



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

Reference

RES/STQ-253

Keywords

handsfree, loudspeaking, narrowband, quality,
speech, 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

The present document can be downloaded from:

<http://www.etsi.org/standards-search>

The present document may be made available in electronic versions and/or in print. The content of any electronic and/or print versions of the present document shall not be modified without the prior written authorization of ETSI. In case of any existing or perceived difference in contents between such versions and/or in print, the only prevailing document is the print of the Portable Document Format (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

<https://portal.etsi.org/TB/ETSIDeliverableStatus.aspx>

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

<https://portal.etsi.org/People/CommitteeSupportStaff.aspx>

Copyright Notification

No part may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying and microfilm except as authorized by written permission of ETSI.

The content of the PDF version shall not be modified without the written authorization of ETSI.

The copyright and the foregoing restriction extend to reproduction in all media.

© European Telecommunications Standards Institute 2017.

All rights reserved.

DECT™, **PLUGTESTS™**, **UMTS™** and the ETSI logo are Trade Marks of ETSI registered for the benefit of its Members.
3GPP™ and **LTE™** are Trade Marks of ETSI registered for the benefit of its Members and
of the 3GPP Organizational Partners.
GSM® and the GSM logo are Trade Marks registered and owned by the GSM Association.

Contents

Intellectual Property Rights	5
Foreword.....	5
Modal verbs terminology.....	5
Introduction	5
1 Scope	6
2 References	6
2.1 Normative references	6
2.2 Informative references.....	7
3 Definitions and abbreviations.....	8
3.1 Definitions.....	8
3.2 Abbreviations	9
4 General considerations	10
4.1 Default Coding Algorithm.....	10
4.2 End-to-end considerations	10
5 Test equipment	10
5.1 IP half channel measurement adaptor.....	10
5.2 Environmental conditions for tests	10
5.3 Accuracy of measurements and test signal generation	11
5.4 Network impairment simulation.....	11
5.5 Acoustic environment.....	12
5.6 Influence of terminal delay on measurements	12
6 Requirements and associated measurement methodologies	13
6.1 Notes	13
6.2 Test setup.....	13
6.2.1 General.....	13
6.2.2 Setup for terminal	14
6.2.2.1 Hands-free measurements	14
6.2.2.2 Measurements in loudspeaking mode	19
6.2.3 Test signal levels.....	19
6.2.3.1 Send.....	19
6.2.3.2 Receive.....	20
6.2.4 Setup of background noise simulation.....	20
6.3 Coding independent parameters	21
6.3.1 Send frequency response	21
6.3.2 Send Loudness Rating (SLR).....	22
6.3.3 Mic mute.....	22
6.3.4 Send distortion	22
6.3.5 Out-of-band signals in send direction	23
6.3.6 Send noise.....	24
6.3.7 Terminal Coupling Loss weighted (TCLw).....	24
6.3.8 Stability loss.....	24
6.3.9 Receive frequency response.....	25
6.3.10 Receive Loudness Rating (RLR)	27
6.3.11 Receive Distortion	28
6.3.12 Out-of-band signals in receive direction	29
6.3.13 Receive noise	29
6.3.14 Double talk performance	30
6.3.14.1 General	30
6.3.14.2 Attenuation range in send direction during double talk $A_{H,S,dt}$	30
6.3.14.3 Attenuation range in receive direction during double talk $A_{H,R,dt}$	31
6.3.14.4 Detection of echo components during double talk	32
6.3.14.5 Minimum activation level and sensitivity of double talk detection.....	33

6.3.15	Switching characteristics	33
6.3.15.1	Note.....	33
6.3.15.2	Activation in send direction	33
6.3.15.3	Silence suppression and comfort noise generation.....	34
6.3.16	Background noise performance	34
6.3.16.1	Performance in send direction in the presence of background noise.....	34
6.3.16.2	Speech quality in the presence of background noise	35
6.3.16.3	Quality of background noise transmission (with far end speech).....	36
6.3.17	Quality of echo cancellation	36
6.3.17.1	Temporal echo effects	36
6.3.17.2	Spectral echo attenuation	37
6.3.17.3	Occurrence of artefacts	37
6.3.17.4	Variable echo path.....	37
6.3.18	Variant impairments; network dependant	38
6.3.18.1	Clock accuracy send.....	38
6.3.18.2	Clock accuracy receive	39
6.3.18.3	Send delay variation	39
6.3.19	Send and receive delay - round trip delay	40
6.4	Codec specific requirements.....	42
6.4.1	Objective listening speech quality MOS-LQO in send direction.....	42
6.4.2	Objective listening quality MOS-LQO in receive direction	42
6.4.3	Quality of jitter buffer adjustment	44
Annex A (informative):	Processing delays in VoIP terminals	46
Annex B (informative):	Bibliography.....	49
History		50

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 (<https://ipr.etsi.org/>).

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 and multimedia Transmission Quality (STQ).

Modal verbs terminology

In the present document "**shall**", "**shall not**", "**should**", "**should not**", "**may**", "**need not**", "**will**", "**will not**", "**can**" and "**cannot**" are to be interpreted as described in clause 3.2 of the [ETSI Drafting Rules](#) (Verbal forms for the expression of provisions).

"**must**" and "**must not**" are **NOT** allowed in ETSI deliverables except when used in direct citation.

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 harmonised 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 (Recommendation ITU-T G.711 [7]) speech algorithms.

The advanced signal processing of terminals is targeted to speech signals. Therefore, wherever possible speech signals are used for testing in order to achieve mostly realistic test conditions and meaningful results.

The present document provides 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 audio terminals.

In contrast to other standards which define minimum performance requirements it is the intention of the present document to specify terminal equipment requirements which enable manufacturers and service providers to enable good quality end-to-end speech performance as perceived by the user.

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

2.1 Normative references

References are either specific (identified by date of publication and/or edition number or version number) or non-specific. For specific references, only the cited version applies. For non-specific references, the latest version of the referenced document (including any amendments) applies.

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.

The following referenced documents are necessary for the application of the present document.

- [1] 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".
- [2] ETSI EN 300 726: "Digital cellular telecommunications system (Phase 2+) (GSM); Enhanced Full Rate (EFR) speech transcoding (GSM 06.60)".
- [3] ETSI TS 126 171: "Digital cellular telecommunications system (Phase 2+); Universal Mobile Telecommunications System (UMTS); LTE; Speech codec speech processing functions; Adaptive Multi-Rate - Wideband (AMR-WB) speech codec; General description (3GPP TS 26.171)".
- [4] Recommendation ITU-T G.108: "Application of the E-model: A planning guide".
- [5] Recommendation ITU-T G.109: "Definition of categories of speech transmission quality".
- [6] Recommendation ITU-T G.122: "Influence of national systems on stability and talker echo in international connections".
- [7] Recommendation ITU-T G.711: "Pulse code modulation (PCM) of voice frequencies".
- [8] Recommendation ITU-T G.723.1: "Dual rate speech coder for multimedia communications transmitting at 5.3 and 6.3 kbit/s".
- [9] Recommendation ITU-T G.726: "40, 32, 24, 16 kbit/s Adaptive Differential Pulse Code Modulation (ADPCM)".
- [10] Recommendation ITU-T G.729: "Coding of speech at 8 kbit/s using conjugate-structure algebraic-code-excited linear prediction (CS-ACELP)".

- [11] Recommendation ITU-T G.729.1: "G.729-based embedded variable bit-rate coder: An 8-32 kbit/s scalable wideband coder bitstream interoperable with G.729".
- [12] Recommendation ITU-T O.41: "Psophometer for use on telephone-type circuits".
- [13] Recommendation ITU-T P.50: "Artificial voices".
- [14] Recommendation ITU-T P.56: "Objective measurement of active speech level".
- [15] Recommendation ITU-T P.58: "Head and torso simulator for telephony".
- [16] Recommendation ITU-T P.79: "Calculation of loudness ratings for telephone sets".
- [17] Recommendation ITU-T P.310: "Transmission characteristics for narrow-band digital handset and headset telephones".
- [18] Recommendation ITU-T P.340: "Transmission characteristics and speech quality parameters of hands-free terminals".
- [19] Recommendation ITU-T P.342: "Transmission characteristics for narrow-band digital loudspeaking and hands-free telephony terminals".
- [20] Recommendation ITU-T P.501: "Test signals for use in telephony".
- [21] Recommendation ITU-T P.502: "Objective test methods for speech communication systems using complex test signals".
- [22] Recommendation ITU-T P.581: "Use of head and torso simulator for hands-free and handset terminal testing".
- [23] IEC 61260-1: "Electroacoustics - Octave-band and fractional-octave-band filters - Part 1: Specifications".
- [24] Recommendation ITU-T P.800.1: "Mean Opinion Score (MOS) terminology".
- [25] ETSI TS 103 224: "Speech and multimedia Transmission Quality (STQ); A sound field reproduction method for terminal testing including a background noise database".
- [26] Recommendation ITU-T P.863.1: "Application guide for Recommendation ITU-T P.863".
- [27] Recommendation ITU-T P.863: "Perceptual objective listening quality assessment".
- [28] ETSI ES 202 737: "Speech and multimedia Transmission Quality (STQ); Transmission requirements for narrowband VoIP terminals (handset and headset) from a QoS perspective as perceived by the user".
- [29] Recommendation ITU-T P.1010: "Fundamental voice transmission objectives for VoIP terminals and gateways".
- [30] IETF RFC 3550: "RTP: A Transport Protocol for Real-Time Applications".

2.2 Informative references

References are either specific (identified by date of publication and/or edition number or version number) or non-specific. For specific references, only the cited version applies. For non-specific references, the latest version of the referenced document (including any amendments) applies.

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

The following referenced documents are not necessary for the application of the present document but they assist the user with regard to a particular subject area.

- [i.1] ETSI EG 202 425: "Speech Processing, Transmission and Quality Aspects (STQ); Definition and implementation of VoIP reference point".

- [i.2] ETSI EG 202 396-3: "Speech and multimedia Transmission Quality (STQ); Speech Quality performance in the presence of background noise; Part 3: Background noise transmission - Objective test methods".
- [i.3] NIST Net.
- NOTE: Available at <http://snad.ncsl.nist.gov/itg/nistnet/>.
- [i.4] Netem.
- NOTE: Available at <http://www.linuxfoundation.org/en/Net:Netem>.
- [i.5] ETSI EG 201 377-1: "Speech and multimedia Transmission Quality (STQ); Specification and measurement of speech transmission quality; Part 1: Introduction to objective comparison measurement methods for one-way speech quality across networks".

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

NOTE: In case of a head and torso simulator (HATS) in the centre 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 Recommendation ITU-T P.58 [15]: "n" is one of the points numbered from 11 to 17 and defined in table 6a of Recommendation ITU-T P.58 [15] (coordinates of far field front point)

NOTE: The HATS HFRP depends on the location(s) of the microphones of the terminal under test: the appropriate axis lip-ring/HATS HFRP is to 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): point 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:

AM-FM	Amplitude Modulation - Frequency Modulation
AMR	Adaptative Multi-Rate
CS	Composite Source
CSS	Composite Source Signal
DRP	ear Drum Reference Point
DUT	Device Under Test
EC	Echo Cancellor
EFR	Enhanced Full Rate
EL	Echo Loss
ERP	Ear Reference Point
ETH	Eidgenössische Technische Hochschule
FFT	Fast Fourier Transform
GSM	Global System for Mobile communications
HATS	Head And Torso Simulator
HFRP	Hands Free Reference Point
HFT	HandsFree Terminal
IEC	International Electrotechnical Commission
IP	Internet Protocol
IPDV	IP Packet Delay Variation
ITU-T	International Telecommunication Union - Telecommunication standardization sector
LE	Earphone coupling Loss
MOS	Mean Opinion Score
MOS-LQOy	Mean Opinion Score - Listening Quality Objective

NOTE: y being N for narrow-band, W for wideband, M for mixed and S for superwideband. See Recommendation ITU-T P.800.1 [24].

