# ITU-T

TELECOMMUNICATION
STANDARDIZATION SECTOR
OF ITU

# G.1030

(02/2014)

SERIES G: TRANSMISSION SYSTEMS AND MEDIA, DIGITAL SYSTEMS AND NETWORKS

Multimedia Quality of Service and performance – Generic and user-related aspects

## Estimating end-to-end performance in IP networks for data applications

Recommendation ITU-T G.1030

# Recommendation ITU-T G.1030

## Estimating end-to-end performance in IP networks for data applications

**Summary**

Recommendation ITU-T G.1030 provides a framework of tools to obtain IP network performance, estimate the performance of user applications, and apply perceptual models to gauge user satisfaction with the end-to-end performance.

The user-perceived performance of data applications on packet networks is dependent on many factors, including the end-to-end performance of the packet network, the application's dependency on the communications network, the performance of the terminals and other devices beyond the purview of the network operator(s), and the user's task and the extent of user interaction with the application. Network designers take these factors into account to assure user satisfaction. Once the application performance has been estimated, perceptual models can be applied to interpret the level of end-to-end performance attained.

This Recommendation assumes that the reader will be able to provide at least some level of detail about each of the key factors above, and will then use the framework of tools to estimate end-to-end performance.

**History**

| Edition | Recommendation | Approval | Study Group | Unique ID* |
|---------|----------------|----------|-------------|------------|
| 1.0 | ITU-T G.1030 | 2005-11-29 | 12 | 11.1002/1000/8672 |
| 2.0 | ITU-T G.1030 | 2014-02-13 | 12 | 11.1002/1000/12122 |

_____

\* To access the Recommendation, type the URL http://handle.itu.int/ in the address field of your web browser, followed by the Recommendation's unique ID. For example, http://handle.itu.int/11.1002/1000/11830-en.

FOREWORD

The International Telecommunication Union (ITU) is the United Nations specialized agency in the field of telecommunications, information and communication technologies (ICTs). The ITU Telecommunication Standardization Sector (ITU-T) is a permanent organ of ITU. ITU-T is responsible for studying technical, operating and tariff questions and issuing Recommendations on them with a view to standardizing telecommunications on a worldwide basis.

The World Telecommunication Standardization Assembly (WTSA), which meets every four years, establishes the topics for study by the ITU-T study groups which, in turn, produce Recommendations on these topics.

The approval of ITU-T Recommendations is covered by the procedure laid down in WTSA Resolution 1.

In some areas of information technology which fall within ITU-T's purview, the necessary standards are prepared on a collaborative basis with ISO and IEC.

NOTE

In this Recommendation, the expression "Administration" is used for conciseness to indicate both a telecommunication administration and a recognized operating agency.

Compliance with this Recommendation is voluntary. However, the Recommendation may contain certain mandatory provisions (to ensure, e.g., interoperability or applicability) and compliance with the Recommendation is achieved when all of these mandatory provisions are met. The words "shall" or some other obligatory language such as "must" and the negative equivalents are used to express requirements. The use of such words does not suggest that compliance with the Recommendation is required of any party.

INTELLECTUAL PROPERTY RIGHTS

ITU draws attention to the possibility that the practice or implementation of this Recommendation may involve the use of a claimed Intellectual Property Right. ITU takes no position concerning the evidence, validity or applicability of claimed Intellectual Property Rights, whether asserted by ITU members or others outside of the Recommendation development process.

As of the date of approval of this Recommendation, ITU had not received notice of intellectual property, protected by patents, which may be required to implement this Recommendation. However, implementers are cautioned that this may not represent the latest information and are therefore strongly urged to consult the TSB patent database at http://www.itu.int/ITU-T/ipr/.

## Table of Contents

# Recommendation ITU-T G.1030

## Estimating end-to-end performance in IP networks for data applications

## 1        Scope

This Recommendation provides a framework of tools and the process for estimating end-to-end performance of applications operating on IP networks.

## 2        References

The following ITU-T Recommendations and other references contain provisions which, through reference in this text, constitute provisions of this Recommendation. At the time of publication, the editions indicated were valid. All Recommendations and other references are subject to revision; users of this Recommendation are therefore encouraged to investigate the possibility of applying the most recent edition of the Recommendations and other references listed below. A list of the currently valid ITU-T Recommendations is regularly published. The reference to a document within this Recommendation does not give it, as a stand-alone document, the status of a Recommendation.

[ITU-T G.1010]    Recommendation ITU-T G.1010 (2001), *End-user multimedia QoS categories*.

[ITU-T G.1031]    Recommendation ITU-T G.1031 (2014), *QoE factors in web-browsing*.

[ITU-T G.1040]    Recommendation ITU-T G.1040 (2006), *Network contribution to transaction time*.

[ITU-T Y.1540]    Recommendation ITU-T Y.1540 (2011), *Internet protocol data communication service – IP packet transfer and availability performance parameters*.

## 3        Definitions

### 3.1      Terms defined elsewhere

None.

### 3.2      Terms defined in this Recommendation

None.

## 4        Abbreviations and acronyms

This Recommendation uses the following abbreviations and acronyms:

HPT          Host Processing Time

IPLR         IP Packet Loss Ratio

MOS          Mean Opinion Score

MSS          Maximum Segment Size

NAT          Network Address Translation

QoS          Quality of Service

RTT          Round-Trip Time

TCP          Transmission Control Protocol

UNI          User Network Interface

# 5 Conventions

None.

# 6 End-to-end performance estimation

This Recommendation covers the process of estimating end-to-end performance of applications operating on IP networks, using:

• performance of the IP network of interest, based on relevant measurements or network modelling results;

• specifications of the application of interest, in terms of its governing protocol(s) with specified options, or a model of the application using network performance and customer appliance performance as input and producing a key metric of application performance (e.g., file download time) as a result;

• a perceptual model intended for the applications of interest, to interpret the application performance as an estimate of the quality experienced by a typical population of users.

Figure 6-1 illustrates the general process used to develop an end-to-end performance estimate



**Figure 6-1 Process to obtain end-to-end performance estimate**

The steps to obtain network and application performance may be combined in some cases, such as when a simulation provides a means to measure the performance of a particular session, or set of sessions.

## 6.1 Network performance assessment

Network performance may be assessed in terms of the packet transfer performance parameters defined in [ITU-T Y.1540] and other relevant standards (e.g., RFCs developed by the IETF IP Performance Metrics Working Group). Besides the usual one-way transfer performance metrics, the most straightforward assessment of application performance will sometimes be possible with round-trip metrics instead.

There are two principal sources of network performance information: measurement and modelling.

