

**Speech Processing, Transmission and Quality Aspects (STQ);  
Transmission requirements for wideband  
VoIP terminals (handset and headset)  
from a QoS perspective as perceived by the user**

---



---

Reference

DES/STQ-00091

---

Keywords

quality, telephony, terminal, VoIP

**ETSI**

650 Route des Lucioles  
F-06921 Sophia Antipolis Cedex - FRANCE

Tel.: +33 4 92 94 42 00 Fax: +33 4 93 65 47 16

Siret N° 348 623 562 00017 - NAF 742 C  
Association à but non lucratif enregistrée à la  
Sous-Préfecture de Grasse (06) N° 7803/88

---

**Important notice**

Individual copies of the present document can be downloaded from:

<http://www.etsi.org>

The present document may be made available in more than one electronic version or in print. In any case of existing or perceived difference in contents between such versions, the reference version is the Portable Document Format (PDF). In case of dispute, the reference shall be the printing on ETSI printers of the PDF version kept on a specific network drive within ETSI Secretariat.

Users of the present document should be aware that the document may be subject to revision or change of status. Information on the current status of this and other ETSI documents is available at

<http://portal.etsi.org/tb/status/status.asp>

If you find errors in the present document, please send your comment to one of the following services:

[http://portal.etsi.org/chaicor/ETSI\\_support.asp](http://portal.etsi.org/chaicor/ETSI_support.asp)

---

**Copyright Notification**

No part may be reproduced except as authorized by written permission.  
The copyright and the foregoing restriction extend to reproduction in all media.

© European Telecommunications Standards Institute 2007.  
All rights reserved.

**DECT**<sup>TM</sup>, **PLUGTESTS**<sup>TM</sup> and **UMTS**<sup>TM</sup> are Trade Marks of ETSI registered for the benefit of its Members.  
**TIPHON**<sup>TM</sup> and the **TIPHON logo** are Trade Marks currently being registered by ETSI for the benefit of its Members.  
**3GPP**<sup>TM</sup> is a Trade Mark of ETSI registered for the benefit of its Members and of the 3GPP Organizational Partners.

# Contents

Intellectual Property Rights .....	5
Foreword.....	5
Introduction .....	5
1 Scope .....	6
2 References .....	6
3 Definitions and abbreviations.....	7
3.1 Definitions .....	7
3.2 Abbreviations .....	8
4 General considerations .....	8
4.1 Coding Algorithm.....	8
4.2 End-to-end considerations .....	8
4.3 Parameters to be investigated .....	9
4.3.1 Basic parameters.....	9
4.3.2 Further Parameters with respect to Speech Processing Devices .....	9
5 Test equipment .....	9
5.1 IP half channel measurement adaptor.....	9
5.2 Network impairment simulation.....	10
6 Acoustic environment.....	10
7 Requirements and associated Measurement Methodologies .....	11
7.1 Test setup.....	11
7.2 Coding independent Parameters .....	13
7.2.1 Send Frequency response.....	13
7.2.2 Send Loudness Rating .....	14
7.2.3 D- Factor.....	15
7.2.4 Linearity Range for SLR.....	15
7.2.5 Send Distortion .....	16
7.2.6 Send Noise .....	16
7.2.7 Sidetone Masking Rating STMR (Mouth to ear).....	17
7.2.8 Sidetone delay.....	17
7.2.9 Terminal Coupling Loss weighted (TCLw).....	18
7.2.10 Stability Loss .....	18
7.2.11 Receive Frequency response.....	20
7.2.12 Receive Loudness Rating.....	21
7.2.13 Receiving Distortion .....	22
7.2.14 Minimum activation level and sensitivity in Receive direction .....	22
7.2.15 Receive Noise .....	22
7.2.16 Automatic Gain Control in Receiving .....	23
7.2.17 Double talk Performance .....	23
7.2.17.1 Attenuation Range in Sending Direction during Double Talk $A_{H,S,dt}$ .....	23
7.2.17.2 Attenuation Range in Receiving Direction during Double Talk $A_{H,R,dt}$ .....	25
7.2.17.3 Detection of Echo Components during Double Talk .....	25
7.2.17.4 Minimum activation level and sensitivity of double talk detection.....	27
7.2.18 Switching characteristics .....	27
7.2.18.1 Activation in Sending Direction.....	27
7.2.18.2 Silence Suppression and Comfort Noise Generation .....	28
7.2.18.3 Performance in Sending in the Presence of Background Noise .....	28
7.2.18.4 Speech Quality in the Presence of Background Noise .....	29
7.2.18.5 Quality of Background Noise Transmission (with Far End Speech) .....	29
7.2.18.6 Quality of background noise transmission (with Near End Speech).....	30
7.2.19 Quality of echo cancellation .....	30
7.2.19.1 Temporal echo effects.....	30

7.2.19.2	Spectral Echo Attenuation.....	31
7.2.19.3	Occurrence of Artefacts .....	31
7.2.20	Variant Impairments; Network Dependant .....	31
7.2.20.1	Delay versus Time Send.....	31
7.2.20.2	Delay versus Time Receive.....	31
7.2.20.3	Quality of Jitter buffer adjustment .....	31
7.3	Codec Specific Requirements.....	32
7.3.1	Send Delay .....	32
7.3.2	Receive delay.....	33
7.3.3	Objective Listening Speech Quality MOS-LQOM in Send direction.....	34
7.3.4	Objective Listening Quality MOS-LQOM in Receive direction .....	35
7.3.4.1	Efficiency of Packet Loss Concealment (PLC).....	36
7.3.4.2	Efficiency of Delay Variation Removal.....	36
<b>Annex A (informative):</b>	<b>Processing delays in VoIP terminals .....</b>	<b>37</b>
<b>Annex B (informative):</b>	<b>Bibliography.....</b>	<b>40</b>
History .....		41

---

## Intellectual Property Rights

IPRs essential or potentially essential to the present document may have been declared to ETSI. The information pertaining to these essential IPRs, if any, is publicly available for **ETSI members and non-members**, and can be found in ETSI SR 000 314: "*Intellectual Property Rights (IPRs); Essential, or potentially Essential, IPRs notified to ETSI in respect of ETSI standards*", which is available from the ETSI Secretariat. Latest updates are available on the ETSI Web server (<http://webapp.etsi.org/IPR/home.asp>).

Pursuant to the ETSI IPR Policy, no investigation, including IPR searches, has been carried out by ETSI. No guarantee can be given as to the existence of other IPRs not referenced in ETSI SR 000 314 (or the updates on the ETSI Web server) which are, or may be, or may become, essential to the present document.

---

## Foreword

This ETSI Standard (ES) has been produced by ETSI Technical Committee Speech Processing, Transmission and Quality Aspects (STQ).

---

## Introduction

Traditionally, the analogue and digital telephones were interfacing switched-circuit 64 kbit/s PCM networks. With the fast growth of IP networks, wideband terminals providing higher audio-bandwidth and directly interfacing packet-switched networks (VoIP) are being rapidly introduced. Such IP network edge devices may include gateways, specifically designed IP phones, soft phones or other devices connected to the IP based networks and providing telephony service. Since the IP networks will be in many cases interworking with the traditional PSTN and private networks, many of the basic transmission requirements have to be harmonized with specifications for traditional digital terminals. However, due to the unique characteristics of the IP networks including packet loss, delay, etc. new performance specification, as well as appropriate measuring methods, will have to be developed. Terminals are getting increasingly complex, advanced signal processing is used to address the IP specific issues.

Note Requirement limits are given in tables, the associated curve when provided is given for illustration.

---

# 1 Scope

The present document provides speech transmission performance requirements for 8 kHz wideband VoIP terminals; it addresses all types of IP based terminals, including wireless and soft phones.

The intention of the present document is to specify equipment requirements which enable manufacturers and service providers to enable good quality end-to-end speech performance.

The objective measurement methodologies and requirements given by the present document. are for handset and headset operation only. Hands-free terminals are outside the scope of the present document.

In addition to basic testing procedures, the present document describes advanced testing procedures taking into account further quality parameters as perceived by the user.

---

# 2 References

The following documents contain provisions which, through reference in this text, constitute provisions of the present document.

- References are either specific (identified by date of publication and/or edition number or version number) or non-specific.
- For a specific reference, subsequent revisions do not apply.
- For a non-specific reference, the latest version applies including subsequent corrigendums and amendments.

Referenced documents which are not found to be publicly available in the expected location might be found at <http://docbox.etsi.org/Reference>.

NOTE: While any hyperlinks included in this clause were valid at the time of publication ETSI cannot guarantee their long term validity.

- [1] ETSI EG 201 377-2: "Speech Processing, Transmission and Quality Aspects (STQ); Specification and measurement of speech transmission quality; Part 2: Mouth-to-ear speech transmission quality including terminals".
- [2] ETSI EG 202 396-1: "Speech Processing, Transmission and Quality Aspects (STQ); Speech quality performance in the presence of background noise; Part 1: Background noise simulation technique and background noise database".
- [3] ETSI EG 202 425: "Speech Processing, Transmission and Quality Aspects (STQ); Definition and implementation of VoIP reference point".
- [4] ETSI I-ETS 300 245-5: "Integrated Services Digital Network (ISDN); Technical characteristics of telephony terminals; Part 5: Wideband (7 kHz) handset telephony".
- [5] ITU-T Recommendation G.107: "The E-model, a computational model for use in transmission planning".
- [6] ITU-T Recommendation G.108: "Application of the E-model: A planning guide".
- [7] ITU-T Recommendation G.109: "Definition of categories of speech transmission quality".
- [8] ITU-T Recommendation G.122: "Influence of national systems on stability and talker echo in international connections".
- [9] ITU-T Recommendation G.711: "Pulse code modulation (PCM) of voice frequencies".
- [10] ITU-T Recommendation G.722: "7 kHz audio-coding within 64 kbit/s".
- [11] ITU-T Recommendation G.722.1: "Low-complexity coding at 24 and 32 kbit/s for hands-free operation in systems with low frame loss".

- [12] ITU-T Recommendation G.729.1: "G.729 based Embedded Variable bit-rate coder: An 8-32 kbit/s scalable wideband coder bitstream interoperable with G.729".
- [13] ITU-T Recommendation G.1020: "Performance parameter definitions for quality of speech and other voiceband applications utilizing IP networks".
- [14] ITU-T Recommendation P.50: "Artificial voices".
- [15] ITU-T Recommendation P.56: "Objective measurement of active speech level".
- [16] ITU-T Recommendation P.57: "Artificial ears".
- [17] ITU-T Recommendation P.58: "Head and torso simulator for telephonometry".
- [18] ITU-T Recommendation P.64: "Determination of sensitivity/frequency characteristics of local telephone systems".
- [19] ITU-T Recommendation P.79: "Calculation of loudness ratings for telephone sets".
- [20] ITU-T Recommendation P.340: "Transmission characteristics and speech quality parameters of hands-free terminals".
- [21] ITU-T Recommendation P.380: "Electro-acoustic measurements on headsets".
- [22] ITU-T Recommendation P.501: "Test signals for use in telephonometry".
- [23] ITU-T Recommendation P.502: "Objective test methods for speech communication systems using complex test signals".
- [24] ITU-T Recommendation P.581: "Use of head and torso simulator (HATS) for hands-free terminal testing".
- [25] ITU-T Recommendation P.862: "Perceptual evaluation of speech quality (PESQ): An objective method for end-to-end speech quality assessment of narrow-band telephone networks and speech codecs".
- [26] IEC 61260: "Electroacoustics - Octave-band and fractional-octave-band filters".
- [27] ISO 3 (1973): "Preferred numbers - Series of preferred numbers".

---

## 3 Definitions and abbreviations

### 3.1 Definitions

For the purposes of the present document, the following terms and definitions apply:

**artificial ear:** device for the calibration of earphones incorporating an acoustic coupler and a calibrated microphone for the measurement of the sound pressure and having an overall acoustic impedance similar to that of the median adult human ear over a given frequency band

**codec:** combination of an analogue-to-digital encoder and a digital-to-analogue decoder operating in opposite directions of transmission in the same equipment

**ear-Drum Reference Point (DRP):** point located at the end of the ear canal, corresponding to the ear-drum position

**freefield reference point:** point located in the free sound field, at least in 1,5 m distance from a sound source radiating in free air (in case of a head and torso simulator [HATS] in the center of the artificial head with no artificial head present)

**freefield equalization:** artificial head is equalized in such a way that for frontal sound incidence in anechoic conditions the frequency response of the artificial head is flat

**Head And Torso Simulator (HATS) for telephonometry:** manikin extending downward from the top of the head to the waist, designed to simulate the sound pick-up characteristics and the acoustic diffraction produced by a median human adult and to reproduce the acoustic field generated by the human mouth

**Mouth Reference Point (MRP):** is located on axis and 25 mm in front of the lip plane of a mouth simulator

**nominal setting of the volume control:** when a receive volume control is provided, the setting which is closest to the nominal RLR of 2 dB

## 3.2 Abbreviations

For the purposes of the present document, the following abbreviations apply:

CSS	Composite Source Signal
D	D-value of terminal
DRP	ear Drum Reference Point
EL	Echo Loss
ELR	Echo Loudness Rating
HATS	Head And Torso Simulator
MOS-LQOy	Mean Opinion Score - Listening Quality Objective

NOTE: See ITU-T Recommendation P.800.1.