MRP	Mouth Reference Point
NIST	National Institute of Standards and Technology
NIST Net	Network Simulation Tool from National Institute of Standards and Technology
NLP	Non Linear Processor
PBX	Private Branch eXchange
PC	Personal Computer
PCM	Pulse Code Modulation
PDA	Personal Digital Assistance
PMRP	Sound Pressure at the Mouth Reference Point
PN	PseudoNoise
POI	Point Of Interconnect
PSTN	Public Switched Telephone Network
QoS	Quality of Service
RLR	Receive Loudness Rating
RLR max	Receive Loudness Rating corresponding to the maximum setting of the volume control
RLR min	Receive Loudness Rating corresponding to the minimum setting of the volume control
RMS	Root Mean Square
RTP	Real Time Protocol
SLR	Send Loudness Rating
TCL _w	Terminal Coupling Loss (weighted)
TCN	Trace Control for Netem
TDM	Time Division Multiplex
TOSQA	Telecommunication Objective Speech Quality Assessment
VAD	Voice Activity Detector
VoIP	Voice over Internet Protocol
xDSL	any Digital Subscriber Line technology

4 General considerations

4.1 Default Coding Algorithm

VoIP terminals shall support the coding algorithm according to Recommendation ITU-T G.711 [7] (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 Recommendation ITU-T G.711 [7] 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 Recommendation ITU-T G.109 [5].

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 Recommendation ITU-T G.108 [4] 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 document is acceptable or representative for the specific configuration is the responsibility of the individual transmission planner.

Recommendation ITU-T G.108 [4] 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 (ETSI I-ETS 300 245-3 [1]) parameters are applicable, if Recommendation ITU-T G.711 [7] is used.

5 Test equipment

5.1 IP half channel measurement adaptor

The IP half channel measurement adaptor is described in ETSI EG 202 425 [i.1].

5.2 Environmental conditions for tests

The following conditions shall apply for the testing environment:

- a) Ambient temperature: 15 °C to 35 °C (inclusive).
- b) Relative humidity: 5 % to 85 %.
- c) Air pressure: 86 kPa to 106 kPa (860 mbar to 1 060 mbar).

5.3 Accuracy of measurements and test signal generation

Unless specified otherwise, the accuracy of measurements made by test equipment shall be equal to or better than:

Table 1: Measurement Accuracy

Item	Accuracy
Electrical signal level	$\pm 0,2$ dB for levels ≥ -50 dBV $\pm 0,4$ dB for levels < -50 dBV
Sound pressure	$\pm 0,7$ dB
Frequency	$\pm 0,2$ %
Time	$\pm 0,2$ %
Application force	± 2 N
Measured maximum frequency	20 kHz
NOTE: The measured maximum frequency is due to P.58 limitations.	

Unless specified otherwise, the accuracy of the signals generated by the test equipment shall be better than:

Table 2: Accuracy of test signal generation

Quantity	Accuracy
Sound pressure level at Mouth Reference Point (MRP)	± 3 dB for frequencies from 100 Hz to 200 Hz ± 1 dB for frequencies from 200 Hz to 4 000 Hz ± 3 dB for frequencies from 4 000 Hz to 8 000 Hz
Electrical excitation levels	$\pm 0,4$ dB across the whole frequency range
Frequency generation	± 2 % (see note)
Time	$\pm 0,2$ %
Specified component values	± 1 %
NOTE: This tolerance may be used to avoid measurements at critical frequencies, e.g. those due to sampling operations within the terminal under test.	

For terminal equipment which is directly powered from the mains supply, all tests shall be carried out within ± 5 % of the rated voltage of that supply. If the equipment is powered by other means and those means are not supplied as part of the apparatus, all tests shall be carried out within the power supply limit declared by the supplier. If the power supply is a.c., the test shall be conducted within ± 4 % of the rated frequency.

5.4 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 NIST Net [i.3] or Netem [i.4].

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 easy to adapt devices to these 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 the exact same distortions with Netem.

NOTE: NIST Net™, NETEM™, Linux™ and X Window System™ are examples of suitable products available commercially. This information is given for the convenience of users of the present document and does not constitute an endorsement by ETSI of these product(s).

5.5 Acoustic environment

Unless stated otherwise 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 shall 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, the test laboratory, in the case where its test room does not conform to anechoic conditions as given in Recommendation ITU-T P.342 [19], has to present difference in results for measurements due to its test room.

Standardized measurement methods for measurements with variable echo paths are for further study.

5.6 Influence of terminal delay on measurements

As delay is introduced by the terminal, care shall be taken for all measurements where exact position of the analysis window is required. It shall be checked that the test is performed on the test signal and not on any other signal.

6 Requirements and associated measurement methodologies

6.1 Notes

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. This will be indicated in the test report.

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.

6.2 Test setup

6.2.1 General

In order to use a compatible test system for all types of speech terminals a HATS (Head and Torso Simulator) will be used instead of freefield microphone (for receive measurement) and artificial mouth (for Send measurement). HATS is described in Recommendation ITU-T P.58 [15].

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, it should be adopted, for testing terminal electroacoustical parameters, the bit rate recognized as giving the best characteristics is selected, e.g.:

- ETSI TS 126 171 [3]: 12,2 kbit/s.

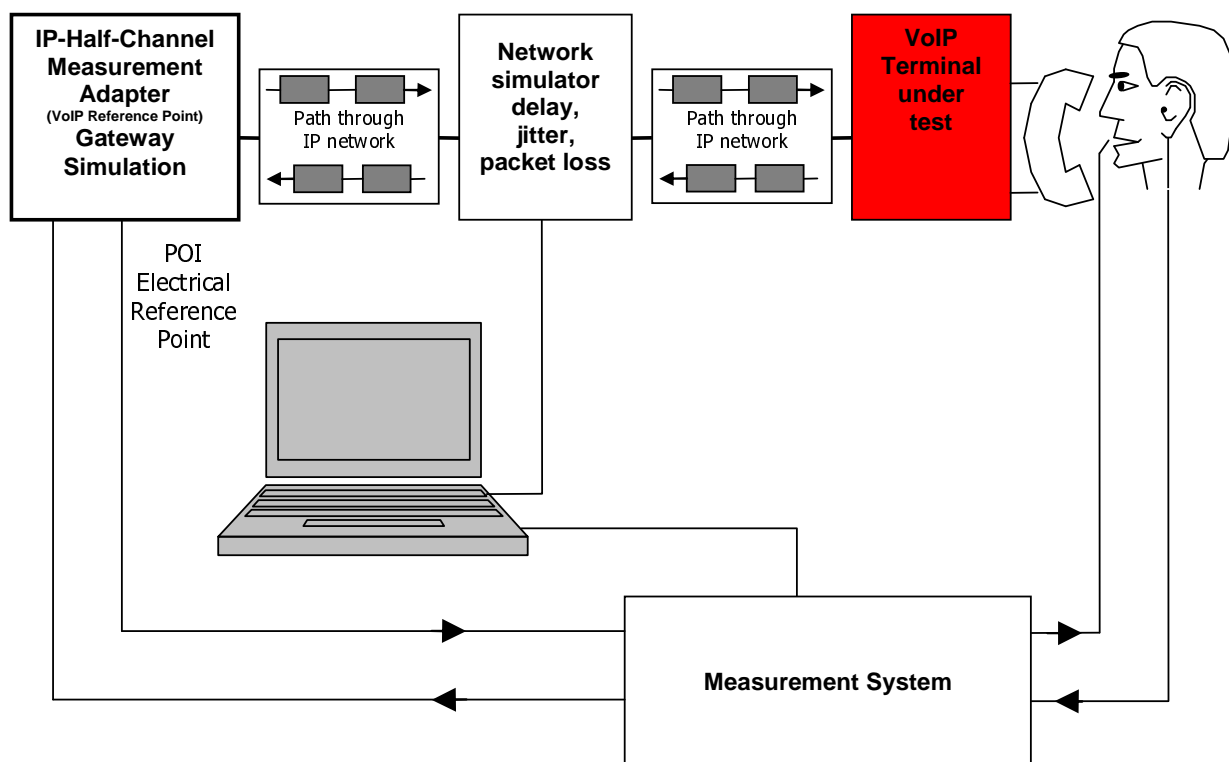


Figure 1: Half channel terminal measurement

6.2.2 Setup for terminal

6.2.2.1 Hands-free measurements

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

Desktop operated hands-free terminal

For HATS test equipment, the definition of hands-free terminals and setups for hands-free terminal can be found in Recommendation ITU-T P.581 [22].

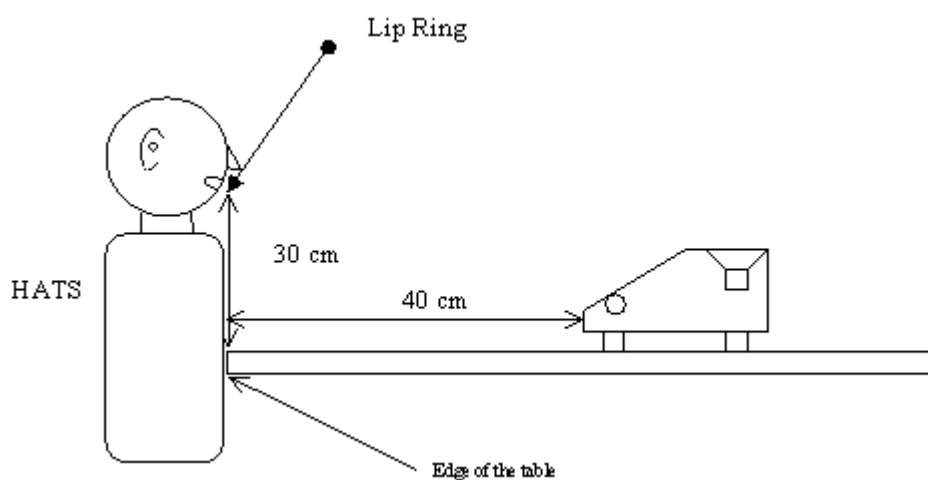


Figure 2: Position for test of desktop hands free terminal side view

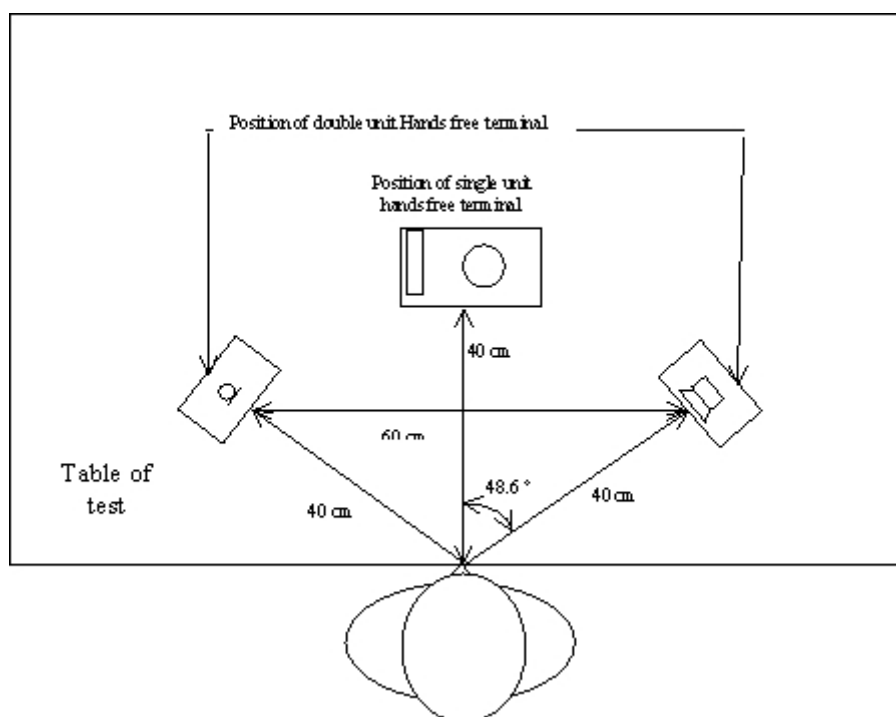


Figure 3: Position for test of desktop hands free terminal top sight

Handheld hands-free terminal

It should be placed in accordance 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.

