

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



Reference

DES/STQ-00103

Keywords

narrowband, terminal, handsfree, loudspeaking,
VoIP, quality

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

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.

DECTTM, **PLUGTESTS**TM and **UMTS**TM are Trade Marks of ETSI registered for the benefit of its Members.
TIPHONTM and the **TIPHON logo** are Trade Marks currently being registered by ETSI for the benefit of its Members.
3GPPTM is a Trade Mark of ETSI registered for the benefit of its Members and of the 3GPP Organizational Partners.

Contents

Intellectual Property Rights	6
Foreword.....	6
Introduction	6
1 Scope	7
2 References	7
3 Definitions and abbreviations.....	8
3.1 Definitions	8
3.2 Abbreviations	9
4 General considerations	9
4.1 Default Coding Algorithm.....	9
4.2 End-to-end considerations	10
4.3 Parameters to be investigated	10
4.3.1 Basic parameters.....	10
4.3.2 Further Parameters with respect to Speech Processing Devices	10
5 Test equipment	11
5.1 IP half channel measurement adaptor.....	11
5.2 Network impairment simulation.....	11
5.3 Acoustic environment.....	12
6 Test setup.....	12
6.1 Setup for terminal.....	13
6.1.1 Hands-free measurements.....	13
6.1.2 Measurements in loudspeaking mode	17
6.2 Test signal levels	17
6.2.1 Sending	17
6.2.2 Receiving.....	18
6.3 Setup of background noise simulation.....	18
7 Measurements and Requirements for Basic Parameters	20
7.1 Coding independent parameters	20
7.1.1 Sending sensitivity/frequency response	20
7.1.1.1 Requirement	20
7.1.1.2 Measurement method	21
7.1.2 Sending loudness rating	21
7.1.2.1 Requirement	21
7.1.2.2 Measurement method	21
7.1.3 Sending distortion	21
7.1.3.1 Requirement	21
7.1.3.2 Measurement method	22
7.1.4 Out-of-band signals in sending direction	22
7.1.4.1 Requirement	22
7.1.4.2 Measurement method	22
7.1.5 Sending noise.....	23
7.1.5.1 Requirement	23
7.1.5.2 Measurement method	23
7.1.6 Receive sensitivity/frequency response	23
7.1.6.1 Requirement	23
7.1.6.2 Measurement method	25
7.1.7 Receive loudness rating	26
7.1.7.1 Requirement	26
7.1.7.2 Measurement method	26
7.1.8 Receiving distortion.....	27
7.1.8.1 Requirement	27

7.1.8.2	Measurement method	27
7.1.9	Out-of-band signals in receiving direction.....	28
7.1.9.1	Requirement	28
7.1.9.2	Measurement Method.....	28
7.1.10	Receiving noise.....	28
7.1.10.1	Requirement	28
7.1.10.2	Measurement method	28
7.1.11	Terminal Coupling Loss	29
7.1.11.1	Requirement	29
7.1.11.2	Measurement method	29
7.1.12	Stability Loss	29
7.1.12.1	Requirement	29
7.1.12.2	Measurement method	29
7.2	Codec Specific Requirements.....	30
7.2.1	Send Delay.....	30
7.2.1.1	Requirement	31
7.2.1.2	Measurement Method.....	31
7.2.2	Receive delay.....	31
7.2.2.1	Requirement	32
7.2.2.2	Measurement Method.....	32
8	Measurements and Requirements for Parameters with respect to Speech Processing Devices	33
8.1	Objective Listening Speech Quality MOS-LQO in Send direction.....	33
8.2	Objective Listening Quality MOS-LQO in Receive direction	33
8.3	Minimum activation level and sensitivity in Receive direction	33
8.4	Automatic Level Control in Receiving.....	33
8.5	Double Talk Performance.....	33
8.5.1	Attenuation Range in Sending Direction during Double Talk $A_{H,S,dt}$	34
8.5.1.1	Requirement	34
8.5.1.2	Measurement Method.....	34
8.5.2	Attenuation Range in Receiving Direction during Double Talk $A_{H,R,dt}$	35
8.5.2.1	Requirement	35
8.5.2.2	Measurement Method.....	35
8.5.3	Detection of Echo Components during Double Talk.....	36
8.5.3.1	Requirement	36
8.5.3.2	Measurement Method.....	36
8.5.4	Minimum activation level and sensitivity of double talk detection	37
8.5.5	Switching characteristics	38
8.5.5.1	Activation in Sending Direction.....	38
8.5.5.1.1	Requirements	38
8.5.5.1.2	Measurement Method	38
8.5.5.2	Silence Suppression and Comfort Noise Generation	39
8.5.5.3	Performance in sending direction in the presence of background noise	39
8.5.5.3.1	Requirement	39
8.5.5.3.2	Measurement Method.....	39
8.5.5.4	Speech Quality in the Presence of Background Noise	40
8.5.5.5	Quality of Background Noise Transmission (with Far End Speech)	40
8.5.5.5.1	Requirements	40
8.5.5.5.2	Measurement Method.....	40
8.5.5.6	Quality of Background Noise Transmission (with Near End Speech).....	40
8.5.5.6.1	Requirement	40
8.5.5.6.2	Measurement Method.....	40
8.5.6	Quality of echo cancellation	41
8.5.6.1	Temporal echo effects	41
8.5.6.1.1	Requirement	41
8.5.6.1.2	Measurement Method.....	41
8.5.6.2	Spectral Echo Attenuation.....	41
8.5.6.2.1	Requirement	41
8.5.6.2.2	Measurement Method.....	42
8.5.6.3	Occurrence of Artifacts	42
8.5.7	Variant Impairments; Network dependant.....	42
8.5.7.1	Delay versus Time Send.....	42

8.5.7.2	Delay versus Time Receive.....	42
8.5.7.3	Quality of Jitter buffer adjustment	42
Annex A (informative):	Processing delays in VoIP terminals	43
Annex B (informative):	Bibliography.....	46
History		47

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, terminals directly interfacing packet-switched networks (VoIP) are being rapidly introduced. Such IP network edge devices may include 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. Also, the VoIP terminals may use other than 64 kbit/s PCM (ITU-T Recommendation G.711 [9]) speech algorithms.

The present document will provide speech transmission performance requirements for narrowband VoIP loudspeaking and hands-free terminals.

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

1 Scope

The present document will provide speech transmission performance requirements for narrowband VoIP loudspeaking and hands-free terminals; it addresses all types of IP based terminals, including wireless, softphones and group terminals.

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.

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

NOTE: The present document does not concern headset terminals.

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 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".
- [2] ETSI EG 202 425: "Speech Processing, Transmission and Quality Aspects (STQ); Definition and implementation of VoIP reference point".
- [3] ETSI I-ETS 300 245-3: "Integrated Services Digital Network (ISDN); Technical characteristics of telephony terminals; Part 3: Pulse Code Modulation (PCM) A-law, loudspeaking and handsfree telephony".
- [4] ETSI TS 126 171: "Digital cellular telecommunications system (Phase 2+); Universal Mobile Telecommunications System (UMTS); AMR speech codec, wideband; General description (3GPP TS 26.171 version 6.0.0 Release 6)".
- [5] ITU-T Recommendation G.108: "Application of the E-model: A planning guide".
- [6] ITU-T Recommendation G.109: "Definition of categories of speech transmission quality".
- [7] ITU-T Recommendation G.122: "Influence of national systems on stability and talker echo in international connections".
- [8] ITU-T Recommendation G.131: "Talker echo and its control".
- [9] ITU-T Recommendation G.711: "Pulse code modulation (PCM) of voice frequencies".
- [10] ITU-T Recommendation G.723.1: "Dual rate speech coder for multimedia communications transmitting at 5.3 and 6.3 kbit/s".

- [11] ITU-T Recommendation G.729: "Coding of speech at 8 kbit/s using conjugate-structure algebraic-code-excited linear prediction (CS-ACELP)".
- [12] ITU-T Recommendation G.1020: "Performance parameter definitions for quality of speech and other voiceband applications utilizing IP networks".
- [13] ITU-T Recommendation O.41: "Psophometer for use on telephone-type circuits".
- [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.58: "Head and torso simulator for telephonometry".
- [17] ITU-T Recommendation P.79: "Calculation of loudness ratings for telephone sets".
- [18] ITU-T Recommendation P.310: "Transmission characteristics for telephone band (300-3400 Hz) digital telephones".
- [19] ITU-T Recommendation P.340: "Transmission characteristics and speech quality parameters of hands-free terminals".
- [20] ITU-T Recommendation P.342: "Transmission characteristics for telephone band (300-3400 Hz) digital loudspeaking and hands-free telephony terminals".
- [21] ITU-T Recommendation P.501: "Test signals for use in telephonometry".
- [22] ITU-T Recommendation P.502: "Objective test methods for speech communication systems using complex test signals".
- [23] ITU-T Recommendation P.581: "Use of head and torso simulator (HATS) for hands-free terminal testing".
- [24] 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".
- [25] 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 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

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)

group-audio terminal: handsfree terminal primarily designed for use by several users which will not be equipped with a handset

handsfree telephony terminal: telephony terminal using a loudspeaker associated with an amplifier as a telephone receiver and which can be used without a handset

HATS Hands-Free Reference Point (HATS HFRP): corresponds to a reference point "n" from ITU-T Recommendation P.58 [16]: "n" is one of the points numbered from 11 to 17 and defined in table 6a of ITU-T Recommendation P.58 [16] (coordinates of far field front point). The HATS HFRP depends on the location(s) of the microphones of the terminal under test: the appropriate axis lip-ring/HATS HFRP shall be as close as possible to the axis lip-ring/HFT microphone under test

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

loudspeaking function: function of a handset telephone using a loudspeaker associated with an amplifier as a telephone receiver

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: setting which is closest to the nominal RLR

softphone: speech communication system based upon a computer

3.2 Abbreviations

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

CSS	Composite Source Signal
DRP	ear Drum Reference Point
EL	Echo Loss
HATS	Head And Torso Simulator
HFRP	Hands Free Reference Point
LAN	Local Area Network
L_E	Earphone coupling Loss
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
PLC	Packet Loss Concealment
PSTN	Public Switched Telephone Network
QoS	Quality of Service
RLR	Receive Loudness Rating
RLRmax	Receive Loudness Rating corresponding to the maximum setting of the volume control
RLRmin	Receive Loudness Rating corresponding to the minimum setting of the volume control
SLR	Send Loudness Rating
TCLw	Terminal Coupling Loss (weighted)
TCN	Trace Control for Netem

4 General considerations

4.1 Default Coding Algorithm

VoIP terminals shall support the coding algorithm according to ITU-T Recommendation G.711 [9] (both μ -law and A-law). VoIP terminals may support other coding algorithms.