MRP	Mouth Reference Point
NLP	Non Linear Processor
PCM	Pulse Code Modulation
PESQ™	Perceptual Evaluation of Speech Quality™
PLC	Packet Loss Concealment
POI	Point Of Interconnect
PSTN	Public Switched Telephone Network
QoS	Quality of Service
RLR	Receive Loudness Rating
SLR	Send Loudness Rating
STMR	SideTone Masking Rating
TCLw	Terminal Coupling Loss (weighted)
TOSQA	Telecommunication Objective Speech Quality Assessment
TCN	Trace Control for Netem

---

## 4 General considerations

### 4.1 Coding Algorithm

The assumed coding algorithm is according to ITU-T Recommendation G.722 [10]. VoIP terminals may support other coding algorithms.

NOTE: Associated Packet Loss Concealment, e.g. as defined in ITU-T Recommendation G.722 [10] Appendix 3 and 4 should be used.

### 4.2 End-to-end considerations

In order to achieve a desired end-to-end speech transmission performance (mouth-to-ear) it is recommended that the general rules of transmission planning are carried out with the E-model of Recommendation G.107 [5] taking into account that the E-model does not yet address wideband transmission planning; this includes the a-priori determination of the desired category of speech transmission quality as defined in ITU-T Recommendation G.109 [7].



While, in general, the transmission characteristics of single circuit-oriented network elements, such as switches or terminals can be assumed to have a single input value for the planning tasks of ITU-T Recommendation G.108 [6], this approach is not applicable in packet based systems and thus there is a need for the transmission planner's specific attention.

In particular the decision as to which delay measured according to the present standard should be acceptable or representative for the specific configuration is the responsibility of the individual transmission planner.

ITU-T Recommendation G.108 [6] with its Amendments provides further guidance on this important issue.

The following optimum terminal parameters from a users' perspective need to be considered:

- Minimized delay in send and receive direction.
- Optimum loudness Rating (RLR, SLR).
- Compensation for network delay variation.
- Packet loss recovery performance.
- Maximized terminal coupling loss.

## 4.3 Parameters to be investigated

### 4.3.1 Basic parameters

The basic parameters are based on I-ETS 300 245-5 [4].

### 4.3.2 Further Parameters with respect to Speech Processing Devices

For VoIP terminals that contain non-linear speech processing devices, the following parameters require additional attention in the context of the present Standard:

- Objective evaluation of speech quality for VoIP terminals.
- Doubletalk capability.
- Time-variant impairments:
  - Switching behaviour.
  - Partial echo effects.
  - Occurrence of artefacts.
  - Clock accuracy.
- Background noise performance of the terminal.
- etc.

The measurements of these further parameters with respect to speech processing devices which are a novelty to terminal requirement standards have been successfully used in the ETSI VoIP speech quality test events TR 102 648-1 (see bibliography).

---

## 5 Test equipment

### 5.1 IP half channel measurement adaptor

The IP half channel measurement adaptor is described in EG 202 425 [3].

## 5.2 Network impairment simulation

At least one set of requirements is based on the assumption of an error free packet network, and at least one other set of requirements is based on a defined simulated malperformance of the packet network.

An appropriate network simulator has to be used, for example NISTnet [<http://snad.ncsl.nist.gov/itg/nistnet/>] or Netem [[tcn.hypert.net](http://tcn.hypert.net)].

Based on the positive experience, STQ have made during the ETSI Speech Quality Test Events with "NIST Net" this will be taken as a basis to express and describe the variations of packet network parameters for the appropriate tests.

Here is a brief blurb about NIST Net:

The NIST Net network emulator is a general-purpose tool for emulating performance dynamics in IP networks. The tool is designed to allow controlled, reproducible experiments with network performance sensitive/adaptive applications and control protocols in a simple laboratory setting. By operating at the IP level, NIST Net can emulate the critical end-to-end performance characteristics imposed by various wide area network situations (e.g. congestion loss) or by various underlying subnetwork technologies (e.g. asymmetric bandwidth situations of xDSL and cable modems).

NIST Net is implemented as a kernel module extension to the Linux operating system and an X Window System-based user interface application. In use, the tool allows an inexpensive PC-based router to emulate numerous complex performance scenarios, including: tunable packet delay distributions, congestion and background loss, bandwidth limitation, and packet reordering / duplication. The X interface allows the user to select and monitor specific traffic streams passing through the router and to apply selected performance "effects" to the IP packets of the stream. In addition to the interactive interface, NIST Net can be driven by traces produced from measurements of actual network conditions. NIST Net also provides support for user defined packet handlers to be added to the system. Examples of the use of such packet handlers include: time stamping / data collection, interception and diversion of selected flows, generation of protocol responses from emulated clients.

The key points of Netem can be summarized as follows:

Netem is nowadays part of most Linux distributions, it only has to be switched on, when compiling a kernel. With netem, there are the same possibilities as with nistnet, there can be generated loss, duplication, delay and jitter (and the distribution can be chosen during runtime). Netem can be run on a Linux-PC running as a bridge or a router (Nistnet only runs on routers).

With an amendment of netem, Trace Control for Netem (TCN) which was developed by ETH Zurich, it is even possible, to control the behaviour of single packets via a trace file. So it is for example possible to generate a single packet loss, or a specific delay pattern. This amendment is planned to be included in new Linux kernels, nowadays it is available as a patch to a specific kernel and to the iproute2 tool (iproute2 contains netem).

It is not advised to define specific distortion patterns for testing in standards, because it will be to easy adapt devices to this patterns (as it is already done for test signals). But if a pattern is unknown to a manufacturer, the same pattern can be used by a test lab for different devices and gives comparable results. It is also possible to take a trace of Nistnet distortions, generate a file out of this and playback exact the same distortions with Netem.

---

## 6 Acoustic environment

In general two possible approaches need to be taken into account: EITHER room noise and background noise are an inherent part of the test environment OR room noise and background noise shall be eliminated to such an extent that their influence on the test results can be neglected.

All measurements shall be conducted under quiet and "anechoic" conditions. Depending on the distance of the transducers from mouth and ear a quiet office room may be sufficient e.g. for handsets where artificial mouth and artificial ear are located close to the acoustical transducers.

However, for some headsets or handset terminals with smaller dimension an anechoic room will be required.

In cases where real or simulated background noise is used as part of the testing environment, the original background noise must not be noticeably influenced by the acoustical properties of the room.

In all cases where the performance of acoustic echo cancellers shall be tested a realistic room which represents the typical user environment for the terminal shall be used.

## 7 Requirements and associated Measurement Methodologies

NOTE 1: In general the test methods as described in the present document apply. If alternative methods exist they may be used if they have been proven to give the same result as the method described in the present document.

NOTE 2: Due to the time variant nature of IP connections delay variation may impair the measurements. In such cases the measurement has to be repeated until a valid measurement result is achieved.

### 7.1 Test setup

The preferred acoustical access to terminals is the most realistic simulation of the "average" subscriber. This can be made by using Head And Torso Simulator (HATS) with appropriate ear simulation and appropriate means to fix handset and headset terminals in a realistic and reproducible way to the HATS. HATS is described in ITU-T Recommendation P.58 [17], appropriate ears are described in ITU-T Recommendation P.57 [16] (type 3.3 and type 3.4 ear), a proper positioning of handsets under realistic conditions is to be found in ITU-T Recommendation P.64 [18].

The preferred way of testing a terminal is to connect it to a network simulator with exact defined settings and access points. The test sequences are fed in either electrically, using a reference codec or using the direct signal processing approach or acoustically using ITU-T specified devices.

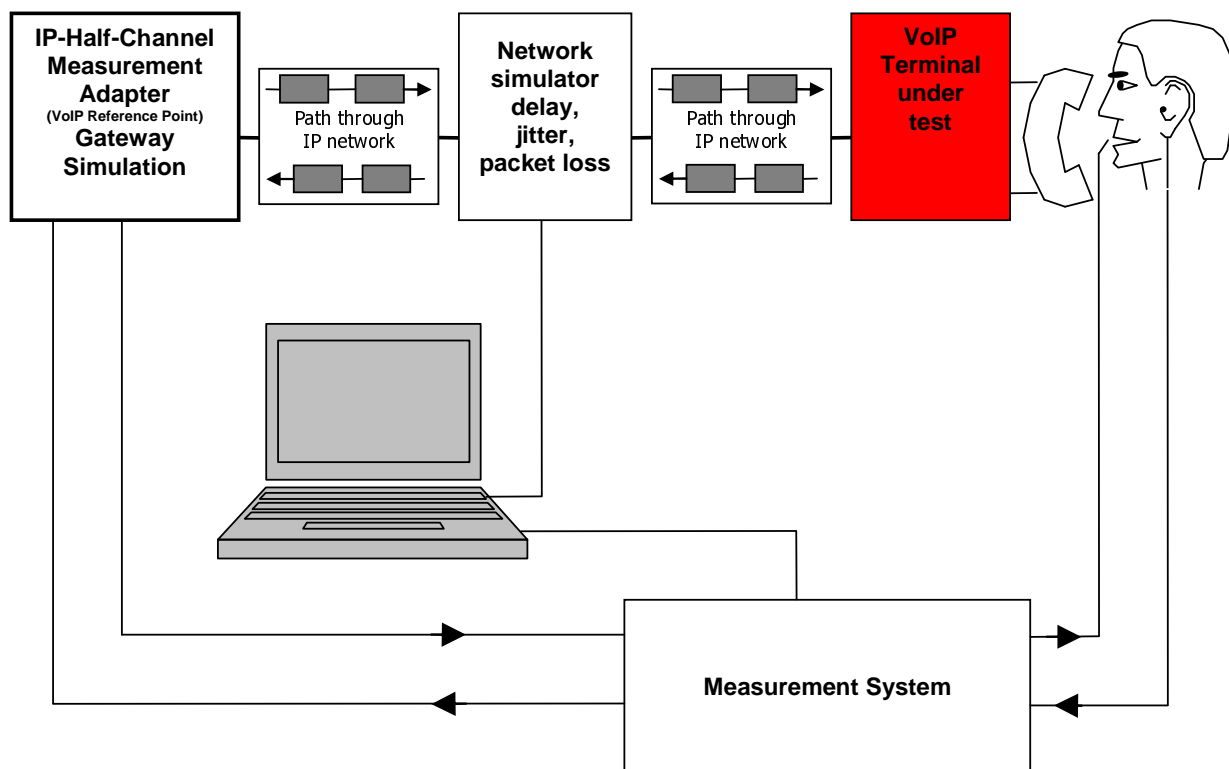


Figure 1: Half channel terminal measurement

### Setup for handsets and headsets

When using a handset telephone the handset is placed in the HATS position as described in ITU-T Recommendation P.64 [18]. The artificial mouth shall conform with ITU-T Recommendation P.58 [17]. The artificial ear shall conform with ITU-T Recommendation P.57 [16], type 3.3 or type 3.4 ears shall be used.

Recommendations for positioning headsets are given in ITU-T Recommendation P.380 [21]. If not stated otherwise headsets shall be placed in their recommended wearing position. Further information about setup and the use of HATS can be found in ITU-T Recommendation P.380 [21].

Unless stated otherwise if a volume control is provided the setting is chosen such that the nominal RLR is met as close as possible.

### Position and calibration of HATS

All the sending and receiving characteristics shall be tested with the HATS, it shall be indicated what type of ear was used at what application force. For handsets if not stated otherwise 8N application force shall be used

The horizontal positioning of the HATS reference plane shall be guaranteed within  $\pm 2^\circ$ .

The HATS shall be equipped with two type 3.3 or type 3.4 artificial ears. For binaural headsets two artificial ears are required. The type 3.3 or type 3.4 artificial ears as specified in ITU-T Recommendation P.57 [16] shall be used. The artificial ear shall be positioned on HATS according to ITU-T Recommendation P.58 [17].

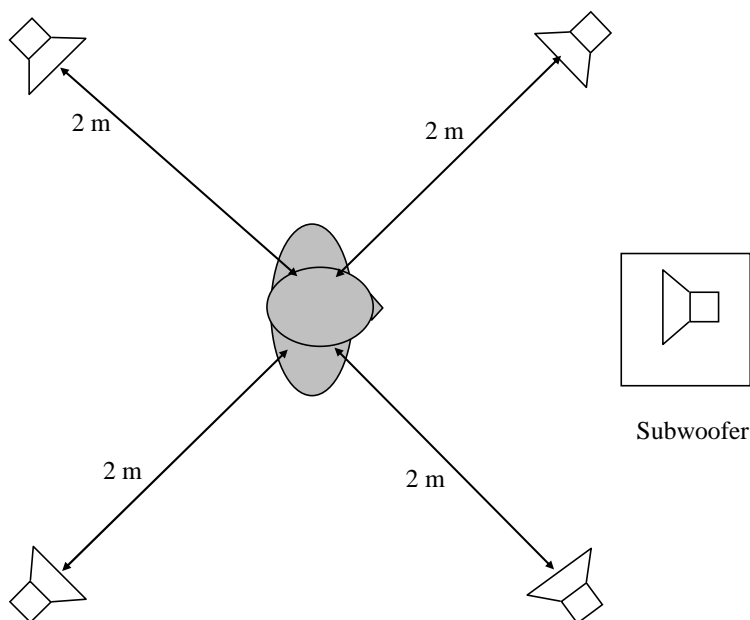
The exact calibration and equalization can be found in ITU-T Recommendation P.581 [24].

