



INTERNATIONAL TELECOMMUNICATION UNION

ITU-T

TELECOMMUNICATION
STANDARDIZATION SECTOR
OF ITU

P.800

(08/96)

SERIES P: TELEPHONE TRANSMISSION QUALITY

Methods for objective and subjective assessment of
quality

**Methods for subjective determination of
transmission quality**

ITU-T Recommendation P.800

(Previously CCITT Recommendation)

ITU-T P-SERIES RECOMMENDATIONS
TELEPHONE TRANSMISSION QUALITY

Vocabulary and effects of transmission parameters on customer opinion of transmission quality	Series P.10
Subscribers' lines and sets	Series P.30 P.300
Transmission standards	Series P.40
Objective measuring apparatus	Series P.50 P.500
Objective electro-acoustical measurements	Series P.60
Measurements related to speech loudness	Series P.70
Methods for objective and subjective assessment of quality	Series P.80 P.800
Audiovisual quality in multimedia services	Series P.900

For further details, please refer to ITU-T List of Recommendations.

ITU-T RECOMMENDATION P.800

METHODS FOR SUBJECTIVE DETERMINATION OF TRANSMISSION QUALITY

Summary

This Recommendation describes methods and procedures for conducting subjective evaluations of transmission quality. The main revision encompassed by this version of this Recommendation is the addition of an annex describing the Comparison Category Rating (CCR) procedure. Other modifications have been made to align this Recommendation with recent revision of Recommendation P.830.

Source

ITU-T Recommendation P.800 was revised by ITU-T Study Group 12 (1993-1996) and was approved under the WTSC Resolution No. 1 procedure on the 30th of August 1996.

Keywords

Absolute Category Rating, Comparison Category Rating, conversational test, Degradation Category Rating, listening test, subjective evaluation, Subjective testing

FOREWORD

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

The World Telecommunication Standardization Conference (WTSC), which meets every four years, establishes the topics for study by the ITU-T Study Groups which, in their turn, produce Recommendations on these topics.

The approval of Recommendations by the Members of the ITU-T is covered by the procedure laid down in WTSC Resolution No. 1 (Helsinki, March 1-12, 1993).

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

NOTE

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

© ITU 1996

All rights reserved. No part of this publication may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying and microfilm, without permission in writing from the ITU.

CONTENTS

	Page
1 Scope.....	1
2 References.....	1
3 Definitions	2
4 Abbreviations.....	3
5 Conventions	3
6 Recommended methods.....	3
6.1 Conversation-opinion tests	3
6.2 Listening-opinion tests.....	3
6.3 Interview and survey tests.....	5
6.4 Other tests	5
Annex A – Conversation-opinion tests	5
A.1 Test facilities.....	5
A.1.1 Physical conditions.....	5
A.1.2 Establishing the connection.....	10
A.1.3 Monitoring.....	10
A.2 Experiment design	10
A.3 Conversation task.....	10
A.4 Test procedure.....	11
A.4.1 Eligibility of subjects.....	11
A.4.2 Opinion scale	11
A.4.3 Instructions to subjects	12
A.4.4 Data collection.....	13
A.4.5 Treatment of results.....	13
Annex B – Listening tests – Absolute Category Rating (ACR).....	14
B.1 Source recordings.....	14
B.1.1 Recording environment	14
B.1.2 Sending system.....	14
B.1.3 Recording system.....	14
B.1.4 Speech material.....	14
B.1.5 Recording procedure.....	15
B.1.6 Talkers	16
B.1.7 Speech levels	16
B.1.8 Calibration signal.....	16

	Page
B.2 Selection of circuit conditions	16
B.2.1 Speech input and listening levels.....	16
B.2.2 Talkers	16
B.2.3 Reference conditions	17
B.2.4 Other conditions	17
B.3 Design of experiment.....	17
B.4 Listening test procedure.....	17
B.4.1 Listening environment.....	17
B.4.2 Listening system.....	17
B.4.3 Listening level	18
B.4.4 Listeners.....	18
B.4.5 Opinion scales recommended by the ITU-T.....	18
B.4.6 Instructions to subjects	19
B.4.7 Statistical analysis and reporting of results.....	20
Annex C – Quantal-Response Detectability Tests.....	20
Annex D –Degradation Category Rating (DCR) method	22
D.1 Introduction.....	22
D.2 Degradation Category Rating (DCR) procedure.....	22
D.2.1 Speech samples.....	22
D.2.2 Reference conditions	22
D.2.3 Stimulus presentation	22
D.2.4 Test instructions.....	23
D.3 Statistical analysis.....	23
Annex E – Comparison Category Rating (CCR) method.....	23
E.1 Introduction.....	23
E.2 Quality reference.....	24
E.3 MNRU references	24
E.4 Presentation to listeners	24
E.5 Data analysis	24
Annex F – The threshold method for comparison of transmission systems with a reference system.....	25
F.1 Introduction.....	25
F.2 Testing procedure.....	26

	Page
F.3 Presentation of signals	26
F.4 Speech sources	27
F.5 Listening environment	27
F.6 Listeners.....	27
F.7 Reliability.....	27
Bibliography.....	28

Introduction

Modern telecommunication networks provide a wide array of voice services using many transmission systems. In particular, the rapid deployment of digital technologies has led to an increased need for evaluating the transmission characteristics of new transmission equipment. In many circumstances, it is necessary to determine the subjective effects of some new transmission equipment or modification to the transmission characteristics of a telephone network. This Recommendation describes methods for obtaining subjective evaluations of transmission systems and components. Recommendation G.113 contains useful information on the impairments that can occur. Recommendation P.11 discusses the effects that transmission impairments may have on the users of telecommunication networks and services. The methods described in this Recommendation may be used to estimate the equipment impairment factors (eifs) or quantization distortion units (qdus) that are described in Recommendation G.113.

Recommendation P.800¹

METHODS FOR SUBJECTIVE DETERMINATION OF TRANSMISSION QUALITY

(Amended at Helsinki, 1993; revised in Geneva, 1996)

1 Scope

This Recommendation contains advice to Administrations on conducting subjective tests of transmission quality in their own laboratories. It does not however deal with types of tests described in detail in other ITU–T Recommendations and documentation, namely:

- a) determination of Reference and Relative Equivalents – see *Handbook on Telephony*, Geneva, 1993;
- b) determination of Loudness Ratings – see Recommendation P.78;
- c) determination of Articulation Ratings (A.E.N. values) – see *Handbook on Telephony*, Geneva, 1993.

Neither does it deal with the various kinds of specialized tests used in the course of developing items of telephone equipment, for the purpose of diagnosing faults and shortcomings, such as Diagnostic Rhyme Tests [1] and other tests dedicated to the study of specific aspects of speech output.

This Recommendation gives the approved methods which are considered to be suitable for determining how satisfactorily given telephone connections may be expected to perform.

The methods indicated here are intended to be generally applicable whatever the form of degradation factors present. Examples of degrading factors include: loss (often frequency dependent); circuit noise; transmission errors (random bit errors as well as erased frames that occur in systems such as mobile communications); environmental noise; sidetone; talker echo; non-linear distortion of various kinds including low bit-rate encoding; propagation time; harmful effects of voice-operated devices; distortions of the time scale arising from packet switching; and time-varying degradations of the communication channel, including those arising in loudspeaking sets. Combinations of two or more of such factors also have to be catered for. Further guidance for specific applications is available in Recommendations P.830 (digital speech codecs), P.84 (DCME/PCME), and P.85 (speech output devices).

2 References

The following Recommendations and other references contain provisions that, through reference in this text, constitute provisions of this Recommendation. At the time of publication, the editions indicated are valid. All Recommendations and other references are subject to revision; all users of this Recommendation are therefore encouraged to investigate the possibility of applying the most recent edition of the Recommendations listed below. A list of the currently valid ITU–T Recommendations is regularly published.

- IEC Publication 1260: 1995, *Electroacoustics – Octave-band and fractional – Octave-band filters*.
- IEC Publication 581-5: 1981, *High fidelity audio equipment and systems; Minimum performance requirements – Part 5: Microphones*.

¹ Formerly Recommendation P.80.