Network measurements allow the assessor to treat the network as a black-box, and produce information that may be useful in the remaining steps of the modelling chain. However, there are several important considerations for the measurement design, including:

1) The sending discipline needs to match the application of interest in some cases. For example, transmission control protocol (TCP) flow control responds to network conditions and tends to fill the queue at the bottleneck, increasing delay beyond that which would otherwise be measured.

2) The non-measurement load must be similar in size and character to the conditions where the performance estimates are to be applicable. Measurements on an unloaded network are not particularly useful.

Network modelling can provide the needed performance characterization when the network is not yet fully constructed (e.g., nodes and links are in place, but not running a critical protocol), or when the key considerations for measurement are not attainable. There are many choices of modelling tools, including commercial products and public domain research tools. Modelling tools require a substantial degree of expertise and information about the network of interest to be used effectively. As an alternative to this fairly precise network modelling, Appendix I provides simplified methods which may be used for estimating end-to-end performance in an IP network. However, the accuracy of this method will depend greatly on the accuracy of the information provided.

## 6.2 Application performance assessment

Application models take the estimates of network performance and information describing application device performance as inputs, and produce one or more key metrics of application performance as outputs.

One such application performance model has been specified in [ITU-T G.1040] for conversational packet exchanges typical of credit card and other point-of-sale transactions.

For long-lived file transfers using TCP's reliable byte transfer service, the models described in Appendix I provide accurate results, assuming the inputs are correct.

## 6.3 Perceptual models

This Recommendation includes a limited perceptual model for a simple web search task in Appendix II.

It should be noted that this model cannot be applied to general web browsing as experienced by the end-user as it does not consider the dynamics of real web browsing. For the more complex case of real web browsing, [ITU-T G.1031] defines perceptual events that have to be considered for creating a perceptual model and goes beyond the events considered in the model described in Appendix II.

## 6.4 Framework for models in the end-to-end performance assessment process

Figure 6-2 illustrates the various alternatives in the process to estimate the end-to-end performance of applications on IP networks. As indicated in this figure, there are many options available to complete the process, although in practice the assessor must combine options that are consistent with the goal of an end-to-end estimate (and with one another).

| Inputs (beyond application of interest) | | Network performance assessment | | Application performance assessment | | Perceptual performance assessment | | End-to-end performance estimate |
|---|---|---|---|---|---|---|---|---|
| Access to live network for testing | ▶ | Measurements<br>• Ping<br>• Unbiased sample<br>• Application stream<br>• Bulk transfer tool | ▶ | Model for long-lived TCP flows (see Appendix I)<br><br>Direct assessment with application tool | ▶ | Model for web browsing (see Appendix II) | ▶ | End-to-end performance estimate |
| Network design (detailed) | ▶ | Modelling tools<br>• Commercial products<br>• Research tools | | Embedded in network modelling tool | | Future model for video teleconferencing | | |
| Network design (limited info) | ▶ | Simplified estimation (see Appendix I) | | ITU-T G.1040 (Note) | | Other models | | |

NOTE – [ITU-T G.1040] provides the network contribution to transaction application performance, and does not supply the full application performance estimation, yet it gives an application-oriented view of network performance.

**Figure 6.2 – Framework for developing end-to-end IP performance estimate**

## 7 Evaluation of end-user QoS

After evaluating end-to-end physical performance, it is essential for network planners to evaluate the end-user quality of service (QoS) taking into account perceptual aspects.

[ITU-T G.1010] provides QoS categories for data applications as well as voice and video applications. In addition, Appendix II provides an opinion model for estimating users' perceived quality of a simple web search task. This enables diagnostic evaluation based on measured/planned end-to-end delay values in terms of customers' opinion.

# Appendix I

# Estimates of network performance with limited information

(This appendix does not form an integral part of this Recommendation.)

## I.1 Introduction

This appendix provides information on simplified methods that may be used for estimating end-to-end performance in an IP network when limited information is available. The accuracy of estimates derived using these methods are highly dependent on the quality of the input information. The material in this appendix is subject to change following further study and evaluation.

## I.2 Reference connection

It is necessary to account for firewalls, proxies (for web and/or performance enhancement), network address translation (NAT) devices, and load balancing switches in the reference connection, as these devices are present in many homes, enterprises, and managed IP networks. These "middleboxes" provide various functions, but remove the end-to-end transparency that is a desirable aspect of an IP network architecture. Figure I.1 shows an example reference connection.
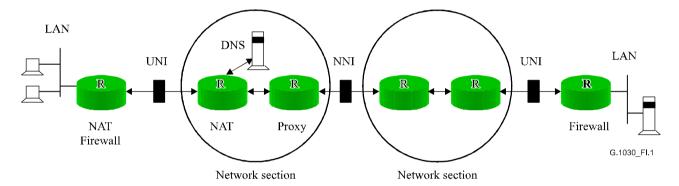


**Figure I.1 – Reference connection with example middleboxes**

This reference connection allows for evaluation of the performance of signalling protocols such as session initiation protocol (SIP) using the same general framework as other IP applications.

## I.3 Concatenation of packet transfer performance values

The introduction of middleboxes means that there will be more individual sections in the end-to-end path, creating a need for accurate accumulation formulae.

### I.3.1 Delay

The mean delay for individual sections are additive. However, it may be noted that the mean delay represents the centre of gravity of the delay variation distribution often observed with packet transport. When packet transfer delay across a network section (or processing time in a host) is represented by the mean, a single sample from the delay distribution is being replaced with its expected value. The time averages of delay are appropriate here because each transaction will sample the underlying delay distribution at many different times over the transaction interval.

### I.3.2 Loss

The Internet packet loss ratio (IPLR) performance is essentially the user network interface (UNI)-UNI performance. It can be assumed that the loss probabilities of end terminals and hosts are negligible.

The method proposed for IPLR concatenation is to invert the probability of successful packet transfer across n network sections, as follows:

$$IPLR_{UNI-UNI} = 1 - \{(1 - IPLR_{NS1}) \times (1 - IPLR_{NS2}) \times (1 - IPLR_{NS3}) \times \cdots \times (1 - IPLR_{NSn})\}$$

This equation relies on the theory of conditional probabilities and assumes that the loss probabilities in each network section are independent. For a UNI-UNI path with two networks, A and B, with loss probabilities pA and pB:

$$\text{Prob}\{\text{success on both networks}\} \quad = \text{Prob}\{\text{success on B}|\text{success on A}\} \times \text{Prob}\{\text{success on A}\}$$

$$= (1 - p_B) \times (1 - p_A)$$

Empirically-derived loss probabilities for individual networks (pA and pB) are equivalent to the conditional probabilities for success on any preceding networks; they require a packet entry event to commence the measurement operation.

### I.3.2.1 Additivity of bursty packet loss

For further study.