### Setup of background noise simulation

A setup for simulating realistic background noises in a lab-type environment is described in EG 202 396-1 [2].

The EG 202 396-1 [2] contains a description of the recording arrangement for realistic background noises, a description of the setup for a loudspeaker arrangement suitable to simulate a background noise field in a lab-type environment and a database of realistic background noises, which can be used for testing the terminal performance with a variety of different background noises.

The principle loudspeaker setup for the simulation arrangement is shown in figure 2.



**Figure 2: Loudspeaker arrangement for background noise simulation**

The equalization and calibration procedure for the setup is described in detail in EG 202 396-1 [2].

If not stated otherwise this setup is used in all measurements where background noise simulation is required.

The following noises of EG 202 396-1 [2] shall be used:

Recording in pub	Pub_Noise_binaural	30 s	L: 77,8 dB(A) R: 78,9 dB(A)	binaural
Recording at sales counter	Cafeteria_Noise_binaural	30 s	L: 68,4 dB(A) R: 67,3 dB(A)	Binaural
Recording in business office	Work_Noise_Office_Callcenter_binaural	30 s	L: 56,6 dB(A) R: 57,8 dB(A)	Binaural

## 7.2 Coding independent Parameters

### 7.2.1 Send Frequency response

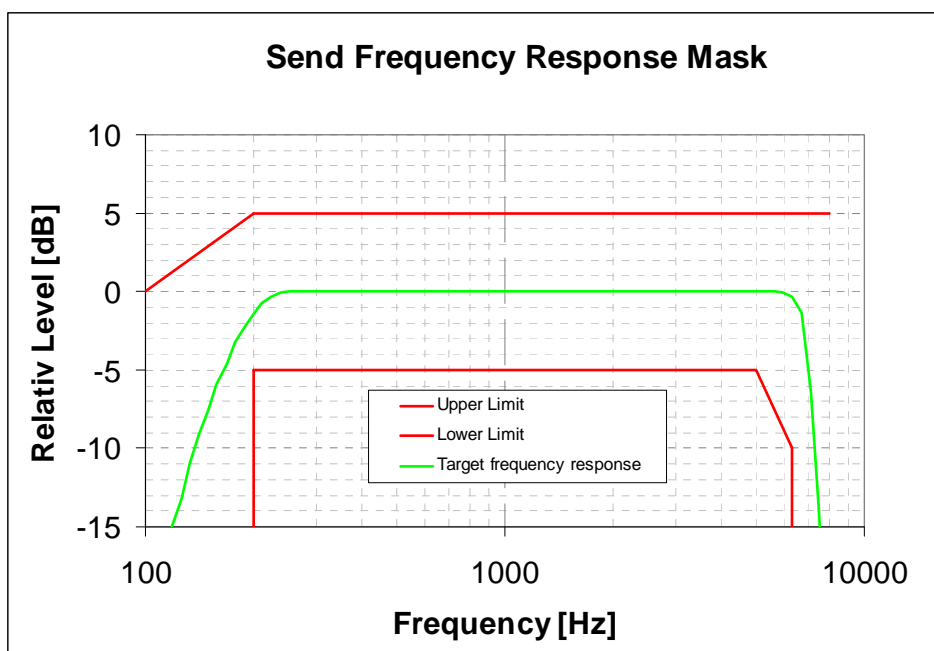
#### Requirement

The send frequency response of the handset or the headset shall be within a mask as defined in table 1 and shown in figure 3. This mask shall be applicable for all types of handsets and headsets.

**Table 1**

Frequency	Upper Limit	Lower Limit
100 Hz	0 dB	
200 Hz	5 dB	-5 dB
5 000 Hz	5 dB	-5 dB
6 300 Hz	5 dB	-10 dB
8 000 Hz	5 dB	

NOTE: The limits for intermediate frequencies lie on a straight line drawn between the given values on a linear (dB) – logarithmic (Hz) scale.



**Figure 3: Send frequency response mask**

NOTE 1: The basis for the target frequency responses in sending and receiving is the orthotelephonic reference response which is measured between 2 subjects in 1 m distance under free field conditions and is assuming an ideal receive characteristic. Under these conditions the overall frequency response shows a rising slope. In opposite to other standards the present document no longer uses the ERP as the reference point for receiving but the free-field. With the concept of free-field based receive measurements a rising slope for the overall frequency response is achieved by a flat target frequency response in sending and a free field based receiving frequency response.

NOTE 2: A "balanced" frequency response is preferable from the perception point of view. If frequency components in the low frequency domain are attenuated in a similar way frequency components in the high frequency domain should be attenuated.

### Measurement Method

The test signal to be used for the measurements shall be the artificial voice according to ITU-T Recommendation P.50 [14]. If the signal to noise ratio in the high frequency domain is not sufficient Composite Source Signal (CSS) as defined in ITU-T Recommendation P.501 [22] shall be used. The spectrum of acoustic signal produced by the artificial mouth is calibrated under free field conditions at the MRP. The test signal level shall be -4,7 dBPa, duration 20 s (10 s female, 10 s male voice), measured at the MRP. The test signal level is averaged over the complete test signal sequence.

The handset terminal is setup as described in clause 7.1. The handset is mounted in the HATS position (see ITU-T Recommendation P.64 [18]). The application force used to apply the handset against the artificial ear is noted in the test report.

In case of headset measurements the tests are repeated 5 times, in conformance with ITU-T Recommendation P.380 [21]. The results are averaged (averaged value in dB, for each frequency).

Measurements shall be made at one twelfth-octave intervals as given by the R.40 series of preferred numbers in ISO 3 [27] for frequencies from 100 Hz to 8 kHz inclusive. For the calculation the averaged measured level at the electrical reference point for each frequency band is referred to the averaged test signal level measured in each frequency band at the MRP.

The sensitivity is expressed in terms of dBV/Pa.

## 7.2.2 Send Loudness Rating

### Requirement

The nominal value of Send Loudness Rating (SLR) shall be:

- $SLR(\text{set}) = 8 \text{ dB} \pm 3 \text{ dB}$

### Measurement Method

The test signal to be used for the measurements shall be the artificial voice according to ITU-T Recommendation P.50 [14], duration 20 s (10 s female, 10 s male voice). If the signal to noise ratio in the high frequency domain is not sufficient CSS as defined in ITU-T Recommendation P.501 [22] shall be used. The spectrum of acoustic signal produced by the artificial mouth is calibrated under free field conditions at the MRP. The test signal level shall be -4,7 dBPa, measured at the MRP. The test signal level is averaged over the complete test signal sequence.

The handset or headset terminal is setup as described in clause 7.1. The handset is mounted in the HATS position (see ITU-T Recommendation P.64 [18]). The application force used to apply the handset against the artificial ear is noted in the test report.

In case of headset measurements the tests are repeated 5 times, in conformance with ITU-T Recommendation P.380 [21]. The results are averaged (averaged value in dB, for each frequency).

The sending sensitivity shall be calculated from each band of the 20 frequencies given in table 1 of ITU-T Recommendation P.79 [19], bands 1 to 20. For the calculation the averaged measured level at the electrical reference point for each frequency band is referred to the averaged test signal level measured in each frequency band at the MRP.

The sensitivity is expressed in terms of dBV/Pa and the SLR shall be calculated according to ITU-T Recommendation P.79 [19], see annex A.

## 7.2.3 D- Factor

### Requirement

For VoIP terminals the D- factor shall be:

- D-factor (DelSM)  $\geq 2$  dB

NOTE: Wideband calculation is for further study, provisionally the measurement is based on narrowband.

### Measurement Method

The background noise simulation as described in clause 7.1 is used.

Handset or headset terminals are mounted as described in clause 7.1. Measurements are made on one-third octave bands according to IEC 61260 [26] for the 14 bands centered at 200 Hz to 4 kHz (bands 4 to 17). For each band the diffuse sound sensitivity  $S_{si}(\text{diff})$  is measured. The sensitivity shall be expressed in terms of dBV/Pa.

The direct sound field sensitivity  $S_{si}(\text{direct})$  is measured as described in clause 7.2.2 (SLR).

The D value according to ITU-T Recommendation P.79 [19], annex E, formula E2 and E3 is calculated in bands 4 to 17. The coefficients  $K_i$  as described in table E1 are used.

The direct sound sensitivity shall be measured using the test set-up specified in clause 7.1 and a speech like test signal as defined in ITU-T Recommendation P.50 [14] or P.501 [22]. The type of test signal used shall be stated in the test report. The direct sound sensitivity is measured in one-third octave bands according to IEC 61260 [26] for the 14 bands centered at 200 Hz to 4 kHz (bands 4 to 17). For each band the direct sound sensitivity  $S_{si}(\text{direct})$  is measured. The sensitivity shall be expressed in terms of dBV/Pa.

The value of the D-factor shall be calculated according to ITU-T Recommendation P.79 [19], annex E, formulas E2 and E3, over the bands from 4 to 17, using the coefficients  $K_i$  from table E1 of ITU-T Recommendation P.79 [19].

## 7.2.4 Linearity Range for SLR

### Requirement

The sensitivity determined with input sound pressure levels between -24,7 dBPa and 5,3 dBPa shall not differ by more than  $\pm 2$  dB from the sensitivity determined with an input sound pressure level of -4,7 dBPa. For the input sound pressure level of 5,3 dBPa a limit of +4/-2 dB applies.

Table 2

Linearity range of SLR: $\Delta\text{SLR} = \text{SLR} - \text{SLR}@-4,7 \text{ dBPa}$			
Input Level	Target $\Delta\text{SLR}$	Upper limit	Lower limit
-24,7 dBPa	0	2 dB	-2 dB
-19,7 dBPa	0	2 dB	-2 dB
-14,7 dBPa	0	2 dB	-2 dB
-9,7 dBPa	0	2 dB	-2 dB
-4,9 dBPa	0	2 dB	-2 dB
-4,7 dBPa	0	0 dB	0 dB
-4,5 dBPa	0	2 dB	-2 dB
0,3 dBPa	0	2 dB	-2 dB
5,3 dBPa	0	4 dB	-4 dB

NOTE: It is assumed that the variation of gain is mostly codec independent. In case codec specific requirements are needed this is found in the codec specific section.

### Measurement Method

The test signal to be used for the measurements shall be the artificial voice according to ITU-T Recommendation P.50 [14]. If the signal to noise ratio in the high frequency domain is not sufficient CSS as defined in ITU-T Recommendation P.501 [22] shall be used. The spectrum of acoustic signal produced by the artificial mouth is calibrated under free field conditions at the MRP. The test signal levels shall be -24,7 dBPa up to 5,3 dBPa in steps of 5 dB, measured at the MRP. The test signal level is averaged over the complete test signal sequence.

The handset terminal is setup as described in clause 7.1. The handset is mounted in the HATS position (see ITU-T Recommendation P.64 [18]). The application force used to apply the handset against the artificial ear is noted in the test report.

The sending sensitivity shall be calculated from each band of the 20 frequencies given in table 1 of ITU-T Recommendation P.79 [19], bands 1 to 20. For the calculation the averaged measured level at the electrical reference point for each frequency band is referred to the averaged test signal level measured in each frequency band at the MRP.

The sensitivity is expressed in terms of dBV/Pa and the SLR shall be calculated according to ITU-T Recommendation P.79 [19], Annex G.

## 7.2.5 Send Distortion

### Requirement

The terminal will be positioned as described in clause 7.1.

The ratio of signal to harmonic distortion shall be above the following mask:

Frequency	Ratio
315 Hz	26 dB
400 Hz	30 dB
1 kHz	30 dB
2 kHz	30 dB
NOTE: Limits at intermediate frequencies lie on a straight line drawn between the given values on a linear (dB ratio) - logarithmic (frequency) scale.	

### Measurement method

The terminal will be positioned as described in clause 7.1.

The signal used is an activation signal followed by a sine-wave signal with a frequency at 315 Hz, 400 Hz, 500 Hz, 630 Hz, 800 Hz, 1 000 Hz and 2 000 Hz. The duration of the sine wave shall be less than 1 s. The sinusoidal signal level shall be calibrated to -4,7 dBPa at the MRP.

The signal to harmonic distortion ratio is measured selectively up to 6,3 kHz.

An artificial voice according to ITU-Recommendation P.50 [14] or a speech like test signal as described in ITU-T Recommendation P.501 [22] can be used for activation. Level of this activation signal will be -4,7 dBPa at the MRP.

NOTE: Depending on the type of codec the test signal used may need to be adapted.

## 7.2.6 Send Noise

### Requirement

The maximum noise level produced by the VoIP terminal at the POI under silent conditions in the sending direction shall not exceed -68 dBm0 (C).

No peaks in the frequency domain higher than 10 dB above the average noise spectrum shall occur.