- IEC Publication 651: 1979, *Sound level meters. (Amendment 1-1993) (Corrigendum March 1994)*.
- ISO 266: 1975, *Acoustics – Preferred frequencies for measurements*.
- ISO 1996-1: 1982, *Acoustics – Description and measurement of environmental noise – Part 1: Basic quantities and procedures*.
- ISO 1996-2: *Acoustics – Description and measurement of environmental noise – Part 2: Acquisition of data pertinent to land use*.
- ISO 1996-3: 1987, *Acoustics – Description and measurement of environmental noise – Part 3: Application to noise limits*.
- ITU-T Recommendation G.113 (1996), *Transmission impairments*.
- CCITT Recommendation G.722 (1988), *7 kHz audio-coding within 64 kbit/s*.
- CCITT Recommendation G.726 (1990), *40, 32, 24 and 16 kbit/s Adaptive Differential Pulse Code Modulation (ADPCM)*.
- CCITT Recommendation G.728 (1992), *Coding of speech at 16 kbit/s using low-delay code excited linear prediction*.
- ITU-T Recommendation G.729 (1996), *Coding of speech at 8 kbit/s using Conjugate-Structure Algebraic-Code-Excited Linear-Prediction (CS-ACELP)*.
- ITU-T Recommendation P.10 (1993), *Vocabulary of terms on telephone transmission quality and telephone sets*.
- ITU-T Recommendation P.11 (1993), *Effect of transmission impairments*.
- CCITT Recommendation P.48 (1988), *Specification for an intermediate reference system*.
- ITU-T Recommendation P.56 (1993), *Objective measurement of active speech level*.
- ITU-T Recommendation P.78 (1993), *Subjective testing method for determination of loudness ratings in accordance with Recommendation P.76*.
- ITU-T Recommendation P.810 (1996), *Modulated Noise Reference Unit (MNRU)*.
- CCITT Recommendation P.82 (1984), *Method for evaluation of service from the standpoint of speech transmission quality*.
- ITU-T Recommendation P.830 (1996), *Subjective performance assessment of telephone-band and wideband digital codecs*.
- ITU-T Recommendation P.84 (1993), *Subjective listening test method for evaluating digital circuit multiplication and packetized voice systems*.
- ITU-T Recommendation P.85 (1994), *A method for subjective performance assessment of the quality of speech voice output devices*.

3 Definitions

For the purposes of this Recommendation, the following definitions apply:

3.1 dBov: dB relative to the overload of a digital system.

3.2 Q : The ratio, in dB, of speech power to modulated noise power in the Modulated Noise Reference Unit, as described in Recommendation P.810.

4 Abbreviations

For the purposes of this Recommendation, the following abbreviations are used:

ACR	Absolute Category Rating
ADPCM	Adaptive Differential Pulse Code Modulation
BER	Bit Error Rate
CCR	Comparison Category Rating
CMOS	Comparison Mean Opinion Score
DCR	Degradation Category Rating
DMOS	Degradation Mean Opinion Score
FER	Frame Erasure Rate
IRS	Intermediate Reference System (Recommendation P.48)
MNRU	Modulated Noise Reference Unit (Recommendation P.810)
MOS	Mean Opinion Score
PCM	Pulse Code Modulation
SNR	Signal-to-Noise Ratio

5 Conventions

Subjective evaluation of telecommunications equipment and systems may, in principle, be conducted using listening-only or conversational methods of subjective testing. As a practical matter, listening-only tests may be the only feasible method of subjective testing during the development of new transmission equipment or telecommunication services. This Recommendation describes recommended procedures for conversational and listening-only methods of subjective evaluation.

6 Recommended methods

6.1 Conversation-opinion tests

Laboratory conversation tests are intended – as far as possible – to reproduce, in the laboratory situation, the actual service conditions experienced by telephone customers. To this end it is necessary to choose the circuit conditions and subjects suitably, and to administer the tests in an appropriate manner.

It is important that the conditions simulated in the test are correctly specified and set up, and measured accurately before and after each experiment; that auxiliary facilities such as dialling and ringing are provided; and that faithful records of the output of each test are kept. Detailed description of the method, considerations and precautions are found in Annex A.

6.2 Listening-opinion tests

Listening-opinion tests are not expected to reach the same standard of realism as conversation tests, and the restrictions are therefore less severe in some respects; but the artificiality that has to be accepted brings with it a necessity for strict control of many things which in conversation tests are allowed to find their own equilibrium.

The recommended test method for listening-only tests is the "Absolute Category Rating" (ACR) method described in Annex B, which is in conformance with the Category Judgement method recommended for conversation tests (see Annex A), and adopted partly for the same reasons.

Category ratings are applied to short groups of unrelated sentences, each of which has been passed through a number of standard processes as well as the processes under test. This method is well-established, and has been applied to analogue and digital telephone connections and to telecommunications devices, such as digital codecs. In the work leading to Recommendations G.726 32 kbit/s ADPCM, G.728, G.729, and G.722, for example, laboratories in different countries performed subjective tests by the same method on the same physical conditions and on identical transmission systems, and the results showed a high degree of consistency.

Other methods commonly used are the Quantal-Response Detectability Method, Degradation Category Rating (DCR), Comparison Category Rating (CCR) and the Threshold Method.

Annex C describes Quantal-Response Detectability Tests, which are suitable for evaluating threshold values of certain quantities and their associated probabilities. For example, the level above which single-frequency interference has a given probability of being objectionable or detectable, or the probability that crosstalk in a given range of levels is intelligible, can best be determined by this method.

An alternative to the Absolute Category Rating method is the Degradation Category Rating (DCR) method which is described in detail in Annex D. The DCR method compares the system under test with a high quality fixed reference and the degradation (from "Inaudible" to "Very annoying") is rated on a five-point scale. This method is suitable when the impairment (especially digital impairments) is small. It may therefore be particularly useful for evaluating similar digital speech processing algorithms. Thus, the DCR method may serve as a means for system optimization once it has been shown by the methods of Annexes A and B that the worst-case connection incorporating the degradation in question is within acceptable limits.

Annex E describes a variation of the DCR procedure called the Comparison Category Rating (CCR) method. As in the DCR, the CCR method compares the system under test with a high quality fixed reference (in the CCR case on a scale from "Much Better" to "Much Worse"). This procedure may be particularly suitable for systems that improve the quality of the input speech (e.g. noise cancellation systems).

The Threshold method, also suitable for system optimization, is described in Annex F. By direct comparison of the system under test with a reference system, such as the Modulated Noise Reference Unit (MNRU, as described in Recommendation P.810), it is possible to equate the value of the reference condition (Q for digital processes) which equals the performance of the system under test.

Information on other types of subjective test methods, which include scaling methods, can be found in 2.6 of the *Handbook on Telephony*.

Listening tests have direct applications in the assessment of physical transmission systems which are essentially unidirectional. Examples include broadcast circuits, public address systems and recorded announcement systems in which listening degradations such as loss, noise and distortion may be present.

Results of listening-only tests can be applied, but only with certain reservations, to the prediction of the assessment for conversation conducted over a two-way system, such as a connection in a public switched telephone network. The provisos are that the effects of the following additional factors are duly taken into account:

- talking degradations (e.g. sidetone and echo);
- conversation degradations (e.g. propagation time and mutilation of speech by the action of voice-operated devices).

The annexes to this Recommendation provide information on preparation of speech material, processing of speech material, experiment philosophy (including choice of circuit conditions), listening test procedure and treatment of results.

6.3 Interview and survey tests

If the rather large amount of effort needed is available and the importance of the study warrants it, transmission quality can be determined by "service observations". Recommended ways of performing these, including the questions to be asked when interviewing customers, are given in Recommendation P.82. To maintain a high degree of precision a total of at least 100 interviews per condition is required.

A disadvantage of the service-observation method for many purposes is that little control is possible over the detailed characteristics of the telephone connections being tested. However, this method does afford a global appreciation of how the "equipment" performs in the real environment.

Further information can be found in 2.5.8.3 of the *Handbook on Telephony*.

6.4 Other tests

Reference [2] gives information of a method that largely overcomes the disadvantages of the interview technique of 6.3, yet retains many of the advantages. This method, termed SIBYL, allows a small proportion of a user's ordinary calls to be passed through special arrangements which modify the normal quality of transmission according to a test programme. If a particular call has been so treated, the volunteer is asked to vote by dialling one of a set of digits to indicate his opinion. In this way, all results are recorded by the controlling computer and complete privacy is maintained.

Annex A

Conversation-opinion tests

A.1 Test facilities

A.1.1 Physical conditions

A.1.1.1 Test cabinets

The two subjects are seated in separate sound-proof cabinets near the point from which the experiment is controlled. The volume of the room is not less than 20 m³, with a reverberation time less than 500 ms (normally in the range of 200-300 ms), for handheld systems such as telephone handsets, or for headset systems; and not less than 30 m³ for handsfree systems (extra care is exercised if reverberation time is an experimental variable).

The internal dimensions of the cabinet are such that standing-wave pattern effects are kept to a minimum. A typical ratio is 5:4:3.

The physical construction of the rooms should be such that sufficient sound attenuation of the outside noise environment is achieved so that the requirements of A.1.1.2.1 are met.

The cabinets are favourably decorated to recreate a natural environment.

A.1.1.2 Noise

A.1.1.2.1 Noise floor

The ambient noise level (when no environmental noise is deliberately introduced) is kept as low as possible. For practical reasons, such as regular changes of fresh air in the cabinet, the target is an upper limit of NC25 [3] or NR25 (see ISO 1996). These values approximate the noise level in homes (sleeping areas), hospitals and libraries.

A.1.1.2.2 Environmental noise

Environmental noise is fed in with the required spectrum (e.g. Hoth spectrum to represent typical room noise – see A.1.1.2.2.1) at the required level (e.g. 50 dBA) measured with a Precision Sound Level Meter conforming to IEC Publication 651, used with the "A weighting" and the "fast" meter characteristic. If different conversations in the same experiment require different room noise levels, then care is taken to prevent the transitions from being too obvious to the subjects. Ideally, room noise should be changed only when subjects are out of the sound-proof rooms. If this is not possible, then changes of level are carried out gradually (at a rate not exceeding 4 dB per second), at a time when no experimental conversation is in progress and when the subjects' attention is otherwise occupied – by communicating with the operator, for example.