### I.3.3 Delay variation

Traditional IP and many data applications are less sensitive to delay variation than applications requiring continuous play-out rates (isochronous or real-time applications, such as VoIP). For this reason, traditional IP applications are sometimes categorized by their "elastic" packet streams, which can be compressed or stretched during transfer with little effect. When VoIP packet spacing changes in transport, the application requires additional buffering to restore the continuous play-out capability and the buffer itself is a source of delay with these "non-elastic" packet streams.

It is difficult to concatenate delay variation of various network sections and processors in an accurate way, since the complete delay distribution is seldom known. However, a reasonable approximation method is given in clause 8 of [b-ITU-T Y.1541].

## I.4 Bottleneck bandwidth

Using the access bandwidth as the limiting factor when calculating end-to-end data transfer time can lead to an overly optimistic view of throughput.

### I.4.1 Key considerations

A more accurate approximation of the actual data transfer time of a fixed-size file or message can be made by including the following considerations in the model:

1) The transfer direction must be given. Access rates are asymmetrical for many new technologies such as DSL and cable modems, and the difference may be a factor of ten or more.

2) Use the real bottleneck bandwidth. Access rate is not necessarily the bottleneck. One of the hosts may throttle the transmission rate (by failing to pass sufficient data to the TCP process or with window limits), host TCP parameters could limit throughput, and middleboxes such as performance enhancing proxies may change the TCP parameters to improve throughput, etc.

3) The transaction may encounter contention for shared resources. Use the effective bottleneck bandwidth implied by the design objectives to support some number of simultaneous users on a link or other shared resource.

4) The bottleneck bandwidth may be time-varying. The bandwidth limits will produce a range of data transfer times.

5) Overhead must be calculated at the point of the bottleneck. Each layer below IP adds some header overhead. Assume a payload size and calculate the overhead percentage; this further increases the data transfer time.

6)     TCP does not achieve transfer capacity equal to the bottleneck bandwidth. Mathis' equation ([b-Mathis]) for TCP gives a good approximation of the steady state transfer capacity, yet the factors above may have overriding importance.

Most of these items require more preparation on the part of the model user or network planner. Details of the TCP capacity approximation are given in clause I.4.2.

## I.4.2     TCP capacity

The overwhelming majority of traditional IP applications use the reliable byte stream transfer services of TCP. A model for behaviour of TCP flow control [b-Mathis], is described as:

$$BW < \frac{MSS}{RTT} \frac{C}{\sqrt{p}}$$

where:

$BW$  is the bandwidth (BW), i.e., data transferred per unit time (cycle time)

$MSS$  is the TCP maximum segment size (MSS)

$RTT$  is the average round-trip time (RTT)

$C$  is a constant that accounts for the effects of random/periodic loss and ACK strategy

$p$  is the packet loss probability

This equation should be viewed as an upper bound on information flow, even for recent TCP enhancements, such as selective acknowledgments (SACK) and TCPReno (fast recovery). It assumes that the TCP connection is sufficiently long-lived to reach equilibrium in the congestion avoidance state. When the total bytes of the data transfer is small, TCP flow control may not reach equilibrium and slow-start behaviour dominates the calculation. A more flexible relationship for TCPReno capacity may be found in [b-Padhye], including the limiting effects of maximum window size.

$$B(p) \approx \min\left(\frac{W_{\max}}{RTT}, \frac{1}{RTT\sqrt{\frac{2bp}{3}} + T_0 \min\left(1, 3\sqrt{\frac{3bp}{8}}\right)p\left(1 + 32p^2\right)}\right)$$

where:

$B(p)$  is the approximate model of TCP throughput [packet/s]

$W_{\max}$  is the maximum window buffer size of receiver [packets]

$RTT$  is the round-trip time [seconds]

$b$  is the number of packets that are acknowledged by a received ACK

$p$  is the probability that a packet is lost

$T_0$  is the time-out for retransmitting an unacknowledged (lost) packet [seconds]

TCP's flow control attempts to grow its sending rate (window size) until it encounters congestion (or a bottleneck link), which it infers from packet loss. Thus, some packet loss is inherent to bottleneck probing. The congestion avoidance flow control of halving the window when a loss occurs, then growing the window one packet at a time until another loss occurs, effectively constrains the throughput to about 75% of the peak window size or sending rate. The equation shown below has been arranged to solve for the packet loss due to TCP's bottleneck probing alone.

$$\sqrt{p_{\text{Probing}}} \approx \frac{MSS}{RTT} \frac{C}{0.75 \times BW(bottleneck\_link)}$$

Table I.1 lists the inherent packet loss associated with a particular bottleneck link speed.

**Table I.1 – TCP packet loss due to bottleneck probing**

| C | MSS | RTT | BW (link) | 75% BW | Delay*BW | Window | p |
|---|---|---|---|---|---|---|---|
| 0.866 | 12 000 | 0.08 | 10 000 000 | 7 500 000 | 800 000 | 66.666667 | 3.00E-04 |
| 0.866 | 12 000 | 0.08 | 1 536 000 | 1 152 000 | 122 880 | 10.24 | 1.27E-02 |
| 0.866 | 12 000 | 0.08 | 768 000 | 576 000 | 61 440 | 5.12 | 5.09E-02 |
| 0.866 | 12 000 | 0.08 | 384 000 | 288 000 | 30 720 | 2.56 | 2.03E-01 |
| 0.866 | 12 000 | 0.08 | 128 000 | 96 000 | 10 240 | **0.8533333** | 1.83E+00 |
| 0.866 | 12 000 | 0.08 | 64 000 | 48 000 | 5 120 | **0.4266667** | 7.32E+00 |
| 0.866 | **2 048** | 0.08 | 128 000 | 96 000 | 10 240 | 5 | 5.33E-02 |
| 0.866 | **2 048** | 0.08 | 64 000 | 48 000 | 5 120 | 2.5 | 2.13E-01 |

Several points may be noted in Table I.1:

1)    Inherent packet loss is very low when the bottleneck link BW is 10 Mbit/s (Ethernet or higher). If network packet loss is $10^{-4}$ or higher, it will tend to reduce the throughput in accordance with [b-Mathis].

2)    For the given parameters (80 ms RTT, 12 000 bits (1 500 bytes) MSS, etc.) about 1 in 100 packets are lost on a $T_1$ link when a single TCP flow probes the bottleneck. Network packet loss at $10^{-3}$ or lower would be almost inconsequential to the resulting throughput.

3)    The Delay*BW product and the optimum window size (product/MSS) for the path are shown, and note that there are two link BW (128 kbit/s and 64 kbit/s) where the window is less than 1 and the calculated loss ratio is meaningless (>1). TCP tuning is warranted in these cases; thus the MSS is reduced to 2 048 bits in order to produce a more reasonable result.