### Measurement Method

For the actual measurement no test signal is used. In order to reliably activate the terminal an activation signal is introduced before the actual measurement. The activation signal shall be a sequence of 4 composite source signals (CSS) as described in ITU-T Recommendation P.501 [22]. The spectrum of the acoustic signal produced by the artificial mouth is calibrated under free field conditions at the MRP. The activation signal level shall be -4,7 dBPa, measured at the MRP. The activation signal level is averaged over the complete activation signal sequence. Alternatively other speech like test signals (e.g. artificial voice) with the same signal level can be used for activation.

The handset terminal is set-up as described in clause 7.1. The handset is mounted at the HATS position (see ITU-T Recommendation P.64 [18]).



The send noise is measured at the POI in the frequency range from 100 Hz to 8 kHz. The analysis window is applied directly after stopping the activation signal but taking into account the influence of all acoustical components (reverberations). The averaging time is 1 s. The test house has to ensure (e.g. by monitoring the time signal) that during the test the terminal remains in activated condition. If the terminal is deactivated during the measurement, the measurement time has to be reduced to the period where the terminal remains in activated condition.

The noise level is measured in dBm0(C).

## 7.2.7 Sidetone Masking Rating STMR (Mouth to ear)

### Requirement

The STMR shall be 16 dB  $\pm$  4 dB for nominal setting of the volume control.

For all other positions of the volume control, the STMR must not be below 8 dB.

NOTE: It is preferable to have a constant STMR independent of the volume control setting.

### Measurement Method

The test signal to be used for the measurements shall be the artificial voice according to ITU-T Recommendation P.50 [14]. The spectrum of the acoustic signal produced by the artificial mouth is calibrated under free field conditions at the MRP. The test signal level shall be -4,7 dBPa, measured at the MRP. The test signal level is averaged over the complete test signal sequence.

The handset or the headset terminal is setup as described in clause 7.1. The handset is mounted in the HATS position (see ITU-T Recommendation P.64 [18]) and the application force shall be 13N on the artificial ear type 3.3 or type 3.4.

Where a user operated volume control is provided, the measurements shall be carried out the nominal setting of the volume control. In addition the measurement is repeated at the maximum volume control setting.

Measurements shall be made at one twelfth-octave intervals as given by the R.40 series of preferred numbers in ISO 3 [27] for frequencies from 100 Hz to 8 kHz inclusive. For the calculation the averaged measured level at each frequency band (ITU-T Recommendation P.79 [19], table 3, bands 1 to 20) is referred to the averaged test signal level measured in each frequency band.

The Sidetone path loss (LmeST), as expressed in dB, and the SideTone Masking Rate (STMR) (in dB) shall be calculated from the formula 5-1 of ITU-T Recommendation P.79 [19], using  $m = 0,225$  and the weighting factors of in table 3 of ITU-T Recommendation P.79 [19].

## 7.2.8 Sidetone delay

### Requirement

The maximum sidetone-round-trip delay shall be  $\leq 5$  ms, measured in an echo-free setup.

### Measurement Method

The handset or the headset terminal is setup as described in clause 7.1. The handset is mounted in the HATS position (see ITU-T Recommendation P.64 [18]).

The test signal is a CS-signal complying with ITU-T Recommendation P.501 [22] using a pn sequence with a length of 4 096 points (for the 48 kHz sampling rate) which equals to the period T. The duration of the complete test signal is as specified in ITU-T Recommendation P.501 [22]. The level of the signal shall be -4,7 dBPa at the MRP.

The cross-correlation function  $\Phi_{xy}(\tau)$  between the input signal  $S_x(t)$  generated by the test system in send direction and the output signal  $S_y(t)$  measured at the artificial ear is calculated in the time domain:

$$\Phi_{xy}(\tau) = \lim_{T \rightarrow \infty} \sum_{t=-T/2}^{T/2} S_x(t) S_y(t + \tau) \quad (1)$$

The measurement window T shall be exactly identical with the time period T of the test signal, the measurement window is positioned to the pn-sequence of the test signal.

The sidetone delay is calculated from the envelope  $E(\tau)$  of the cross-correlation function  $\Phi_{xy}(\tau)$ . The first maximum of the envelope function occurs in correspondence with the direct sound produced by the artificial mouth, the second one occurs with a possible delayed sidetone signal. The difference between the two maxima corresponds to the sidetone delay. The envelope  $E(\tau)$  is calculated by the Hilbert transformation  $H\{xy(\tau)\}$  of the cross-correlation:

$$H\{xy(\tau)\} = \sum_{-\infty}^{\infty} \frac{\Phi_{xy}(u)}{\Pi(\tau - u)} \quad (2)$$

$$E(\tau) = \sqrt{[\Phi_{xy}(\tau)]^2 + \{H[\Phi_{xy}(\tau)]\}^2} \quad (3)$$

It is assumed that the measured sidetone delay is less than  $T/2$ .

## 7.2.9 Terminal Coupling Loss weighted (TCLw)

### Requirement

The TCLw shall be  $\geq 55$  dB.

With the volume control set to maximum TCLw shall be  $\geq 46$  dB. The volume control shall be set back to nominal after each call unless  $\text{TCLw} \geq 55$  dB can be maintained also with maximum volume setting.

### Measurement Method

The handset or headset terminal is setup as described in clause 7.1. The handset is mounted in the HATS position (see ITU-T Recommendation P.64 [18]) and the application force shall be 2N on the artificial ear type 3.3 or type 3.4 as specified in ITU-T Recommendation P.57 [16]. The ambient noise level shall be less than -64 dBPa(A) for handset and headset terminals. The attenuation from electrical reference point input to electrical reference point output shall be measured using a speech like test signal.

Before the actual test a training sequence consisting of 10 s artificial voice male and 10 s artificial voice female according to ITU-T Recommendation P.50 [14] is altered. The training sequence level shall be -16 dBm0 in order not to overload the codec.

The test signal is a PN-sequence complying with ITU-T Recommendation P.501 [22] with a length of 4 096 points (for the 48 kHz sampling rate) and a crest factor of 6 dB. The length of the complete test signal composed of at least four sequences of CSS shall be at least one second (1,0 s). The test signal level is -3 dBm0 (from 50 Hz to 7 kHz). The low crest factor is achieved by random alternation of the phase between  $-180^\circ$  and  $180^\circ$ .

The TCLw is calculated according to ITU-T Recommendation G.122 [8], clause B.4 (trapezoidal rule) but using the frequency range of 300 to 6 700 Hz (instead of 300 to 3 400 Hz). For the calculation the averaged measured echo level at each frequency band is referred to the averaged test signal level measured in each frequency band. For the measurement a time window has to be applied adapted to the duration of the actual test signal (200 ms).

NOTE: The extension of the frequency range is for further study.

## 7.2.10 Stability Loss

### Requirement

With the handset lying on and the transducers facing a hard surface, the attenuation from the digital input to the digital output shall be at least 6 dB at all frequencies in the range of 100 Hz to 8 kHz. In case of headsets the requirement applies for the closest possible position between microphone and headset receiver.

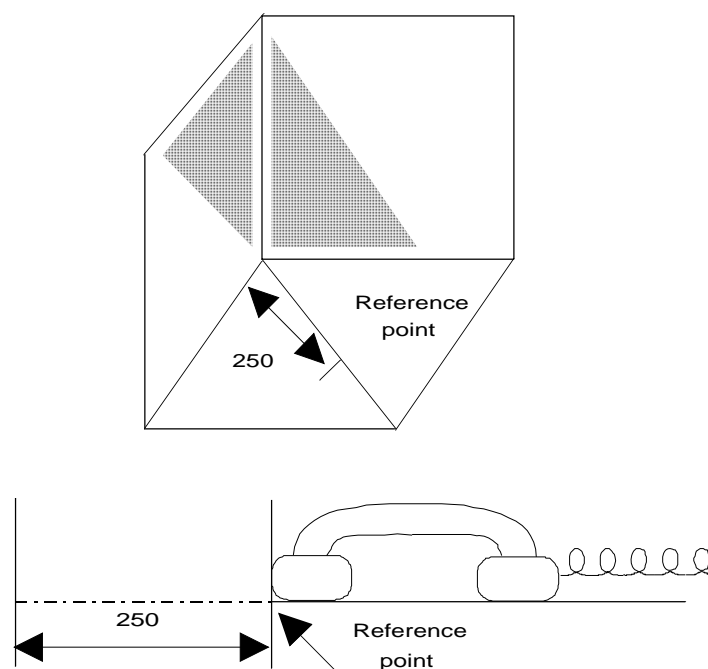
NOTE: Depending on the type of headset it may be necessary to repeat the measurement in different positions.

### Measurement Method

Before the actual test a training sequence consisting of 10 s artificial voice male and 10 s artificial voice female according to ITU-T Recommendation P.50 [14] is altered. The training sequence level shall be -16 dBm0 in order not to overload the codec.

The test signal is a PN sequence complying with ITU-T Recommendation P.501 [22] with a length of 4 096 points (for the 48 kHz sampling rate) and a crest factor of 6 dB. The duration of the test signal is 250 ms. With an input signal of -3 dBm0, the attenuation from digital input to digital output shall be measured for frequencies from 100 Hz to 8 kHz under the following conditions:

- a) the handset or the headset, with the transmission circuit fully active, shall be positioned on one inside surface that is of three perpendicular plane, smooth, hard surfaces forming a corner. Each surface shall extend 0,5 m from the apex of the corner. One surface shall be marked with a diagonal line, extending from the corner formed by the three surfaces, and a reference position 250 mm from the corner, as shown in figure 4;
- b1) the handset, with the transmission circuit fully active, shall be positioned on the defined surface as follows:
  - 1) the mouthpiece and ear cup shall face towards the surface;
  - 2) the handset shall be placed centrally, the diagonal line with the ear cup nearer to the apex of the corner;
  - 3) the extremity of the handset shall coincide with the normal to the reference point, as shown in figure 4.
- b2) the headset, with the transmission circuit fully active, shall be positioned on the defined surface as follows:
  - 1) the microphone and the receiver shall face towards the surface;
  - 2) the headset receiver shall be placed centrally at the reference point as shown in figure 4;
  - 3) the headset microphone is positioned as close as possible to the receiver.



NOTE: All dimensions in mm.

**Figure 4**

## 7.2.11 Receive Frequency response

### Requirement

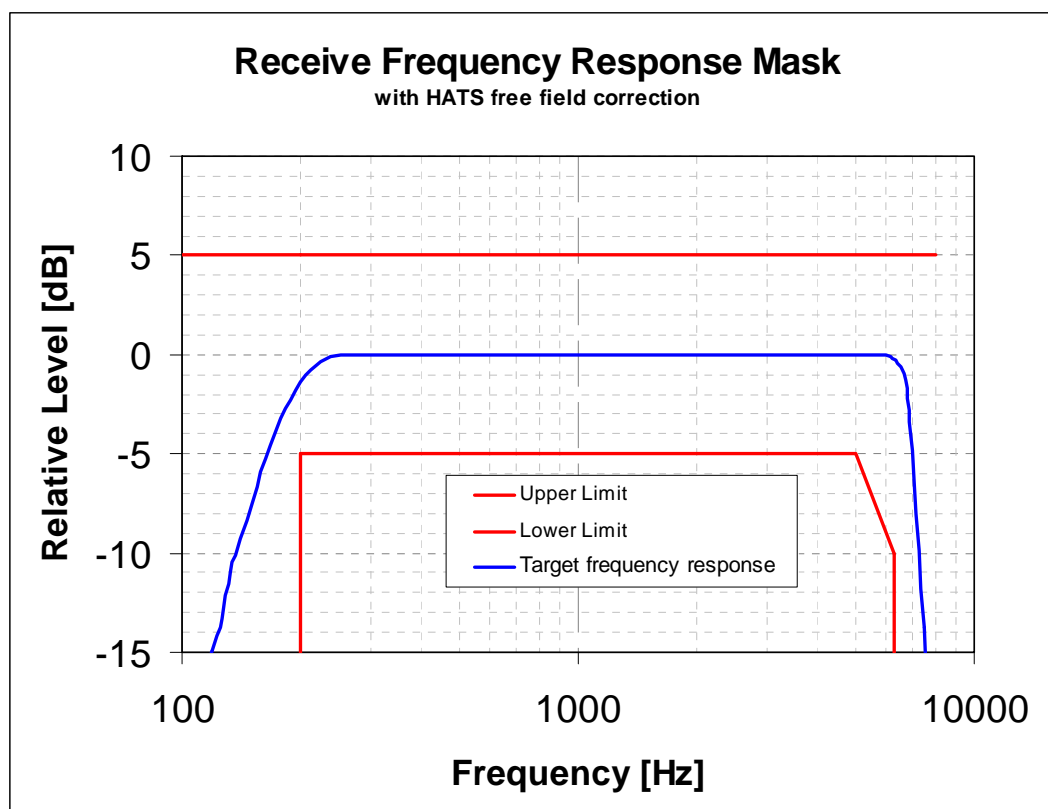
The receive frequency response of the handset or the headset shall be within a mask as defined in table 3 and shown in figure 5. The application force for handsets is 2N, 8N and 13N. This mask defined for 8 N application force shall be applicable for all types of headsets.