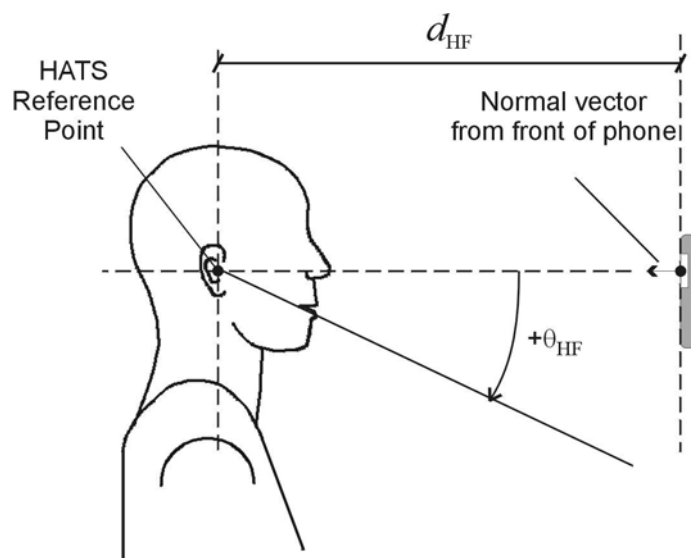


Figure 4: Configuration of Hand-Held loudspeaker relative to the HATS side view

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 receive measurement, d_{HFS} is the distance for Send measurement, and d_{EM} is the distance from ERP to MRP.

When no operating distance is specified by manufacturer, value for d_{HFS} 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 to the following conditions:

Softphone including speakers and microphone

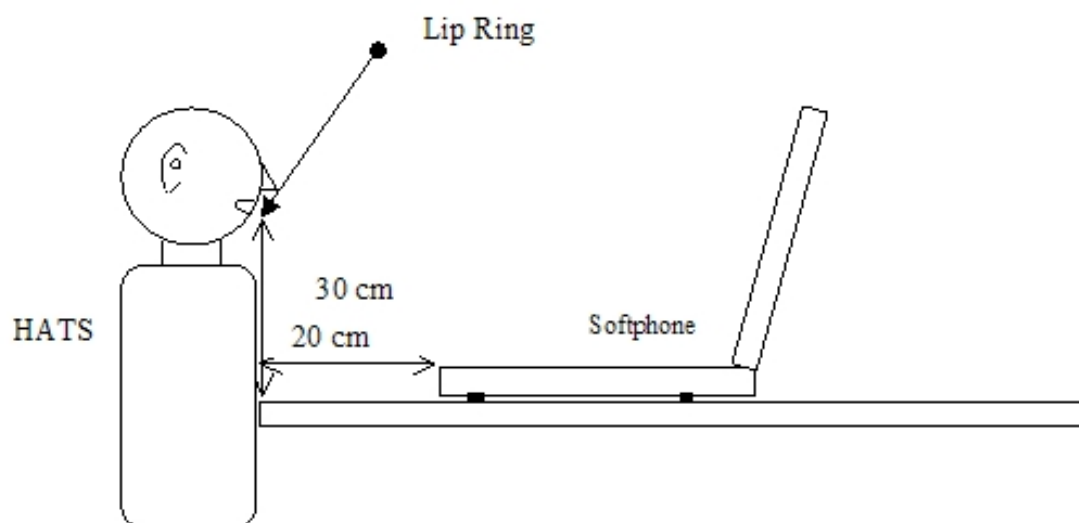


Figure 5: Configuration of softphone relative to the HATS side view

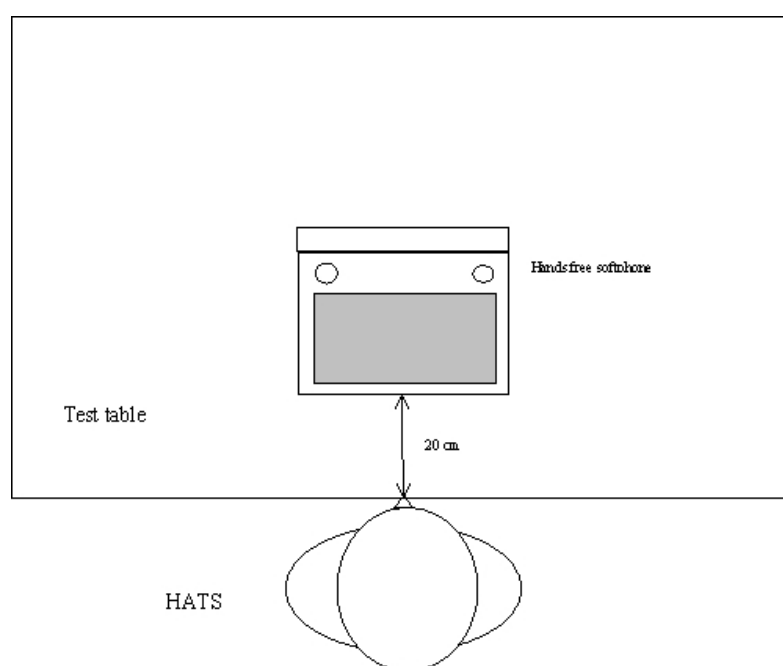


Figure 6: Configuration of softphone relative to the HATS top sight

Softphone with separate speakers

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

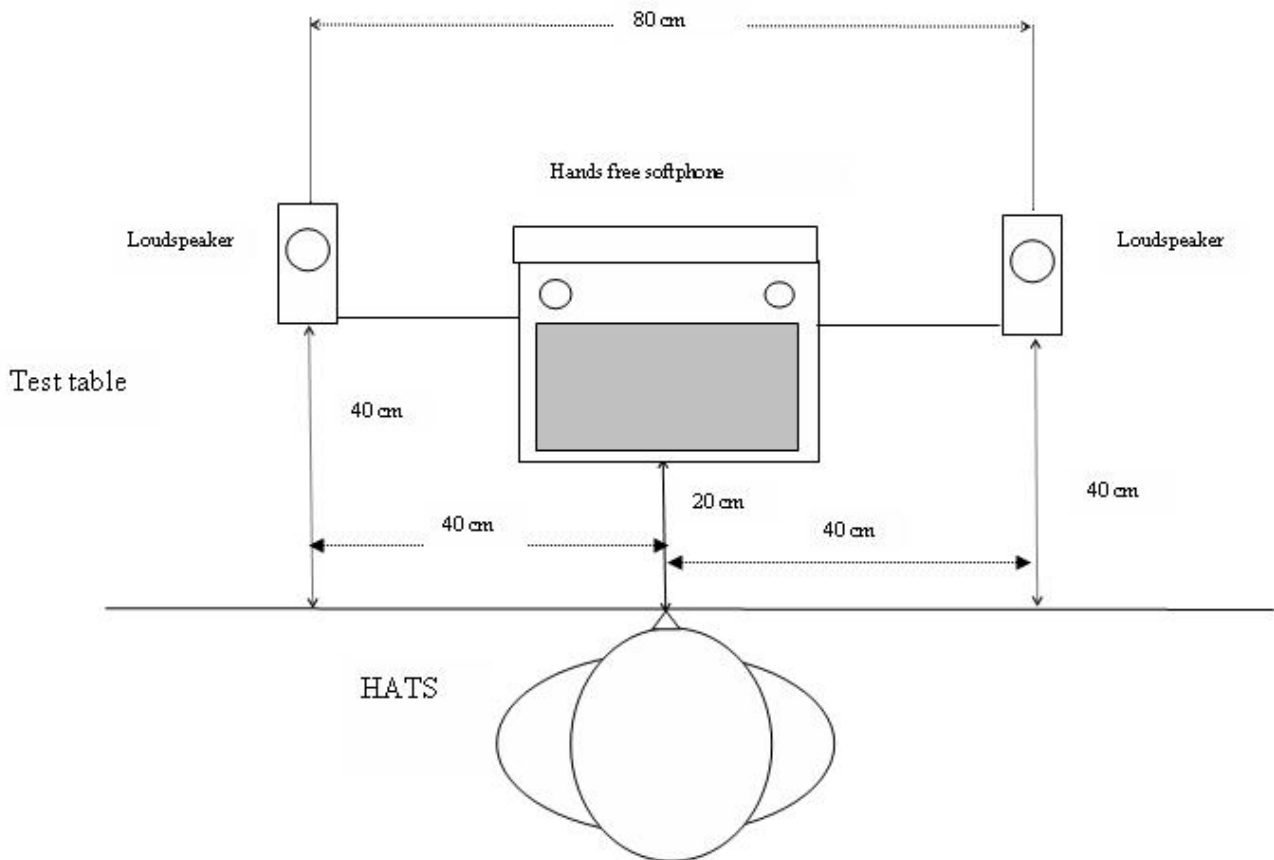


Figure 7: Configuration of softphone using external speakers relative to the HATS top sight

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

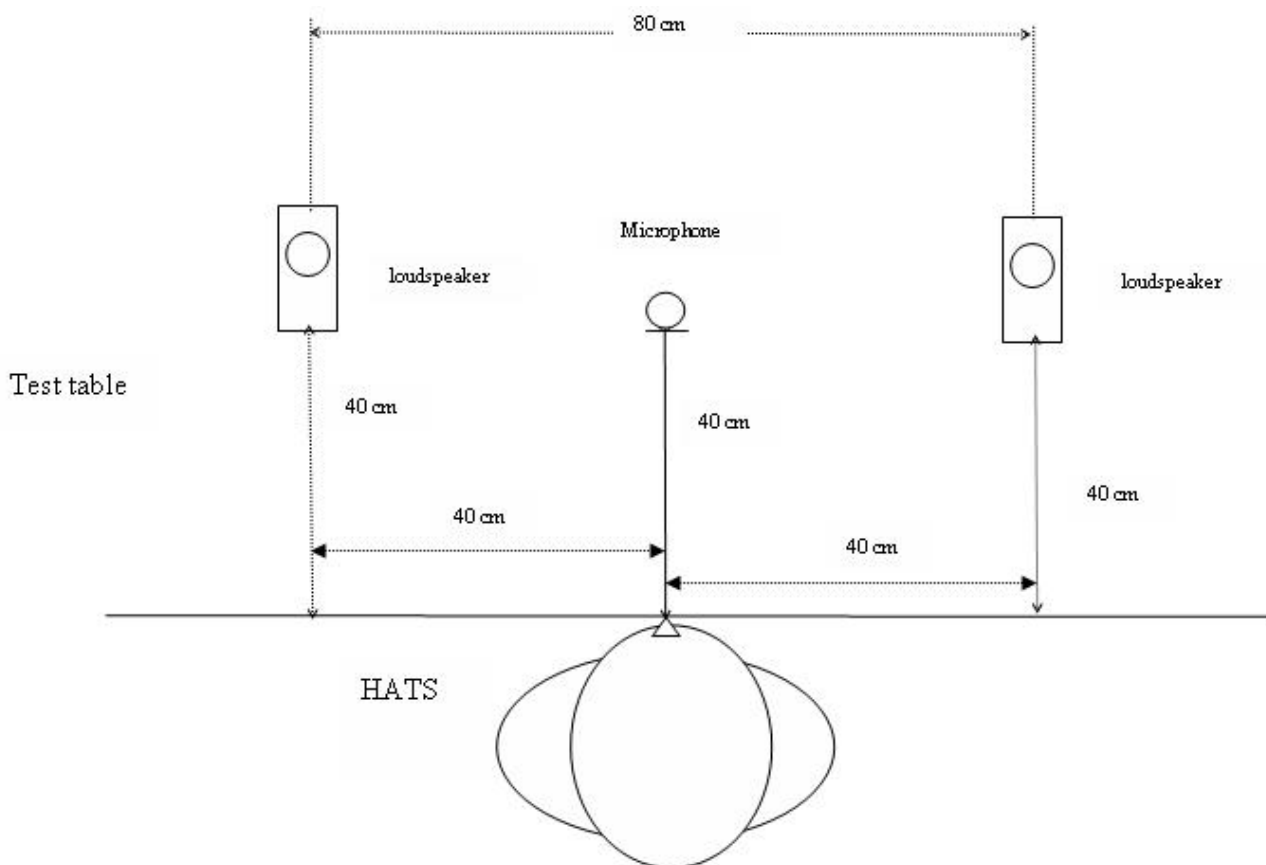


Figure 8: Configuration of softphone using external speakers and microphone relative to the HATS top sight

Group audio 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.