4)    Note that the BW reduction due to lower layer overhead has not been addressed here.

The estimated network loss ratio is combined with the probing loss ratio as follows:

$$p_{Total} = 1 - \{(1 - p_{Network}) \times (1 - p_{Probing})\}$$

When the network and probing packet loss ratios are of the same order of magnitude, then it is somewhat pessimistic to combine them and compute the BW because some of the probing losses will not occur.

Finally, it may be noted that TCP parameter tuning is not addressed in any detail here, except to say that parameters can limit throughput in some cases. If the maximum window size is insufficient to fill the round-trip path, then the throughput will be limited to window × MSS/RTT [b-Padhye]. Many TCP connections never experience packet loss throughout their lifetime, owing to small windows and TCP's tendency to fill the queue in front of the bottleneck (and increasing the RTT for some packets – which is why an average RTT is used).

## I.5    Handshaking time

Packet exchange may be categorized in two phases:

1)    Handshaking; and

2)    Data transfer. (The traditional $3 \times 3$ matrix treats disconnection as a separate phase. Some applications do not have this phase, for example point-of-sale terminals for credit cards and VoIP signalling. The final packet exchange to close the connection in this phase is included, as a simplification.)

If a packet or its response is lost during initial handshaking, the sender usually waits a specified period of time before retransmitting the message. The waiting time, or retransmission time-out, can be on the order of one to three seconds, and significantly extends the initial handshake time when a loss occurs.

If all network-related times (including those caused by loss) are separated from the host processing times, there are two metrics that align with separable administrative responsibilities: network and CPE/hosts.

This decomposition is fairly straightforward for the initial handshaking phase. Packet transfer times and time-outs due to loss are attributed to the networks.

For data transfer time, it may be sufficient to give the time and indicate whether the limiting factor is attributable to network loss, network delay, or host processing/settings.

### I.5.1 Effect of loss during handshaking

The IP network contribution to initial handshaking time can be determined as illustrated in Figure I.2 below. First, take the case where eight messages and responses must be exchanged and all packet transfers are successful. The probability of this case is shown below:
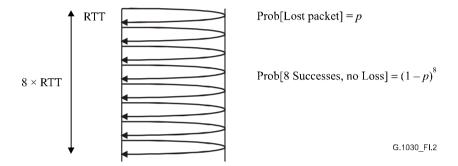


**Figure I.2 – Handshake with 8 round-trip turns and no packets lost**

This is a simple case, but it represents the IP network contribution to almost all of the customer experience when round-trip packet loss probability, $p$, is sufficiently low. Round-trip loss probability is given as:

$$p = p_{RT} = 1 - \{(1 - p_{1-way}) \times (1 - p_{other-way})\}$$

The time to accomplish all eight exchanges is 8 times the average RTT, plus the remote host processing time (HPT), and this total is the handshaking time. These components may be calculated separately, as:

$$Handshake\_Time = NCTT + Total\_HPT$$

where network contribution to transaction time (NCTT) is defined in [ITU-T G.1040].

Note that as an example, the post office protocol version 3 (POP3) handshake has eight request/response exchanges. Assuming that p = $10^{-3}$, RTT + HPT = 0.080 seconds, retransmit time-out (RTO) = 1 second, and that there are 350 000 transaction attempts, the probabilities for each of the cases of loss are shown in Table I.2.

**Table I.2 – Handshake time for cases with 0 to 3 packets lost**

| Losses | Handshake time [s] | Probability of occurrence | Transactions |
|--------|---------------------|----------------------------|--------------|
| 0 | 0.64 | 0.99202794 | 347 210 |
| 1 | 1.64 | 0.00793622 | 2 778 |
| 2 | 2.64 | 3.5713E-05 | 12 |
| 3 | 3.64 | 1.1904E-07 | 0 |

This example shows that the IP network contributes less than a second to average, 95th percentile and 99th percentile of handshake time for a POP3 transaction with eight RT exchanges.

### I.5.2 Effect of packet loss during data transfer

Packet loss will require retransmissions, and with TCP, the congestion avoidance flow control reduces the sending rate, as described in clause I.4.

### I.6 Application example: HTTP transaction

Figure I.3 shows an HTTP transaction identifying handshake time. The remaining packet exchange constitutes the data transfer time from the client perspective.
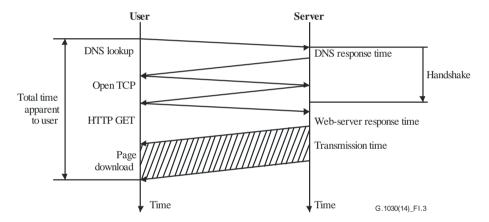


**Figure I.3 – HTTP protocol exchange**

The reference connection of Figure I.1 is used in the examples below, showing two network sections, with NATs, firewalls, and a proxy as example middleboxes. Note that the proxy does not have the requested objects stored in a cache, and makes a request to the remote web-server.

### I.6.1 Handshake time

There are two request-responses conducted over different paths, the DNS look-up and the TCP 3-way handshake (where the SYN and SYN-ACK are timed, the last ACK is assumed to proceed along with the HTTP GET and does not add significant time). Although the proxy splits the TCP connection, it is treated as additional delay on a single connection in this example.

Since most of the accumulation is simple addition or multiplication, the values are simply tabulated and the totals are shown (rather than writing out equations with many new parameters).

First, the components of DNS performance are shown in Table I.3 (and referring to the reference connection of Figure I.1).

**Table I.3 – DNS look-up**

| Client to DNS | RTT (net) [ms] | HPT (Proc time) [ms] | RT (1-p) | per flow BW |
|---|---|---|---|---|
| LAN | 2 | | 0.99999 | 5.0E+6 |
| NAT/firewall | | 4 | 0.99999 | |
| Link 1 | 10 | | 0.9999 | 5.0E+6 |
| NAT/edge router | | | 0.99999 | |
| NAT to DNS link | 4 | | 0.99999 | 10.0E+6 |
| DNS | | 6 | 0.99999 | |
| Totals (inverted loss) | 16 | 10 | 0.00015 | |

Last, the components of TCP connection set-up are shown in Table I.4 (and referring to the reference connection of Figure I.1).