NOTE: Packet Loss Concealment as defined in e.g. appendix I of ITU-T Recommendation G.711 [9] 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 general rules of transmission planning tasks are carried out with the E-model taking into account that E model does not directly address handsfree or loudspeaking terminals; this includes the a-priori determination of the desired category of speech transmission quality as defined in ITU-T Recommendation G.109 [6].

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 [5] 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 is acceptable or representative for the specific configuration is the responsibility of the individual transmission planner.

ITU-T Recommendation G.108 [5] 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.
- Some more basic (I-ETS 300 245-3 [3]) parameters are applicable, if ITU-T Recommendation G.711 [9] is used.

4.3 Parameters to be investigated

4.3.1 Basic parameters

The basic parameters are given in ETS 300 245-3 [3], ITU-T Recommendation P.342 [20] and ITU-T Recommendation P.340 [19].

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 document.

The measurements for further parameters with respect to speech processing devices which are a novelty to terminal requirement standards, have been successfully used in TR 102 648-1 (see bibliography).

- Objective evaluation of speech quality for VoIP terminals.
- Minimum activation level and sensitivity in Receive direction.
- Automatic Level Control in Receiving.
- Double Talk Performance.
- Minimum activation level and sensitivity of double talk detection.
- Switching characteristics.
- Quality of echo cancellation.
- Variant Impairments; Network dependant.
- etc.

5 Test equipment

5.1 IP half channel measurement adaptor

The IP half channel measurement adaptor is described in EG 202 425 [2].

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 loss of performance 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].

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.

5.3 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. But this is not applicable for handsfree and loudspeaking terminals.

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.

In case where an anechoic room is not available the test room has to be an acoustically treated room with few reflections and a low noise level.

Considering this, test laboratory, in the case where its test room does not conform to anechoic conditions as given in ITU-T Recommendation P.342 [20], has to present difference in results for measurements due to its test room.

6 Test setup

In order to use a compatible test system with handset measurements a HATS (Head and Torso Simulator) will be used instead of free field microphone (for receiving measurement) and artificial mouth (for sending measurement). HATS is described in ITU-T Recommendation P.58 [16].

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.

When, a coder with variable bite rate is used, we should adopt, for testing terminal electroacoustical parameters, the bit rate recognized as giving the best characteristics is selected, e g.:

- TS 126 171 [4]: 12,2 kbit/s.

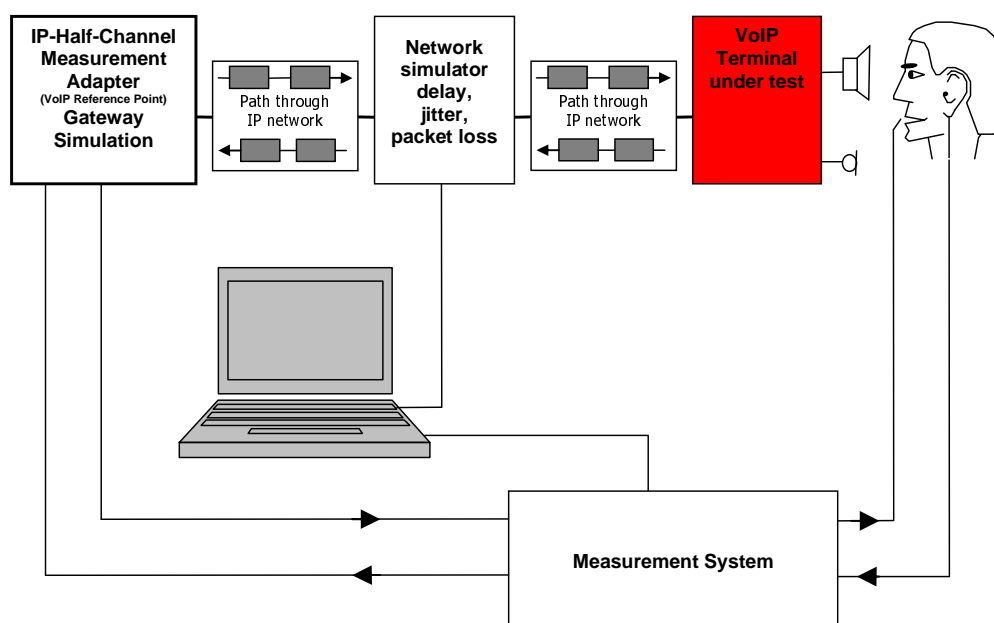


Figure 1: Half channel terminal measurement

6.1 Setup for terminal

6.1.1 Hands-free measurements

The ear used for measurement shall be indicated in the test report.

Desktop operated loudspeaker terminal

For HATS test equipment, definition of loudspeaker terminal and setups for loudspeaker terminal can be found in ITU-T Recommendation P.581 [23].

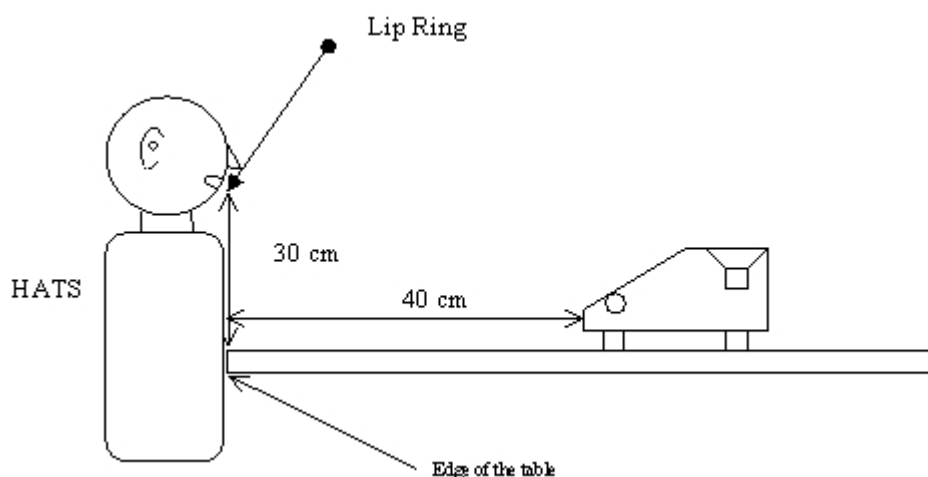


Figure 2: Position for test of desktop hands free terminal side view

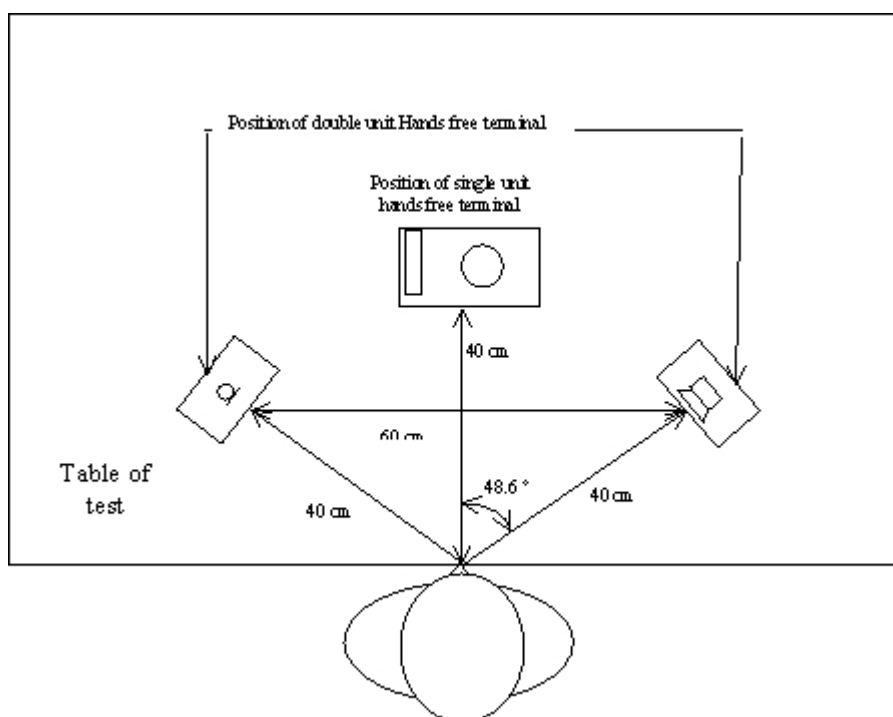


Figure 3: Position for test of desktop hands free terminal top sight

Handheld loudspeaker terminal

It should be placed in according to figure 4. The HATS should be positioned so that the HATS Reference Point is at a distance d_{HF} from the centre point of the visual display of the Mobile Station. The distance d_{HF} is specified by the manufacturer. A vertical angle θ_{HF} may be specified by the manufacturer.