Spectra with appropriate long-term characteristics are given in A.1.1.2.2.1 and A.1.1.2.2.2.

For some purposes it is necessary to use noise that fluctuates in level or spectrum, such as tape recordings of actual office noise or traffic noise. In such cases it should be ensured that the statistical characteristics are stable when averaged over a reasonably short period of time such as one minute.

It is recommended that the noise level and spectrum are measured at least twice; at the beginning and end of the experiment. Any significant variation in the two measurements, when compared with each other, must be assessed by the experimenter as it may cast doubt on the validity of the experiment.

It is essential to ensure that the loudspeakers and amplifiers are capable of faithfully reproducing the required noise.

A.1.1.2.2.1 Room noise

The room noise shall have a power density spectrum corresponding to that published by Hoth [4]. Table A.1 gives the spectrum density adjusted in level to produce a reading of 50 dBA on a sound level meter conforming to IEC Publication 651. This is produced in Figure A.1. This spectrum is independent of level, i.e. for 40 dBA the level in each band shall be 10 dB less than that shown in Table A.1. Additional information on the power in each one-third octave band is also given in Table A.1.

TABLE A.1/P.800

Room noise spectrum

Frequency (Hz)	Spectrum density (dB SPL/Hz)	Bandwidth $10 \log_{10} \Delta f$ (dB)	Total power in each 1/3rd octave band (dB SPL)	Tolerance (dB)
100	32.4	13.5	45.9	±3
125	30.9	14.7	45.4	
160	29.1	15.7	44.9	
200	27.6	16.5	44.1	
250	26.0	17.6	43.6	
315	24.4	18.7	43.1	
400	22.7	19.7	42.3	
500	21.1	20.6	41.7	
630	19.5	21.7	41.2	
800	17.8	22.7	40.4	
1000	16.2	23.5	39.7	
1250	14.6	24.7	39.3	
1600	12.9	25.7	38.7	
2000	11.3	26.5	37.8	
2500	9.6	27.6	37.2	
3150	7.8	28.7	36.5	
4000	5.4	29.7	34.8	
5000	2.6	30.6	33.2	
6300	-1.3	31.7	30.4	
8000	-6.6	32.7	26.0	

NOTES

- 1 The electrical input signal, e.g. white noise, shall be band-limited to the 1/3rd octave bands centred on the ISO preferred frequencies (ISO 266) between 100 Hz and 8000 Hz with the band edges conforming to the filters described in IEC 1260.
- 2 The acoustical room noise is difficult to control at low frequencies, especially in the unspecified region below 100 Hz because of the dimensions of typical test cabinets, poor attenuation of such cabinets and the influence of extraneous noises, e.g. air-conditioning plant. It is therefore desirable to select a test cabinet that keeps these unwanted low frequency sound pressure levels to a minimum.

A.1.1.2.2 Internal vehicle noise

Two spectra representing internal vehicle noise [5], [6] are recommended. They are adequately represented by simplified curves [7]: one spectrum for moving vehicles and the other for stationary vehicles. Table A.2 gives the spectrum densities together with additional information on the power in each one-third octave band. The spectrum density for moving vehicles is shown in Figure A.2 a) and for stationary vehicles in Figure A.2 b). These spectra are independent of level.

NOTE – The noise spectra in Table A.2 should be considered provisional. More detailed specifications are under study.

Table A.3 gives the computed values of the unweighted sound pressure levels for various speeds calculated over the ISO one-third octave frequency bands centred on 63 Hz to 8000 Hz.

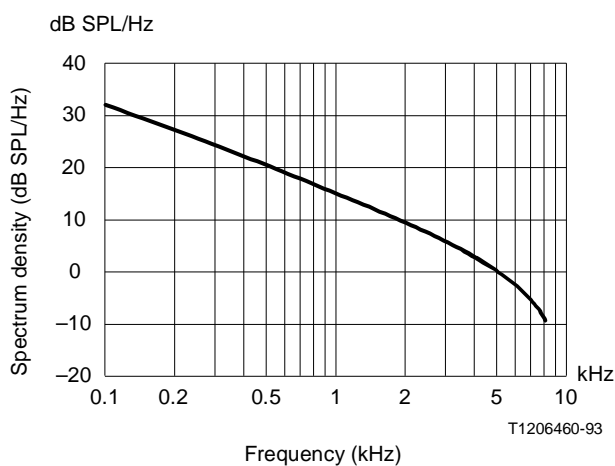


FIGURE A.1/P.800
Room noise spectral density

TABLE A.2/P.800
Internal vehicle noise spectra

Frequency (Hz)	Spectrum density (dB SPL/Hz)		Bandwidth $10 \log_{10} \Delta f$ (dB)	Total power in each 1/3rd octave band (dB SPL)		Tolerance (dB)
	Moving	Stationary		Moving	Stationary	
63	72.3	58.3	11.7	84.0	70.0	±3
80	69.3	55.0	12.7	82.0	66.7	
100	66.5	49.8	13.5	80.0	63.3	
125	63.3	45.1	14.7	78.0	60.0	
160	60.3	42.0	15.7	76.0	56.7	
200	57.5	36.8	16.5	74.0	53.3	
250	54.4	34.7	17.6	72.0	52.3	
315	51.3	32.6	18.7	70.0	51.3	
400	48.3	30.6	19.7	68.0	50.3	
500	45.4	28.7	20.6	66.0	49.3	
630	42.3	26.6	21.7	64.0	48.3	
800	39.3	24.6	22.7	62.0	47.3	
1000	36.5	22.8	23.5	60.0	46.3	
1250	33.3	20.6	24.7	58.0	45.3	
1600	30.3	18.6	25.7	56.0	44.3	
2000	27.5	16.8	26.5	54.0	43.3	
2500	24.4	14.7	27.6	52.0	42.3	
3150	21.3	12.6	28.7	50.0	41.3	
4000	18.3	10.6	29.7	48.0	40.3	
5000	15.4	8.7	30.6	46.0	39.3	
6300	12.3	6.6	31.7	44.0	38.3	
8000	9.3	4.6	32.7	42.0	37.3	

TABLE A.3/P.800

Computed sound pressure levels of spectra

Spectra		Sound pressure level, unweighted (dB SPL)
Moving	30 km/h	80
	80 km/h	85
	110 km/h	90
Stationary		75

NOTES to Tables A.2 and A.3:

- 1 These values apply for typical vehicles. Discretion may be used to adjust the levels downwards for luxury vehicles and upwards for noisier vehicles.
- 2 Because of the practical difficulty of generating such high sound pressure levels at low frequencies, and because normal speech contains no apparent energy below about 63 Hz in which range of frequencies the ear is also comparatively insensitive it is probably advisable to restrict the recommended noise spectrum to frequencies above 63 Hz. However, it should be borne in mind that low and medium frequency vibrations have important physiological and psychological effects which should be studied in their own right.
- 3 The electrical input signal, e.g. white noise, shall be band-limited to the 1/3rd octave bands centred on the ISO preferred frequencies (ISO 266) between 63 Hz and 8000 Hz with the band edges conforming to the filters described in IEC 1260.
- 4 The acoustical room noise is difficult to control at low frequencies especially in the unspecified region below 63 Hz because of the dimensions of typical test cabinets, poor attenuation of such cabinets and the influence of extraneous noises, e.g. air-conditioning plant. It is therefore desirable to select a test cabinet that keeps these unwanted low frequency sound pressure levels to a minimum.

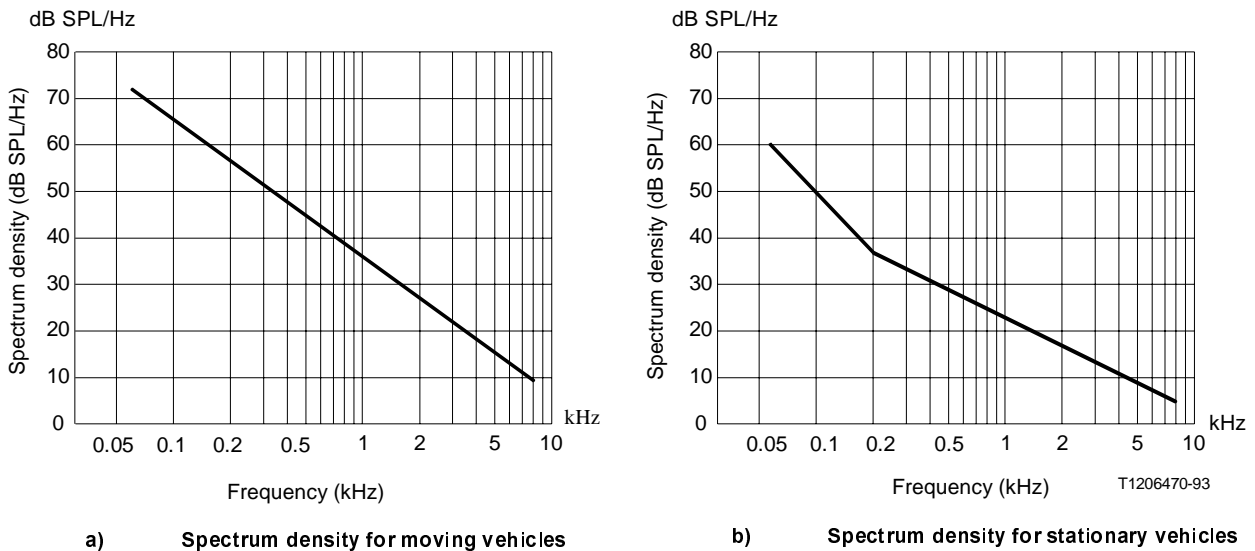


FIGURE A.2/P.800

Vehicle noise spectral density

A.1.1.3 Noise measurement position

It is recommended that the measurement of Sound Pressure Level (SPL) in the test cabinets (see A.1.1.1) shall be made as follows:

- furniture in position;
- no subject or test personnel present;
- the SPL shall be measured at a vertical distance of 740 mm above the centre of the seat of the subject's chair with a meter conforming to Recommendation P.54 using "A" weighting;
- the spectrum of the environmental noise shall be measured in one-third octaves, centred at the preferred frequencies as defined in ISO 266, and must stay within the specified tolerances, e.g. ± 3 dB for Hoth noise (see A.1.1.2.2.1);
- in rooms where more than one subject is to be tested the difference in dBA, for all subject positions, shall not vary by more than ± 2 dB.

NOTE – It is suggested that the minimum distance between each loudspeaker and the measurement position should be 1.5 m.

A.1.2 Establishing the connection

When establishing the laboratory connection, the following must be taken into consideration:

- telephone sets;
- initial call set-up;
- laboratory representation of telephone connections.

It is recommended that the sensitivity/frequency characteristic of the connection is measured at least twice; at the beginning and end of the experiment. Any significant variation in the two measurements, when compared with each other, must be assessed by the experimenter as it may cast doubt on the validity of the experiment.

Detailed information on this aspect can be found in 2.5.8.2 of the *Handbook on Telephonometry*.

A.1.3 Monitoring

Monitoring can take many forms but the three most commonly used are:

- *Intercommunication system* – Essential to allow the subject and experimenter to communicate with each other.
- *Visual monitoring* – This has two purposes; the first being safety and the second to observe the peculiarities of the subject, e.g. how they hold the telephone handset.
- *Tape recordings and recording system* – This facilitates other useful information to be gleaned, e.g. duration of the call, speech voltage, speech activity.

Detailed information on this aspect can be found in 2.5.8.2 of the *Handbook on Telephonometry*.

A.2 Experiment design

There are a number of methods suitable for use in designing experiments, e.g. Latin Squares, Youden Squares, Balanced Incomplete Blocks and Randomization with Replication to name but a few. The experimenter must decide for himself the method to be used taking into account the number of test conditions, accuracy of results and the ability to make sound judgements from the findings.

Suitable designs include those of the $n \times n$ graeco-latin square type. A detailed description of these designs can be found in 2.5.8.2 of the *Handbook on Telephonometry*.

A.3 Conversation task

Every effort is to be made to ensure that conversations are purposeful, and that subjects have full opportunity to exploit the transmission capabilities of the test circuit.

The general rule is that every conversation should have a natural beginning and a natural ending. Unless it is absolutely necessary, the conversation must never be terminated in the middle of the task (with the exception of Simplified Conversation Tests which are described in 2.5.8.2 of the *Handbook on Telephonometry*).

An example of a conversational task can be found in 2.5.8.2d of the *Handbook on Telephonometry*.

A.4 Test procedure

A.4.1 Eligibility of subjects

Subjects taking part in the conversation tests are chosen at random from the normal telephone using population, with the provisos that:

- a) they have not been directly involved in work connected with assessment of the performance of telephone circuits, or related work such as speech coding; and
- b) they have not participated in any subjective test whatever for at least the previous six months, and not in a conversation test for at least one year.

If the available population is unduly restricted, then allowance must be made for this fact in drawing conclusions from the results.

No steps are taken to balance the numbers of male and female subjects unless the design of the experiment requires it. Subjects are arbitrarily paired in the experimental design prior to the test and remain thus paired for its duration.

A.4.2 Opinion scale

The following opinion scales are those recommended by the ITU-T.

A.4.2.1 Conversation opinion scale

Various five-point category-judgement scales may be used for different purposes. The layout and wording of opinion scales, as seen by subjects in experiments, is very important, and should follow the standard arrived at through years of experience. The following opinion scale is the most frequently used for ITU-T applications and equivalent wording should be used depending on language which might result in small variations to the original English text.

This is a category rating obtained from each subject at the end of each conversation.

Opinion of the connection you have just been using

Excellent

Good

Fair

Poor

Bad

The experimenter allocates the following values to the scores:

Excellent = 5; Good = 4; Fair = 3; Poor = 2; Bad = 1

and all further statistical processing is performed in terms of these numbers. The arithmetic mean of any collection of these opinion scores is called the mean conversation-opinion score, and is represented by the symbol MOS_C (or, where suffix notation is not available, the symbol MOS_c).

NOTE – In the past, the equivalences Excellent = 4, Good = 3, Fair = 2, Poor = 1, Bad = 0 have often been used. Anyone using results from earlier experiments must be aware that the mean scores must all be increased by one to be comparable with those now obtained; otherwise there is no difference in the numerical processing that may be applied.

A.4.2.2 Difficulty scale

This is a binary response obtained from each subject at the end of each conversation.

Did you or your partner have any difficulty in talking or hearing over the connection?

Yes

No

The experimenter allocates the following values to the scores:

Yes = 1 No = 0

The quantity evaluated (percentage of "yes" responses) is called percentage Difficulty or per cent "Difficult", and is denoted by the symbol %D. The corresponding simple proportion is denoted by the symbol d; in other words, $%D = 100d$.

NOTE – It is often the case that the nature of the difficulty is required and then it is usual for the experimenter to ask the subject to describe in his/her own words their perception of the difficulty.

The layout and wording of the opinion scale, as seen by subjects in experiments, are very important, and should follow the standard arrived at through years of experience: see A.4.3.

A.4.2.3 Other opinion scales

Other opinion scales that may be suitable are variants of the methods of "magnitude estimation" and "cross-modality matching" [8]. The responses on these scales may be one of the following:

- a) one of a numerical series of categories labelled 1, 2, 3, 4, 5 (and denoted as such to the subject), but with descriptions attached only to the first and the last, to identify the subjective dimension;
- b) a numerical mark on a scale from one to a number much greater than five – say 10 or 100; or
- c) a length proportional to some property (e.g. quality), marked manually along a given straight line.

A survey of experimental methods can be found in 2.6.2 of the *Handbook on Telephonometry*.

A.4.3 Instructions to subjects

Instructions are given to subjects on arrival for their first visit. It is normal for the subject to receive a letter prior to arrival which contains non-technical information on the experiment and what will be expected of them. An example of such a letter can be found in 2.5.8.2, Table 3/2.5 of the *Handbook on Telephonometry*.

They are asked whether they have read and understood the letter. Any obscurities are clarified, and opportunity is given for asking questions. The sound-proof rooms and their facilities are demonstrated. Subjects are informed how many calls will be comprised in this visit. On subsequent visits the subjects are merely informed that the procedure will be the same as before, with possibly a different number of calls. An example of some operational details required by an experiment can be found in 2.5.8.2, Table 4/2.5 of the *Handbook on Telephonometry*.

A.4.4 Data collection

Speech levels, and related data such as durations and activity factors, may be derived from the tape recordings, but are now normally measured on-line, by computer-controlled meters, and stored directly into computer files for subsequent analysis.

Two subjective responses are collected per conversation per subject by the experimenter. The essential data consists of the conversation-opinion score and the Difficulty decision. These responses may be collected using any suitable means, including pencil and paper, electronic buttons, keyboards, keypads, or computer touch-screen terminals. A sample response form can be found in 2.5.8.2, Table 5/2.5 of the *Handbook on Telephony*.

A.4.5 Treatment of results

This is a very extensive subject, and only a brief outline can be given here.

Each conversation gives rise to two conversation opinions on the scale: Excellent – Good – Fair – Poor – Bad (scored respectively 5, 4, 3, 2, 1), two votes on the Difficulty scale (scored 1 = Yes, 0 = No), two measured active speech levels and one value of duration. In particular cases information may be collected about other variables; e.g. video recordings may be made in order to observe how subjects hold their handsets, or other data may be derived from the opinion forms or from the audio recordings.

The average of the opinion scores should be calculated for each test condition. Confidence limits should be evaluated and significance tests performed by conventional analysis-of-variance techniques.

The usual assumptions underlying the analysis of variance are sufficiently nearly satisfied in the case of opinion score, active speech level and most other variables of interest; but they are not satisfied – particularly the assumption of constant residual variance – in the case of a binary variant like Difficulty score. In spite of this, experience confirms the observation made in other fields [9] that the analysis of variance technique is robust enough to give reasonable results even with such extreme departures from the statistically ideal conditions. The results from the first stage of the analysis of variance of the Difficulty scores should be regarded with some reserve; but once it is established that there are no unexplained abnormalities in the results and no unexplained conflicts with the outcome of the corresponding analysis of the MOS_C results, then the second stage (detailed analysis of the averages from the End-Condition combinations) can be confidently undertaken with the aid of a mathematical transformation.