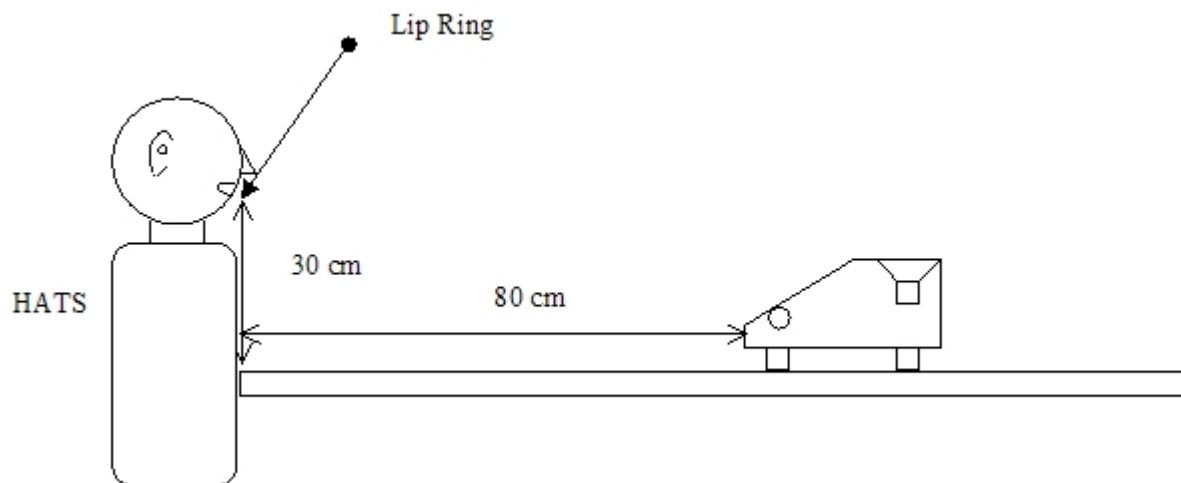


Figure 9: Configuration of group audio terminal relative to the HATS side view

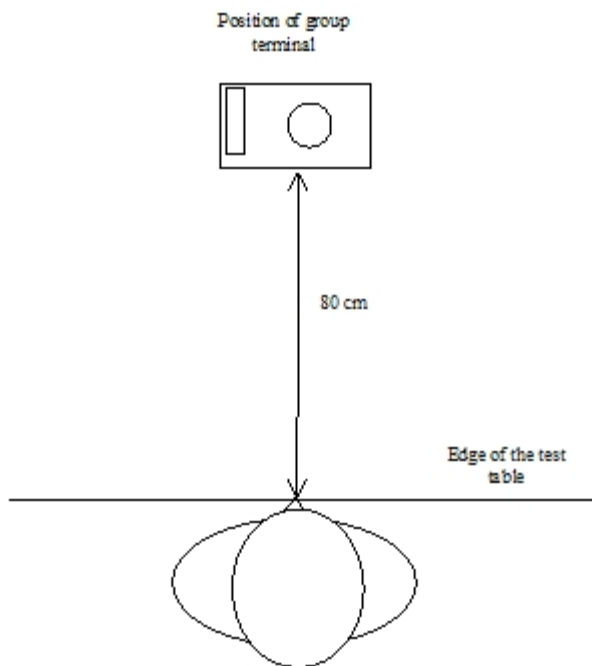


Figure 10: Configuration of group audio terminal relative to the HATS top sight

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.

6.2.2.2 Measurements in loudspeaking mode

For those measurements HATS will be used.

It will be positioned as defined in clause 6.2.2 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. For the handset 8N application force shall be used.

NOTE: Only desktop terminals are concerned by loudspeaking measurement.

6.2.3 Test signal levels

6.2.3.1 Send

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 freefield conditions at the MRP. The total level on the frequency range is set to -4,7 dBPa.
- The spectrum at MRP is recorded. This spectrum is used as the reference for the send characteristics.
- The spectrum at the HATS HFRP is recorded.
- The spectrum at HATS HFRP is calibrated to the nominal spectrum given for the relevant HATS HFRP in Recommendation ITU-T P.58 [15], table 7d and table 7e.
- 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 Send 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 shall be calculated as follows:

- **SmJ** = $20 \log V_s - 20 \log \text{PMRP}$.

where:

- **V_s** is the measured voltage across the appropriate termination (unless stated otherwise, a 600 Ω termination).
- **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.

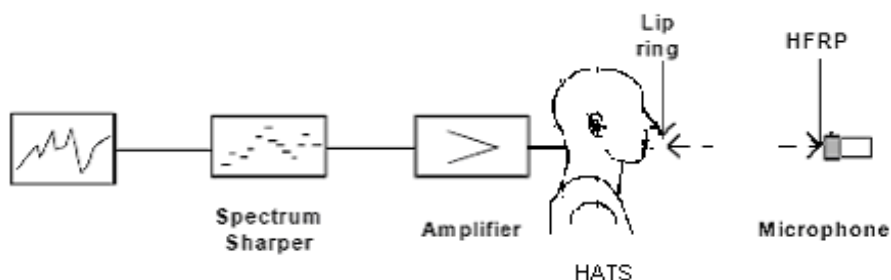


Figure 11: Calibration at HFRP

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 audio terminal: 85 cm with -33,3 dBPa.

6.2.3.2 Receive

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 freefield equalized.

6.2.4 Setup of background noise simulation

A setup for simulating realistic background noises in a lab-type environment is described in ETSI TS 103 224 [25].

If not stated otherwise this setup is used in all measurements where background noise simulation is required. The following noises of ETSI TS 103 224 [25] shall be used.

Table 2a

Pub Noise (Pub)	HATS and microphone array in a pub	30 seconds	1: 75,2 dB 2: 75,1 dB 3: 74,9 dB 4: 75,1 dB 5: 74,8 dB 6: 74,8 dB 7: 74,8 dB 8: 75,0 dB
Sales Counter (SalesCounter)	HATS and microphone array in a supermarket	30 seconds	1: 65,5 dB 2: 65,3 dB 3: 65,2 dB 4: 65,5 dB 5: 65,6 dB 6: 65,3 dB 7: 65,2 dB 8: 65,3 dB
Callcenter 2 (Callcenter)	HATS and microphone array in business office	30 seconds	1: 59,3 dB 2: 59,3 dB 3: 59,5 dB 4: 59,6 dB 5: 59,4 dB 6: 59,3 dB 7: 59,3 dB 8: 59,5 dB

6.3 Coding independent parameters

6.3.1 Send frequency response

Requirement

The Send sensitivity/frequency response shall be within the limits given in table 3.

Table 3

Frequency	Upper limit	Lower limit
100 Hz	0 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	

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

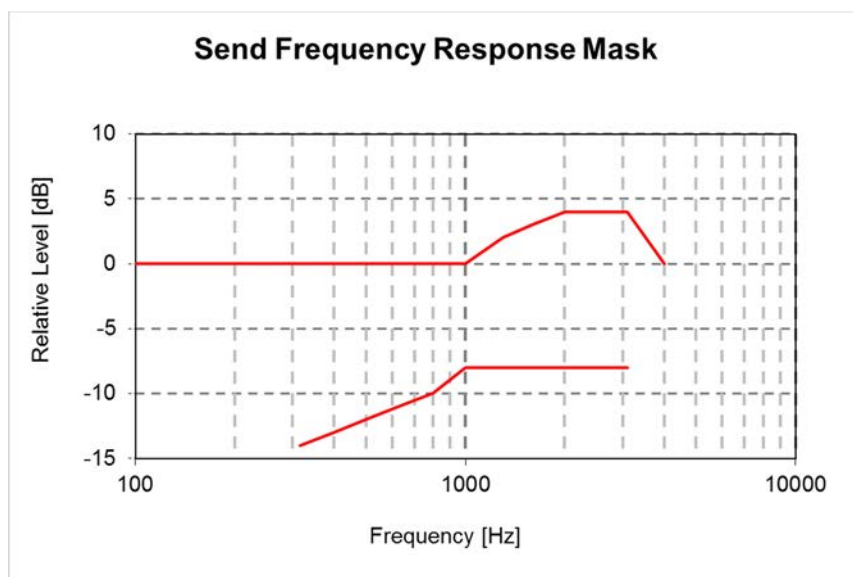


Figure 12: Send frequency response mask for HFT

Measurement method

The terminal will be positioned as described in clause 6.2.

The test signal to be used for the measurements shall be the British-English single talk sequence described in clause 7.3.2 of Recommendation ITU-T P.501 [20]. The spectrum of acoustic signal produced by the artificial mouth is calibrated under freefield conditions at the MRP. The signal level is adjusted according to clause 6.2.3.1.

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

Measurements shall be made at one third-octave bands as given by the IEC 61260-1 [23] 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 dBV/Pa.

6.3.2 Send Loudness Rating (SLR)

Requirement

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

This value is derived from Recommendation ITU-T P.310 [17]. According to Recommendation ITU-T P.340 [18] 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 6.2).

Measurement method

The terminal will be positioned as described in clause 6.2.

For a correct activation of the system, the test signal to be used for the measurements shall be the British-English single talk sequence described in clause 7.3.2 of Recommendation ITU-T P.501 [20]. The spectrum of acoustic signal produced by the artificial mouth is calibrated under freefield 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.

Calibration is realized as explained in clause 6.2.3.1.

SLR shall be calculated according Recommendation ITU-T P.79 [16].

6.3.3 Mic mute

Requirement

The SLR (Send Loudness Rating) with mic mute on shall be 50 dB higher than with mic mute off.

Measurement method

The terminal will be positioned as described in clause 6.2.

For a correct activation of the system, the test signal to be used for the measurements shall be the British-English single talk sequence described in clause 7.3.2 of Recommendation ITU-T P.501 [20]. The spectrum of acoustic signal produced by the artificial mouth is calibrated under freefield 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.

Calibration is realized as explained in clause 6.2.3.1.

The send sensitivity shall be calculated from each band of the 20 frequencies given in table 1 of Recommendation ITU-T P.79 [16], 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 Recommendation ITU-T P.79 [16], annex A.

6.3.4 Send distortion

Requirement

The terminal will be positioned as described in clause 6.2.

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

Table 4

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.	

Measurement method

The terminal will be positioned as described in clause 6.2.

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 second. 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.

The female speaker signal of the short conditioning sequence described in clause 7.3.7 of Recommendation ITU-T P.501 [20] shall be used for activation. The level of this activation signal shall be -4,7 dBPa at the MRP.

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

6.3.5 Out-of-band signals in send direction

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 5.

Table 5: Out-of-band signal limit, send

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.	

Measurement method

The terminal will be positioned as described in clause 6.2.

The female speaker of the short conditioning sequence described in clause 7.3.7 of Recommendation ITU-T P.501 [20] shall be used for activation. The 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.

6.3.6 Send noise

Requirement

The limit for the Send 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 dependent of activity of system.

Measurement method

The terminal will be positioned as described in clause 6.2.

The female speaker of the short conditioning sequence described in clause 7.3.7 of Recommendation ITU-T P.501 [20] shall be used for activation. The 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 Recommendation ITU-T O.41 [12].

Spectral peaks are measured in the frequency domain in the frequency range from 100 Hz to 3,4 kHz. The frequency spectrum of the idle channel noise is measured by a spectral analysis having a noise bandwidth of 8,79 Hz (determined using FFT 8 k samples/48 kHz sampling rate with Hanning window or equivalent). The idle channel noise spectrum is stated in dB. A smoothed average idle channel noise spectrum is calculated by a moving average (arithmetic mean) 1/3rd octave wide across the idle noise channel spectrum stated in dB (linear average in dB of all FFT bins in the range from $2^{-(1/6)}f$ to $2^{+(1/6)}f$). Peaks in the idle channel noise spectrum are compared against a smoothed average idle channel noise spectrum.

6.3.7 Terminal Coupling Loss weighted (TCLw)

Requirement

TCLw shall be greater than 46 dB.

TCLw shall be not less than 40 dB for any setting of the volume control.

NOTE: A $TCLw \geq 50$ dB is recommended as a performance objective. Depending on the idle channel noise in the sending direction, it may not always be possible to measure an echo loss ≥ 50 dB.

Measurement method

The setup for terminal is described in clause 6.2.

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

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

The test signal is the compressed real speech signal described in clause 7.3.3 of Recommendation ITU-T P.501 [20].

The TCLw is calculated according to Recommendation ITU-T G.122 [6], 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.

6.3.8 Stability loss

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

Measurement method

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

For loudspeaking mode handset is placed at 50 cm beside terminal with transducers facing the table as in figure 13.