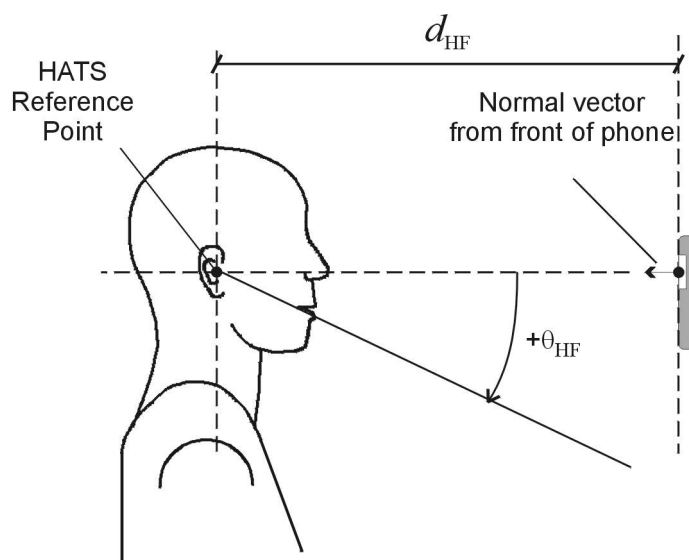


Figure 4: Configuration of Hand-Held loudspeaker relative to the HATS side view

The HATS reference point should be located at a distance d_{HF} from the centre of the visual display of the Mobile Station. The distance d_{HF} is specified by the manufacturer, $d_{HFR}=d_{HF}$, $d_{HFS}=d_{HF}-d_{EM}$, where d_{HFR} is the distance for receiving measurement, d_{HFS} is the distance for sending measurement, and d_{EM} is the distance from ERP to MRP.

When no operating distance is specified by manufacturer, value for d_{HFR} will be 30 cm. A calculation of d_{EM} for HATS gives 12 cm.

A value of 42 cm will be taken for d_{HF} .

Softphone (computer-based terminals)

Two types of softphones are to be considered:

- Type 1 is to be used as a desktop type (e.g. notebook).
- Type 2 is to be used as a handheld type (e.g. PDA).

When manufacturer gives conditions of use, they will apply for test.

If no other requirement is given by manufacturer softphone will be positioned according the following conditions:

Softphone including speakers and microphone

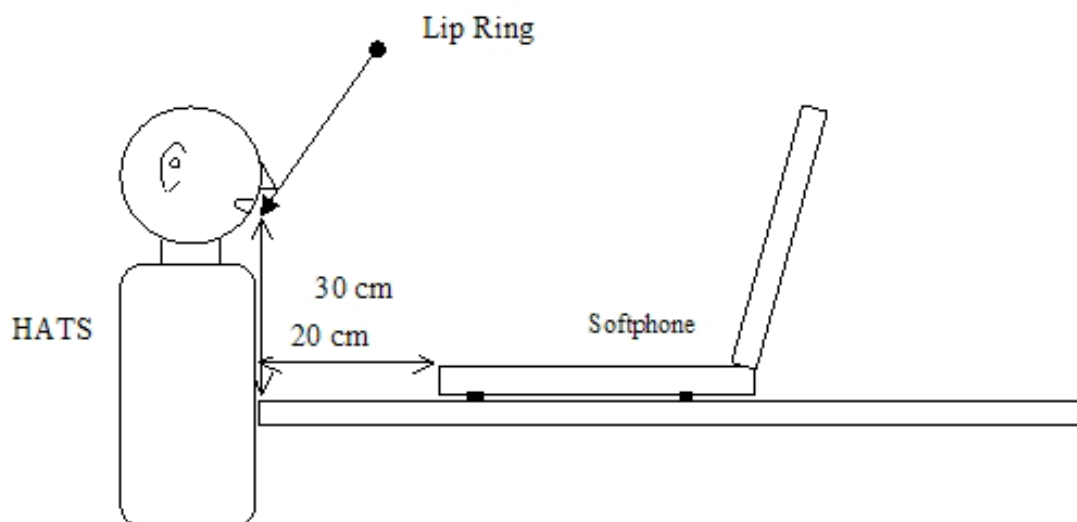


Figure 5: Configuration of softphone relative to the HATS side view

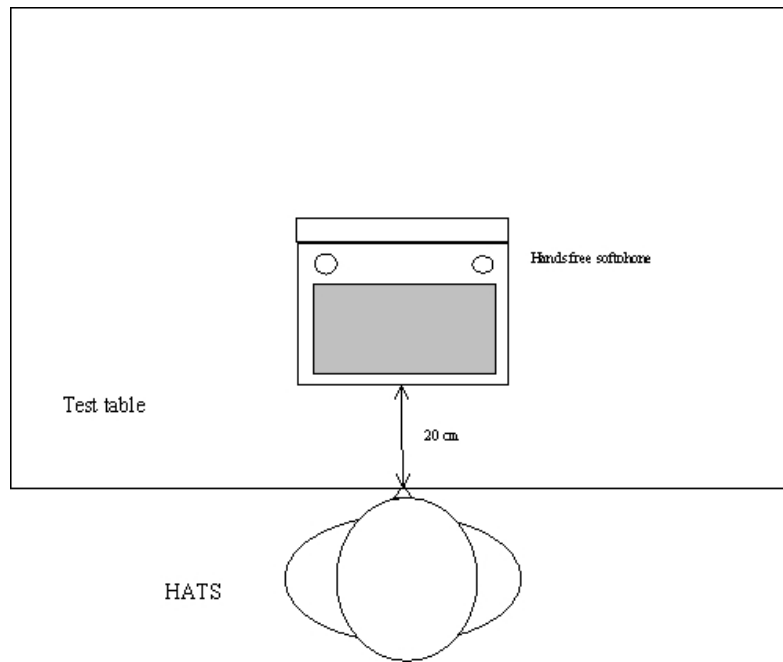


Figure 6: Configuration of softphone relative to the HATS top sight

Softphone with separate speakers

When separate loudspeakers are used, system will be positioned as in figure 7.

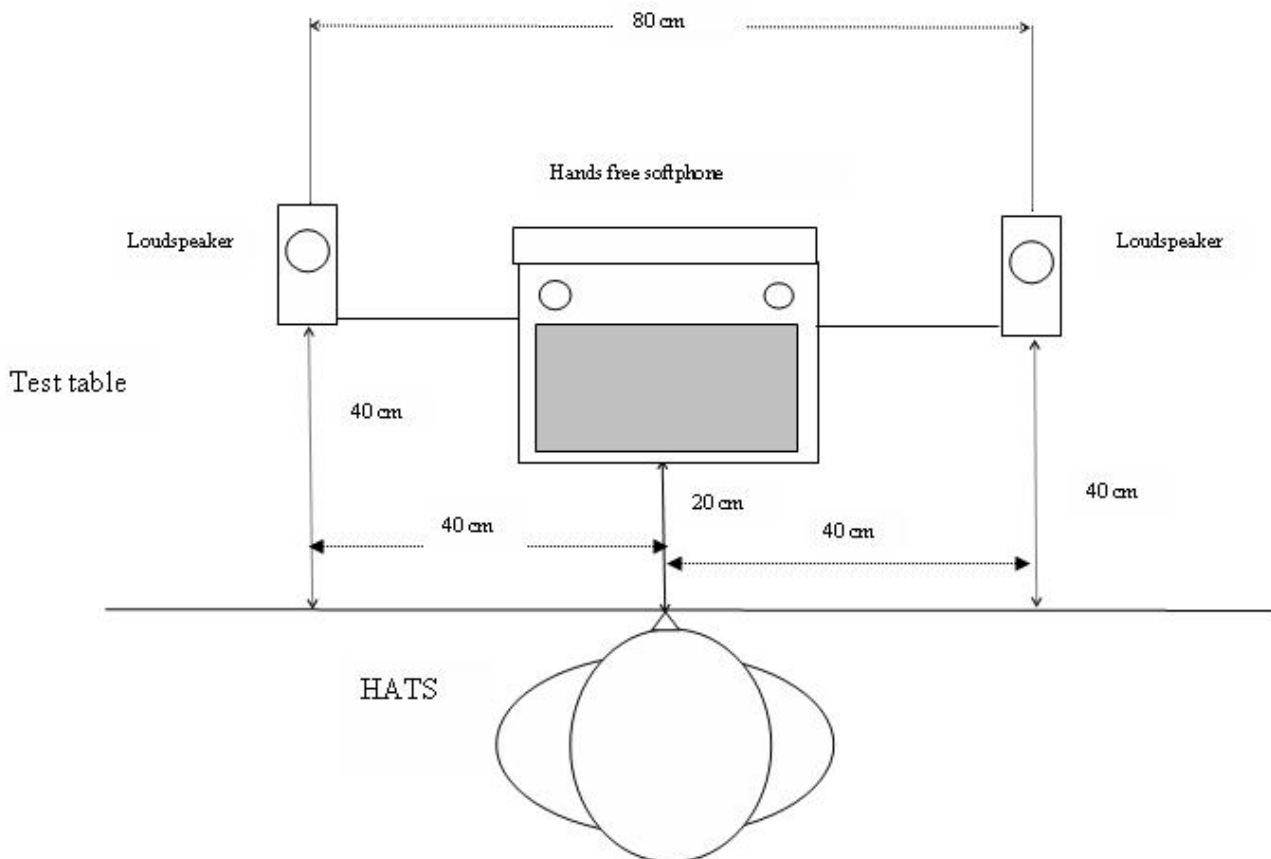


Figure 7: Configuration of softphone using external speakers relative to the HATS top sight

When external microphone and speakers are used, system will be positioned as in figure 8.