**Table 3: Receive Frequency Response Mask**

Frequency	Upper Limit 8N	Lower Limit 8N	Upper Limit 2N and 13N	Lower Limit 2N and 13N
100 Hz	4 dB		6 dB	
200 Hz	4 dB	-8 dB	6 dB	-10 dB
300 Hz	4 dB	-4 dB	6 dB	-6 dB
5 000 Hz	4 dB	-4 dB	6 dB	-6 dB
6 300 Hz	4 dB		6 dB	
8 000 Hz	4 dB			

NOTE 1 : The limit curves shall be determined by straight lines joining successive co-ordinates given in the table, where frequency response is plotted on a linear dB scale against frequency on a logarithmic scale. is a floating or "best fit" mask

NOTE 2: The basis for the target frequency responses in sending and receiving is the orthotelephonic reference response which is measured between 2 subjects in 1 m distance under free field conditions and is assuming an ideal receive characteristic. Under these conditions the overall frequency response shows a rising slope. In opposite to other standards the present document no longer uses the ERP as the reference point for receiving but the free-field. With the concept of free-field based receive measurements a rising slope for the overall frequency response is achieved by a flat target frequency response in sending and a freefield based receiving frequency response.



**Figure 5: Receive frequency response mask for 8N application force**

NOTE: A "balanced" frequency response is preferable from the perception point of view. If frequency components in the low frequency domain are attenuated in a similar way frequency components in the high frequency domain should be attenuated.

### Measurement Method

Receive frequency response is the ratio of the measured sound pressure and the input level.  
(dB relative Pa/V)

$$S_{\text{Jeff}} = 20 \log (p_{e\text{ff}} / v_{\text{RCV}}) \text{ dB rel 1 Pa / V} \quad (4)$$

$S_{\text{Jeff}}$	Receive Sensitivity; Junction to HATS Ear with free field correction.
$p_{e\text{ff}}$	DRP Sound pressure measured by ear simulator Measurement data are converted from the Drum Reference Point to free field.
$v_{\text{RCV}}$	Equivalent RMS input voltage.

The test signal to be used for the measurements shall be the artificial voice according to ITU-T Recommendation P.50 [14], duration 20 s (10 s female, 10 s male voice). If the signal to noise ratio in the high frequency domain is not sufficient CSS as defined in ITU-T Recommendation P.501 [22] shall be used. The test signal level shall be -16 dBm<sub>0</sub>, measured according to ITU-T Recommendation P.56 [15] at the digital reference point or the equivalent analogue point.

The handset terminal or the headset terminal is setup as described in clause 7.1. The handset is mounted in the HATS position (see ITU-T Recommendation P.64 [18]). The application forces used to apply the handset against the artificial ear is 2N, 8N and 13N.

In case of headset measurements the tests are repeated 5 times, in conformance with ITU-T Recommendation P.380 [21]. The results are averaged (averaged value in dB, for each frequency).

The HATS is free field equalized as described in ITU-T Recommendation P.581 [24]. The equalized output signal is power-averaged on the total time of analysis. The 1/12 octave band data are considered as the input signal to be used for calculations or measurements.

Measurements shall be made at one twelfth-octave intervals as given by the R.40 series of preferred numbers in ISO 3 [27] for frequencies from 100 Hz to 8 kHz inclusive. For the calculation the averaged measured level at each frequency band is referred to the averaged test signal level measured in each frequency band.

The sensitivity is expressed in terms of dBPa/V.

## 7.2.12 Receive Loudness Rating

### Requirement

The nominal value of Receive Loudness Rating (RLR) shall be:

- RLR(set) = 2 dB ± 3 dB
- RLR (binaural headset) = 8 dB ± 3 dB for each earphone

### Measurement Method

The test signal to be used for the measurements shall be the artificial voice according to ITU-T Recommendation P.50 [14], duration 20 s (10 s female, 10 s male voice). If the signal to noise ratio in the high frequency domain is not sufficient CSS as defined in ITU-T Recommendation P.501 [22] shall be used. The test signal level shall be -16 dBm<sub>0</sub>, measured at the digital reference point or the equivalent analogue point. The test signal level is averaged over the complete test signal sequence.

The handset terminal or the headset terminal is setup as described in clause 7.1. The handset is mounted in the HATS position (see ITU-T Recommendation P.64 [18]). The application force used to apply the handset against the artificial ear is noted in the test report. The HATS is **NOT** freefield equalized as described in ITU-T Recommendation P.581 [24]. The DRP-ERP correction as defined in ITU-T Recommendation P.57 [16] is applied. The application force used to apply the handset against the artificial ear is noted in the test report. By default, 8N will be used.

In case of headset measurements the tests are repeated 5 times, in conformance with ITU-T Recommendation P.380 [21]. The results are averaged (averaged value in dB, for each frequency).

The receiving sensitivity shall be calculated from each band of the 20 frequencies given in table 1 of ITU-T Recommendation P.79 [19], bands 1 to 20. For the calculation the averaged measured level at each frequency band is referred to the averaged test signal level measured in each frequency band.

The sensitivity is expressed in terms of dBPa/V and the RLR shall be calculated according to ITU-T Recommendation P.79 [19], annex A. No leakage correction shall be applied for the measurement.

## 7.2.13 Receiving Distortion

### Requirement

The ratio of signal to harmonic distortion shall be above the following mask:

**Table 4**

Frequency	Signal to distortion ratio limit, receiving
315 Hz	26 dB
400 Hz	30 dB
500 Hz	30 dB
800 Hz	30 dB
1 kHz	30 dB
2 kHz	30 dB
NOTE:	Limits at intermediate frequencies lie on a straight line drawn between the given values on a linear (dB ratio) - logarithmic (frequency) scale.

### Measurement Method

The handset terminal or the headset terminal is positioned as described in clause 7.1.

The signal used is an activation signal followed by a sine-wave signal with a frequency at 315 Hz, 400 Hz, 500 Hz, 630 Hz, 800 Hz, 1 000 Hz and 2 000Hz.

An artificial voice according to ITU-Recommendation P.50 [14] or a speech like test signal as described in ITU-T Recommendation P.501 [22] can be used for activation.

The signal level shall be -16 dBm0.

Measurement are made at 315 Hz, 400 Hz, 500 Hz, 630 Hz, 800 Hz, 1 000 Hz and 2 000 Hz.

The signal to harmonic distortion ratio is measured selectively up to 6,3 kHz.

The ratio of signal to harmonic distortion shall be measured at the DRP of the artificial ear with the free field equalization active.

NOTE: Depending on the type of codec the test signal used may need to be adapted.

## 7.2.14 Minimum activation level and sensitivity in Receive direction

For further study.

## 7.2.15 Receive Noise

### Requirement

Telephone sets with adjustable receive levels shall be adjusted so that the RLR is as close as possible to the nominal RLR.

The receive noise shall be less than -57 dBPa(A).

Where a volume control is provided, the measured noise shall not be greater than -54 dBPa(A) at the maximum setting of the volume control.

### Measurement Method

The handset terminal or the headset terminal is setup as described in clause 7.1.

An artificial voice according to ITU-Recommendation P.50 [14] or a speech like test signal as described in ITU-T Recommendation P.501 [22] can be used for activation. The activation signal level shall be -16 dBm0.

The A-weighted noise level shall be measured at DRP of the artificial ear with the free field equalization active.

## 7.2.16 Automatic Gain Control in Receiving

For further study.

## 7.2.17 Double talk Performance

During double talk the speech is mainly determined by 2 parameters: impairment caused by echo during double talk and level variation between single and double talk (attenuation range).

In order to guarantee sufficient quality under double talk conditions the talker Echo Loudness Rating (ELR) should be high and the attenuation inserted should be as low as possible. Terminals which do not allow double talk in any case should provide a good echo attenuation which is realized by a high attenuation range in this case.

The most important parameters determining the speech quality during double talk are (see ITU-T Recommendations P.340 [20] and P.502 [23]):

- Attenuation range in sending direction during double talk  $A_{H,S,dt}$ .
- Attenuation range in receiving direction during double talk  $A_{H,R,dt}$ .
- Echo attenuation during double talk.

### 7.2.17.1 Attenuation Range in Sending Direction during Double Talk $A_{H,S,dt}$

#### Requirement

Based on the level variation in sending direction during double talk  $A_{H,S,dt}$  the behaviour of the terminal can be classified according to table 5.

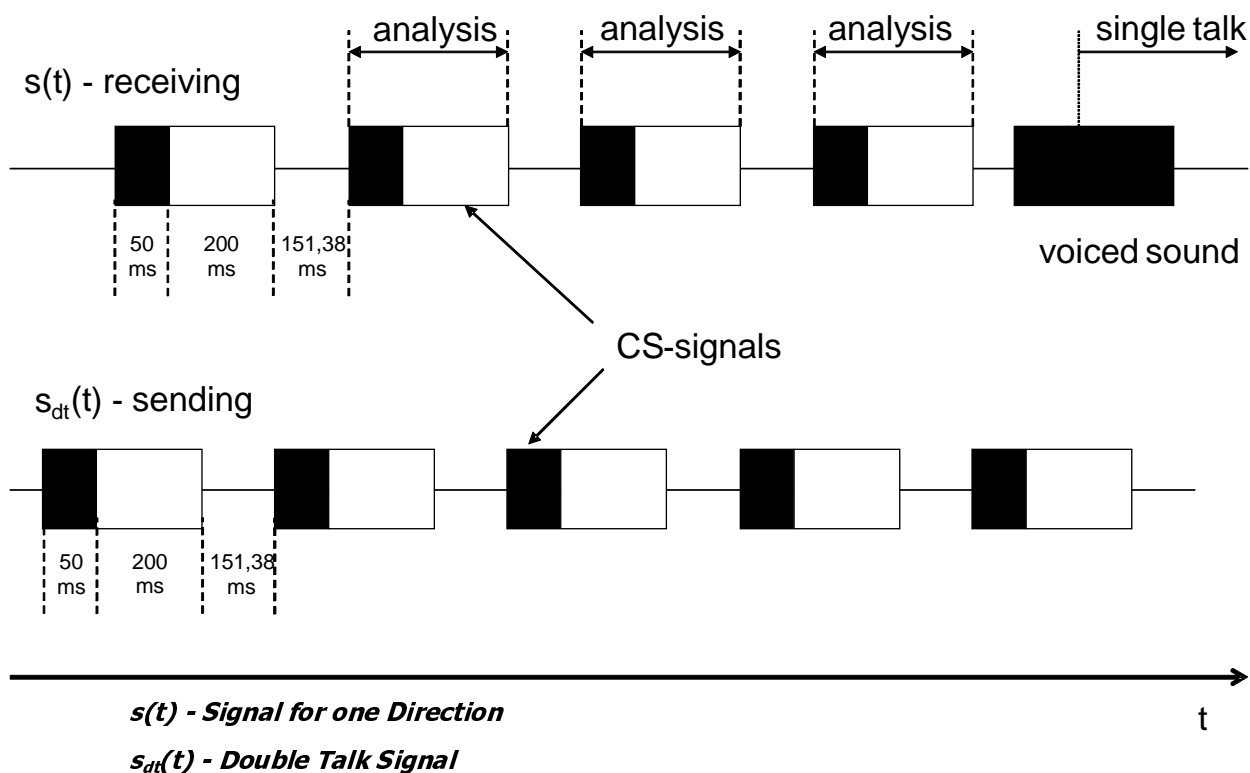
**Table 5**

Category (according to ITU-T Rec. P.340 [20])	1	2a	2b	2c	3
	<i>Full Duplex Capability</i>	<i>Partial Duplex Capability</i>			<i>No Duplex Capability</i>
$A_{H,S,dt}$ [dB]	$\leq 3$	$\leq 6$	$\leq 9$	$\leq 12$	$> 12$

In general this table provides a quality classification of terminals regarding double talk performance. However, this does not mean that a terminal which is category 1 based on the double talk performance is of high quality concerning the overall quality as well.

#### Measurement Method

The test signal to determine the attenuation range during double talk is shown in figure 6. A sequence of uncorrelated CS signals is used which is inserted in parallel in sending and receiving direction.



**Figure 6: Double Talk Test Sequence with overlapping CS signals in sending and receiving direction**

Figure 6 indicates that the sequences overlap partially. The beginning of the CS sequence (voiced sound, black) is overlapped by the end of the pn-sequence (white) of the opposite direction. During the active signal parts of one signal the analysis can be conducted in sending and receiving direction. The analysis times are shown in figure 6 as well. The test signals are synchronized in time at the acoustical interface. The delay of the test arrangement should be constant during the measurement.

The settings for the test signals are as follows:

**Table 6**

	Receiving Direction	Sending Direction
<b>Pause Length between two Signal Bursts</b>	151,38 ms	151,38 ms
<b>Average Signal Level (Assuming an Original Pause length of 101,38 ms)</b>	-16 dBm0	-4,7 dBPa
<b>Active Signal Parts</b>	-14,7 dBm0	-3 dBPa

The test arrangement is according to clause 7.

When determining the attenuation range in sending direction the signal measured at the electrical reference point is referred to the test signal inserted.