**Table I.4 – Client-server path for TCP connection time and data transfer time**

| Client to server | RTT (net) [ms] | HPT (Proc time) [ms] | RT (1-p) | per flow BW [Mbits/s] |
|---|---|---|---|---|
| LAN | 2 | | 0.99999 | 5.0E+6 |
| NAT/firewall | | 4 | 0.99999 | |
| UNI link 1 | 10 | | 0.9999 | 5.0E+6 |
| NAT/edge router | | | 0.99999 | |
| NAT to proxy link | 4 | | 0.99999 | 10.0E+6 |
| Proxy | | 4 | 0.99999 | |
| NNI link | 20 | | 0.99999 | 10.0E+6 |
| Network section 2 | 14 | | 0.99999 | 10.0E+6 |
| UNI link 2 | 10 | | 0.99999 | 3.0E+6 |
| Firewall | | 2 | 0.99999 | |
| LAN | 2 | | 0.99999 | 5.0E+6 |
| Server (TCP proc) | | 1 | 0.99999 | |
| Totals (inverted loss) | 62 | 11 | 0.00021 | |

With a DNS lookup time = 16 ms + 10 ms = 26 ms, and TCP open time of 62 ms + 11 ms = 73 ms, the handshake time is 99 ms. This time appears to be valid for more than 99.9% of HTTP transactions.

### I.6.2 Data transfer time

For data transfer time calculations, the first steps are to determine whether the client-server loss will play a role.

As shown in Table I.4 (and referring to the reference connection of Figure I.1), the lowest per-flow BW occurs on UNI link 2, at 3 Mbit/s. The delay × BW product is 0.073 × 3 Mbit/s = 219 kbit/s – indicating an optimum window size of 18.25 packets at 12 000 bits (1 500 bytes) MSS. TCP can be tuned to accommodate this path, so there is no constraint associated with window size.

Using the bottleneck link BW, TCP's inherent packet loss ratio due to probing is calculated as:

$$\sqrt{p_{\text{Probing}}} \approx \frac{MSS}{RTT} \frac{C}{0.75 \times BW} \text{, so } p = \left( \frac{12\,000}{0.073} \frac{0.866}{0.75 \times 3M} \right)^2 = 0.004003$$

Since the calculated network loss (0.00021) is approximately *P*probing/20, it is considered negligible, and TCP's byte stream throughput will be 0.75 × 3 Mbit/s = 2.25 Mbit/s.

If, on the other hand, the estimated network loss ratio was on the same order as the probing packet loss ratio, then it would be somewhat pessimistic to combine them and re-compute the BW:

$$p_{Total} = 1 - \{(1 - p_{Network}) \times (1 - p_{Probing})\}$$

Assuming a webpage with complex graphics, the transfer time for a 1 Mbit webpage is 0.444 seconds. The time for the server to process the HTTP GET (0.01 seconds) can be added for a total of 0.454 seconds.

### I.6.3    Total time apparent to user

The total time that is apparent to users is the sum of the handshake time and the data transfer time: 0.099 + 0.454 = 0.553 seconds.

### I.7    Summary

This appendix describes a methodology for estimating end-to-end performance in IP networks. The results obtained may be used in comparison with user-centric performance targets, e.g., [ITU-T G.1010], to estimate overall end-user satisfaction with different multimedia applications.

# Appendix II

## Opinion model for a simple web search task

(This appendix does not form an integral part of this Recommendation.)

### II.1 Scope

In this appendix, a model is provided for mapping the response and download times, as measured in the network or calculated from the HTTP transaction times, to the perceived quality of a simple web search session. This model is based on experiments where the response and download times in a simple web search session were manipulated [b-Beerends], [b-Van]. The scope of the model is currently limited to web search sessions that consist of two steps: a first step in which a search request is made, and a second step in which a results page is shown. A simple extension towards single-timing events, where the impact of waiting for a single page is modelled, is also provided.

It is pointed out that this simple web search task does not reflect the dynamics of real web browsing such as flow experience, interaction with the web site and consecutive load time patterns throughout a real web browsing session. Information on these influence factors is given in [ITU-T G.1031].

### II.2 Introduction

An important observation in modelling the perceived quality in web browsing is the fact that the expected maximal session time will dominate the perceived quality. If one expects a session time of 100 seconds, the perceived quality of a 10-second session will be much higher than if one expects a session time of 1 second. Therefore, the model takes a context-dependent approach by using three different time-scales: 6, 15 and 60 seconds, corresponding to fast, medium and slow network contexts, respectively.
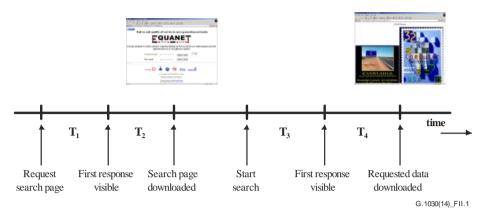
In general, quality perception related to response time can be classified according to the following three perceptual regions [b-Nielsen]:

1) **Instantaneous experience**: A delay of about 0.1 second is the limit for having the feel that the system is reacting instantaneously, an important limit for conversational services (e.g., chatting).

2) **Uninterrupted experience**: A delay of about 1.0 second is the limit for the user's flow of thought to remain uninterrupted, even though the user does lose the feeling that the service is operating directly, an important limit for interactive services (e.g., gaming).

3) **Loss of attention**: A delay of about 10 seconds is the limit for keeping the user's attention focused on the dialogue. For longer delays, users will want to perform other tasks while waiting for the computer to finish, so they should be given feedback indicating when the computer expects to be done. Feedback during the delay is especially important if the response time is likely to be highly variable since users will then not know what to expect.

Regarding download times, subjects tend to adapt their quality judgment towards the expected download time [b-Dellaert]. When subjects are informed about the expected download time, they are willing to accept longer download times.

The model in this Recommendation describes the relationship between the different response and download times within web browsing sessions and the corresponding perceived web browsing quality for a given *maximum* session time within a certain network and system configuration. The model is applicable to a wide range of network and system configurations, as well as to web browsing services for a wide variety of users. The subjective experiments, which served as the basis for the model,

mimicked a real-life web browse experience as closely as possible. Three subjective web browse experiments with time-scales of approximately 6, 15 and 60 seconds, representing fast, medium, and slow network contexts, respectively, were used in the model development. In each session, a subject first retrieves a search page and then a page that shows the search results. In Figure II.1, the time-line of such a session is shown. The first two time-intervals, $T_1$ and $T_2$, represent the non-interactive response and download times of the search page. The second two time-intervals, $T_3$ and $T_4$, represent the interactive response and download times of the result page.



T_1 is the non-interactive response time and was manipulated using Java scripting. $T_2$ is the non-interactive download time and was manipulated using the network manipulator. $T_3$ and $T_4$ are the equivalents for the interactive part. The sum $T_1+T_2+T_3+T_4$ represents the session time.

**Figure II.1 – Experiment timers**

Due to the known differences in behaviour between trained experts and untrained, so-called naïve, users, a distinction is made between these groups in the development of the model. Separating these two groups allows for a model to be developed that predicts the quality of web browsing for a large population of users.