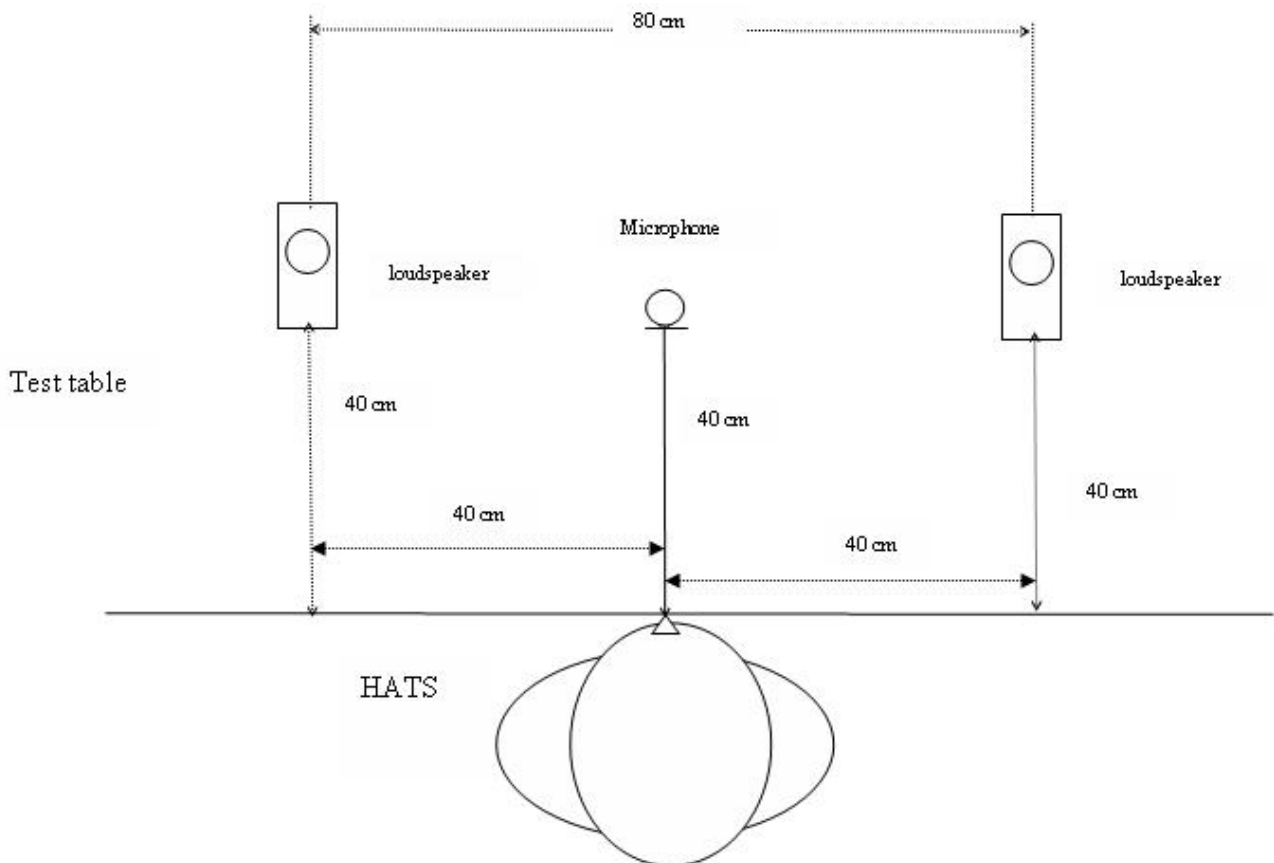


Figure 8: Configuration of softphone using external speakers and microphone relative to the HATS top sight

Group terminal

When manufacturer gives conditions of use, they will apply for test.

When no requirement from manufacturer is given, the following conditions will be used by test laboratory.

Measurement will be conducted by using a HATS test equipment.

The following test position will be used.

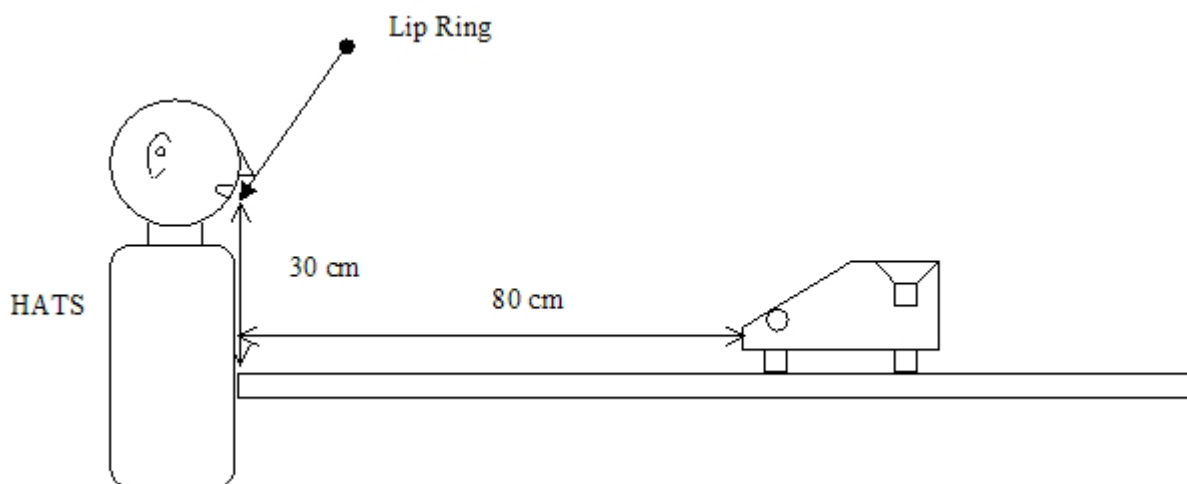


Figure 9: Configuration of group terminal relative to the HATS side view

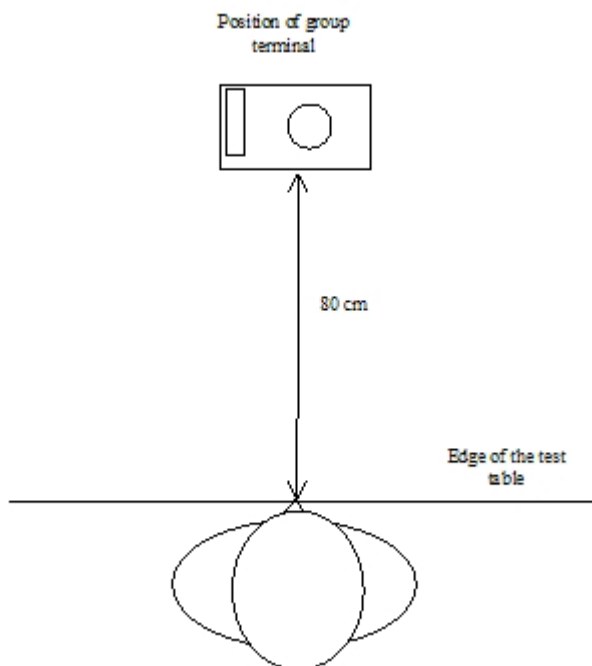


Figure 10: Configuration of group terminal relative to the HATS top sight

NOTE: In case of special casing where those conditions are not realistic, test laboratory can use a different position more representative of real use. The conditions of test will be given in the test report.

6.1.2 Measurements in loudspeaking mode

For those measurements HATS will be used.

It will be positioned as defined in clause 6.1.1 measurement will be performed on one ear and handset will be placed on the other ear. The ear used for measurement will be specified in test report.

NOTE: Only desktop terminals are concerned by loudspeaking measurement.

6.2 Test signal levels

6.2.1 Sending

Unless specified otherwise, the test signal level shall be -4,7 dBPa at the MRP.

The following procedure shall be used to perform the calibration of the artificial mouth of the HATS:

- The input signal from the artificial mouth is first calibrated under free-field conditions at the MRP. The total level on the frequency range is set to -4,7 dBPa.
- The spectrum at MRP is recorded.
- Then the level is adjusted to the level given further in this text (depending of type of terminal tested (for example -24,3 dBPa at 30 cm for handheld terminal)).
- The level at MRP (measured in third octave bands) adjusted at the first step (with total level of -4,7 dBPa) is used as the reference for sending characteristics.

The test setup shall be in conformance with, figure 11 but, depending on the type of terminal, the appropriate distance and level will be used. When using this calibration method, send sensitivity must be calculated as follows:

$$\mathbf{SmJ} = 20 \log \mathbf{Vs} - 20 \log \mathbf{PMRP}$$

where:

\mathbf{Vs} is the measured voltage across the appropriate termination (unless stated otherwise, a 600 Ω termination).

\mathbf{PMRP} is the applied sound pressure at the MRP during the first step of calibration.

NOTE: Reason for this procedure of calibration in two steps is to take into account the different variation of signal with distance by using different implementations of HATS.

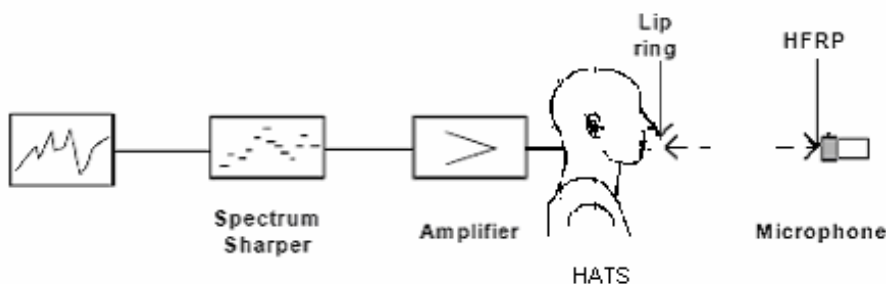


Figure 11: Calibration at HFRP

The distance used for level calibration corresponds to the following values:

Desktop terminal: 50 cm and level to adjust - 28,7 dBPa

Handheld terminal: 30 cm with - 24,3 dBPa

Softphone: 36 cm with - 25,8 dBPa

Group terminal: 85 cm with - 33,3 dBPa

6.2.2 Receiving

Unless specified otherwise, the applied test signal level at the digital input shall be -16 dBm0.

All measurement values produced by HATS are intended to be free-field equalized.

6.3 Setup of background noise simulation

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

The general procedure for setup a background noise simulation arrangement is described in EG 202 396-1 [1]. The EG 202 396-1 [1] 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 12.