The level is determined as level vs. time from the time domain. The integration time of the level analysis is 5 ms. The attenuation is determined from the level difference measured at the beginning of the double talk always with the beginning of the CS-signal in sending direction until its complete activation (during the pause in the receiving channel). The analysis is performed over the complete signal starting with the second CS-signal. The first CS-signal is not used for the analysis.



### 7.2.17.2 Attenuation Range in Receiving Direction during Double Talk $A_{H,R,dt}$

#### Requirement

Based on the level variation in receiving direction during double talk  $A_{H,R,dt}$  the behaviour of the terminal can be classified according to table 7.

**Table 7**

Category (according to ITU-T Rec. P.340 [20])	1	2a	2b	2c	3
	Full Duplex Capability	Partial Duplex Capability			No Duplex Capability
$A_{H,R,dt}$ [dB]	$\leq 3$	$\leq 5$	$\leq 8$	$\leq 10$	$> 10$

In general this table provides a quality classification of terminals regarding double talk performance. However, this does not mean that a terminal which is category 1 based on the double talk performance is of high quality concerning the overall quality as well.

#### Measurement Method

The test signal to determine the attenuation range during double talk is shown in figure 6. A sequence of uncorrelated CS signals is used which is inserted in parallel in sending and receiving direction. The test signals are synchronized in time at the acoustical interface. The delay of the test arrangement should be constant during the measurement.

The settings for the test signals are as follows:

**Table 8**

	Receiving Direction	Sending Direction
<b>Pause Length between two Signal Bursts</b>	151,38 ms	151,38 ms
<b>Average Signal Level (Assuming an Original pause Length of 101,38 ms)</b>	-16 dBm0	-4,7 dBPa
<b>Active Signal Parts</b>	-14,7 dBm0	-3 dBPa

The test arrangement is according to clause 7.

When determining the attenuation range in receiving direction the signal measured at the artificial ear referred to the test signal inserted.

The level is determined as level vs. time from the time domain. The integration time of the level analysis is 5 ms. The attenuation is determined from the level difference measured at the beginning of the double talk always with the beginning of the CS-signal in receiving direction until its complete activation (during the pause in the sending channel). The analysis is performed over the complete signal starting with the second CS-signal. The first CS-signal is not used for the analysis.

### 7.2.17.3 Detection of Echo Components during Double Talk

#### Requirement

Echo Loss during double talk is the echo suppression provided by the terminal during double talk measured at the electrical reference point.

NOTE: The echo attenuation during double talk is based on the parameter Talker Echo Loudness Rating (TELRdt). It is assumed that the terminal at the opposite end of the connection provides nominal Loudness Rating (SLR + RLR = 10 dB).

Under these conditions the requirements given in the table below are applicable (more information can be found in annex A of the ITU-T Recommendation P.340 [20]).

Table 9

Category (according to ITU-T Rec. P.340 [20])	1	2a	2b	2c	3
	Full Duplex Capability	Partial Duplex Capability			No Duplex Capability
<b>Echo Loss [dB]</b>	$\geq 27$	$\geq 23$	$\geq 17$	$\geq 11$	$< 11$

### Measurement Method

The test arrangement is according to clause 7.1.

The double talk signal consists of a sequence of orthogonal signals which are realized by voice-like modulated sine waves spectrally shaped similar to speech. The measurement signals used are shown in figure 7. A detailed description can be found in ITU-T Recommendation P.501 [22].

The signals are fed simultaneously in sending and receiving direction. The level in sending direction is -4,7 dBPa at the MRP (nominal level), the level in receiving direction is -16 dBm0 at the electrical reference point (nominal level).

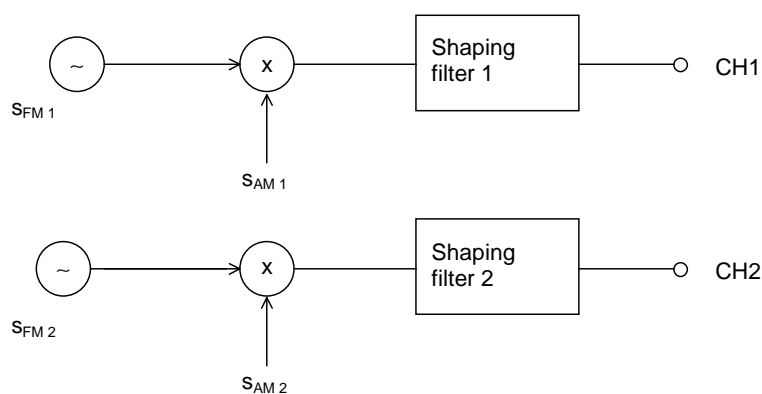


Figure 7: Measurement signals

$$s_{FM1,2}(t) = \sum A_{FM1,2} * \cos(2\pi t n * F_{01,2}) ; n= 1,2,\dots \quad (5)$$

$$s_{AM1,2}(t) = A_{AM1,2} * \cos(2\pi t F_{AM1,2}); \quad (6)$$

The settings for the signals are as follows:

**Table 10: Parameters of the two Test Signals for Double Talk Measurement based on AM-FM modulated sine waves**

Receiving Direction			Sending Direction			
$f_m$ [Hz]	$f_{\text{mod}(fm)}$ [Hz]	$F_{\text{am}}$ [Hz]		$f_m$ [Hz]	$f_{\text{mod}(fm)}$ [Hz]	$F_{\text{am}}$ [Hz]
250	±5	3		270	±5	3
500	±10	3		540	±10	3
750	±15	3		810	±15	3
1 000	±20	3		1 080	±20	3
1 250	±25	3		1 350	±25	3
1 500	±30	3		1 620	±30	3
1 750	±35	3		1 890	±35	3
2 000	±40	3		2 160	±35	3
2 250	±40	3		2 400	±35	3
2 500	±40	3		2 900	±35	3
2 750	±40	3		3 150	±35	3
3 000	±40	3		3 400	±35	3
3 250	±40	3		3 650	±35	3
3 500	±40	3		3 900	±35	3
3 750	±40	3				

NOTE: Parameters of the Shaping Filter: Low Pass Filter, 5 dB/oct.

The test signal is measured at the electrical reference point (sending direction). The measured signal consists of the double talk signal which was fed in by the artificial mouth and the echo signal. The echo signal is filtered by comb filter using mid-frequencies and bandwidth according to the signal components of the signal in receiving direction (see ITU-T Recommendation P.501 [22]). The filter will suppress frequency components of the double talk signal.

In each frequency band which is used in receiving direction the echo attenuation can be measured separately. The requirement for category 1 is fulfilled if in any frequency band the echo signal is either below the signal noise or below the required limit. If echo components are detectable, the classification is based on the table above. The echo attenuation is to be achieved for **each individual frequency band** according to the different categories.

#### 7.2.17.4 Minimum activation level and sensitivity of double talk detection

For further study.

### 7.2.18 Switching characteristics

NOTE: Additional requirements may be needed in order to further investigate the effect of NLP implementations on the users' perception of speech quality.

#### 7.2.18.1 Activation in Sending Direction

The activation in sending direction is mainly determined by the built-up time  $T_{r,S,\min}$  and the minimum activation level ( $L_{S,\min}$ ). The minimum activation level is the level required to remove the inserted attenuation in sending direction during idle mode. The built-up time is determined for the test signal burst which is applied with the minimum activation level.

The activation level described in the following is always referred to the test signal level at the Mouth Reference Point (MRP).

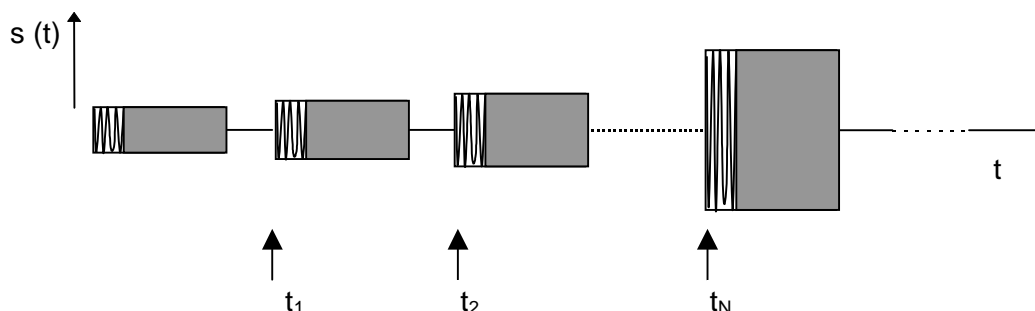
##### Requirement

The minimum activation level  $L_{S,\min}$  shall be  $\leq -20$  dBPa.

The built-up time  $T_{r,S,\min}$  (measured with minimum activation level) should be  $\leq 15$  ms.

### Measurement Method

The structure of the test signal is shown in figure 8. The test signal consists of CSS components according to ITU-T Recommendation P.501 [22] with increasing level for each CSS burst.



**Figure 8: Test Signal to Determine the Minimum Activation Level and the Built-up Time**

The settings of the test signal are as follows:

**Table 11**

	CSS Duration/ Pause Duration	Level of the first CS Signal (active Signal Part at the MRP)	Level Difference between two Periods of the Test Signal
CSS to Determine Switching Characteristic in Sending Direction	~250 ms / ~450 ms	-23 dBPa (see note)	1 dB
NOTE: The level of the active signal part corresponds to an average level of -24,7 dBPa at the MRP for the CSS according to ITU-T Recommendation P.501 [22] assuming a pause of about 100 ms.			

It is assumed that the pause length of about 450 ms is longer than the hang-over time so that the test object is back to idle mode after each CSS burst.

The test arrangement is described in clause 7.1.

The level of the transmitted signal is measured at the electrical reference point. The measured signal level is referred to the test signal level and displayed vs. time. The levels are calculated from the time domain using an integration time of 5 ms.

The minimum activation level is determined from the CSS burst which indicates the first activation of the test object. The time between the beginning of the CSS burst and the complete activation of the test object is measured.

NOTE: If the measurement using the CS-Signal does not allow to clearly identify the minimum activation level, the measurement may be repeated by using a one syllable word instead of the CS-Signal. The word used should be of similar duration, the average level of the word should be adapted to the CS-signal level of the according CS-burst.

### 7.2.18.2 Silence Suppression and Comfort Noise Generation

For further study.

### 7.2.18.3 Performance in Sending in the Presence of Background Noise

#### Requirement

The level of comfort noise shall be within in a range of +2 and -5 dB compared to the original (transmitted) background noise. The noise level is calculated with psophometric weighting.

NOTE 1: It is advisable that the comfort noise matches the original signal as good as possible (from a perceptual point of view).

NOTE 2: Input for further specification necessary (e.g. on temporal matching).

The spectral difference between comfort noise and original (transmitted) background noise shall be within the mask given through straight lines between the breaking points on a logarithmic (frequency) - linear (dB sensitivity) scale as given in table 12.

**Table 12: Requirements for Spectral Adjustment of Comfort Noise (Mask)**

Frequency	Upper Limit	Lower Limit
200 Hz	12 dB	-12 dB
800 Hz	12 dB	-12 dB
800 Hz	10 dB	-10 dB
2 000 Hz	10 dB	-10 dB
2 000 Hz	6 dB	-6 dB
4 000 Hz	6 dB	-6 dB
8 000 Hz	6 dB	-6 dB
NOTE: All sensitivity values are expressed in dB on an arbitrary scale.		

### Measurement Method

The background noise simulation as described in clause 7.1 is used.

The handset terminal is set-up as described in clause 7.1. The handset is mounted at the HATS position (see ITU-T Recommendation P.64 [18]).

First the background noise transmitted in send is recorded at the POI for a period of at least 20 s.

In a second step a test signal is applied in receiving direction consisting of an initial pause of 10 s and a periodical repetition of the Composite Source Signal in receiving direction (duration 10 s) with nominal level to enable comfort noise injection simultaneously with the background noise. For the measurement the background noise sequence has to be started at the same point as it was started in the previous measurement. Alternatively other speech like test signals (e.g. artificial voice) with the same signal level can be used.

The transmitted signal is recorded in sending direction at the POI.

The power density spectra measured in sending direction without far end speech simulation averaged between 10 s and 20 s is referred to the power density spectrum measured in sending direction determined during the period with far end speech simulation in receiving direction averaged between 10 s and 20 s. Level and spectral differences between both power density spectra are analysed and compared to the requirements.

### 7.2.18.4 Speech Quality in the Presence of Background Noise

For further study, taking into account EG 202 396-3.

### 7.2.18.5 Quality of Background Noise Transmission (with Far End Speech)

#### Requirement

The test is carried out applying the Composite Source Signal in receiving direction. During and after the end of Composite Source Signal bursts (representing the end of far end speech simulation) the signal level in sending direction should not vary more than 10 dB (during transition to transmission of background noise without far end speech). The measurement is conducted for all types of background noise as defined in clause 7.1.

#### Measurement Method

The test arrangement is according to clause 7.1.

The background noises are generated as described in clause 7.1.

First the measurement is conducted without inserting the signal at the far end. At least 10 s of noise is analysed. The background signal level versus time is calculated using a time constant of 35 ms. This is the reference signal.