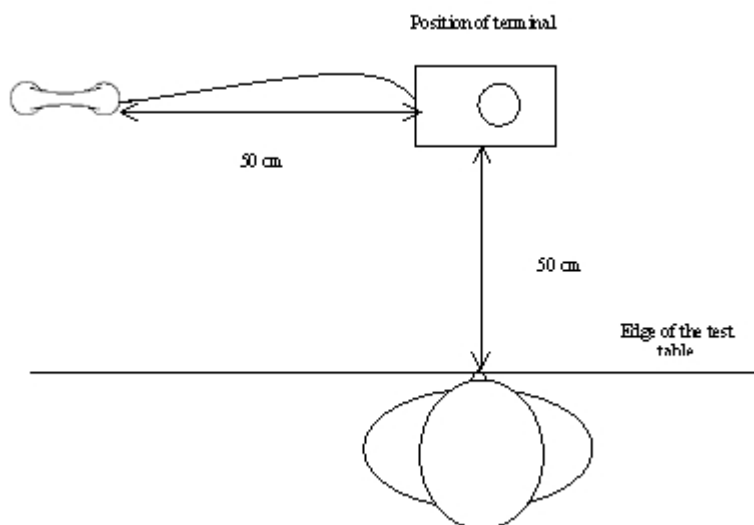


Figure 13: Stability loss position for loudspeaking function

Before the actual test a training sequence consisting of the British-English single talk sequence described in clause 7.3.2 of Recommendation ITU-T P.501 [20] shall be applied. The training sequence level shall be -16 dBm0 in order not to overload the codec.

The test signal is a PN sequence complying with Recommendation ITU-T P.501 [20] 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 200 Hz to 4 kHz.

6.3.9 Receive frequency response

Requirement

The following masks are required for handsfree and loudspeaking terminals.

Desktop operated loudspeaker

Table 6: Receive frequency response mask-desktop

Frequency	Upper limit	Lower limit
100 Hz	6 dB	
315 Hz	6 dB	-9 dB
400 Hz	6 dB	-6 dB
3 150 Hz	6 dB	-6 dB
4 000 Hz	6 dB	
NOTE: The limits for intermediate frequencies lie on a straight line drawn between the given values on a linear (dB) - logarithmic (kHz) scale.		

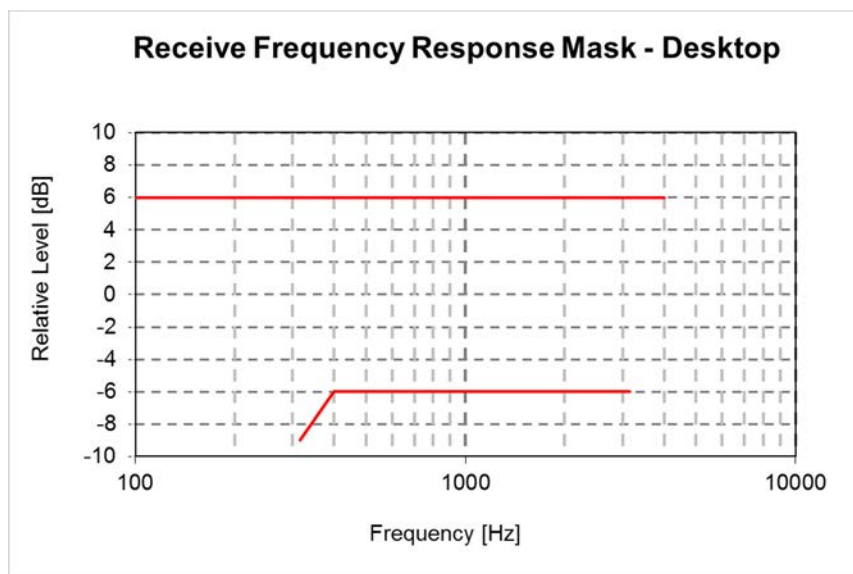


Figure 14: Receive frequency response mask for Desktop hands free terminals

Handheld terminal

Table 7: Receive frequency response mask-handheld

Frequency	Upper limit	Lower limit
100 Hz	6 dB	
500 Hz	6 dB	-9 dB
630 Hz	6 dB	-6 dB
3 150 Hz	6 dB	-6 dB
4 000 Hz	6 dB	
NOTE: The limits for intermediate frequencies lie on a straight line drawn between the given values on a linear (dB) - logarithmic (kHz) scale.		

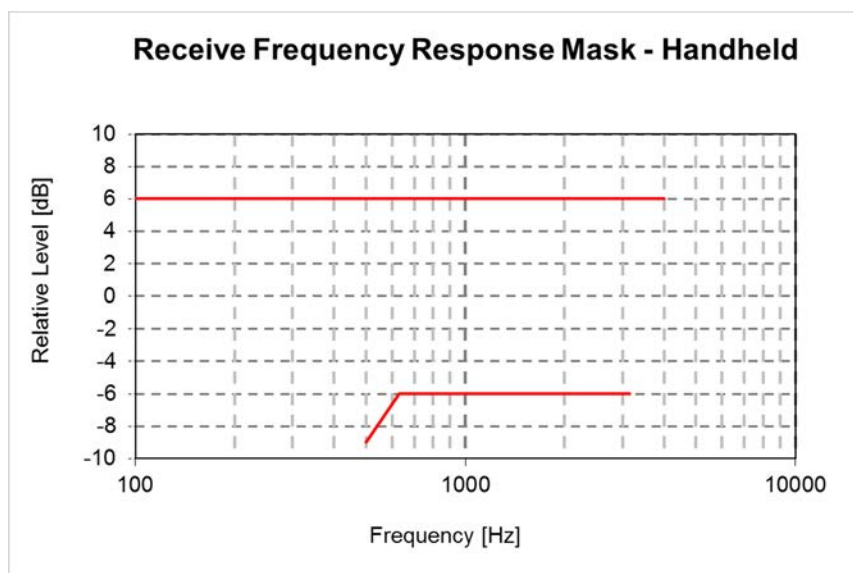


Figure 15: Receive frequency response mask for Hand-held HFTs

Softphone (computer-based terminals)

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

Type 2 requirement as for handheld terminal.

Group audio terminal

Same requirement as desktop terminals.

Measurement method

Test setup is described in clause 6.2.

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).

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

Where:

- S_{Jeff} Receive Sensitivity; Junction to HATS Ear with freefield correction.
- p_{eff} DRP Sound pressure measured by ear simulator Measurement data are converted from the Drum Reference Point to freefield.
- v_{RCV} Equivalent RMS input voltage.

The test signal to be used for the measurements shall be British-English single talk sequence described in clause 7.3.2 of Recommendation ITU-T P.501 [20]. The test signal level shall be -20 dBm0, measured according to Recommendation ITU-T P.56 [14] at the digital reference point or the equivalent analogue point.

The HATS is freefield equalized as described in Recommendation ITU-T P.581 [22]. The equalized output signal is power-averaged on the total time of analysis. The 1/3 octave band data are considered as the input signal to be used for calculations or measurements.

Measurements shall be made at one third-octave bands as given by the IEC 61260-1 [23] 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.

6.3.10 Receive Loudness Rating (RLR)

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 shall be less than (louder) or equal to -2 dB: $\text{RLR}_{\text{max}} \leq -2$ dB.

Range of volume control shall be equal or exceed 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 shall be less than (louder) or equal to 5 dB: $\text{RLR}_{\text{max}} \leq 5$ dB.

Range of volume control shall be equal or exceed 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 audio 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 shall be less than (louder) or equal to -6 dB: $\text{RLR max} \leq -6$ dB.

Range of volume control shall be equal or exceed 19 dB.

Measurement method

Test setup is described in clause 6.2.

The test signal to be used for the measurements shall be the British-English single talk sequence described in clause 7.3.2 of Recommendation ITU-T P.501 [20]. The test signal level shall be -20 dBm0, measured according to Recommendation ITU-T P.56 [14] at the digital reference point or the equivalent analogue point.

The receive sensitivity shall be calculated from each band of the 14 frequencies given in table 1 of Recommendation ITU-T P.79 [16], 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 shall be calculated according to the formula 5-1 of Recommendation ITU-T P.79 [16], using the receive weighting factors from table 1 and according to clause 6, of Recommendation ITU-T P.79 [16]. The RLR shall then be corrected as RLR minus 14 dB according to Recommendation ITU-T P.340 [18], and without the LE factor.

6.3.11 Receive Distortion

Requirement

Desktop operated loudspeaker

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

Table 8

Frequency	Signal to distortion ratio limit, receive for desktop terminal at nominal volume	Signal to distortion ratio limit, receive for handheld terminal at nominal volume	Signal to distortion ratio limit, receive 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 6.2.

The ratio of signal to harmonic distortion is given in table 8.

Softphone (computer-based terminal)

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

Type 2 requirement as for handheld terminal.

Group audio terminal

Same requirement as for desktop terminal.

Measurement method

Test setup is described in clause 6.2.

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 second. Appropriate signals for activation and signal combinations can be found in Recommendation ITU-T P.501 [20]. The sinusoidal signal level shall be calibrated to -16 dBm0.

The female speaker signal of the short conditioning sequence described in clause 7.3.7 of Recommendation ITU-T P.501 [20] shall be used for activation. Level of this activation signal shall 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.

6.3.12 Out-of-band signals in receive direction

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 9.

Table 9: Out-of-band signal limit, receive

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.	

Measurement method

Test setup is described in clause 6.2.

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.

The female speaker signal of the short conditioning sequence described in clause 7.3.7 of Recommendation ITU-T P.501 [20] shall be used for activation. Level of this activation signal shall be -16 dBm0.

6.3.13 Receive noise

Requirement

The noise level measured until 10 kHz shall not exceed -54 dBPa(A) at **nominal setting of the volume control**.

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

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

Measurement method

The test setup is described in clause 6.2.

The A-weighted noise level shall be measured at DRP of the artificial ear with the freefield equalization active. The noise level is measured until 10 kHz.

The female speaker signal of the short conditioning sequence described in clause 7.3.7 of Recommendation ITU-T P.501 [20] shall be used for activation. Level of this activation signal shall be -16 dBm0.

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

Spectral peaks are measured in the frequency domain in the frequency range from 100 Hz to 3,4 kHz. The frequency spectrum of the idle channel noise is measured by a spectral analysis having a noise bandwidth of 8,79 Hz (determined using FFT 8 k samples/48 kHz sampling rate with Hanning window or equivalent). The idle channel noise spectrum is stated in dB. A smoothed average idle channel noise spectrum is calculated by a moving average (arithmetic mean) 1/3rd octave wide across the idle noise channel spectrum stated in dB (linear average in dB of all FFT bins in the range from $2^{(-1/6)}f$ to $2^{(+1/6)}f$). Peaks in the idle channel noise spectrum are compared against a smoothed average idle channel noise spectrum.

6.3.14 Double talk performance

6.3.14.1 General

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 Recommendations ITU-T P.340 [18] and P.502 [21]):

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

6.3.14.2 Attenuation range in send direction during double talk $A_{H,S,dt}$

Requirement

Based on the level variation in Send direction during double talk $A_{H,S,dt}$ the behaviour of the terminal can be classified according to table 10.

Table 10

Category (according to Recommendation ITU-T P.340 [18])	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]	≤ 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.

Measurement method

The test signal to determine the attenuation range during double talk is the double talk speech sequence as defined in clause 7.3.5 of Recommendation ITU-T P.501 [20] as shown in figure 16. The competing speaker is always inserted as the double talk sequence sdt(t) either in send or receive and is used for analysis.