Detailed descriptions of the analysis can be found in 2.5.9 of the *Handbook on Telephony*.

As a further aid to the review of the data, graphs if appropriate should be plotted showing the mean opinion score as a function of the parameter under test, e.g. MOS_C versus circuit attenuation. On the graph the vertical axis should always be MOS_C .

Annex B

Listening tests – Absolute Category Rating (ACR)

B.1 Source recordings

In order to eliminate unwanted variability in the speech source, samples of speech having the desired standardized properties should first be prepared in recorded or stored form, as follows.

B.1.1 Recording environment

The talker should be seated in a quiet room with volume between 30 and 120 m³ and a reverberation time less than 500 ms (preferably in the range 200-300 ms). The room noise level must be below 30 dBA with no dominant peaks in the spectrum.

The room noise characteristic should be reported in as complete a form as possible, e.g. dBA, long-term spectrum, and amplitude-time distribution. It is desirable to record a 30-second sample of the room noise for detailed investigation if this proves to be necessary.

B.1.2 Sending system

Whatever sending system is chosen, e.g. local telephone or Intermediate Reference System (IRS) as specified in Recommendation P.48, the system should be calibrated according to the relevant Recommendation (e.g. Recommendation P.64), and the sending sensitivity-frequency characteristic should be reported in full. Annex D/P.830 describes the "modified IRS" that has been deemed appropriate for evaluation of all-digital connections using speech codecs.

It is recommended that the sending sensitivity characteristic of the connection is measured at least twice; at the beginning and end of the experiment. Any significant variation in the two measurements, when compared with each other, must be assessed by the experimenter as it may cast doubt on the validity of the experiment.

B.1.3 Recording system

The recording system must be of high (studio) quality and can be any of the following:

- a) A conventional two-track tape recorder. The type of equalization must be stated, but IEC is recommended. High grade tape (low print-through, low noise) should be used at all times.
- b) A two-channel digital audio processor with a high quality video cassette recorder or Digital Audio Tape (DAT) machine.
- c) A computer-controlled digital storage system.

The third system is the best and most versatile, but practical reasons often dictate the choice of one of the other systems. In these, one of the two tracks should be used for recording the speech and the other for inserting control signals at a level and a frequency chosen to avoid crosstalk problems.

B.1.4 Speech material

The speech material should consist of simple, meaningful, short sentences, chosen at random as being easy to understand (from current non-technical literature or newspapers, for example). These sentences should be made up into lists in random order in such a way that there is no obvious connection of meaning between one sentence and the next. Very short and very long sentences should be avoided, the aim being that each sentence when spoken should fit into a time-slot of 2–3 seconds. Examples of sentences are shown in Table B.1.

The experimenter must decide how many sentences are required in each group to constitute a speech sample. A minimum of two and a maximum of five are recommended. The time interval between sentences, during which circuit noise may be heard and adaptive processes settle into new states, is

also important. It is advisable to record the longest groups that may be needed, as it is always possible to obtain shorter groups by copying or replaying parts of longer ones.

Groups are combined into lists consisting of five or ten groups each, so that a complete list can be used as a series of samples subjected to the same treatment but with listening level or some other parameter varied when the list is reproduced.

TABLE B.1/P.800

Examples of speech material

You will have to be very quiet.
There was nothing to be seen.
They worshipped wooden idols.
I want a minute with the inspector.
Did he need any money?

B.1.5 Recording procedure

The following recording scheme has been extensively used and is recommended.

Speech is recorded from a linear microphone and low-noise amplifier with a flat frequency response as specified in IEC Publication 581-5. The microphone is positioned between 140 mm and 200 mm from the talker's lips. In some applications, it may be necessary to use a windscreen, it is used if breath puffs from the talker are noticed.

The same speech may be recorded simultaneously from the sending output of an Intermediate Reference System (IRS, see Recommendation P.48), with the handset held in the normal manner. If the investigation in view specifically requires it, another telephone instrument may be used in place of the IRS.

Two separate recording systems are used simultaneously: one for recording the wideband speech in one channel, and the other for recording the telephone speech in the corresponding channel. The other channel of each recording system is used for recording control signals as explained in 2.5 of the *Handbook on Telephonometry*.

This dual recording system ensures that the same speech is recorded in two forms (telephone speech and wideband speech). Normally only one of these is required in any one experiment, but there are occasions when it is necessary to use both, and it is an advantage in any case to be able to make comparative measurements on the two versions.

The active speech level, as defined in Recommendation P.56, is observed during recording. Care is taken during the recording process that the active speech level in both recording systems is between 20 and 30 dB below the overload point of the recording system for each sentence measured separately. Any group of sentences for which this does not hold is re-recorded.

It is recommended that the ratio of the active speech level to psophometrically weighted noise level (for definition see 8.2.3/P.830), SNR(p), on the recording media should be > 40 dB with an objective of 50 dB.

All speech samples used in one experiment may be different: this is essential for Listening-Effort tests and desirable for other types.

B.1.6 Talkers

There must be as many talkers as required in the design of the experiment (see B.3).

Talkers should pronounce the sentences fluently but not dramatically, and have no speech deficiencies such as stutter; they should adopt a speaking level that is comfortable to them as individuals and which they can maintain fairly constantly.

B.1.7 Speech levels

The recordings when completed are played back, and the active speech level of each sentence is measured with a meter conforming to Recommendation P.56. The lists (announcements, sentences and control tones) are then re-recorded on to a second system with the necessary gain adjustments, so as to bring each group of sentences to the standardized active speech level specified below, and still preserve the proper time relationships between the sentences and the tone signals in the other channel.

For the narrow-band speech, the standardized level is derived by measuring and adjusting the narrow-band recorded signal directly; the recommended target is -26 dB ($+0.5$ dB) relative to the peak overload level of the recording system. The calibration tone has its r.m.s. level equal to the mean active level of the re-recorded speech.

For the wideband speech, account must be taken of the intended use of the recordings. It is sometimes appropriate to adopt the same levels as for telephone speech, but if the recording is intended for playback through a loudspeaker or artificial mouth, then the individual target speech levels should be such that equality is maintained at the output of the whole electric-acoustic replay chain.

B.1.8 Calibration signal

At the beginning of each resulting recording, 20 seconds or more of tone are inserted at the re-recording stage (for calibration purposes) at a level that is in a known relationship to the mean active speech level (most conveniently, equal to it). This calibration tone is normally at 1000 Hz, but may be at some other frequency if the recordings are intended for playing through systems (such as certain types of sub-band coder) that respond to 1000 Hz in a special manner.

This tone can then be used later to adjust the mean input speech levels (see B.4.3).

B.2 Selection of circuit conditions

B.2.1 Speech input and listening levels

In the selection of circuit conditions particular attention should be applied to:

- range of input levels;
- range of listening levels:
 - there is no universal optimum listening level;
 - a variety of listening levels will occur in practice;
 - comparability considerations;
 - interactions may occur.

A detailed explanation of these aspects can be found in 2.5.8.1 of the *Handbook on Telephonometry*.

B.2.2 Talkers

Since sophisticated processes often affect male and female voices differently, the experimental design should provide for two types of voice as a balanced factor; scores for male and female speech

should be evaluated separately, only to be averaged if they yield main effects and interactions that are not statistically different.

Moreover, to reduce the danger that the results may depend heavily on peculiarities of the voices chosen, it is essential for more than one male and more than one female voice to be used in a balanced design.

B.2.3 Reference conditions

Every experiment should include reference conditions so that experiments made in different laboratories or at a different time in the same laboratory can be sensibly compared. Such reference conditions will depend on what is being assessed. For a digital system the reference conditions may include the Modulated Noise Reference Unit (MNRU) conforming to Recommendation P.810; other controlled degradations are appropriate in other cases (e.g. signal-to-noise ratio, see 8.2.3/P.830).

B.2.4 Other conditions

Besides the requirements of B.2.1 to B.2.3 inclusive, other conditions will be included depending on the purpose of the test. For instance, room noise might be a variable as well as bit-error rate for a digital system or Rayleigh fading for a radio system.

B.3 Design of experiment

The design of the experiment uses the same principles as given in A.2.

In addition, the experiment design must cater for the following:

- a) requirements of B.2;
- b) order-of-presentation effect.

For a given sample of subjects the test is limited in size by the maximum length of session possible without fatigue. If the experiment is too large to be catered for in one session then it is prudent to sub-divide into two or more sessions. Ideally no session should last for more than 20 minutes and in no case should a session exceed 45 minutes.

B.4 Listening test procedure

B.4.1 Listening environment

The listening room should meet the same conditions as the recording room (see B.1.1) with the exception that the environmental noise (see A.1.1.2.2) should be set to the appropriate level. See A.1.1.2.2.1 and A.1.1.2.2.2 for examples of noise spectra.

It is recommended that the noise level and spectrum are measured at least twice; at the beginning and end of the experiment. Any significant variation in the two measurements, when compared with each other, must be assessed by the experimenter as it may cast doubt on the validity of the experiment.