In a second step the same measurement is conducted but with inserting the CS-signal at the far end. The exactly identical background noise signal is applied. The background noise signal must start at the same point in time which was used for the measurement without far end signal. The background noise should be applied for at least 5 s in order to allow adaptation of the noise reduction algorithms. After at least 5 s a Composite Source Signal according to ITU-T Recommendation P.501 [22] is applied in receiving direction with a duration of  $\geq 2$  CSS periods. The test signal level is -16 dBm0 at the electrical reference point.

The sending signal is recorded at the electrical reference point. The test signal level versus time is calculated using a time constant of 35 ms.

The level variation in sending direction is determined during the time interval when the CS-signal is applied and after it stops. The level difference is determined from the difference of the recorded signal levels vs. time between reference signal and the signal measured with far end signal.

### 7.2.18.6 Quality of background noise transmission (with Near End Speech)

#### **Requirement**

The test is carried out applying a simulated speech signal in sending direction. During and after the end of the simulated speech signal (Composite Source Signal bursts) the signal level in sending direction should not vary more than 10 dB.

#### **Measurement Method**

The test arrangement is according to clause 7.1.

The background noises are generated as described in clause 7.1. The background noise should be applied for at least 5 s in order to allow adaptation of the noise reduction algorithms.

The near end speech is simulated using the Composite Source Signal according to ITU-T Recommendation P.501 [22] with a duration of  $\geq 2$  CSS periods. The test signal level is -4,7 dBPa at the MRP.

The sending signal is recorded at the electrical reference point. The test signal level versus time is calculated using a time constant of 35 ms.

First the measurement is conducted without inserting the signal at the near end. The signal level is analysed vs. time. In a second step the same measurement is conducted but with inserting the CS-signal at the near end. The level variation is determined by the difference between the background noise signal level without inserting the CS-signal and the maximum level of the noise signal during and after the CS-bursts in sending direction.

### 7.2.19 Quality of echo cancellation

#### 7.2.19.1 Temporal echo effects

##### **Requirement**

This test is intended to verify that the system will maintain sufficient echo attenuation during single talk. The measured echo attenuation during single talk should not decrease by more than 6 dB from the maximum measured during the TCLw test.

##### **Measurement Method**

The test arrangement is according to clause 7.1.

The test signal consists of periodically repeated Composite Source Signal according to ITU-T Recommendation P.501 [22] with an average level of -5 dBm0 as well as an average level of -25 dBm0. The echo signal is analysed during a period of at least 2,8 s which represents 8 periods of the CS signal. The integration time for the level analysis shall be 35 ms, the analysis is referred to the level analysis of the reference signal.

The measurement result is displayed as attenuation vs. time. The exact synchronization between input and output signal has to be guaranteed.

NOTE 1: In addition tests with more speech like signals should be made, e.g. ITU-T Recommendation P.50 [14] to see time variant behaviour of EC.

NOTE 2: The analysis is conducted only during the active signal part, the pauses between the Composite Source Signals are not analysed. The analysis time is reduced by the integration time of the level analysis (35 ms).

## 7.2.19.2 Spectral Echo Attenuation

### Requirement

The echo attenuation vs. frequency shall be below the tolerance mask given in table 13.

**Table 13: Echo attenuation limits**

Frequency	Limit
100 Hz	-20 dB
200 Hz	-30 dB
300 Hz	-38 dB
800 Hz	-34 dB
1 500 Hz	-33 dB
2 600 Hz	-24 dB
4 000 Hz	-24 dB
8 000 Hz	-24 dB

NOTE 1: All sensitivity values are expressed in dB on an arbitrary scale.  
 NOTE 2: The limit at intermediate frequencies lies on a straight line drawn between the given values on a log (frequency) - linear (dB) scale.

During the measurement it should be ensured that the measured signal is really the echo signal and not the Comfort Noise which possibly may be inserted in sending direction in order to mask the echo signal.

### Measurement Method

The test arrangement is according to clause 7.1.

Before the actual measurement a training sequence is fed in consisting of 10 s CS signal according to ITU-T Recommendation P.501 [22]. The level of the training sequence is -16 dBm0.

The test signal consists of a periodically repeated Composite Source Signal. The measurement is carried out under steady-state conditions. The average test signal level is -16 dBm0, averaged over the complete test signal. 4 CS signals including the pauses are used for the measurement which results in a test sequence length of 1,4 s. The power density spectrum of the measured echo signal is referred to the power density spectrum of the original test signal. The analysis is conducted using FFT analysis with 8 k points (48 kHz sampling rate, Hanning window).

The spectral echo attenuation is analysed in the frequency domain in dB.

## 7.2.19.3 Occurrence of Artefacts

For further study.

## 7.2.20 Variant Impairments; Network Dependant

For further study.

### 7.2.20.1 Delay versus Time Send

For further study.

### 7.2.20.2 Delay versus Time Receive

For further study.

### 7.2.20.3 Quality of Jitter buffer adjustment

For further study.

## 7.3 Codec Specific Requirements

### 7.3.1 Send Delay

For a VoIP terminal, send delay is defined as the one-way delay from the acoustical input (mouthpiece) of this VoIP terminal to its interface to the packet based network. The total send delay is the upper bound on the mean delay and takes into account the delay contributions of all of the elements shown in figures 2 and A.1 in ITU-T Recommendation G.1020 [13], respectively.

The sending delay  $T(s)$  is defined as follows:

$$T(s) = T(ps) + T(la) + T(rif) + T(asp) \quad (\text{Formula 1})$$

Where:

$T(ps)$  = packet size =  $N * T(fs)$

$N$  = number of frames per packet

$T(fs)$  = frame size of encoder

$T(la)$  = look-ahead of encoder

$T(aif)$  = air interface framing

$T(asp)$  = allowance for signal processing

The additional delay required for IP packet assembly and presentation to the underlying link layer will depend on the link layer. When the link layer is a LAN (e.g. Ethernet), this additional time will usually be quite small. For the purposes of the present document it is assumed that in the test setup this delay can be neglected.

NOTE 1: The size of  $T(aif)$  is for further study.

#### Requirement

The allowance for signal processing shall be  $T(asp) < 10$  ms

NOTE 2: With the knowledge of the codec specific values for  $T(fs)$  and  $T(la)$  the requirements for send delay for any type of coder and any frame size  $T(fs)$  can easily be calculated by formula 1. Table 14 provides requirements calculated accordingly for frequently used codecs and packet sizes.

**Table 14**

Codec	N	T(fs) in ms	T(ps) in ms	T(la) in ms	T(aif) in ms	T(asp) in ms	T(s) Requirement in ms
<b>G.722 [10]</b>	80	0,0625	10	0	0	10	< 20,0625
	160	0,0625	20	0	0	10	< 30,0625
<b>G.722.1 [11]</b>	1	20	10	5	0	10	To be completed
	2	20	20	5	0	10	To be completed
<b>L16-256</b>	160	0,0625	10	0	0	10	< 20,0625

Further information about the different sources of delay for different codecs can be found in annex A.

#### Measurement Method

The test signal to be used for the measurements shall be a Composite Source Signal (CSS) as described in ITU-T Recommendation P.501 [22]. The test signal consists of the voiced part as described in ITU-T Recommendation P.501 [22] followed by a pseudo random noise sequence with a periodicity of minimum 500 ms. The spectrum of acoustic signal produced by the artificial mouth is calibrated under free field conditions at the MRP. The test signal level shall be -4,7 dBPa, measured at the MRP. The test signal level is averaged over the complete test signal sequence.

NOTE 3: If the expected delay is higher than 500 ms a pseudo random sequence with a higher periodicity should be used.



The handset terminal is setup as described in clause 7.1. The handset is mounted in the HATS position (see ITU-T Recommendation P.64 [18]). The application force used to apply the handset against the artificial ear shall be stated in the test report.

The delay is calculated using the cross correlation function between the signal at the electrical test point and the signal at the MRP. The cross correlation analysis has to be chosen in such a way that the maximum delay of 500 ms can be analysed. The measurement is corrected by the delay introduced by the test equipment.

The delay is expressed in ms, determined from the maximum of the cross correlation function.

NOTE 4: Delay may be time variant. Therefore constant monitoring of the actual delay may be required when evaluating the range of delay which can be observed in a given connection. The test setup should take into account either real network conditions or the tools needed to simulate typical causes for time variant delay (e.g. packet loss) during the measurement period. Other methods like running cross correlation or delay estimation procedures e.g. used in PESQ (ITU-T Recommendation P.862 [25]) may be used.

### 7.3.2 Receive delay

For a VoIP terminal, receive delay is defined as the one-way delay from the interface to the packet based network of this VoIP terminal to its acoustical output (earpiece). The total receive delay is the upper bound on the mean delay and takes into account the delay contributions of all of the elements shown in figures 3, A.1 and A.2 of ITU-T Recommendation G.1020 [13], respectively.

The receiving delay  $T(r)$  is defined as follows:

$$T(r) = T(fs) + T(fi) + T(aif) + T(jb) + T(plc) + T(asp) \quad (\text{Formula 2})$$

Where:

$T(fs)$  = frame size of encoder

$T(fi)$  = filter processing delay

$T(aif)$  = air interface framing

$T(jb)$  = jitter buffer size

$T(plc)$  = PLC buffer size

$T(asp)$  = allowance for signal processing

The additional delay required for IP packet dis-assembly and presentation from the underlying link layer will depend on the link layer. When the link layer is a LAN (e.g. Ethernet), this additional time will usually be quite small. For the purposes of the present document it is assumed that in the test setup this delay can be neglected.

NOTE 1: The size of  $T(aif)$  is for further study.

#### Requirements

The allowance for signal processing shall be  $T(asp) < 10$  ms.

The additional delay introduced by the jitter buffer shall be  $T(jb) \leq 10$  ms.

For Coders without integrated PLC the additional PLC buffer size shall be  $T(plc) < 10$  ms.

For Coders with integrated PLC the additional PLC buffer size shall be  $T(plc) = 0$  ms.

NOTE 2: With the knowledge of the codec specific values for  $T(fs)$  and  $T(la)$  the requirements for receive delay for any type of coder and any frame size  $T(fs)$  can easily be calculated by formula 2. Table 15 provides requirements calculated accordingly for some frequently used codecs and packet sizes as an example.

Table 15

Codec	N	T(fs)	T(fi)	T(aif)	T(jb)	T(plc)	T(asp)	T(r) Requirement
G.722 [10]	80	0,0625	0	0	10	10	10	< 30,0625
G.722 [10]	160	0,0625	0	0	10	10	10	< 30,0625
G.722.1 [11]	1	20	0	0	10	0	10	< 40
G.722.1 [11]	2	20	0	0	10	0	10	< 40
L 16-256	160	0,0625	0	0	10	10	10	< 30,0625

NOTE 1: T(ps) = packet size = N \* T(fs).  
NOTE 2: N = number of frames per packet.

NOTE 3: These requirements are based on the lowest possible delay values which can be expected under ideal network conditions. Caution must be exercised to ensure that the terminal is operated under optimum conditions in order to avoid adverse effects, e.g. network conditions, settings and memory effects of the terminal jitter buffer.

### Measurement Method

The test signal to be used for the measurements shall be a Composite Source Signal (CSS) as described in ITU-T Recommendation P.501 [22]. The test signal consists of the voiced part as described in ITU-T Recommendation P.501 [22] followed by a pseudo random noise sequence with a periodicity of minimum 500 ms. The test signal level shall be -16 dBm0, measured at the electrical test point. The test signal level is averaged over the complete test signal sequence.

The handset terminal is setup as described in clause 7.1. The handset is mounted in the HATS position (see ITU-T Recommendation P.64 [18]). The application force used to apply the handset against the artificial ear shall be stated in the test report.

The delay is calculated using the cross correlation function between the signal at the electrical test point and the signal at the DRP. The cross correlation analysis has to be chosen in such a way that the maximum delay of 500 ms can be analysed. The measurement is corrected by the delay introduced by the test equipment.

The delay is expressed in ms, determined from the maximum of the cross correlation function.

NOTE 4: Delay may be time variant. Therefore constant monitoring of the actual delay may be required when evaluating the range of delay which can be observed in a given connection. The test setup should take into account either real network conditions or the tools needed to simulate typical causes for time variant delay (e.g. packet loss) during the measurement period. Other methods like running cross correlation or delay estimation procedures e.g. used in PESQ (ITU-T Recommendation P.862 [25]) may be used.

## 7.3.3 Objective Listening Speech Quality MOS-LQOM in Send direction

The listening speech quality tests are conducted under clean network conditions.

### Requirements

The requirements for the listening speech quality are as follows:

Speech coder	MOS-LQOM
G.722 [10]	> 4,2
G.729.1 @ 32 kbit/s [12]	> 4,5
G.722.1 [11]	> 4,2
L16-256	> 4,5

NOTE: Currently no test method is available for terminals, TOSQA 2001 is one method (EG 201 377-2 [1]), which may be used in half-channel scenarios.

### Measurement method

For further study.

### 7.3.4 Objective Listening Quality MOS-LQOM in Receive direction

The listening speech quality tests are conducted under clean network conditions as well as with network impairments simulated. In addition to the listening speech quality tests the delay is measured.

