

*NRSC  
REPORT*

# NATIONAL RADIO SYSTEMS COMMITTEE

**NRSC-R55  
EIA/NRSC DAR Systems  
Subjective Tests of Audio Quality  
and Transmission Impairments  
Final Report  
July 21, 1995**



NAB: 1771 N Street, N.W.  
Washington, DC 20036  
Tel: (202) 429-5356 Fax: (202) 775-4981



CEA: 1919 South Eads Street  
Arlington, VA 22202  
Tel: (703) 907-7660 Fax: (703) 907-8113

Co-sponsored by the Consumer Electronics Association and the National Association of Broadcasters  
<http://www.nrscstandards.org>

## NRSC-R55

### NOTICE

NRSC Standards, Guidelines, Reports and other technical publications are designed to serve the public interest through eliminating misunderstandings between manufacturers and purchasers, facilitating interchangeability and improvement of products, and assisting the purchaser in selecting and obtaining with minimum delay the proper product for his particular need. Existence of such Standards, Guidelines, Reports and other technical publications shall not in any respect preclude any member or nonmember of the Consumer Electronics Association (CEA) or the National Association of Broadcasters (NAB) from manufacturing or selling products not conforming to such Standards, Guidelines, Reports and other technical publications, nor shall the existence of such Standards, Guidelines, Reports and other technical publications preclude their voluntary use by those other than CEA or NAB members, whether to be used either domestically or internationally.

Standards, Guidelines, Reports and other technical publications are adopted by the NRSC in accordance with the NRSC patent policy. By such action, CEA and NAB do not assume any liability to any patent owner, nor do they assume any obligation whatever to parties adopting the Standard, Guideline, Report or other technical publication.

This Guideline does not purport to address all safety problems associated with its use or all applicable regulatory requirements. It is the responsibility of the user of this Guideline to establish appropriate safety and health practices and to determine the applicability of regulatory limitations before its use.

Published by  
CONSUMER ELECTRONICS ASSOCIATION  
Technology & Standards Department  
1919 S. Eads St.  
Arlington, VA 22202

NATIONAL ASSOCIATION OF BROADCASTERS  
Science and Technology Department  
1771 N Street, NW  
Washington, DC 20036

©2009 CEA & NAB. All rights reserved.

*This document is available free of charge via the NRSC website at [www.nrscstandards.org](http://www.nrscstandards.org). Republication or further distribution of this document, in whole or in part, requires prior permission of CEA or NAB.*

## FOREWORD

NRSC-R55, EIA/NRSC DAR Systems – Subjective Tests of Audio Quality and Transmission Impairments – Final Report, documents the results of subjective tests conducted at the Communications Research Center (CRC) from June 1994 to March 1995. These tests were performed to assess the audio quality of Digital Audio Radio (DAR) systems submitted to the DAR Subcommittee of the Electronics Industries Association (precursor to CEA) and the DAB Subcommittee of the National Radio Systems Committee.

An eight-page summary of this work that was included with Comments submitted to the FCC by the Consumer Electronics Manufacturers Association (CEMA) on July 13, 1999 (as part of MM Docket No. 99-25, In the Matter of Creation of a Low Power Radio Service) is also provided. In this summary the systems that were tested are identified in Table 3 which is excerpted here:

<b>Designation</b>	<b>System</b>	<b>Audio coding</b>	<b>Bit rate (kbps)</b>
<i>a</i>	Eureka-147	Musicam	224
<i>b</i>	Eureka-147	Musicam	192
<i>c</i>	AT&T/Lucent	PAC	160
<i>d</i>	AT&T/Amati, DSB	PAC	160
<i>e</i>	AT&T/Amati, LSB	PAC	160
<i>f</i>	VOA/LPL	PAC	160
<i>g</i>	USADR FM-2	Musicam	256
<i>h</i>	USADR FM-1	Musicam	256
<i>i</i>	USADR AM	Musicam	96 kbps (32 kHz ref.)
<i>j</i>	USADR AM	Musicam	96 kbps (48 kHz ref.)

The NRSC is jointly sponsored by the Consumer Electronics Association and the National Association of Broadcasters. It serves as an industry-wide standards-setting body for technical aspects of terrestrial over-the-air radio broadcasting systems in the United States.

## **Summary of CRC Subjective Test Program**

Submitted by Consumer Electronics Manufacturers Association (CEMA) to the FCC on July 13, 1999 (as part of MM Docket No. 99-25 (In the Matter of Creation of a Low Power Radio Service))

## Appendix 2

### Subjective Assessments of Audio Quality of DAR Systems

#### I. Introduction

This document describes the procedures and results of subjective tests conducted at the Communications Research Centre (CRC), Ottawa, Ontario, Canada, performed to assess the audio quality of digital audio radio (DAR) systems submitted to the Electronic Industries Association's Digital Audio Radio Subcommittee.