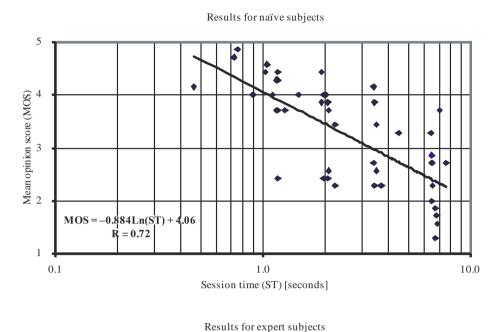## II.3    Subjective web browse quality experiment and results

In the experiments, and consequently in the model, the ITU-T absolute category rating scale [b-ITU-T P.800] was used (i.e., the five-point scale where 5: excellent, 4: good, 3: fair, 2: poor, 1: bad). In each experiment, 49 sessions were presented, each consisting of:
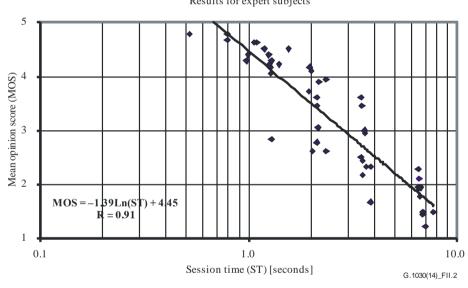
•        requesting, retrieving, and displaying of a search page;

•        typing and submitting a search term on this page;

•        retrieving and displaying of the results page.

To obtain consistent data, each session used exactly the same results pages and subjects were asked to type in the same search query in every session. From the perspective of the subjects, the search engine first has to find the results page, which then has to be downloaded. For each of the 49 sessions, different combinations of $T_1$ through $T_4$ are configured, varying the sum $T_1+T_2+T_3+T_4$, i.e., the session time, from 0 to the time-scale for this set of experiments.

Overall results are given in Figures II.2-II.4 where the session time (i.e., $T_1+T_2+T_3+T_4$) is plotted versus the mean opinion score (MOS) value for all three experiments. For the long duration context, see Figure II.4. In this context, the results for naïve and expert subjects were about the same and MOS values were calculated over the whole population. For the two shorter duration contexts, see Figures II.2 and II.3. In these contexts, naïve and expert subjects behaved differently, and the correlations between session times and perceived quality are significantly different and are thus presented separately. All results show the same behaviour, i.e., that the perceived quality goes down linearly with the logarithm of the session time. The correlation in the long duration experimental context is sufficiently high (>0.9) to make reliable quality predictions for both naïve and expert subjects. In
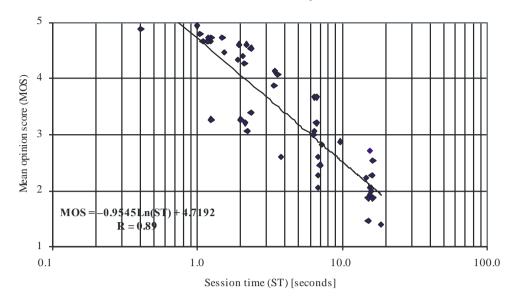
general, correlations above 0.9 are strived for in psychophysical modelling of quality perception [b-ITU-T P.862]. The results also show that for the 6-second experiment with naïve subjects the correlation between session time and perceived quality is far too low (0.72) in order to allow for a simple model based on session time only.

Results for naïve subjects



$$MOS = -0.884Ln(ST) + 4.06$$
$$R = 0.72$$

Results for expert subjects



$$MOS = -1.39Ln(ST) + 4.45$$
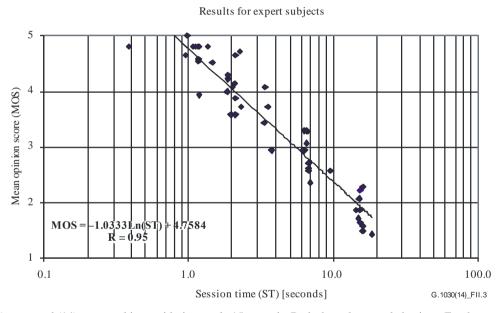$$R = 0.91$$

G.1030(14)_FII.2

Results for (7) naïve and (18) expert subjects with time-scale 6 seconds. For the naïve subjects, the correlation is too low to allow for accurate MOS predictions. For the expert subjects, the MOS can be predicted from a logarithmic interpolation of the session time between 0.67 and 12 seconds.
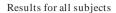
**Figure II.2 – Results for subjects with time-scale 6 seconds**

Results for naïve subjects



$$MOS = -0.9545 Ln(ST) + 4.7192$$
$$R = 0.89$$

Results for expert subjects



$$MOS = -1.0333 Ln(ST) + 4.7584$$
$$R = 0.95$$

G.1030(14)_FII.3

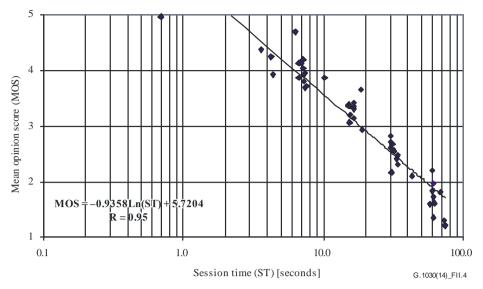Results for (15) naïve and (14) expert subjects with time-scale 15 seconds. Both show the same behaviour. For the naïve subjects, the correlation is slightly too low to allow for accurate MOS predictions. For the expert subjects, the MOS can be predicted from a logarithmic interpolation of the session time between about 0.79 and 38 seconds.

**Figure II.3 – Results for subjects with time-scale 15 seconds**

Results for all subjects

MOS = −0.9358Ln(ST) + 5.7204
R = 0.95

Session time (ST) [seconds]

G.1030(14)_FII.4

Results for all (12+17=29) subjects with time-scale 60 seconds. The MOS can be predicted from a logarithmic interpolation of the session time between 2.16 and 155 seconds.

**Figure II.4 – Results for subjects with time-scale 60 seconds**

## II.4 Modelling results

### II.4.1 60-second context

For the 60-second context, the correlations between session time and subjective quality are very good (0.95) for both naïve and expert subjects and the mapping from session time to subjective quality can directly be used in the objective model. The regression shows that session times below about 2 seconds lead to the maximum MOS value of 5 (see Figure II.4), while the minimum MOS value of 1 is obtained for session times that are larger than the longest session time in the actual experiment. A general mapping from session time to web browsing quality for the long duration context is constructed by defining a minimum (*Min*) and a maximum (*Max*) session time and using a logarithmic interpolation between these extreme session times. When writing $MOS = a - b \cdot \ln(SessionTime)$, and substituting *MOS* = 5 for *SessionTime = Min,* and *MOS* = 1 for *SessionTime = Max*, session times between *Min* and *Max are obtained*:

$$MOS = \frac{4}{\ln(Min \, / \, Max)} \cdot \left( \ln(SessionTime) - \ln(Min) \right) + 5 \qquad \text{(II-1)}$$

For the long duration experiment (see Figure II.4) the regression is:

$$MOS = 5.72 - 0.936 \cdot \ln(SessionTime), \qquad \text{(II-2)}$$

clipped between MOS 1.0 and 5.0.