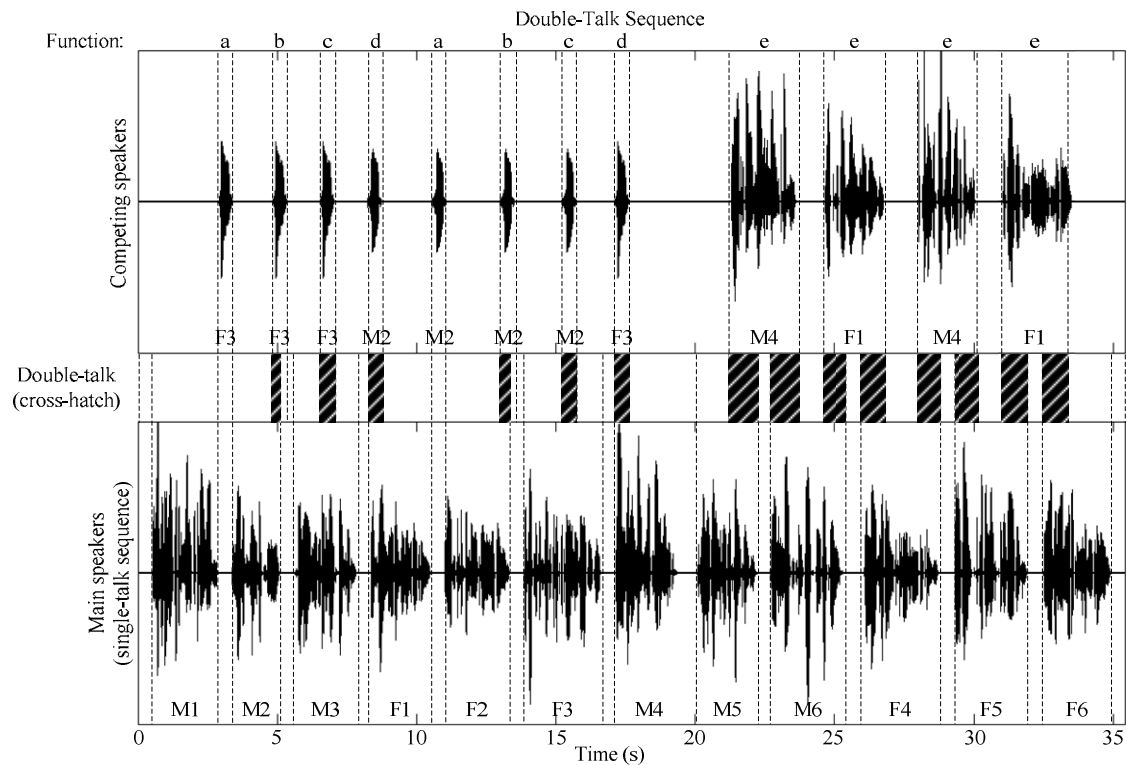


Figure 16: Double talk test sequence with overlapping speech sequences in send and receive direction

The attenuation range during double talk is determined as described in Appendix III of Recommendation ITU-T P.502 [21]. The double talk performance is analysed for each word and sentence produced by the competing speaker. The requirement has to be met for each word and sentence produced by the competing speaker.

6.3.14.3 Attenuation range in receive direction during double talk $A_{H,R,dt}$

Requirement

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

Table 11

Category (according to Recommendation ITU-T P.340 [18])	1	2a	2b	2c	3
	Full Duplex Capability	Partial Duplex Capability			No Duplex Capability
$A_{H,R,dt}$ [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.

Measurement method

Test setup is described in clause 6.2.

The test signal to determine the attenuation range during double talk is shown in figure 16. A sequence of speech signals is used which is inserted in parallel in Send and receive 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 attenuation range during double talk is determined as described in Appendix III of Recommendation ITU-T P.502 [21]. The double talk performance is analysed for each word and sentence produced by the competing speaker. The requirement has to be met for each word and sentence produced by the competing speaker.

6.3.14.4 Detection of echo components during double talk

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 12 are applicable (more information can be found in annex A of the Recommendation ITU-T P.340 [18]).

Table 12

Category (according to Recommendation ITU-T P.340 [18])	1	2a	2b	2c	3
	<i>Full Duplex Capability</i>	<i>Partial Duplex Capability</i>			<i>No Duplex Capability</i>
Echo Loss [dB]	≥ 27	≥ 23	≥ 17	≥ 11	< 11

NOTE 1: The echo attenuation during double talk is based on the parameter Talker Echo Loudness Rating (TEL_{dt}). It is assumed that the terminal at the opposite end of the connection provides nominal Loudness Rating ($SLR + RLR = 10$ dB).

Measurement method

Test setup is described in clause 6.2.

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 17. A detailed description can be found in Recommendation ITU-T P.501 [20].

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

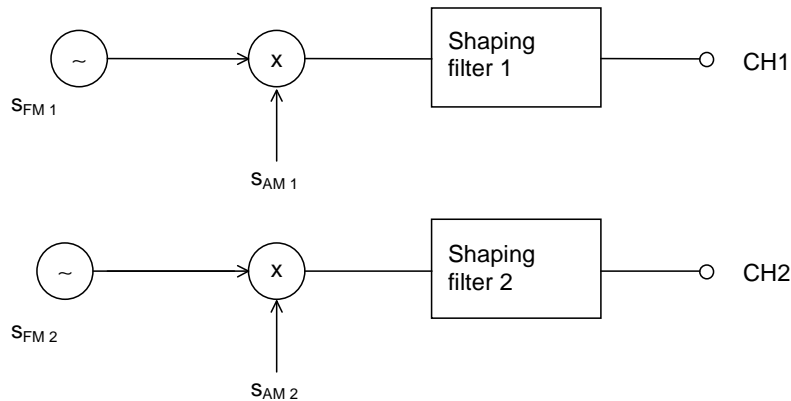


Figure 17: Measurement signals

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

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

NOTE 2: A is determined by the required test signal level as found in the individual test cases.

The settings for the signals are as follows.

Table 13: Parameters of the two Test Signals for Double Talk Measurement based on AM-FM modulated sine waves

Send Direction				Receive Direction		
f_m [Hz]	$f_{mod(fm)}$ [Hz]	F_{am} [Hz]		f_m [Hz]	$f_{mod(fm)}$ [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		2 650	± 35	3
2 750	± 40	3		2 900	± 35	3
3 000	± 40	3		3 150	± 35	3
3 250	± 40	3		3 400	± 35	3
3 500	± 40	3		3 650	± 35	3
3 750	± 40	3		3 900	± 35	3

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

The test signal is measured at the electrical reference point (Send 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 receive direction (see Recommendation ITU-T P.501 [20]). The filter will suppress frequency components of the double talk signal.

In each frequency band which is used in receive 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 table 13. The echo attenuation is to be achieved for **each individual frequency band** according to the different categories.

6.3.14.5 Minimum activation level and sensitivity of double talk detection

For further study.

6.3.15 Switching characteristics

6.3.15.1 Note

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

6.3.15.2 Activation in send direction

The activation in Send 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 Send 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).

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.

Measurement method

Test setup is described in clause 6.2.

The test signal is the "short words for activation" sequence described in clause 7.3.4 of Recommendation ITU-T P.501 [20] with increasing level for each single word.

The settings of the test signal are as follows.

Table 14

	Single word/ Pause Duration	Level of the first single word (active Signal Part at the MRP)	Level Difference between two Periods of the Test Signal
Single word to Determine Switching Characteristic in Send Direction	~600 ms/ ~500 ms	-24 dBP _a (see note)	1 dB
NOTE: The signal level is determined for each utterance individually according to Recommendation ITU-T P.56 [14].			

It is assumed that the pause length of about 400 ms is longer than the hang-over time so that the test object is back to idle mode after each single word.

The level of the transmitted signal is measured at the electrical reference point. The test signal is filtered by the transfer function of the test object. The measured signal level is referred to the filtered 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 single word which indicates the first activation of the test object. The time between the beginning of the single word and the complete activation of the test object is measured.

6.3.15.3 Silence suppression and comfort noise generation

For further study.

6.3.16 Background noise performance

6.3.16.1 Performance in send direction in the presence of background noise

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 15.

Table 15: 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.		

Measurement method

Test setup is described in clause 6.2.

The background noise simulation as described in clause 6.2 is used.

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

In a second step a test signal is applied in receive direction consisting of an initial pause of 10 seconds and a periodical repetition of the Composite Source Signal (CSS) in receive direction (duration 10 seconds) 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 Send direction at the POI.

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

6.3.16.2 Speech quality in the presence of background noise

Requirement

Speech Quality for wideband systems can be tested based on ETSI EG 202 396-3 [i.2]. The test method is applicable for narrowband (100 Hz to 4 kHz) and wideband (100 Hz to 8 kHz) transmission systems. LQOn is used for narrowband and LQOw is used for wideband systems. The test method described leads to three MOS-LQO quality numbers:

- N-MOS-LQOn: Transmission quality of the background noise.
- S-MOS-LQOn: Transmission quality of the speech.
- G-MOS-LQOn: Overall transmission quality.

For the background noises defined in clause 6.3, the following requirements apply:

- N-MOS-LQOn $\geq 3,0$.
- S-MOS-LQOn $\geq 3,0$.
- G-MOS-LQOn $\geq 3,0$.

NOTE: It is recommended to test the terminal performance with other types of background noises if the terminal is likely to be exposed to other noises than specified in clause 6.2.

Measurement method

The background noise simulation as described in clause 6.2 is used. The terminal is set-up as described in clause 6.2.

The background noise should be applied for at least 5 seconds in order to adapt noise reduction algorithms in advance the test.

The near end speech signal consists of 8 sentences of speech (2 male and 2 female talkers, 2 sentences each). Appropriate speech samples can be found in Recommendation ITU-T P.501 [20]. The preferred language is English since the objective method was validated with English language in narrowband. The test signal level is +1,3 dBPa at the MRP.

Three signals are required for the tests:

- 1) The clean speech signal is used as the undisturbed reference (see ETSI EG 202 396-3 [i.2]).
- 2) The speech plus undisturbed background noise signal is recorded at the terminal's microphone position using an omni directional measurement microphone with a linear frequency response between 50 Hz and 6 kHz.
- 3) The Send signal is recorded at the electrical reference point.

N-MOS-LQOn, S-MOS LQOn and G-MOS LQOn are calculated as described in ETSI EG 202 396-3 [i.2].

6.3.16.3 Quality of background noise transmission (with far end speech)

Requirements

The test is carried out applying the Composite Source Signal in receive direction. During and after the end of Composite Source Signal bursts (representing the end of far end speech simulation) the signal level in Send 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.

NOTE: The intention of this measurement is to detect impairments (modulations, switching and others) influencing the background noise transmitted from the terminal under test when a signal from the distant end (receiving side of the terminal under test) is present. Under these test conditions no modulation of the transmitted signal should occur. Modulation, switching or other type of impairments might be caused by an improper behaviour of a nonlinear processor working in conjunction with the echo canceller and erroneously switching or modulating the transmitted background noise.

Measurement method

Test setup is described in clause 6.2.

The background noises are generated as described in clause 6.2.

First the measurement is conducted without inserting the signal at the far end. At least 10 seconds 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 shall 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 Recommendation ITU-T P.501 [20] is applied in receive direction with a duration of ≥ 2 CSS periods. The test signal level is -16 dBm0 at the electrical reference point.

The Send 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 Send 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.

6.3.17 Quality of echo cancellation

6.3.17.1 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 echo attenuation.

Measurement method

Test setup is described in clause 6.2.

The test signal consists of periodically repeated Composite Source Signal according to Recommendation ITU-T P.501 [20] 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 seconds 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.

The difference between the maximum attenuation and the minimum attenuation is measured.

NOTE 1: In addition tests with more speech like signals should be made, e.g. Recommendation ITU-T P.50 [13] to see time variant behaviour of EC. However for such tests the simple broadband attenuation based test principle as described above cannot be applied due to the time varying spectral content of the speech like signals.

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 (35 ms) of the level analysis taking into account the exponential character of the integration time in any tolerance scheme.

NOTE 3: Care should be taken not to confuse noise or comfort noise with residual echo. In cases of doubt the measured echo signal should be compared to the residual noise signal measured under the same conditions without inserting the receive signal. If the level vs. time analysis leads to the identical result it can be assumed that no echo but just comfort noise is present.

6.3.17.2 Spectral echo attenuation

Requirement

The echo attenuation vs. frequency shall be below the tolerance mask given in table 16.

Table 16: 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 Send direction in order to mask the echo signal.

Measurement method

Test setup is described in clause 6.2.

Before the actual measurement a training sequence is fed in consisting of the compressed real speech signal described in clause 7.3.3 of Recommendation ITU-T P.501 [20]. The level of the training sequence shall be -16 dBm0.

The test signal is the compressed real speech signal described in clause 7.3.3 of Recommendation ITU-T P.501 [20]. The measurement is carried out under steady-state conditions. The average test signal level shall be -16 dBm0, averaged over the complete test signal. 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.

6.3.17.3 Occurrence of artefacts

For further study.

6.3.17.4 Variable echo path

Requirement

This test is intended to verify that the system will maintain sufficient echo attenuation during single talk with dynamic changing echo paths. The measured echo level over time during single talk should not be more than 10 dB above the minimum noise level during the measurement.

Measurement method

Test setup for desktop hands free terminals: A notebook is positioned at least 20 cm in front of the device (or devices) with the transducers, as shown in figure 18. The notebook lid is moved during the measurement.