B.4.2 Listening system

Whatever listening system is chosen (e.g. local telephone system, Intermediate Reference System as specified in Recommendation P.48, a loudspeaker system), the system should be calibrated according to the relevant Recommendation (e.g. Recommendation P.64) and the receiving sensitivity/frequency characteristic should be reported in full. Annex D/P.830 describes the "modified IRS" that has been deemed appropriate for evaluation of all-digital connections using speech codecs.

It is recommended that the receiving sensitivity/frequency characteristic of the connection is measured at least twice; at the beginning and end of the experiment. Any significant variation in the

two measurements, when compared with each other, must be assessed by the experimenter as it may cast doubt on the validity of the experiment.

B.4.3 Listening level

The gain of the system should be set in such a way that the calibration tone (see B.1.8) played back from the processed tapes produces the required listening level.

Variations in listening level, as required by the experiment design, can either be accommodated by use of:

- a) attenuators/amplifiers in the listening system; or
- b) included in the processing or re-recording stage.

The second method is not recommended because it is difficult to maintain a high enough signal-to-noise ratio at low levels, and because flexibility and variety in the randomization are greatly reduced.

The listening level should always be recorded. Information on this subject can be found in 2.5 of the *Handbook on Telephony*.

B.4.4 Listeners

Subjects taking part in listening tests are chosen at random from the normal telephone using population, with the provisos that:

- a) they have not been directly involved in work connected with assessment of the performance of telephone circuits, or related work such as speech coding;
- b) they have not participated in any subjective test whatever for at least the previous six months, and not in any listening-opinion test for at least one year; and
- c) they have never heard the same sentence lists before.

If the available population is unduly restricted, then allowance must be made for this fact in drawing conclusions from the results.

In some cases screening of subjects may be necessary and a method based on Annex B/P.78 may be applicable.

B.4.5 Opinion scales recommended by the ITU-T

Various five-point category-judgement scales may be used for different purposes. The layout and wording of opinion scales, as seen by subjects in experiments, is very important, and should follow the standard arrived at through years of experience. The following opinion scales are those most frequently used for ITU-T applications and equivalent wording should be used depending on language which might result in small variations to the original English text:

a) Listening-quality scale

<i>Quality of the speech</i>	<i>Score</i>
Excellent	5
Good	4
Fair	3
Poor	2
Bad	1

The quantity evaluated from the scores (mean listening-quality opinion score, or simply mean opinion score) is represented by the symbol MOS.

b) **Listening-effort scale**

The heading of the listening-effort opinion scale is particularly important. Without it, the other descriptions are liable to be seriously misunderstood.

<i>Effort required to understand the meanings of sentences</i>	<i>Score</i>
Complete relaxation possible; no effort required	5
Attention necessary; no appreciable effort required	4
Moderate effort required	3
Considerable effort required	2
No meaning understood with any feasible effort	1

The quantity evaluated from the scores (mean listening-effort opinion score) is represented by the symbol MOS_{LE} but where suffix notation is not available, the symbol MOS_{le} is used.

c) **Loudness-preference scale**

<i>Loudness preference</i>	<i>Score</i>
Much louder than preferred	5
Louder than preferred	4
Preferred	3
Quieter than preferred	2
Much quieter than preferred	1

The quantity evaluated from the scores (mean loudness-preference opinion score) is represented by the symbol MOS_{LP} but where suffix notation is not available, the symbol MOS_{lp} is used.

NOTE – Examples of alternative subjective scales, which should only be used if the above three opinion scales do not meet the needs of the experimenter, are given in 2.6 of the *Handbook on Telephonometry* and CCIR Report 751, Volume VIII.3, 1986.

B.4.6 Instructions to subjects

An example of typical instructions is given in Table B.2. The instructions must be given (verbally as well, if necessary) prior to commencement of the experiment. When the subject has understood the instructions, he/she should listen to the preliminary list and give his opinions. No suggestion should be made to the subjects that the preliminary samples include the best or worst in the range to be covered, or exhaust the range of conditions they can expect to hear. After the preliminary list, there should be sufficient time allowed for answering possible questions by the subjects. Questions about procedure or about the meaning of the instructions should be answered, but any technical questions must be met with the response, "We cannot tell you anything about that until the experiment is finished".

TABLE B.2/P.800

Example of instructions to subjects

LISTENING EXPERIMENT No. ...

In this experiment you will be listening to short groups of sentences via the telephone handset, and giving your opinion of the speech you hear.

On the table in front of you is a box with five illuminated press buttons. When all the lamps go on, you will hear ... sentences. Listen to these, and when the lamps go out, press the appropriate button to indicate your opinion on the following scale.

EFFORT REQUIRED TO UNDERSTAND THE MEANINGS OF SENTENCES

- | | |
|---|--|
| 5 | Complete relaxation possible; no effort required. |
| 4 | Attention necessary; no appreciable effort required. |
| 3 | Moderate effort required. |
| 2 | Considerable effort required. |
| 1 | No meaning understood with any feasible effort. |

The button you have pressed will light up for a short time. Then the lamp will go out, and there will be a brief pause before all the lamps go on again for the next group of ... sentences.

There will be a longer pause after every ... groups (each calling for an opinion). There will be a total of ... groups in this visit, and a similar number in your subsequent visit(s).

Thank you for your help in this experiment.

B.4.7 Statistical analysis and reporting of results

The numerical mean (over subjects) should be calculated for each condition at each listening level, and these means listed for initial inspection (so that effects such as those due to male and female speech can be seen).

Calculation of separate standard deviations for each condition is not recommended. Confidence limits should be evaluated and significance tests performed by conventional analysis-of-variance techniques.

NOTE – In the past, the equivalences, for example, Excellent = 4, Good = 3, Fair = 2, Poor = 1, Bad = 0 have often been used. Anyone using results from earlier experiments must be aware that the mean scores must all be increased by one to be comparable with those now obtained; otherwise there is no difference in the numerical processing that may be applied.

The method of analysis of the opinion scales of B.4.5 follows the principles stated in A.4.5.

As a further aid to the review of the data, graphs if appropriate should be plotted showing the mean opinion score as a function of the parameter under test, e.g. MOS versus circuit attenuation. On the graph the vertical axis should always be MOS.

Averaging the scores of the male and female talkers should be made with care and does not imply that this step would be warranted for a detailed study and interpretation of results unless the significance tests justify it.

Annex C

Quantal-Response Detectability Tests

The best method for obtaining information on the detectability or some analogous property of a sound (such as echo) as a function of some objective quantity (such as listening level), is a quantal-response method similar in principle to that described in 2.2 of the *Handbook on Telephonometry*.

The main difference is that the subject's response is not a decision in the form "Reference" or "Test" (the designation of the louder of the two circuits), but a vote on a scale such as [10]:

Detectability opinion scale

- A Objectionable
- B Detectable
- C Not detectable

where B is understood to mean "Detectable but not objectionable".

Scales of this sort, usually with three points, may be used in a variety of quantal-response tests; for example the scale as shown above may be used where the stimulus is echo, reverberation, sidetone, voice-switching mutilation, or interfering tones, while crosstalk and perhaps echo in some circumstances may be judged on the scale Intelligible – Detectable – Not detectable.

It is sometimes permissible to regard these votes as opinion scores, with values 2, 1, 0 respectively, and treat them in the same sort of way as listening or conversation opinion scores. But this is often unsatisfactory because the decisions on such a scale as detectability (see above) are not really equivalents of responses on a continuous scale – as votes on such scales as "Loudness preference" (see B.4.5), may be legitimately taken to be – but effectively embody two distinct dichotomies (for example detectable/not detectable and objectionable/not objectionable), which though not independent may nevertheless call different psychological processes into action: in other words, Objectionability or Intelligibility differs in kind, not merely in degree, from Detectability. For this reason a more profitable method of analysis is to express the probability of response according to each dichotomy separately, as a function of some objective variable, by fitting probit or logit equations, and then using the quantiles or other parameters as a basis of comparison between circuit conditions, in a manner analogous to that used in applying articulation scores.

The actual conduct of experiments of this type resembles that of listening-effort tests (see Annex B), but there are some differences. In particular, it is advisable that the first presentation of the signal in each run should be at a high listening level, so that the listener is left in no doubt what kind of signal is a candidate for his decisions. Where sidetone or echo is involved, the subject will be required to talk as well as listen.

Simple audiometric measurements, as described in Recommendation P.78, are usually performed on subjects who participate in these experiments, so that results can be expressed relative to their threshold of hearing.

For examples of the application of these techniques, see [11].

Noise, fading and other disturbances are sometimes investigated by means of responses on a scale with many more points; for example [12]:

- A *Inaudible* – Noise completely undetectable.
- B *Just audible* – Noise can just be detected by listening carefully.
- C *Slight* – Noise detectable, but not disturbing.
- D *Moderate* – Noise slightly disturbing.
- E *Rather loud* – Noise causes appreciable disturbance.
- F *Loud* – Noise very disturbing, but call would be continued.
- G *Intolerable* – Noise so loud that the call would be abandoned, or operator asked to change the line.