### II.4.2 6-second and 15-second contexts

For the 6-second and 15-second experimental contexts, the correlations between session time and subjective quality are much lower than for the 60-second context, and a more advanced model for the prediction of the subjective quality for naïve and expert users is constructed using the idea that for shorter duration session times the last download time ($T_4$ in the experiment) has a more severe impact on the final perceived web browsing quality than the other response and download times ($T_1$, $T_2$, $T_3$ in the experiment). Table II.1 gives the weight factors with which $T_1$ through $T_4$ have to be weighted in order to get a quantity that has the highest correlation with the subjectively determined MOS values. This quantity, the weighted session time:

$$WeightedST = WT1 \cdot T1 + WT2 \cdot T2 + WT3 \cdot T3 + WT4 \cdot T4 \tag{II-3}$$

can be mapped to the MOS value using the same logarithmic interpolation between minimum and maximum session times as used in equation II-1:

$$MOS = \frac{4}{\ln(Min / Max)} \cdot \big(\ln(WeightedST) - \ln(Min)\big) + 5 \tag{II-4}$$

Table II.1 shows that for the shortest duration context, the impact of the last download time is more than twice as large as the impact of the other download and response times. It also shows a significantly different behaviour for naïve and expert subjects, with the optimal weighting for naïve subjects showing a larger impact from the last download time than the optimal weights for the expert subjects. For the naïve subjects, the impact of the large download time is more than four times as large as the impact of the other download and response times. For expert subjects, only a factor of about two exists, while the overall best weight shows an impact that is about three times as large.

**Table II.1 – Optimal model weighting for $T_1$, $T_2$, $T_3$ and $T_4$ with the associated model correlations between objective timing and subjective MOS results**

|  | WT1 | WT2 | WT3 | WT4 | *Min* | *Max* | Correlation |
|---|---|---|---|---|---|---|---|
| 6 s expert | 0.56 | 0.84 | 0.80 | 1.80 | | | 0.97 |
| 6 s naïve | 0.37 | 0.40 | 0.60 | 2.63 | | | 0.93 |
| **6 s overall** | **0.47** | **0.60** | **0.71** | **2.22** | **0.62** | **13.5** | **0.95** |
| | | | | | | | |
| 15 s expert | 0.63 | 0.77 | 1.11 | 1.49 | | | 0.98 |
| 15 s naïve | 0.48 | 0.70 | 0.88 | 1.95 | | | 0.96 |
| **15 s overall** | **0.54** | **0.72** | **0.98** | **1.76** | **0.81** | **39** | **0.97** |
| | | | | | | | |
| 60 s expert | 0.84 | 0.77 | 1.22 | 1.18 | | | 0.99 |
| 60 s naïve | 0.64 | 1.01 | 1.12 | 1.24 | | | 0.98 |
| **60 s overall** | **0.73** | **0.90** | **1.16** | **1.22** | **2.22** | **151** | **0.98** |
| ***60 s overall, no weighting, see Figure II.4*** | ***1.00*** | ***1.00*** | ***1.00*** | ***1.00*** | ***2.16*** | ***155*** | ***0.95*** |
| NOTE – The sum of the weighting coefficients is normalized to 4.0 in order to be able to compare normal session times ($T_1+T_2+T_3+T_4$) with weighted session times. The weighting as used in the standardized model together with the Min and Max times used in equation II-4 are shown in bold text. For long session time contexts, >60 seconds, the simple model without weighting may be used (bold italics). | | | | | | | |

Table II.1 also shows that for the medium duration context, the weight factors for naïve and expert subjects as well as the overall weight factors are in-between the weight factors for the short and long duration context experiments. This shows the validity of the weighting approach, thus allowing for an interpolation between the different *experimental session context times* (i.e., the time-scales 6, 15, and 60 seconds) in order to obtain weightings for other context times.

## II.4.3    Summary

This appendix standardizes a simple model for the 60-second duration context that allows assessing of web browse sessions, for which the maximum session time is about 155 seconds, using equation II-2. This appendix further standardizes three advanced models using the best overall weights from Table II.1 in combination with the following mappings from weighted session time to perceived browse quality in terms of mean opinion scores:

$$MOS = 4.38 - 1.30 \cdot \ln(WeightedSessionTime) \tag{II-5}$$

clipped between MOS 1.0 and 5.0 for short duration sessions.

$$MOS = 4.79 - 1.03 \cdot \ln(WeightedSessionTime) \tag{II-6}$$

clipped between MOS 1.0 and 5.0 for medium duration sessions.

$$MOS = 5.76 - 0.948 \cdot \ln(WeightedSessionTime) \tag{II-7}$$

clipped between MOS 1.0 and 5.0 for long duration sessions.

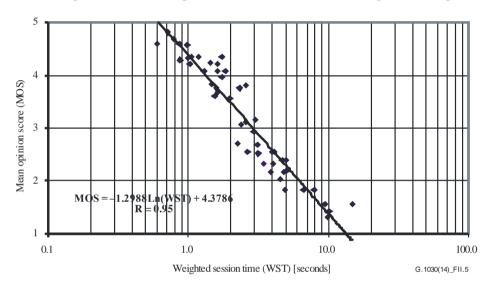The results from the regression fits of equations II-5, II-6 and II-7 are given in Figures II.5 to II.7.



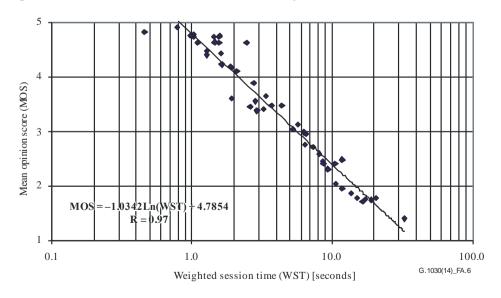**Figure II.5 – Model versus data for all subjects with time-scale 6 seconds**



**Figure II.6 – Model versus data for all subjects with time-scale 15 seconds**
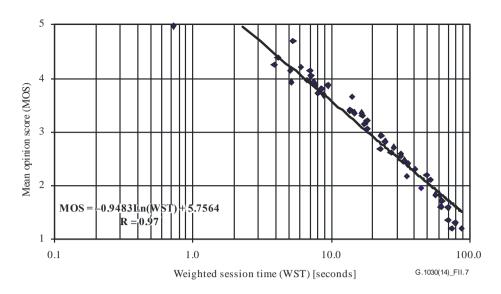
**Figure II.7 – Model versus data for all subjects with time-scale 60 seconds**

## II.5 Perceived quality of one-page web browse sessions and single timing events

Based on the experimental data that was the starting point for deriving equations II-1 and II-4, a relationship is derived between the session time of web browse sessions consisting of a single webpage and the perceived quality. In addition, sessions of a single page for which the download time always equals 0 is of interest. This situation occurs when the information retrieved appears instantaneously to the user after waiting some period time. This last relationship maps a single timing event towards the subjectively perceived quality in terms of a MOS score.