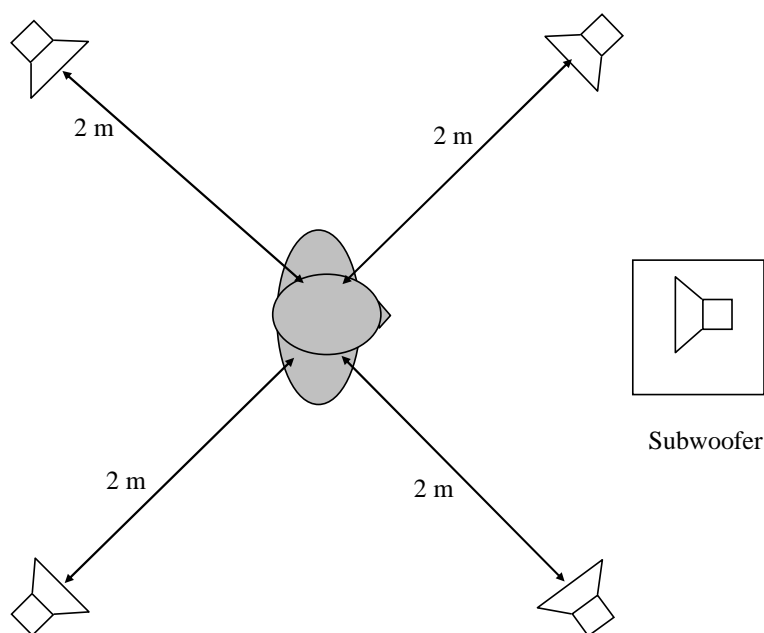


Figure 12: Loudspeaker arrangement for background noise simulation

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

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 standard.

NOTE 2: Due to time variant nature of IP connection, delay variation may impair the measurement. In such case, the measurement has to be repeated until a valid measurement can be achieved.

7.1 Coding independent parameters

7.1.1 Sending sensitivity/frequency response

7.1.1.1 Requirement

The sending sensitivity/frequency response shall be within the limits given in table 1.

Table 1

Frequency	Upper limit	Lower limit
200 Hz	0 dB	-∞ dB
315 Hz	0 dB	-14 dB
400 Hz	0 dB	-13 dB
500 Hz	0 dB	-12 dB
630 Hz	0 dB	-11 dB
800 Hz	0 dB	-10 dB
1 000 Hz	0 dB	-8 dB
1 300 Hz	2 dB	-8 dB
1 600 Hz	3 dB	-8 dB
2 000 Hz	4 dB	-8 dB
3 100 Hz	4 dB	-8 dB
4 000 Hz	0 dB	-∞ dB

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

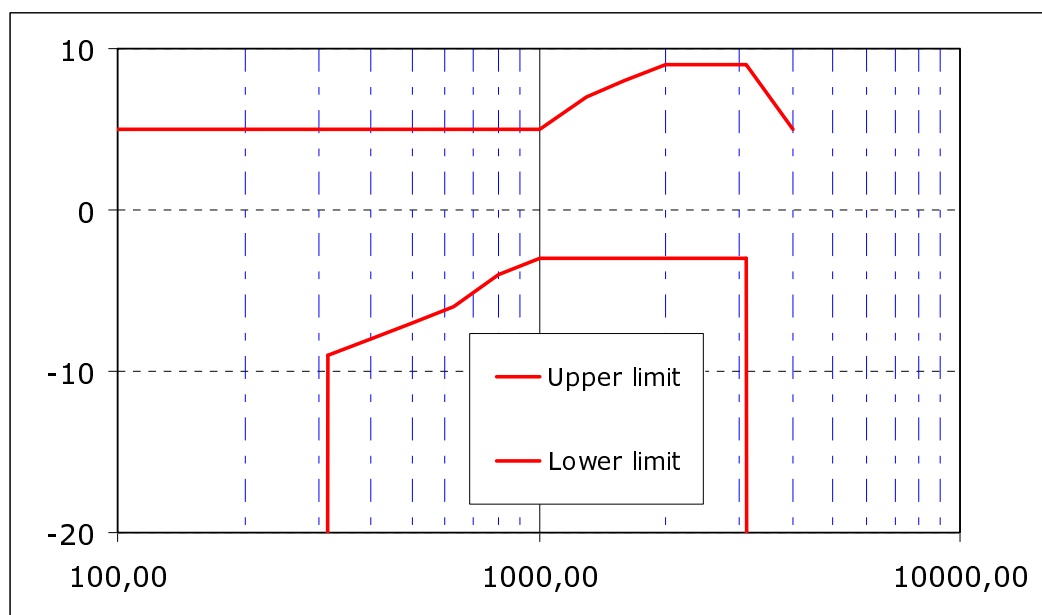


Figure 13: Sending sensitivity/frequency mask for HFT

7.1.1.2 Measurement method

The terminal will be positioned as described in clause 6.1.

An artificial voice according to ITU-T Recommendation P.50 [14] or a speech like test signal as described in ITU-T Recommendation P.501 [21] can be used for test. The type of test signal used shall be stated in the test report. The spectrum of acoustic signal produced by the artificial mouth is calibrated under free field conditions at the MRP. The signal level is adjusted according to clause 6.2.1.

The spectrum at the MRP and the actual level at the MRP (measured in third octaves) is used as reference to determine the sending sensitivity SmJ.

7.1.2 Sending loudness rating

7.1.2.1 Requirement

The value of SLR shall be $+13 \text{ dB} \pm 3 \text{ dB}$.

This value is derived from ITU-T Recommendation P.310 [18]. According to ITU-T Recommendation P.340 [19] the SLR of a hands-free telephone should be about 5 dB higher than the SLR of the corresponding handset telephone.

This value will be identical for all type of terminal (desktop, handheld, etc.) difference in efficiency will be given by conditions for measurement (see clause 6.1).

7.1.2.2 Measurement method

The terminal will be positioned as described in clause 6.1.

An artificial voice according to ITU-T Recommendation P. 50 [14] or a speech like test signal as described in ITU-T Recommendation P.501 [21] can be used to test. The type of test signal used shall be stated in the test report. 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 $-4,7 \text{ dBPa}$, measured at the MRP. The test signal level is averaged over the complete test signal sequence.

Calibration is realized as explained in clause 6.2.1.

SLR shall be calculated according ITU-T Recommendation P.79 [17].

7.1.3 Sending distortion

7.1.3.1 Requirement

The terminal will be positioned as described in clause 6.1.

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

Table 1a

Frequency	Ratio
315 Hz	26 dB
400 Hz	30 dB
1 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.	

7.1.3.2 Measurement method

The terminal will be positioned as described in clause 6.1.

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 and 1 000 Hz, the duration of the sine-wave shall be of 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 3,15 kHz.

An artificial voice according to ITU-T Recommendation P. 50 [14] or a speech like test signal as described in ITU-T Recommendation P.501 [21] 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.

7.1.4 Out-of-band signals in sending direction

7.1.4.1 Requirement

With any signal above 4,6 kHz and up to 8 kHz applied at the MRP at a level of -4,7 dBPa, the level of any image frequency shall be below the level obtained for the reference signal by at least the amount (in dB) specified in table 2.

Table 2: Out-of-band signal limit, sending

Frequency	Limit
4,6 kHz	30 dB
8 kHz	40 dB
NOTE: The limits for intermediate frequencies lie on a straight line drawn between the given values on a linear (dB) - logarithmic (kHz) scale.	

7.1.4.2 Measurement method

The terminal will be positioned as described in clause 6.1.

For a correct activation of the system, an artificial voice according to ITU-T Recommendation P.50 [14] or a speech like test signal as described in ITU-T Recommendation P.501 [21] shall be used for activation. Level of this activation signal shall be -4,7 dBPa at the MRP.

For the test, an out-of-band signal shall be provided as a frequency band signal centred on 4,65 kHz, 5 kHz, 6 kHz, 6,5 kHz, 7 kHz and 7,5 kHz respectively. The level of any image frequencies at the digital interface shall be measured.

The levels of these signals shall be -4,7 dBPa at the MRP.

The complete test signal is constituted by t1 ms of in-band signal (reference signal), t2 ms of out-of-band signal and another time t1 ms of in-band signal (reference signal).

The observation of the output signal on the first and second in-band signals permits control if the set is correctly activated during the out-of-band measurement. This measurement shall be performed during t2 period:

- a value of 250 ms is suggested for t1;
- t2 depends on the integration time of the analyser, typically less than 150 ms.

7.1.5 Sending noise

7.1.5.1 Requirement

The limit for the sending noise is the following:

- send noise level maximum -64 dBm0p.

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

Requirement as for other tests is identical for all types of terminals.

NOTE: Softphones with cooling devices (fans) can produce a rather high level of noise, furthermore largely dependant of activity of system.

7.1.5.2 Measurement method

The terminal will be positioned as described in clause 6.1.

For a correct activation of the system, an artificial voice according to ITU-T Recommendation P. 50 [14] or a speech like test signal as described in ITU-T Recommendation P.501 [21] shall be used for activation. Level of this activation signal shall be -4,7 dBPa at the MRP.

The psophometric noise level at the output of the test setup is measured. The psophometric filter is described in ITU-T Recommendation O.41 [13].

7.1.6 Receive sensitivity/frequency response

7.1.6.1 Requirement

The following masks are required for handsfree and loudspeaking terminals.

Desktop operated loudspeaker

Table 3: Receiving frequency response-desktop

Frequency	Upper limit	Lower limit
200 Hz	6 dB	-∞ dB
250 Hz	6 dB	-∞ dB
315 Hz	6 dB	-9 dB
400 Hz	6 dB	-6 dB
3 100 Hz	6 dB	-6 dB
4 000 Hz	6 dB	-∞ dB
NOTE: The limits for intermediate frequencies lie on a straight line drawn between the given values on a linear (dB) - logarithmic (kHz) scale.		