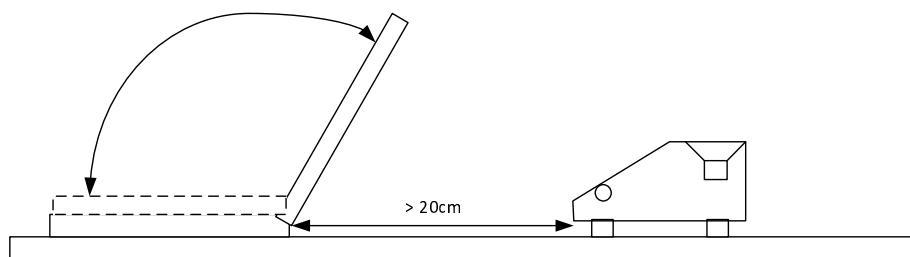


Figure 18: Positioning of DUT

Test setup for softphone: The test setup is described in clause 6.2. The notebook lid is moved during the measurement, as shown in figure 19. This setup is valid for all combinations of notebook with or without external speakers or microphone:

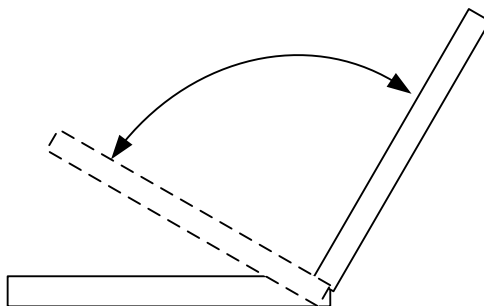


Figure 19: Positioning of DUT

Test setup for other handsfree devices: for further study.

NOTE: Care should be taken to not generate noise during the movement of the notebook lid. Because of this, this measurement is not applicable for a softphone without external microphone.

As test signal the compressed real speech signal described in clause 7.3.3 of Recommendation ITU-T P.501 [20] is used. The signal level shall be -10 dBm0. The terminal volume control is set to nominal RLR. The first 4 sentences of the test signal are used to allow full convergence of the echo canceller. The next 4 sentences (from 10,75 s to 22,5 s) are used for the analysis. The echo signal level is analysed over time. The echo signal level is analysed for 11,75 s, using a time constant of 35 ms.

The measurement result is displayed as echo level versus time.

No level peak should be more than 10 dB above the minimum noise level during the measurement.

6.3.18 Variant impairments; network dependant

6.3.18.1 Clock accuracy send

Requirement

The clock accuracy in send direction between the VoIP-Terminal and the IP reference interface shall be less than 150 ppm under ideal network conditions.

NOTE: The clock accuracy does not cover all possible network configurations. Especially it is not sufficient for data transmission or distributed TDM PBX where synchronization is required.

Measurement method

A sequence of CS signals (active signal length = 250 ms) is repeated for 120 s in order to analyse clock accuracy and any other time-variant delay. The pause length between two CS bursts is 100 ms and 1,2 s after every fourth burst in order to simulate a speech pause, which may lead to buffer adjustments. The test signal level shall be -4,7 dBPa at the MRP.

A cross correlation analysis versus time is carried out over the whole 120 s sequence between the received and the original test signal. The duration of the measurement (120 s) is indicated on the x-axis, the result of the cross correlation analysis (delay) is plotted on the y-axis.

The resulting clock accuracy within an analysis time range of at least 60 s is calculated as follows:

$$clockaccuracy[ppm] = \frac{delaychange[s]}{analysisduration[s]} \cdot 1 \cdot 10^6 \quad (4)$$

6.3.18.2 Clock accuracy receive

Requirement

The clock accuracy in receive direction between the IP reference interface and the VoIP-Terminal shall be less than 150 ppm under ideal network conditions.

Measurement method

A sequence of CS signals (active signal length = 250 ms) is repeated for 120 s in order to analyse clock accuracy and any other time-variant delay. The pause length between two CS bursts is 100 ms and 1,2 s after every fourth burst in order to simulate a speech pause, which may lead to buffer adjustments. The test signal level at the IP reference interface shall be -16 dBm0.

A cross correlation analysis versus time is carried out over the whole 120 s sequence between the received and the original test signal. The duration of the measurement (120 s) is indicated on the x-axis, the result of the cross correlation analysis (delay) is plotted on the y-axis.

The resulting clock accuracy within an analysis time range of at least 60 s is calculated as follows:

$$clockaccuracy[ppm] = \frac{delaychange[s]}{analysisduration[s]} \cdot 1 \cdot 10^6 \quad (5)$$

6.3.18.3 Send delay variation

Requirement

The measured maximum delay variation in send direction of the VoIP-terminal under test should be less than 5 ms.

NOTE: Any delay variation introduced in send direction will lead to potentially increased delay due to increased de-jitter buffer at the far end terminal.

Measurement method

The RTP data stream in send direction should be monitored with a tap or a switch providing a monitoring port, positioned at the location of the network impairment simulator (see clause 6.2). The test arrangement is according to clause 6.2.

The monitoring time should be 60 s. A signal like the British-English single talk sequence described in clause 7.3.2 of Recommendation ITU-T P.501 [20] s played back in send direction using a nominal level of -4.7 dBPa at the MRP. This speech signal is only necessary to make sure, RTP is played out, even in the case VAD is active.

The delay variation for each packet D(i) is evaluated according to IETF RFC 3550 [30]:

$$\begin{aligned} d(i) &= \Delta t_{\text{eff}(i)} - \Delta t_{\text{exp}(i)} \\ D(i) &= (15 * D(i-1) + |d(i)|) / 16 \end{aligned} \quad (6)$$

With:

- $\Delta t_{\text{exp}(i)}$ = the expected time between packet i and packet $i-1$; and
- $\Delta t_{\text{eff}(i)}$ = the effective time between packet i and packet $i-1$.

Maximum delay variation = MAX(D(i))

6.3.19 Send and receive delay - round trip delay

The roundtrip delay of a VoIP-terminal is defined as the sum of send and receive delays. In the following clauses the calculation of the requirements for send and receive delay are explained. For a telecommunication connection, only the roundtrip delay can be experienced. For this reason, also the requirement for VoIP-terminals is given only for the roundtrip delay. As long as the measured roundtrip delay fulfils the requirements, send or receive delays may be above the theoretical requirements.

Requirement

It is recognized that the end to end delay should be as small as possible in order to ensure high quality of the communication.

The roundtrip delay of the VoIP-terminal T_{rtd} (sum of receive and send delay) shall be less than 100 ms. (category B in Recommendation ITU-T P-1010 [29]).

NOTE 1: The limit for the roundtrip delay T_{rtd} of the VoIP-terminal is derived from the sum of the send and receive delay limits.

NOTE 2: This requirement is based on the lowest possible delay values which can be expected under ideal network conditions. Caution should 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

Send direction

The delay in send direction is measured from the MRP to POI. The delay measured in send direction is:

$$T_s + t_{\text{System}} \quad (7)$$

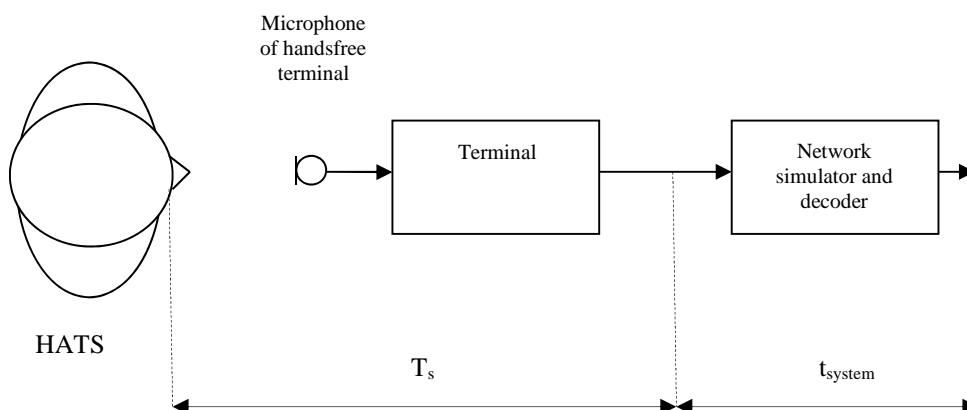


Figure 20: Different blocks contributing to the delay in send direction

The system delay t_{system} is depending on the transmission method used and the network simulator. The delay t_{system} shall be known.

- 1) For the measurements a Composite Source Signal (CSS) according to Recommendation ITU-T P.501 [20] is used. The pseudo random noise (pn)-part of the CSS has to be longer than the maximum expected delay. It is recommended to use a pn sequence of 16 k samples (with 48 kHz sampling rate). The test signal level is -4,7 dBPa at the MRP:
 - The reference signal is the original signal (test signal).
 - The setup of the loudspeaking/handsfree terminal is in correspondence to clause 6.2.
- 2) The delay is determined by cross-correlation analysis between the measured signal at the electrical access point and the original signal. The measurement is corrected by delays which are caused by the test equipment.
- 3) The delay is measured in ms and the maximum of the cross-correlation function is used for the determination.

Receive direction

The delay in receive direction is measured from POI to the Drum Reference Point (DRP). The delay measured in receive direction is:

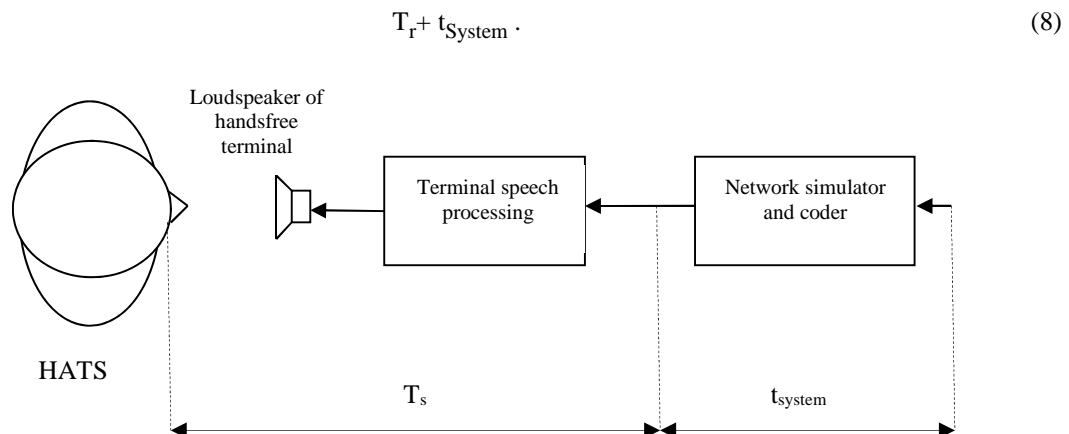


Figure 21: Different blocks contributing to the delay in receive direction

The system delay t_{system} is depending on the transmission system and on the network simulator used. The delay t_{system} shall be known:

- 1) For the measurements a Composite Source Signal (CSS) according to Recommendation ITU-T P.501 [20] is used. The pseudo random noise (pn)-part of the CSS has to be longer than the maximum expected delay. It is recommended to use a pn sequence of 16 k samples (with 48 kHz sampling rate). The test signal level is -16 dBm0 at the electrical interface (POI):
 - The reference signal is the original signal (test signal).
- 2) The test arrangement is according to clause 6.2.
- 3) The delay is determined by cross-correlation analysis between the measured signal at the DRP and the original signal. The measurement is corrected by delays which are caused by the test equipment.
- 4) The delay is measured in ms and the maximum of the cross-correlation function is used for the determination.

6.4 Codec specific requirements

6.4.1 Objective listening speech quality MOS-LQO 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.

Table 17

Speech coder	MOS-LQON (P.863 or TOSQA 2001)	MOS-LQOM (TOSQA 2001)
Recommendation ITU-T G.711 [7]	(ffs)	(ffs)
Recommendation ITU-T G.729 [10]	(ffs)	(ffs)
Recommendation ITU-T G.723.1 [8]	(ffs)	(ffs)
Recommendation ITU-T G.726 @ 32 kbit/s [9]	(ffs)	(ffs)
GSM EFR [2] and AMR @ 12,2 kbit/s [3]	(ffs)	(ffs)
Recommendation ITU-T G.729.1 @ 8 kbit/s [11]	(ffs)	(ffs)

NOTE 1: In narrowband acoustics, Recommendation ITU-T P.863 [27] is recommending using the superwideband mode with a narrowband reference signal, resulting in a prediction on the narrowband scale.