The derivation of the one-page mappings starts with the observation that the perceived quality goes down linearly with the logarithm of the session time between a minimum and maximum session time for which the MOS scores are 5.0 and 1.0 respectively (see Figures II.2 to II.7). The two-page browse data, presented in Figures II.2-II.4, shows that the minimum session time in the three subjective experimental contexts that were used, varies between approximately 0.7 and 2.2 and increases when the maximum duration in the experiment increases (see the results overview in Table II.2).

**Table II.2 – Non-weighted Min and Max times
in the 2-page browse experiment**

|  | *Min* (sec.) | *Max* (sec.) |
|---|---|---|
| 6 s expert | 0.67 | 12 |
| 15 s expert | 0.79 | 38 |
| **60 s overall** | 2.16 | 155 |

Based on these results a minimum session time (*Min*) is defined, for which a MOS score of 5.0 is achieved, as *Min* = 0.011*Max* + 0.47, with *Max* representing the maximum session time that is expected to occur. This allows a general mapping from session time to the MOS score for the two-page experiment for any expected maximum duration (*Max*) of the two-page web search session:

$$MOS_{2-page} = \frac{4}{\ln((0.011Max+0.47)/Max)} \cdot \left(\ln(SessionTime)-\ln(0.011Max+0.47)\right)+5 \qquad \text{(II-8)}$$

From the experimental data that has been presented it is estimated that this relationship is likely to hold for all two-page web browse sessions between about 10 and 200 seconds.

For arbitrary one-page sessions that only have one response time, $T_1$, and one download time, $T_2$, (see Figure II.1) the highest quality is perceived for a minimum value that is about half of the value

obtained in the two-page sessions. This leads to the minimum session time definition $Min = 0.005Max + 0.24$ while the mapping to MOS value is given by equation II-9:

$$MOS_{1-page} = \frac{4}{\ln((0.005Max+0.24)/Max)} \cdot \left(\ln(SessionTime) - \ln(0.005Max+0.24)\right) + 5 \qquad \text{(II-9)}$$

From the experimental data that has been presented it is estimated that this relationship is likely to hold for all one-page web browse sessions between approximately 5 and 100 seconds.

Similarly, for arbitrary one-page sessions in which the download time always equals 0, or for any single timing event, the MOS value is given by equation II-10:

$$MOS_{single\ timing\ event} = \frac{4}{\ln((0.003Max+0.12)/Max)} \cdot \left(\ln(SessionTime) - \ln(0.003Max+0.12)\right) + 5 \quad \text{(II-10)}$$

From the experimental data that has been presented it is estimated that this relationship is likely to hold for all single timing events between about 3 and 50 seconds. Note that for such a single interaction, the minimum time equals 0.12 seconds, corresponding to an instantaneous perception threshold [b-Nielsen].

It should be noted that the above models for the perceived quality of one-page web browse sessions and single timing events should be validated by actual experiments.

# Bibliography

[b-ITU-T P.800]    Recommendation ITU-T P.800 (1996), *Methods for subjective determination of transmission quality*.

[b-ITU-T P.862]    Recommendation ITU-T P.862 (2001), *Perceptual evaluation of speech quality (PESQ): An objective method for end-to-end speech quality assessment of narrow-band telephone networks and speech codecs*.

[b-ITU-T Y.1541]   Recommendation ITU-T Y.1541 (2006), *Network performance objectives for IP-based services*.

[b-Beerends]    Beerends, J.G., Van Der Gaast, S., Ahmed, O.K. (2004), *Web browse quality modelling*, White contribution COM 12-C 3 to ITU-T Study Group 12, November.

[b-Dellaert]    Dellaert, G.C., Kahn, B.E. (1998), *How Tolerable is Delay? Consumers' Evaluations of Internet Websites after Waiting*.
htttp://greywww.kub.nl:2080/greyfiles/center/1998/64.html

[b-Mathis]    Mathis, M., Semke, J., Madavi, J., Ott, T. (1997), *The macroscopic behavior of TCP congestion avoidance algorithm*, Computer communications review, Vol. 27, No. 3, July, ISSN# 0146-4833.
http://www.psc.edu/networking/papers/model_ccr97.ps

[b-Nielsen]    Nielsen, J. (1994), *Response Times: The Three Important Limits*. Available:
http://www.useit.com/papers/responsetime.html

[b-Padhye]    Padhye, J., Firoiu, V., Towsley, D. and Kurose, J. (1998), *Modeling TCP Throughput: a Simple Model and its Empirical Validation*, SIGCOMM.
ftp://gaia.cs.umass.edu/pub/Padhye-Firoiu98:TCP-throughput.ps.Z

[b-Van]    Van Der Gaast, S., Beerends, J.G., Ahmed, O.K. and Meeuwissen, H.B. (2005), *Quantification and prediction of end-user perceived web-browsing quality*, submitted on March 24.

# SERIES OF ITU-T RECOMMENDATIONS

| | |
|---|---|
| Series A | Organization of the work of ITU-T |
| Series D | General tariff principles |
| Series E | Overall network operation, telephone service, service operation and human factors |
| Series F | Non-telephone telecommunication services |
| **Series G** | **Transmission systems and media, digital systems and networks** |
| Series H | Audiovisual and multimedia systems |
| Series I | Integrated services digital network |
| Series J | Cable networks and transmission of television, sound programme and other multimedia signals |
| Series K | Protection against interference |
| Series L | Construction, installation and protection of cables and other elements of outside plant |
| Series M | Telecommunication management, including TMN and network maintenance |
| Series N | Maintenance: international sound programme and television transmission circuits |
| Series O | Specifications of measuring equipment |
| Series P | Terminals and subjective and objective assessment methods |
| Series Q | Switching and signalling |
| Series R | Telegraph transmission |
| Series S | Telegraph services terminal equipment |
| Series T | Terminals for telematic services |
| Series U | Telegraph switching |
| Series V | Data communication over the telephone network |
| Series X | Data networks, open system communications and security |
| Series Y | Global information infrastructure, Internet protocol aspects and next-generation networks |
| Series Z | Languages and general software aspects for telecommunication systems |