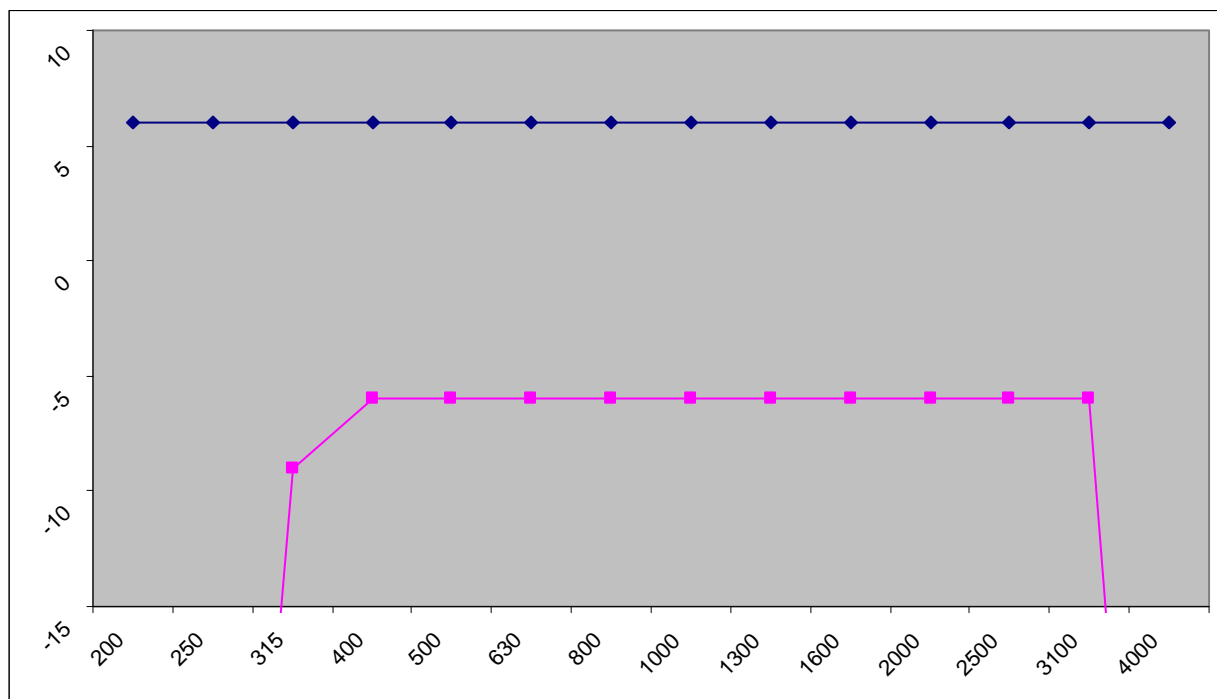


Figure 14: Receiving sensitivity/frequency mask for Desktop hands free terminal

Handheld terminal

Table 4: Receiving frequency response-handheld

Frequency	Upper limit	Lower limit
200 Hz	6 dB	-∞ dB
400 Hz	6 dB	-∞ dB
500 Hz	6 dB	-9 dB
630 Hz	6 dB	-6 dB
3 100 Hz	6 dB	-6 dB
4 000 Hz	6 dB	-∞ dB

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

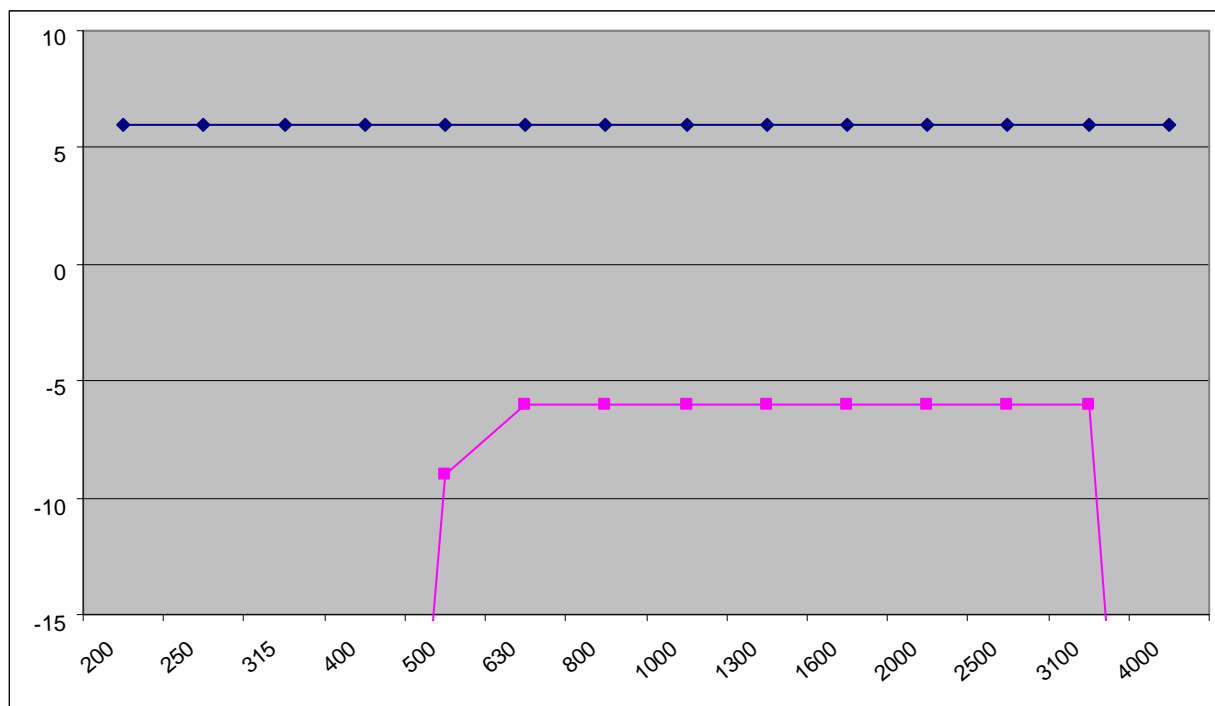


Figure 15: Receiving sensitivity/frequency mask for Hand-held HFT

Softphone (computer-based terminals)

Type 1 or softphone with external speakers: requirement as for desktop terminal

Type 2 requirement as for handheld terminal

Group terminal

Same requirement as desktop terminals.

7.1.6.2 Measurement method

Test setup is described in clause 6.1.

Measurement is operated at nominal value of volume control.

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

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

S_{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_{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]. The test signal level shall be -20 dBm0, measured according to ITU-T Recommendation P.56 [15] at the digital reference point or the equivalent analogue point.

The HATS is free field equalized as described in ITU-T Recommendation P.581 [23]. 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 [25] for frequencies from 100 Hz to 4 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.

7.1.7 Receive loudness rating

7.1.7.1 Requirement

Desktop operated loudspeaker

Nominal value of RLR will be 5 ± 3 dB. This value has to be fulfilled for one position of volume range.

Value of RLR at upper part of volume range must be less than (louder) or equal to -2 dB: $RLR_{max} \leq -2$ dB.

Range of volume control must be equal or exceed 15 dB: $(RLR_{min} - RLR) \geq 15$ dB.

Handheld terminal

Nominal value of RLR will be 9 ± 3 dB. This value has to be fulfilled for one position of volume range.

Value of RLR at upper part of volume range must be less than (louder) or equal to 2 dB: $RLR_{max} \leq 2$ dB.

Range of volume control must be equal or exceed 15 dB: $(RLR_{min} - RLR) \geq 15$ dB.

Softphone (computer-based terminal)

Type 1 or softphone with external speakers: requirement as for desktop terminal.

Type 2 requirement as for handheld terminal.

Group terminal

Nominal value of RLR will be 5 ± 3 dB. This value has to be fulfilled for one position of volume range.

Value of RLR at upper part of volume range must be less than (louder) or equal to -6 dB: $RLR_{max} \leq -6$ dB.

Range of volume control must be equal or exceed 19 dB: $(RLR_{min} - RLR) \geq 19$ dB.

7.1.7.2 Measurement method

Test setup is described in clause 6.1.

The RLR shall be calculated according to ITU-T Recommendation P.79 [17].

The receiving sensitivity shall be calculated from each band of the 14 frequencies given in table 1 of ITU-T Recommendation P.79 [17], bands 4 to 17. 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 dB Pa/V and the RLR(cal) shall be calculated according to the formula 5-1 of ITU-T Recommendation P.79 [17], using the receiving weighting factors from table 1 and according to clause 6, of ITU-T Recommendation P.79 [17]; The RLR shall then be computed as RLR(cal) minus 14 dB according to ITU-T Recommendation P.340 [19], and without the LE factor.

7.1.8 Receiving distortion

7.1.8.1 Requirement

Desktop operated loudspeaker

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

Table 5

Frequency	Signal to distortion ratio limit, receiving for desktop terminal at nominal volume	Signal to distortion ratio limit, receiving for handheld terminal at nominal volume	Signal to distortion ratio limit, receiving for all terminals at maximum volume
315 Hz	26 dB		
400 Hz	30 dB		
500 Hz	30 dB	20 dB	
800 Hz	30 dB	30 dB	20 dB
1 kHz	30 dB	30 dB	
NOTE: The limits for intermediate frequencies lie on a straight line drawn between the given values on a linear (dB) - logarithmic (kHz) scale.			

Handheld terminal

The terminal will be positioned as described in clause 6.1.

The ratio of signal to harmonic distortion is given in table 5.

Softphone (computer-based terminal)

Type 1 or softphone with external speakers: requirement as for desktop terminal.

Type 2 requirement as for handheld terminal.

Group terminal

Same requirement as for desktop terminal.

7.1.8.2 Measurement method

Test setup is described in clause 6.1.

The signal used is an activation signal followed by a series sine-wave signal with a frequency at 315 Hz, 400 Hz, 500 Hz, 630 Hz, 800 Hz and 1 000 Hz, The duration of the sine-wave shall be of less than 1 s. The sinusoidal signal level shall be calibrated to -16 dBm0.