These scales are more nearly of the quantized-continuum type, like the Loudness preference scale, and can be treated similarly.

Annex D

Degradation Category Rating (DCR) method

D.1 Introduction

The Absolute Category Rating (ACR) method described in Annex B tends to lead to low sensitivity in distinguishing among good quality circuits. A modified version of the ACR procedure, called the Degradation Category Rating (DCR) [13] procedure, affords higher sensitivity. This procedure is adapted from the CCIR Recommendation [14] for evaluation of good quality circuits. The DCR procedure, which in particular uses an annoyance scale and a quality reference before each configuration to be evaluated, seems to be suitable for evaluating good quality speech.

D.2 Degradation Category Rating (DCR) procedure

D.2.1 Speech samples

Each configuration is evaluated by means of judgements on speech samples from at least four talkers. Each sample should be composed of two sentences separated by approximately 0.5 s of silence. These two samples (S1, S2), hence four different sentences, should be selected from a wider corpus composed of phonetically balanced sentences so that the mean score obtained in evaluating reference (e.g. MNRU for digital processes) circuits for these sentences is about the same as that obtained for the wider corpus. Therefore the corpus consists of eight samples defined as follows:

- talker T1 reading samples S1, S2;
- talker T2 reading samples S1, S2;
- talker T3 reading samples S1, S2;
- talker T4 reading samples S1, S2;
- etc.

This results in a repetition of the two samples during the test. It is felt that this is not a critical factor for the procedure where a degradation is evaluated with regard to the reference. This is especially true for good telephone quality, where the intelligibility of speech is nearly perfect. The use of different samples for each configuration, as is often done in ACR experiments (where the speaker and the sentence effects are confounded), could be one of the reasons for lack of sensitivity in the ACR method.

Some variations of this basic scheme are allowed: increase the number of talkers, mix sentence and talker effects. However, it is important that all configurations are evaluated on the same corpus.

D.2.2 Reference conditions

Reference conditions shall be included, e.g. for digital processes multiplicative noise with Q values within the range 10 to 30 dB with a minimum of four steps is desirable.

A quality reference should be chosen to be inserted before each judgement. Usually source conditions are used, i.e. samples with no more degradation than those introduced by sending systems and limitations of frequency bandwidth. Thus, the choice of the quality reference depends on the application, i.e. for standard telephony, the source signal is 3.4 kHz bandwidth limited, for wideband telephony it is 7 kHz band limited and for high quality sound, the signal is 15 or 20 kHz band limited.

D.2.3 Stimulus presentation

The stimuli are presented to listeners by pairs (A-B) or repeated pairs (A-B-A-B) where A is the quality reference sample and B the same sample processed by the system under evaluation. The

purpose of the reference sample is to anchor each judgement of the listeners. Some "null pairs" (A-A), at least one for each talker, are included to check the quality of anchoring. Using a reference and subjective judgements with respect to that reference is quite a common procedure in psychoacoustics. It tends to result in a good sensitivity for the overall evaluation by listeners. Samples A and B should be separated by 0.5–1 s. In a repeated pair procedure (A-B-A-B), the separation between the two pairs should be 1–1.5 s.

The order effect observed in a one-sample listening tests (e.g. ACR) is not observed with the DCR procedure. Thus, only one random order of presentation can be used. Therefore the basic test and reference conditions will be eight times (four talkers × two samples) the number of nominal conditions.

D.2.4 Test instructions

The subjects should be instructed to rate the conditions according to the five point degradation category scale as follows:

- 5 Degradation is inaudible.
- 4 Degradation is audible but not annoying.
- 3 Degradation is slightly annoying.
- 2 Degradation is annoying.
- 1 Degradation is very annoying.

The quantity evaluated from the scores (degradation mean opinion score) is represented by the symbol DMOS.

D.3 Statistical analysis

Sensitivities can be quantified by means of a statistical multiple comparison test. When an *a posteriori* comparison of circuits is needed a Tukey [15] Honestly Significant Difference (HSD) test can be applied effectively. The HSD test is designed to make all pair-wise comparisons among the means and to determine the significance of the differences in the mean values.

Annex E

Comparison Category Rating (CCR) method

E.1 Introduction

The Comparison Category Rating (CCR) method is similar to the Degradation Category Rating (DCR) method described in Annex D. Listeners are presented with a pair of speech samples on each trial. In the DCR procedure, a reference (unprocessed) sample is presented first, followed by the same speech sample, which has been processed by some technique. In the DCR method, listeners always rate the amount by which the processed (second) sample is *degraded* relative to the unprocessed (first) sample. In the CCR procedure, the order of the processed and unprocessed samples is chosen at random for each trial. On half of the trials, the unprocessed sample is followed by the processed sample. On the remaining trials, the order is reversed. Listeners use the following scale to judge the quality of the second sample relative to that of the first:

The Quality of the Second Compared to the Quality of the First is:

- 3 Much Better
- 2 Better
- 1 Slightly Better
- 0 About the Same
- 1 Slightly Worse
- 2 Worse
- 3 Much Worse

In effect, listeners provide two judgements with one response: "Which sample has better quality?" and "By how much?" The DCR and the CCR methods are particularly useful for assessing the performance of telecommunications systems when the input has been corrupted by background noise. However, an advantage of the CCR method over the DCR procedure is the possibility to assess speech processing that either degrades or improves the quality of the speech.

The quantity evaluated from the scores (comparison mean opinion score) is represented by the symbol CMOS.

NOTE – Caution should be exercised when using the CCR method. Some laboratories have found the method to be useful in evaluating noise reduction systems. However, when this method was used in the recent subjective evaluations of the G.729 (8 kbit/s) codec, the method was found to be too sensitive when evaluating the performance of the codec for speech embedded in background noise.

E.2 Quality reference

The reference (unprocessed) sample (Quality reference or Direct connection) is presented either before or after the processed or degraded signal. The reference sample is generated using the same talker and speech material as used for the processed sample. This reference sample will be corrupted by the same noise (if any) and processed through the same preliminary processes, such as transmitter characteristic, logarithmic companding, etc. Thus, there will be a different quality reference for each of the test conditions.

E.3 MNRU references

MNRU reference conditions should be included to calibrate the judgement scale. These multiplicative noise references are used without being further mixed with environmental noises.

E.4 Presentation to listeners

Each of the speech samples is presented to the listener through the quality reference condition and through a test codec or reference condition (e.g. Recommendation G.726, MNRU). In addition, a "Null pair" should be included for each of the quality references. On these trials, the quality reference is presented twice.

Listeners should judge the quality of the second sample relative to the quality of the first sample. This judgement is made on the 7-point scale shown in E.1. Sample instructions for the listeners are shown in Table E.1.

E.5 Data analysis

Some care must be exercised when analysing the data from a CCR experiment. As half of the trials for any test condition are presented in the order (unprocessed, processed), and the other half are

presented in the opposite order, simple averaging of the numerical scores should yield a CMOS of approximately 0 for all conditions. It is necessary to recode the raw data. If the order of presentation is (processed, unprocessed), then the sign of the numerical score must be reversed (i.e. $-1 \rightarrow 1$, $-2 \rightarrow 2$, ..., $2 \rightarrow -2$, $1 \rightarrow -1$). The recoded scores may be used to compute CMOS, standard deviations, etc. Thus, results are presented in terms of the (unprocessed, processed) order. Appropriate Analysis of Variance, or other statistical tests, may also be performed on the recoded scores. However, comparison opinion scores may not be presumed to represent a linear interval scale. Therefore, statistics for ordinal scales may need to be applied instead.

TABLE E.1/P.800

Example of instructions to subjects

INSTRUCTIONS TO LISTENERS

Comparison category rating test

**"Evaluation of the influence of various environmental noises
on the quality of different telephone systems"**

In this experiment you will hear pairs of speech samples that have been recorded through various experimental telephone equipment. You will listen to these samples through the telephone handset in front of you.

What you will hear is one pair of sentences, a short period of silence, and another pair of sentences. You will evaluate the quality of the second pair of sentences compared to the quality of the first pair of sentences.

You should listen carefully to each pair of samples. Then, when the green light is on, please record your opinion about the quality of the second sample relative to the quality of the first sample using the following scale:

The Quality of the Second Compared to the Quality of the First is:

3:	Much Better
2:	Better
1:	Slightly Better
0:	About the Same
-1:	Slightly Worse
-2:	Worse
-3:	Much Worse

You will have five seconds to record your answer by pushing the button corresponding to your choice. There will be a short pause before the presentation of next pair of sentences.

We will begin with a short practice session to familiarize you with the test procedure. The actual tests will take place during sessions of 10 to 15 minutes.

Annex F

**The threshold method for comparison of transmission
systems with a reference system**

F.1 Introduction

By direct comparison of a transmission system with a reference system, it is possible to assess the performance of the system under test in terms of a degradation characteristic of the reference system which can be varied and set to defined values. An example of such a characteristic is signal-to-noise ratio (for definition see 8.2.3/P.830), SNR(p). The method described here leads to a threshold of equality defined as 50% preference level between the MNRU and the digital system.

