obtain an error bound of roughly one second in most cases. The formulae will help users and application providers to “choose the right network for the right task”, alternatively to limit data sizes in order to keep download times of an acceptable level. Furthermore, it will help providers to choose the level of provisioning they would like to offer in order to yield a reasonable trade-off between quality and cost.

Besides of refinements of the proposed approaches and formulae, future work will include validation and parameterisation of the download time and throughput approximations for other access network technologies such as WLAN, ADSL and Ethernet. At this point, a deeper analysis of TCP traces and a quantitative comparison to the results obtained from the model proposed in [5] would be of interest. We will also study the impact of data loss and delay variation on the proposed approximations, as they will affect the quasi-stationary throughput. Finally, the practical applicability of the formulae to assess and implement network selection policies for seamless communications will be studied.

## Acknowledgements

The authors would like to thank the Swedish Agency for Innovation Systems, VINNOVA, for sponsoring this work through the EViMonA project (d-nr 2007/02505). They are also grateful to the Indian-Swedish project that supported the collaboration, and in particular Prof. Ashok Jhunjhunvala and Prof. Sara Eriksén for their great effort and commitment.

## References

1. Info24, “Homepage,” <http://www.info24.se/>, (last seen 2010-07-31).
2. ITU-T Recommendation P.10/G.100 (incl. Amendment 2), “Vocabulary for performance and quality of service,” July 2006 (2008).
3. J. Shaikh, M. Fiedler, and D. Collange, “Quality of Experience from user and network perspectives,” *Annals of Telecommunications, Special Issue on Quality of Experience : 1/Metrics and performance evaluation*, vol. 65, no. 1-2, Jan.-Feb. 2010,

- electronically available at <http://www.springerlink.com>, DOI 10.1007/s12243-009-0142-x.
4. T. Gonsalves and A. Bhardwaj, "Comparison of AT-Tester with other popular testers for Quality of Service Experience (QoSE) of an internet connection," <http://www.broadbandasia.info/>, Aug. 2009.
  5. M. Mellia, I. Stoica, and H. Zhang, "TCP model for short lived flows," *IEEE Comm. Letters*, vol. 6, no. 2, pp. 85–87, Feb. 2002.
  6. R. Chakravorty and I. Pratt, "WWW performance over GPRS," in *Proc. 4th Int. Workshop on Mobile and Wireless Communication Networks (MWCN 2002)*, Sept. 2002, pp. 527–531.
  7. Zona Research Inc., "The economic impacts of unacceptable web-site download speeds," Report, 1999.
  8. J. Nielsen, *Usability Engineering*. Morgan Kaufman, 1994.
  9. G. M. Rose, R. Evaristo, and D. Straub, "Culture and consumer responses to web download time: a four-continent study of mono and polochronism," *IEEE Trans. on Engineering Management*, vol. 50, no. 1, pp. 31–44, Feb. 2003.
  10. ITU-T Recommendation G.1030, "Estimating end-to-end performance in IP networks for data applications," November 2005.
  11. M. Fiedler, T. Hofffeld, and P. Tran-Gia, "A generic quantitative relationship between Quality of Experience and Quality of Service," *IEEE Network, Special Issue on Improving Quality of Experience for Network Services*, no. 2, pp. 36–41, March/April 2010.
  12. M. Fiedler and T. Hofffeld, "Quality of Experience-related differential equations and provisioning-delivery hysteresis," in *Proc. 21st ITC Specialist Seminar on Multimedia Applications*, Miyazaki, Japan, March 2010. [Online]. Available: <http://www.ieice.org/proceedings/ITC-SS21/ITC-SS21-Proceedings.pdf>
  13. T. Hofffeld, P. Tran-Gia, and M. Fiedler, "Quantification of Quality of Experience for edge-based applications," in *20th International Teletraffic Congress (ITC20)*, Ottawa, Canada, June 2007.
  14. X. Chen, W. Wang, and J. Nie, "Analysis of web response time in asymmetrical wireless network," in *Proc. 11th IEEE Singapore Int. Conf. on Communication Systems (ICCS 2008)*, Nov. 2008, pp. 1427–1430.
  15. M. Miyagi, K. Ohkubo, M. Kataoka, and S. Yoshizawa, "Performance prediction method for web-access response time distribution using formula," in *Proc. IEEE/IFIP Network Operations and Management Symposium (NOMS 2004)*, vol. 1, Apr. 2004, pp. 905–906.
  16. M. C. Chan and R. Ramjee, "TCP/IP performance over 3G wireless links with rate and delay variation," in *MobiCom '02: Proc. of the 8th annual international conference on mobile computing and networking*. New York, NY, USA: ACM, 2002, pp. 71–82.
  17. S. Voskarides and et al., "Practical evaluation of GPRS use in telemedicine system in Cyprus," in *Proc. 4th Int. IEEE EMBS Special Topic Conference on Information Technology Applications in Biomedicine*, April 2003, pp. 39–42.
  18. E. Baccarelli, M. Biagi, N. Cordeschi, and C. Pelizzoni, "Minimization of download times of large files over wireless channels," *IEEE Trans. on Mobile Computing*, vol. 6, no. 10, pp. 1105–1115, Oct. 2007.
  19. P. Arlos and M. Fiedler, "Influence of the packet size on the one-way delay on the down-link in 3G networks," in *Proc. ISWPC 2010*, Modena, Italy, May 2010.
  20. P. Reichl, S. Egger, R. Schatz, and A. dAlconzo, "The logarithmic nature of QoE and the role of the Weber-Fechner Law in QoE assessment," in *Proc. IEEE ICC'10*, Cape Town, South Africa, May 2010.