An artificial voice according to ITU-T Recommendation P. 50 [14] or a speech like test signal as described in ITU-T Recommendation P.501 [21] can be used for activation. Level of this activation signal will be -16 dBm0.

The signal to harmonic distortion ratio is measured selectively up to 3,15 kHz.

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

7.1.9 Out-of-band signals in receiving direction

7.1.9.1 Requirement

Any spurious out-of-band image signals in the frequency range from 4,6 kHz to 8 kHz measured selectively shall be lower than the in-band level measured with a reference signal. The minimum level difference between the reference signal level and the out-of-band image signal level shall be as given in table 6.

Table 6: Out-of-band signal limit, receiving

Frequency	Signal limit
4,6 kHz	35 dB
8 kHz	45 dB
NOTE: The limits for intermediate frequencies lie on a straight line drawn between the given values on a linear (dB) - logarithmic (kHz) scale.	

7.1.9.2 Measurement Method

Test setup is described in clause 6.1.

Measurement is operated at nominal value of volume control.

The signal used is an activation signal followed by a sine-wave signal. For input signals at the frequencies 500 Hz, 1 000 Hz, 2 000 Hz and 3 150 Hz applied at the level of -16 dBm₀, the level of spurious out-of-band image signals at frequencies up to 8 kHz is measured selectively at measurement point.

An artificial voice according to ITU-Recommendation P. 50 [14] or a speech like test signal as described in ITU-T Recommendation P.50 [14] can be used for activation. Level of this activation signal will be -16 dBm₀.

7.1.10 Receiving noise

7.1.10.1 Requirement

A-weighted

The noise level shall not exceed -54 dBPa(A) at **nominal setting of the volume control**.

Octave band spectrum

The level in any 1/3-octave band, between 100 Hz and 10 kHz shall not exceed a value of -64 dBPa.

NOTE 1: No peaks in the frequency domain higher than 10 dB above the average noise spectrum should occur.

NOTE 2: For softphone fan noise must be avoided in order to fulfil this condition.

7.1.10.2 Measurement method

Test setup is described in clause 6.1.

A signal is applied to input of test system in order to ensure correct activation of receiving state. An artificial voice according to ITU-Recommendation P. 50 [14] or a speech like test signal as described in ITU-T Recommendation P.501 [21] can be used for activation. Level of this activation signal will be -16 dBm₀.

The noise shall be measured just after interrupting the activation signal.

7.1.11 Terminal Coupling Loss

7.1.11.1 Requirement

In order to meet the ITU-T Recommendation G.131 [8] talker echo objective requirements, the recommended weighted terminal coupling loss during single talk (TCL_{wst}) should be greater than 55 dB when measured under free field conditions at nominal setting of volume control.

A TCL_w greater than 46 dB is considered as acceptable.

TCL_{wst} shall be not less than 40 dB for any setting of the volume control.

7.1.11.2 Measurement method

The setup for terminal is described in clause 6.1.

For hands-free measurement, HATS is positioned but not used.

For loudspeaking measurement, handset is positioned on HATS (right ear).

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 [21] 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 4 kHz). The low-crest factor is achieved by random-alternation of the phase between -180° and 180°.

The TCL_w is calculated according to ITU-T Recommendation G.122 [7], clause B.4 (trapezoidal rule). 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 (e.g. 200 ms) has to be applied adapted to the duration of the actual test signal.

7.1.12 Stability Loss

7.1.12.1 Requirement

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. It must exceed 6 dB for all frequencies and for all settings of volume control.

7.1.12.2 Measurement method

For handsfree mode test set-up is identical as for TCL_w.

For loudspeaking mode handset is placed at 50 cm beside terminal with transducers facing the table as in figure 16.

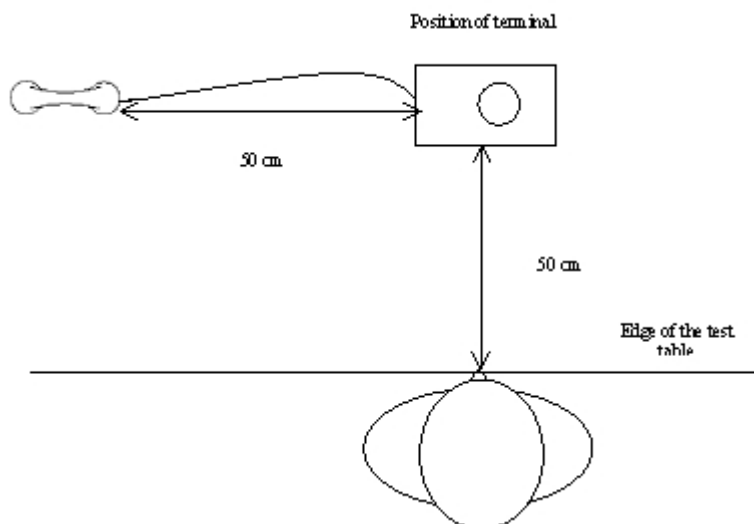


Figure 16: Stability loss position for loudspeaking function

7.2 Codec Specific Requirements

7.2.1 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 figure 2 of ITU-T Recommendation G.1020 [12] and in figure A.1 of ITU-T Recommendation G.1020 [12], respectively.

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

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

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

7.2.1.1 Requirement

The allowance for signal processing shall be $T(\text{asp}) < 40$ ms (including processing time for handsfree).

NOTE: With the knowledge of the codec specific values for $T(\text{fs})$ and $T(\text{la})$ the requirements for send delay for any type of coder and any frame size $T(\text{fs})$ can easily be calculated by formula 1. Table 7 provides requirements calculated accordingly for frequently used codecs and packet sizes.

Table 7

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
G.711 [9]	80	0,125	10	0	0	40	< 50
G.711 [9]	160	0,125	20	0	0	40	< 60
G.729, G729 A and G.729 B [11]	1	10	10	5	0	40	< 55
G.729, G729 A and G.729 B [11]	2	10	20	5	0	40	< 65
ITU-T Rec. G.723.1 (5,3 kbit/s and 6,3 kbit/s) [10]	1	30	30	7,5	0	40	< 77,5

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

7.2.1.2 Measurement Method

Test setup is described in clause 6.1.

The test signal to be used for the measurements shall be a Composite Source Signal (CSS) as described in ITU-T Recommendation P.501 [21]. The test signal consists of the voiced part as described in ITU-T Recommendation P.501 [21] 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 1: If the expected delay is higher than 500 ms a pseudo random sequence with a higher periodicity should be used.

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 2: 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 [24]) may be used.

7.2.2 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 (ears of HATS) 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 figure 3 of ITU-T Recommendation G.1020 [12] and in figure A.2 of ITU-T Recommendation G.1020 [12], respectively.

The receiving delay $T(\text{r})$ is defined as follows:

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

Where:

$T(fs)$ = frame size of encoder

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

7.2.2.1 Requirement

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 1: 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 8 provides requirements calculated accordingly for some frequently used codecs and packet sizes as an example.

Table 8

Codec	N	T(fs) in ms	T(aif) in ms	T(jb) in ms	T(plc) in ms	T(asp) in ms	T(r) Requirement in ms
G.711 [9]	80	0,125	0	10	10	10	< 30,125
G.711 [9]	160	0,125	0	10	10	10	< 30,125
G.729, G.729 A and G.729 B [11]	1	10	0	10	0	10	< 30
G.729, G.729 A and G.729 B [11]	2	10	0	10	0	10	< 30
G.723.1 (5,3 kbit/s and 6,3 kbit/s) [10]	1	30	0	10	0	10	< 50
NOTE 1: $T(ps) = \text{packet size} = N * T(fs)$.							
NOTE 2: N = number of frames per packet.							

NOTE 2: 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.

7.2.2.2 Measurement Method

Test setup is described in clause 6.1.

The test signal to be used for the measurements shall be a Composite Source Signal (CSS) as described in ITU-T Recommendation P.501 [21] 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.

NOTE 1: If the expected delay is higher than 500 ms a pseudo random sequence with a higher periodicity should be used.

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 2: 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 [24]) may be used.

8 Measurements and Requirements for Parameters with respect to Speech Processing Devices

8.1 Objective Listening Speech Quality MOS-LQO in Send direction

For further study.

8.2 Objective Listening Quality MOS-LQO in Receive direction

For further study.

8.3 Minimum activation level and sensitivity in Receive direction

For further study.

8.4 Automatic Level Control in Receiving

For further study.

8.5 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 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 [19] and P.502 [22]):

- 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.

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

8.5.1.1 Requirement

Based on the level variation in sending direction during double talk $A_{H,S,dt}$ the behavior of the terminal can be classified according to table 9.

Table 9

Category (according to ITU-T Rec. P.340 [19])	1	2a	2b	2c	3
	Full Duplex Capability	Partial Duplex Capability			No Duplex Capability
$A_{H,S,dt}$ [dB]	≤ 3	≤ 6	≤ 9	≤ 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.

8.5.1.2 Measurement Method

Test setup is described in clause 6.1.

The test signal to determine the attenuation range during double talk is shown in figure 17. A sequence of uncorrelated CS signals is used which is inserted in parallel in sending and receiving direction.