F.2 Testing procedure

A listening-only test procedure is used. A signal pair consisting of a reference signal and a test signal is presented to listeners, who are then asked to indicate which of the signals in the pair they judge to have the highest quality (preference rating). Subjective equivalence is defined as the reference value corresponding to the intersection point of the regression curve of the preference scores at the 50% preference level. An example of equivalent SNR obtained with hypothetical preference scores is shown in Figure F.1.

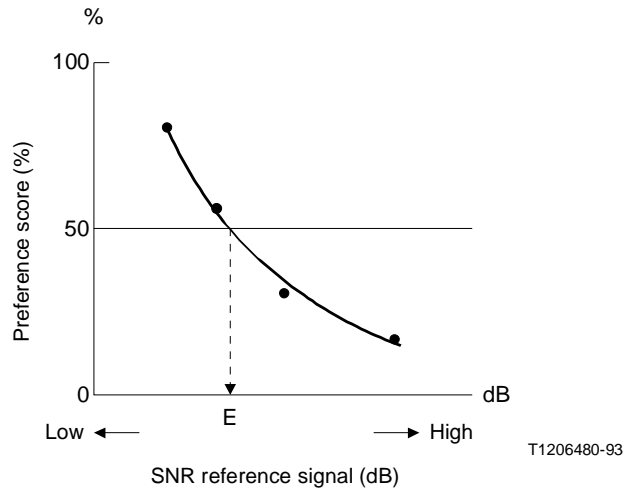
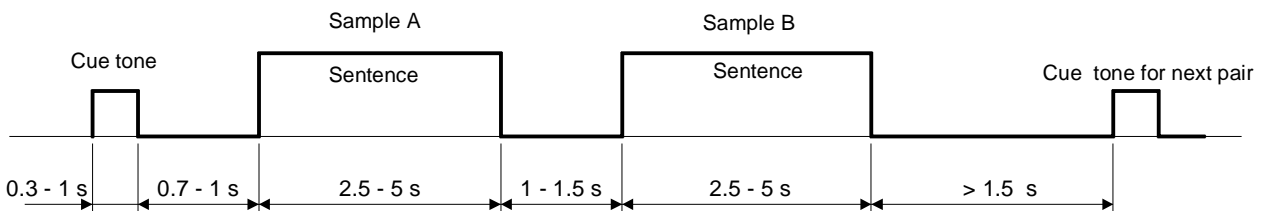


FIGURE F.1/P.800

Example of an equivalence threshold E with hypothetical preference scores

F.3 Presentation of signals

Reference signal A and test signal B are arranged in an equal number of A-B pairs and B-A pairs, and presented in random order. Several degradation levels spaced for example, at 2 dB intervals, are introduced in the reference path so that the range of preference scores extends from 20% to 80%, where the 50% preference lies in the middle of the degradation range. A timing diagram of the presentation is shown in Figure F.2.



T1206490-93

FIGURE F.2/P.800

Timing diagram of the presentation

The subject is required to make a judgement and respond by saying "A is better" or "B is better" (forced choice). The response "A equals B", or "No difference" is forbidden. The duration of the presentation should be limited to about six minutes in order not to tire the listeners. More listening samples may be presented after a suitable rest period. At least two, preferably four or five replications (repetitions of identical presentations) are recommended.

NOTE – If the reference system is available in hardware and the degradation characteristic can easily be changed between presentations, a simplified procedure can be used. In this case the balancing to equally perceived quality is done by the subject. The adjustment is made during the pause between the pairs. The reference signal is always presented first. Presentation continues until the subject reports that the equality threshold has been reached.

F.4 Speech sources

It is necessary to use short sentences spoken by at least two males and two females, preferably four or six of each; different sentences are required for each speaker. The duration should be 2.5-5 seconds for speech and less than 10-15 seconds for music signals. Clicks at the beginning and end of the samples must be avoided. A linear microphone of sufficient bandwidth should be used to record the source signals in a sound-absorbent room having an ambient noise of less than 20 dBA and a reverberation time of less than 0.3 seconds in the band 125 - 8000 Hz. If digital recording equipment is used, the quantizing noise level should be less than the noise level in 14-bit linear PCM.

F.5 Listening environment

A high-fidelity sound reproduction system should be used for the listening test. When listening is carried out with loudspeakers, the reproduction equipment should be studio-quality and the listening room should conform to CCIR Report 797 or IEC 268-13. If headphones are used, diotic (binaural) listening is preferable. The bandwidth shall be at least as wide as that of the system under test.

F.6 Listeners

Although it is preferred that listeners should be selected according to the description in the ACR method (see Annex B), this is not a strict condition in the pair comparison test. If the purpose of the listening test is to obtain the opinions of untrained listeners, untrained subjects are necessary. However, if this is not the purpose of the test, then trained listeners can be used and the reliability of the listening test can be extended by increasing the number of replications for each listener. The minimum number of listeners is six, but should preferably be twelve or more. Several subjects may listen simultaneously but it must be ensured that their responses are obtained independently.

F.7 Reliability

Since variations in preference score in subjective tests are assumed to conform to a t-distribution, the score variation width r which yields 95% reliability at score u ($0 \leq u \leq 1$) over the number (n) of trials (i.e. the number of repetitions for each presentation pair multiplied by the number of subjects and number of source signals) is presented in equation (E-1).

$$r = \pm t(n - 1, 0.05) \cdot \sqrt{u(1 - u) / (n - 1)} \quad (\text{E-1})$$

NOTE – The threshold method is expected to give stable and reliable results even for high quality systems with little degradation.

Degradation can be introduced in the reference system, e.g. by addition of white noise. For digital systems, multiplicative noise as defined in Recommendation P.810 (MNRU) is recommended. For

wideband digital speech coders, the use of a wideband MNRU, as described in Recommendation P.810, is recommended. For some purposes shaped noise instead of white may be appropriate.

Bibliography

- [1] VOIERS (W.D.): Evaluating processed speech using the Diagnostic Rhyme Test, *Speech Technology*, Volume 1, No. 4, pp. 30-39, January-February 1983.
- [2] CCITT Supplement No. 5 to Recommendation P.74, *The SIBYL method of subjective testing*, *Red Book*, Volume V.
- [3] BERANEK (L.L.): Noise and Vibration Control, *McGraw-Hill*, pp. 564-566, 1971.
- [4] HOTH (D.F.): Room noise spectra at subscribers' telephone locations, *J.A.S.A.*, Volume 12, pp. 99-504, April 1941.
- [5] CCITT Question 24/XII, Contribution COM XII-120, *Noise inside light motor vehicles*, study period 1981-1984.
- [6] CCITT Question 24/XII, Contribution COM XII-134, *Internal vehicle noise spectra*, study period 1981-1984.
- [7] CCITT Contribution COM XII-208, *Comparison of the results of vehicle noise submitted by France and BT*, study period 1981-1984.
- [8] STEVENS (S.S.): Psychophysics – Introduction to its perceptual, neural and social prospects, *John Wiley and Sons*, 1975.
- [9] CLARINGBOLD (P.J.): The within-animal bioassay with quantal responses, *Journal of the Royal Statistical Society*, Series B, Volume 18, No. 1, pp. 133-137, 1956.
- [10] RICHARDS (D.L.): Telecommunication by speech, subclause 3.5.2, *Butterworths*, London, 1973.
- [11] *Ibid*, subclauses 3.5.3 and 4.5.1.
- [12] *Ibid*, subclause 4.2.1.6.
- [13] COMBESCURE (P.) *et al*: Quality evaluation of speech coded at 32 kbit/s by means of degradation category ratings, *Proc. ICASSP 82 (International Conference on Acoustics, Speech and Signal Processing)*, Vol. 2, Paris, May 1982.
- [14] CCIR Document 11/17, *Subjective assessment of the quality of television pictures (EBU)*, study period 1978-1982.
- [15] TUKEY (J.W.): The problem of multiple comparisons, *Ditton*, Princeton University, Ed. 1953.
- [16] GABRIELSSON (A.): Statistical treatment of data from listening tests on sound-reproducing systems, Report TA No. 92, *KTH Karolinska Institutet*, Department of Technical Audiology, S-10044 Stockholm, Sweden, November 1979.
- [17] IEC Publication 268-13, Annex 3, subclause 3.3 (a condensed version of [16]).

ITU-T RECOMMENDATIONS SERIES

- Series A Organization of the work of the ITU-T
- Series B Means of expression
- Series C General telecommunication statistics
- Series D General tariff principles
- Series E Telephone network and ISDN
- Series F Non-telephone telecommunication services
- Series G Transmission systems and media
- Series H Transmission of non-telephone signals
- Series I Integrated services digital network
- Series J Transmission of sound-programme and television signals
- Series K Protection against interference
- Series L Construction, installation and protection of cables and other elements of outside plant
- Series M Maintenance: international transmission systems, telephone circuits, telegraphy, facsimile and leased circuits
- Series N Maintenance: international sound-programme and television transmission circuits
- Series O Specifications of measuring equipment
- Series P Telephone transmission quality**
- Series Q Switching and signalling
- Series R Telegraph transmission
- Series S Telegraph services terminal equipment
- Series T Terminal equipment and protocols for telematic services
- Series U Telegraph switching
- Series V Data communication over the telephone network
- Series X Data networks and open system communication
- Series Z Programming languages