#### Requirement

The requirement for the listening speech quality and the delay under clean network conditions are as follows:

Speech coder	MOS-LQOM
G.722 [10]	> 4,0
G.722.1 [11]	> 4,0
G.729.1 @32 kbit/s [12]	> 4,2
L16-256	> 4,2

NOTE: The MOS-LQOM requirements in receiving are lower than the requirements set in sending. This takes into account that in receiving the impairment introduced by a non ideal frequency response characteristics in receiving in addition to the impairment introduced by the codec impairment is more dominant than in sending.

#### Test method

For further study.

For the performance tests with network impairments the following settings are used:

**Table 16: Network Conditions for Electrical-Acoustical Measurements (Speech Samples)**

Condition	Packet Loss (Equal)	Delay Variation
0c (see note 2) (VAD)	0	No
1	0	No
2	0	20 ms (see note 1)
3	1%	No
4	1%	20 ms (see note 1)
5	3%	No

NOTE 1: Delay Variation produced with a Pareto-Distribution and  $r = 0,5$ .  
 NOTE 2: VAD on, all other conditions (1-5) tested with VAD off.  
 NOTE 3: For some network emulation tools, it is necessary to introduce a constant delay to offer the possibility to generate a delay variation distribution. This delay has to be subtracted from the measured delay before interpreting the results.

**Table 17: Requirements for ITU-T Recommendation G.722 [10] speech codecs**

Condition	MOS-LQOM	Delay
1	> 4,0	< 30,0625 ms
2	> 3,8	< 50,0625 ms
3	> 3,8	< 30,0625 ms
4	> 3,8	< 50,0625 ms
5	> 3,6	< 30,0625 ms

NOTE: The settings are derived from the ones used in the ETSI Plugtest VoIP speech quality test events.

**Table 18: Requirements for G.722.1 speech codecs**

<b>Condition</b>	<b>MOS-LQOM</b>	<b>Delay</b>
1	> 4,0	< 40 ms
2	> 3,8	< 60 ms
3	> 3,8	< 40 ms
4	> 3,8	< 60 ms
5	> 3,8	< 40 ms

#### 7.3.4.1 Efficiency of Packet Loss Concealment (PLC)

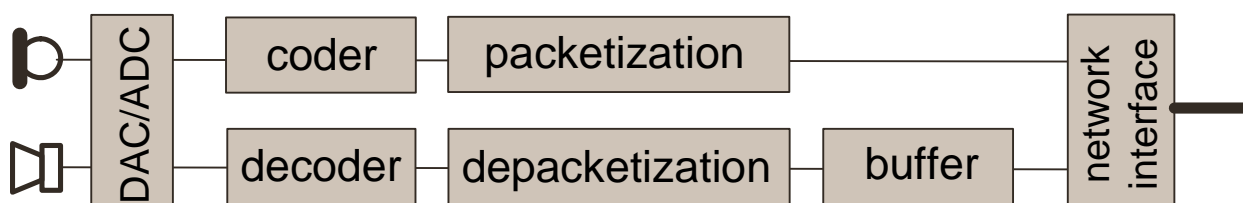
For further study.

#### 7.3.4.2 Efficiency of Delay Variation Removal

For further study.

## Annex A (informative): Processing delays in VoIP terminals

This annex gives some elements about delays generated in VoIP terminals. At first, we consider only wired terminals. These terminals could be schematized as shown in figure A.1.



**Figure A.1: Synoptic of the different functions implemented in a VoIP terminal**

The implemented functions in the sending part of the terminal are:

- The analog-digital conversion.
- The encoding.
- The packetization.
- The interfacing with the network.

The implemented functions in the receiving part of the terminal are:

- The interfacing with the network.
- The depacketization.
- The buffering.
- The decoding.
- The digital-analog conversion.

Let us examine each function's contribution to the processing delay characterizing VoIP terminals.

On the sending part of the terminal, the **network interface** operates the transfer of digital data from IP stack to IP network. At the reception, the network interface operates the transfer of digital data from IP network to IP stack. The network interface has a low contribution to the delay. The contribution is estimated at less than 2 ms per transmission way (sending and receiving direction).

The **packetization** represents the transfer of the audio frames through the IP stack, from the telephony applicative part of the terminal to the transmission network. The packetization consists in adding specific headers (associated to different protocols) to audio frames. The delay associated to the packetization is considered as no significant and included into encoding time.

**Encoding** corresponds to the compression of the speech signal. The delay associated to the encoding process depends on the implemented codec and the payload's length (number of audio frames) inserted into each IP packet. On the sending part of the terminal, encoding is the main contribution to the processing delay. The delay can strongly change according to the codec and the payload's length.

**Analog to digital conversion** consists in transforming speech signal from analog to digital format. The processing delay associated to the conversion is considered as no significant.

**Digital to analog conversion** consists in transforming speech signal from digital to analog format. As analog to digital conversion, the processing delay associated to digital to analog conversion is considered as no significant.

The **depacketization** represents the transfer of the audio frames through the IP stack, from transmission network to the telephony applicative part of the terminal. The depacketization consists in tacking off the headers associated to protocols to get back audio frames after transmission. The delay associated to the depacketization is considered as no significant and included into the decoding processing time.

The first role of the **jitter buffer** is to ensure synchronization between sending and receiving terminals. This synchronization is carried out by buffering the audio frames received from the IP stack before sending them to the decoder. The second role of the jitter buffer is to smooth a possible variation of the transmission time. If synchronization of sending and receiving terminals requires a minimum size of buffer, smoothing transmission delay variation requires a buffer size depending on jitter produced by the network. High variations of transmission time involve an important size of the buffer to smooth jitter. Jitter buffers can be implemented either as buffer with static size(s) (several sizes are possible) or as dynamic buffer. In the last case, size management is carried out according to QoS present on the network interface. Jitter buffer is the main contribution to the processing time on the reception part of VoIP terminal.

**Decoding** corresponds to the rebuilding of speech signal from receiving audio frames. The delay associated to decoding depends on the codec implemented. Decoding contributes in a significant way to the processing time on the reception part of VoIP terminal.

Table A1 presents the processing times of VoIP terminals for different codecs and IP packet payload's lengths.

In this table, x1, x2, x3, x4, y5, x6 and x7 represent the encoding delays according to selected codec. In the same way, y1, y2, y3, y4, y5, y6 and y7 represent the decoding delays according to selected codec.

According to selected codec and payload's length, columns 5 and 6 show overall encoding and decoding delays respectively. Overall encoding time takes into account algorithm, encoding and packetization delays. Overall decoding time takes into account algorithm, decoding and depacketization delays.

Column 7 shows for each codec and payload's length the real time condition. It stands for the maximum duration to encode and decode at the same time. IP terminals have to meet this requirement.

Column 10 shows the minimum delay induced by the jitter buffer. To ensure a correct running of the VoIP terminal, the minimal size of jitter buffer has to correspond to the IP packet payload's length. Furthermore, a double buffering operation induces 10 additional ms in the overall jitter buffer processing.

Column 12 shows the minimum end-to-end delay induced by two terminals connected to a "perfert" network (i.e. with no jitter, no packet loss and with a null transmission delay), with real time condition at the lower limit (i.e. no significant encoding and decoding times).

Column 13 shows the minimum end-to-end delay induced by two terminals connected to a "perfert" network (i.e. with no jitter, no packet loss and with a null transmission delay), with real time condition at the upper limit (i.e. encoding + decoding times very close to the payload size).

Table A.1

Codec	Frame	Lookahead	Payload	Sending processing delay = Algorithm delay + coding and packetization delay	Receiving processing delay = Algorithm delay + coding and packetization delay	Real time condition	Network interface and ADC delay	Network interface and DAC delay	Minimum delay of the jitter buffer	Maximum delay of the jitter buffer	Minimum End to End delay with the lower jitter buffer processing time when real time condition is minimum (x+y=0)	Minimum End to End delay with the lower jitter buffer processing time when real time condition is maximum (x+y=upper limit)	Maximum End to End delay with the higher jitter buffer processing time when real time condition is minimum (x+y=0)	Maximum End to End delay with the higher jitter buffer processing time when real time condition is maximum (x+y=upper limit)
G.711	1	0	10	10+x1	y1	$x1+y1 < 10$ ms	2	2	20	400	34	44	414	424
	1	0	20	$2*(10+x1)$	$2*y1$	$2*(x1+y1) < 20$ ms	2	2	30	400	54	74	424	444
	1	0	30	$3*(10+x1)$	$3*y1$	$3*(x1+y1) < 30$ ms	2	2	40	400	74	104	434	464
	1	0	40	$4*(10+x1)$	$4*y1$	$4*(x1+y1) < 40$ ms	2	2	50	400	94	134	444	484
	1	0	50	$5*(10+x1)$	$5*y1$	$5*(x1+y1) < 50$ ms	2	2	60	400	114	164	454	504
	1	0	60	$6*(10+x1)$	$6*y1$	$6*(x1+y1) < 60$ ms	2	2	70	400	134	194	464	524
G.729	10	5	10	$(10+x2)+5$	y2	$x2+y2 < 10$ ms	2	2	20	400	39	49	419	429
	10	5	20	$(2*(10+x2))+5$	$2*y2$	$2*(x2+y2) < 20$ ms	2	2	30	400	59	79	429	449
	10	5	30	$(3*(10+x2))+5$	$3*y2$	$3*(x2+y2) < 30$ ms	2	2	40	400	79	109	439	469
	10	5	40	$(4*(10+x2))+5$	$4*y2$	$4*(x2+y2) < 40$ ms	2	2	50	400	99	139	449	489
	10	5	50	$(5*(10+x2))+5$	$5*y2$	$5*(x2+y2) < 50$ ms	2	2	60	400	119	169	459	509
	10	5	60	$(6*(10+x2))+5$	$6*y2$	$6*(x2+y2) < 60$ ms	2	2	70	400	139	199	469	529
G.723.1	30	7,5	30	$(30+x3)+7,5$	y3	$x3+y3 < 30$ ms	2	2	40	400	81,5	111,5	441,5	471,5
	30	7,5	60	$(2*(30+x3))+7,5$	$2*y3$	$2*(x3+y3) < 60$ ms	2	2	70	400	141,5	201,5	471,5	531,5
NB-AMR	20	5	20	$(20+x4)+5$	y4	$x4+y4 < 20$ ms	2	2	30	400	59	79	429	449
	20	5	40	$(2*(20+x4))+5$	$2*y4$	$2*(x4+y4) < 40$ ms	2	2	50	400	99	139	449	489
	20	5	60	$(3*(20+x4))+5$	$3*y4$	$3*(x4+y4) < 60$ ms	2	2	70	400	139	199	469	529
G.722	10	1,5	10	$(10+x5)+1,5$	y5	$x5+y5 < 10$ ms	2	2	20	400	35,5	45,5	415,5	425,5
	10	1,5	20	$(2*(10+x5))+1,5$	$2*y5$	$2*(x5+y5) < 20$ ms	2	2	30	400	55,5	75,5	425,5	445,5
	10	1,5	30	$(3*(10+x5))+1,5$	$3*y5$	$3*(x5+y5) < 30$ ms	2	2	40	400	75,5	105,5	435,5	465,5
	10	1,5	40	$(4*(10+x5))+1,5$	$4*y5$	$4*(x5+y5) < 40$ ms	2	2	50	400	95,5	135,5	445,5	485,5
	10	1,5	50	$(5*(10+x5))+1,5$	$5*y5$	$5*(x5+y5) < 50$ ms	2	2	60	400	115,5	165,5	455,5	505,5
	10	1,5	60	$(6*(10+x5))+1,5$	$6*y5$	$6*(x5+y5) < 60$ ms	2	2	70	400	135,5	195,5	465,5	525,5
WB-AMR	20	5	20	$(20+x6)+5$	$y6+0,94$	$x6+y6 < 20$ ms	2	2	30	400	59,94	79,94	429,94	449,94
	20	5	40	$(2*(20+x6))+5$	$2*y6+0,94$	$2*(x6+y6) < 40$ ms	2	2	50	400	99,94	139,94	449,94	489,94
	20	5	60	$(3*(20+x6))+5$	$3*y6+0,94$	$3*(x6+y6) < 60$ ms	2	2	70	400	139,94	199,94	469,94	529,94
G.729.1	20	25	20	$(20+x7)+25+1,97$	$y7+1,97$	$x7+y7 < 20$ ms	2	2	30	400	82,94	102,94	452,94	472,94
	20	25	40	$(2*(20+x7))+25+1,97$	$2*y7+1,97$	$2*(x7+y7) < 40$ ms	2	2	50	400	122,94	162,94	472,94	512,94
	20	25	60	$(3*(20+x7))+25+1,97$	$3*y7+1,97$	$3*(x7+y7) < 60$ ms	2	2	70	400	162,94	222,94	492,94	552,94

---

## Annex B (informative): Bibliography

ETSI TR 102 648-1: "Speech Processing, Transmission and Quality Aspects (STQ); Test Methodologies for ETSI Test Events and Results; Part 1: VoIP Speech Quality Testing".

ITU-T Recommendation P.51: "Artificial mouth".



---

## History

<b>Document history</b>		
V1.1.1	July 2007	Membership Approval Procedure    MV 20070914: 2007-07-17 to 2007-09-14
V1.2.1	October 2007	Publication