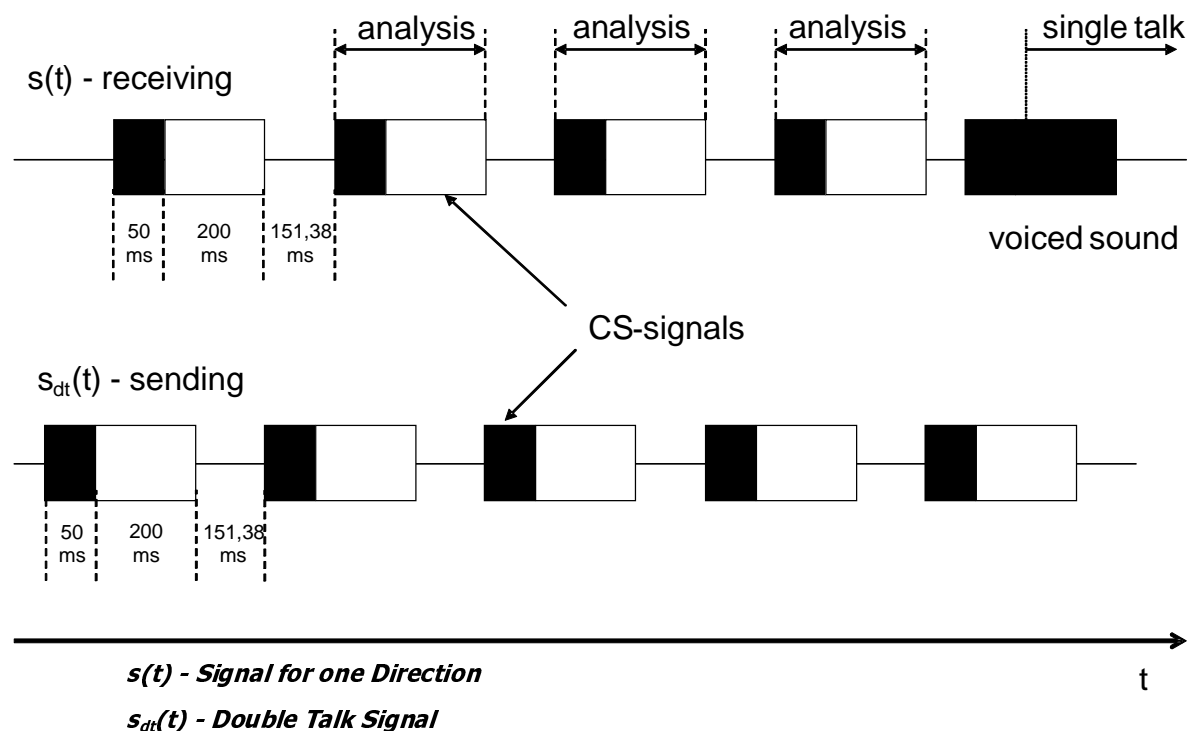


Figure 17: Double Talk Test Sequence with overlapping CS signals in sending and receiving direction

Figure 17 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 17 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 10

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

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.

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

8.5.2.1 Requirement

Based on the level variation in receiving direction during double talk $A_{H,R,dt}$ the behavior of the terminal can be classified according to table 11.

Table 11

Category (according to ITU-T Rec. P.340 [19])	1	2a	2b	2c	3
	<i>Full Duplex Capability</i>	<i>Partial Duplex Capability</i>			<i>No Duplex Capability</i>
$A_{H,R,dt}$ [dB]	≤ 3	≤ 5	≤ 8	≤ 10	> 10

In general table 11 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.

8.5.2.2 Measurement Method

Test setup is described in clause 6.1.

The test signal to determine the attenuation range during double talk is shown in figure 17. 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 12

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

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.

8.5.3 Detection of Echo Components during Double Talk

8.5.3.1 Requirement

"Echo Loss" (EL) is the echo suppression provided by the terminal measured at the electrical reference point. Under these conditions the requirements given in table 13 are applicable (more information can be found in annex A of the ITU-T Recommendation P.340 [19]).

Table 13

Category (according to ITU-T Rec. P.340 [19])	1	2a	2b	2c	3
	Full Duplex Capability	Partial Duplex Capability			No Duplex Capability
Echo Loss [dB]	≥ 27	≥ 23	≥ 17	≥ 11	< 11

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

8.5.3.2 Measurement Method

Test setup is described in clause 6.1.

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 the figure below. A detailed description can be found in ITU-T Recommendation P.501 [21].

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

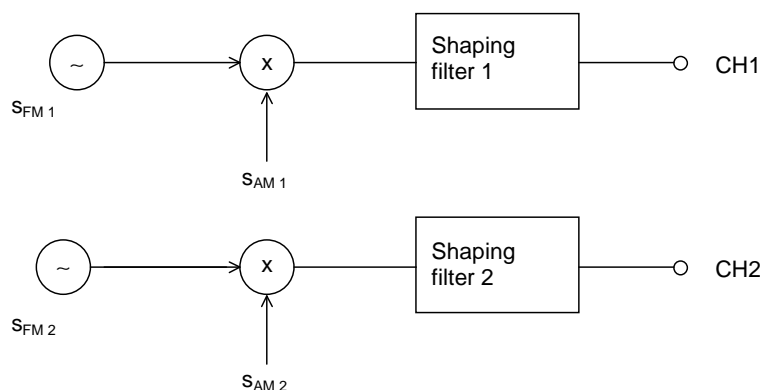


Figure 17a: Measurement signals

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

$$s_{AM1,2}(t) = A_{AM1,2} * \cos(2\pi t F_{AM1,2}); \quad (3)$$

The settings for the signals are as follows.

Table 13a: Settings for the signal

Receiving Direction			Sending Direction			
f_m [Hz]	$f_{\text{mod}(fm)}$ [Hz]	F_{am} [Hz]		f_m [Hz]	$f_{\text{mod}(fm)}$ [Hz]	F_{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.

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

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 [21]). 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.

8.5.4 Minimum activation level and sensitivity of double talk detection

For further study.

8.5.5 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.

8.5.5.1 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).

8.5.5.1.1 Requirements

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

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

8.5.5.1.2 Measurement Method

Test setup is described in clause 6.1.

The structure of the test signal is shown in figure 18. The test signal consists of CSS components according to ITU-T Recommendation P.501 [21] with increasing level for each CSS burst.

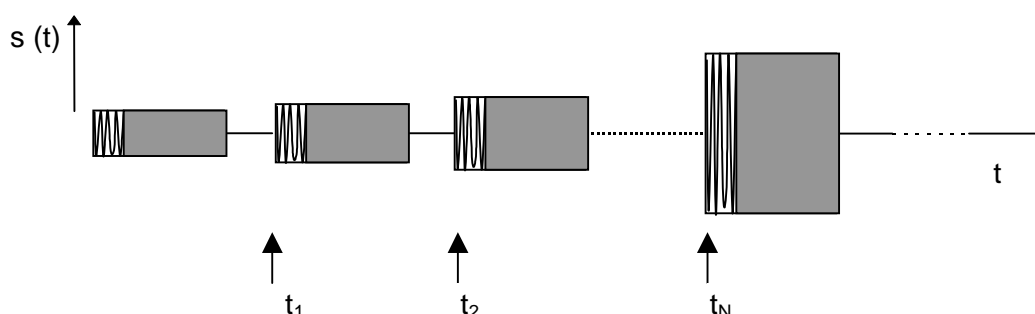


Figure 18: Test Signal to Determine the Minimum Activation Level and the Built-up Time

The settings of the test signal are as follows:

	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 [21] 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 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.

8.5.5.2 Silence Suppression and Comfort Noise Generation

For further study.

8.5.5.3 Performance in sending direction in the presence of background noise

8.5.5.3.1 Requirement

The level of comfort noise, if implemented, shall be within a range of +2 dB 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 14.

Table 14: 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
NOTE: All sensitivity values are expressed in dB on an arbitrary scale.		

8.5.5.3.2 Measurement Method

Test setup is described in clause 6.1.

The background noise simulation as described in clause 6.3 is used.

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 (CSS) 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.

8.5.5.4 Speech Quality in the Presence of Background Noise

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

8.5.5.5 Quality of Background Noise Transmission (with Far End Speech)

8.5.5.5.1 Requirements

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 6.3.

8.5.5.5.2 Measurement Method

Test setup is described in clause 6.1.

The background noises are generated as described in clause 6.3.

First the measurement is conducted without inserting the signal at the far end. At least 10 s of noise are 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 seconds in order to allow adaptation of the noise reduction algorithms. After at least 5 seconds a Composite Source Signal according to ITU-T Recommendation P.501 [21] is applied in receiving direction with a duration of ≥ 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.

8.5.5.6 Quality of Background Noise Transmission (with Near End Speech)

8.5.5.6.1 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.

8.5.5.6.2 Measurement Method

Test setup is described in clause 6.1.

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 [21] with a duration of ≥ 2 CSS periods. The test signal level shall be -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.

8.5.6 Quality of echo cancellation

8.5.6.1 Temporal echo effects

8.5.6.1.1 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.

8.5.6.1.2 Measurement Method

Test setup is described in clause 6.1.

The test signal consists of periodically repeated Composite Source Signal according to ITU-T Recommendation P.501 [21] 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 behavior 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).

8.5.6.2 Spectral Echo Attenuation

8.5.6.2.1 Requirement

The echo attenuation vs. frequency shall be below the tolerance mask given in table 15.

Table 15: Echo attenuation

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
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.

8.5.6.2.2 Measurement Method

Test setup is described in clause 6.1.

Before the actual measurement a training sequence is fed in consisting of 10 seconds CS signal according to ITU-T Recommendation P.501 [21]. The level of the training sequence shall be -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.

8.5.6.3 Occurrence of Artifacts

For further study.

8.5.7 Variant Impairments; Network dependant

8.5.7.1 Delay versus Time Send

For further study.

8.5.7.2 Delay versus Time Receive

For further study.

8.5.7.3 Quality of Jitter buffer adjustment

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.

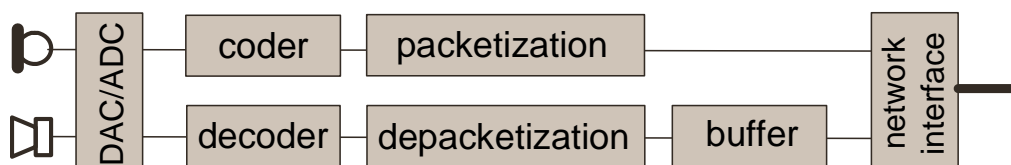


Figure A.1: Synoptic of the different functions implemented in a VoIP terminal

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 A.1 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".

History

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