NOTE 2: Not sufficient experience is available so far with Recommendation ITU-T P.863 [27] and TOSQA 2001 (ETSI EG 201 377-1 [i.5]) for measuring handsfree terminals. Therefore the numbers for MOS-LQOS and MOS-LQON are for further study.

NOTE 3: The use of the codecs G.723.1 [8], G.729 [10] and G.729.1 [11] is not recommended due to low quality.

Measurement method

Objective listening speech quality is measured using Recommendation ITU-T P.863 [27] in superwideband mode.

The test signal to be used for the measurements shall be 4 sentence pairs (male/female) fulfilling the requirements of Recommendation ITU-T P.863.1 [26]. The 4 sentence pairs are taken from Recommendation ITU-T P.501 [20]. It shall be stated, which sentence pairs were used. The test signal level is averaged over all sentence pairs (4 sentence pairs). The measurement is done 4 times, every time using another pair of the speech sentences. The result of the measurement is the averaged value of all 4 measurements.

NOTE 4: With Recommendation ITU-T P.863 [27] narrowband VoIP terminals can be measured in narrowband mode as well as in superwideband mode. If backwards comparability of results is needed (e.g. with older subjective test results), narrowband mode should be chosen.

NOTE 5: For the use of P.863 the following applies (see Recommendation ITU-T P.863.1 [26]):

- Superwideband Context (MOS-LQOS):
 - Reference Signal Superwideband flat filtered 50 Hz to 14 kHz;
 - Test Signal Superwideband flat filtered 50 Hz to 14 kHz.

NOTE 6: An alternative test method is TOSQA 2001 (ETSI EG 201 377-1 [i.5]). With TOSQA, terminals used in narrowband mode only should be measured based on MOS-LQON. Terminals used in narrowband and wideband mode should be measured based on MOS-LQOM.

6.4.2 Objective listening quality MOS-LQO 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.

Requirements

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

Table 18

Speech coder	MOS-LQON (P.863 or TOSQA 2001)	MOS-LQOM (TOSQA 2001)
Recommendation ITU-T G.711 [7]	(ffs)	(ffs)
Recommendation ITU-T G.729 [10]	(ffs)	(ffs)
Recommendation ITU-T G.723.1 [8]	(ffs)	(ffs)
Recommendation ITU-T G.726 @ 32 kbit/s [9]	(ffs)	(ffs)
GSM EFR [2] and AMR @ 12,2 kbit/s [3]	(ffs)	(ffs)
Recommendation ITU-T G.729.1 @ 8 kbit/s [11]	(ffs)	(ffs)

NOTE 1: In narrowband acoustics, Recommendation ITU-T P.863 [27] is recommending using the superwideband mode with a narrowband reference signal, resulting in a prediction on the narrowband scale.

NOTE 2: Not sufficient experience is available so far with Recommendation ITU-T P.863 [27] and TOSQA 2001 (ETSI EG 201 377-1 [i.5]) for measuring handsfree terminals. Therefore the numbers for MOS-LQOS and MOS-LQON are for further study.

NOTE 3: The use of the codecs G.723.1 [8], G.729 [10] and G.729.1 [11] is not recommended due to low quality.

Measurement method

Objective listening speech quality is measured using Recommendation ITU-T P.863 [27] in superwideband mode.

The test signal to be used for the measurements shall be 4 sentence pairs (male/female) fulfilling the requirements of Recommendation ITU-T P.863.1 [26]. The 4 sentence pairs are taken from Recommendation ITU-T P.501 [20]. It shall be stated, which sentence pairs were used. The test signal level is averaged over all sentence pairs (4 sentence pairs). The measurement is done 4 times, every time using another pair of the speech sentences. The result of the measurement is the averaged value of all 4 measurements.

NOTE 4: An alternative test method is TOSQA 2001 (ETSI EG 201 377-1 [i.5]). With TOSQA, terminals used in narrowband mode only should be measured based on MOS-LQON. Terminals used in narrowband and wideband mode should be measured based on MOS-LQOM.

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

Table 19: Network conditions for electrical-acoustical measurements (speech samples)

Condition	Packet Loss (Equal)	Delay Variation
0 (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.		
NOTE 4: The settings are derived from the ones used in the ETSI Plugtest VoIP speech quality test events.		

Table 20: Requirements for G.711 speech codecs

Condition	MOS-LQON (P.863 or TOSQA 2001)	MOS-LQOM (TOSQA 2001)	Delay
0	(ffs)	(ffs)	< 31 ms
1	(ffs)	(ffs)	< 31 ms
2	(ffs)	(ffs)	< 51 ms
3	(ffs)	(ffs)	< 31 ms
4	(ffs)	(ffs)	< 51 ms
5	(ffs)	(ffs)	< 31 ms

Table 21: Requirements for G.729 speech codecs

Condition	MOS-LQON (P.863 or TOSQA 2001)	MOS-LQOM (TOSQA 2001)	Delay
1	(ffs)	(ffs)	< 30 ms
2	(ffs)	(ffs)	< 50 ms
3	(ffs)	(ffs)	< 30 ms
4	(ffs)	(ffs)	< 50 ms
5	(ffs)	(ffs)	< 30 ms

Table 22: Requirements for G.723.1 speech codecs

Condition	MOS-LQON (P.863 or TOSQA 2001)	MOS-LQOM (TOSQA 2001)	Delay
1	(ffs)	(ffs)	< 50 ms
2	(ffs)	(ffs)	< 70 ms
3	(ffs)	(ffs)	< 50 ms
4	(ffs)	(ffs)	< 70 ms
5	(ffs)	(ffs)	< 50 ms

NOTE 5: In narrowband acoustics, Recommendation ITU-T P.863 [27] is recommending using the superwideband mode with a narrowband reference signal, resulting in a prediction on the narrowband scale.

NOTE 6: An alternative test method is TOSQA 2001 (ETSI EG 201 377-1 [i.5]). With TOSQA, terminals used in narrowband mode only should be measured based on MOS-LQON. Terminals used in narrowband and wideband mode should be measured based on MOS-LQOM.

6.4.3 Quality of jitter buffer adjustment

Requirements

The speech quality during and after inserted IP delay variation shall be as follows.

Table 23: Requirements for variant network impairments

Codec	MOS-LQON
G.711	> 3,9
G.729	> 3,4
G.723.1	> 3,1

The delay measured 20 seconds after ending of the IP delay variation shall be max. 10 ms higher than the delay measured before the IP delay variation.

Measurement method

The test signal consists of a CSS-signal, followed by 5 times the same speech sentence, fulfilling the requirements of Recommendation ITU-T P.863.1 [26], then again a CSS signal (20 seconds after the IP delay variation stops). This test is done 8 times with 8 single sentences taken from 4 sentence pairs from Recommendation ITU-T P.501 [20]. The speech signal level is averaged over all used (original) sentences (8 sentences).

NOTE 1: The 8 used sentences consist of the 8 single sentences taken from the 4 sentence pairs used in clauses 6.4.1 and 6.4.2.

NOTE 2: For every new measurement a new call has to be setup to start with an initial delay. Depending on the algorithm used in the variable jitter buffer (e.g. jitter buffer starting with a high fill size), it may be necessary to let some time pass under clean conditions until the measurement is started.

The first CSS signal is used to measure the delay prior to the IP impairment (in clean network conditions). The second CSS signal is used to measure the delay 20 seconds after the IP impairment stops. The difference of the two delays is the measurement result for the variation of the jitter buffer per measurement. The overall result is the average of all 8 measurements.

The first sentence (during which IPDV of 50 ms is applied) is used to measure the speech quality during jitter buffer adaption (low to high). MOS-LQON of the first sentence is measured using Recommendation ITU-T P.863 [27] in superwideband mode. The overall result is the average MOS-LQON of the 8 measurements.

The second to the fifth sentence (every 5 seconds a sentence) are used to measure the speech quality during jitter buffer adaption (high to low). MOS-LQON is measured using Recommendation ITU-T P.863 [27] in superwideband mode for each of these four sentences. The minimum MOS-LQON of these four sentences is used for the averaging over all 8 measurements. The overall result for the speech quality during jitter buffer adaption (high to low) is the average of the minimum MOS-LQON-value of the 8 measurements.

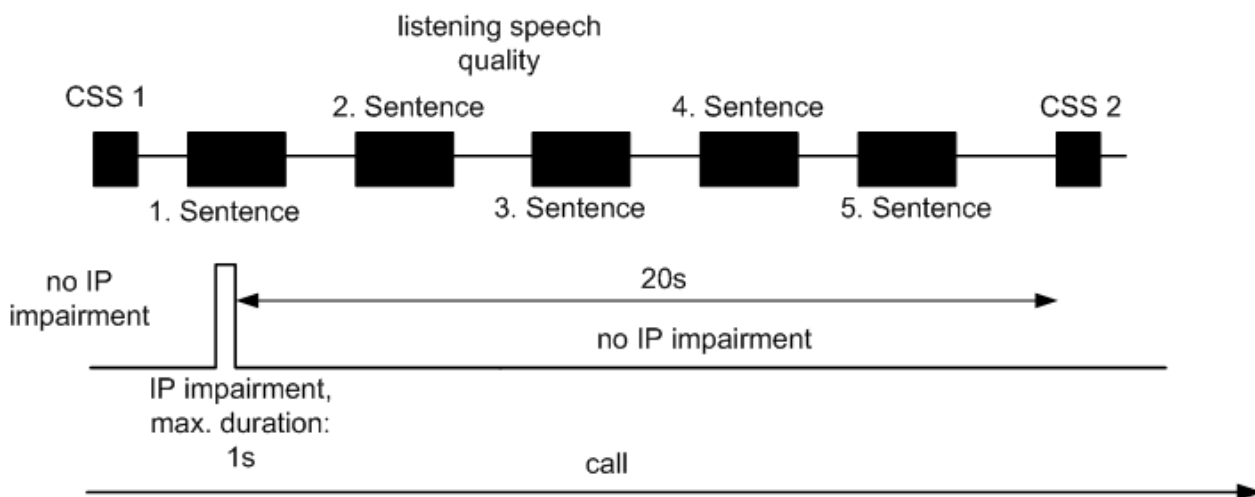


Figure 22: Test Sequence to measure quality of Jitter buffer adjustment (with 1 of 10 sentences)

The IP impairment consists of additional packet delay (IPDV) up to 50 ms, during max. 1 second. The impairment can be in form of jitter, but also with only some single packets delayed. An example for the impairment can be found in annex B of ETSI ES 202 737 [28].

NOTE 3: Care should be given, that no packet reordering occurs (this could happen if e.g. one packet is delayed by 50 ms and the next one is not delayed, they will change order, which will not happen in real networks except in a failover situation or with bad implementations of load balancing).

Annex A (informative): Processing delays in VoIP terminals

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

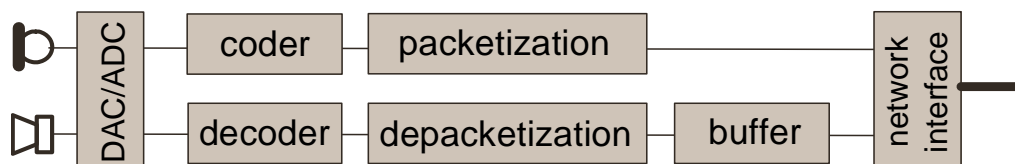


Figure A.1: Synoptic of the different functions implemented in a VoIP terminal

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

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

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

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

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

On the Send 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 (Send and receive 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 Send 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.

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

Digital to analogue conversion consists in transforming speech signal from digital to analogue format. As analogue to digital conversion, the processing delay associated to digital to analogue conversion is considered as not 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 Send and receive terminals. This synchronization is carried out by buffering the audio frames received from the IP stack before Send them to the decoder. The second role of the jitter buffer is to smooth a possible variation of the transmission time. If synchronization of Send and receive 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 receive 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 "perfect" 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 "perfect" 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)	Minimum 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

- Recommendation ITU-T G.131: "Talker echo and its control".
- Recommendation ITU-T G.1020: "Performance parameter definitions for quality of speech and other voiceband applications utilizing IP networks".
- 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		
V1.2.1	October 2007	Publication
V1.3.1	September 2009	Publication
V1.3.2	September 2010	Publication
V1.4.1	March 2015	Publication
V1.6.1	December 2016	Membership Approval Procedure MV 20170131: 2016-12-02 to 2017-01-31
V1.5.1	January 2017	Publication
V1.6.1	February 2017	Publication