A total of nine DAR systems were submitted for testing and are labeled in these results as *a* to *i*. Subjective audio quality was assessed in the absence of any transmission error, thus evaluating the quality of the audio source coding component of each system. One of the nine systems was tested with two different comparison references because the sampling rate for that system was lower than for the other 8 systems, and this report refers to 10 systems noted as *a* to *j*.

#### II. Subjective Assessment Procedures

A panel of three expert listeners selected final test materials from the initial pool of program segments received from the evaluation subcommittees. This panel selected nine materials, two of which were stressful to each system under test. These are listed in Table 1.

A total of 21 listeners went through the test process for two days each, to complete the 90 rating trials (10 systems x 9 materials). The equipment, listening environment and procedures were the standard ones used in subjective tests at the CRC as described in ITU-R Rec. BS.1116 [1]. Statistical evaluations assessed each individual's listening expertise by way of a *t*-test, which showed that no listener who took part in the experiment scored below 2.00. Therefore, they all showed that they were able to discriminate correctly between hidden reference and system versions across all the trials in the experiment.

The actual scale used by the subjects is shown in Figure 1. It is a 5 grade rating scale (1.0 to 5.0) where listeners were instructed to use a single decimal point. In effect, this is a 41 point scale. The subjects were instructed to treat this as a continuous scale but, to facilitate the subjects' orientation, category labels were associated with the scale. Thus, 1.0 to 1.9 is a "very annoying" range; 2.0 to 2.9 is "annoying"; 3.0 to 3.9 is "slightly annoying"; 4.0 to 4.9 is "perceptible but not annoying". Finally, 5.0 is "imperceptible".

The listener's task on a trial is to compare each of two alternative versions of an audio material labeled "B" and "C" with a known Reference version, labeled "A", of the same

material. The subject knows that one of the alternatives ("B" or "C") is a "hidden reference", identical to the Reference, and that the other alternative is one that has been processed through a DAR system. The subject does not know which is which, but must decide this through listening. He or she then assigns a grade to both "B" and "C" alternatives, as compared to the known Reference "A", using the 1.0 to 5.0 scale. A is that the alternative the subject has decided is the "hidden reference" must be graded 5.0. And so, *at least one* of the two grades on each trial must be a 5.0

Thus two totally interdependent scores from the listener are recorded on each trial. This deliberate interdependence is handled by subtracting the score given to the true hidden reference from the score given the true processed version (i.e., DSB System minus reference). so that in a graphical plot of outcomes, the data will fall in the same geometric quadrant as they would if the actual 1.0 to 5.0 scores used by the subjects were plotted. Thus the scores are transformed so that the 1.0 to 5.0 range of the original scale becomes, instead, -4.0 to 0.0 in the analysis and presentation of results. These difference grades or "diffgrades" represent the relative differences between the grades given to the hidden reference and the ones given to the DSB system under test.

### **III. Test Results**

For visual clarity, the average quality diffgrades obtained in the experiment are divided between Figures 2(a) and 2(b) rather than being shown within a single graph. Six of them appear in the first figure, four in the second. In addition to the average score among the listeners for each of the audio materials, the overall average diffgrade (the average across all audio materials for each system) is plotted in the "System Averages" column at the right-hand side of these Figures.

Table 2 shows the overall average diffgrade for each audio material and for each system as well as the overall (average) diffgrade for each system in the right-hand column. This table shows all the numbers that are plotted in Figure 2(a) and 2(b). In Table 2, the average diffgrades across all listeners for each audio material occupy a separate row for each DSB system. The average diffgrades are entered to two decimal figures. Systems are arranged by row in alphabetical order using the letters attributed to the ten systems tested -- part of the "double blind" procedures followed throughout the tests..

### **IV. Overall System Results**

The statistical method used to evaluate the present results is the Analysis of Variance (ANOVA) which has been officially recommended in ITU-R Rec. BS.1116 [1]. The experimental design used for these tests permitted the rigorous application of this analytic method. The first item for discussion is the overall average diffgrade for systems. The ANOVA showed that the overall experimental differences among systems in the tests have a very fine resolution of 0.17 of a grade in the transformed diffgrade scale.

For completeness, however, if a reader is interested in evaluating overall differences among audio materials independent of systems (as shown in the averages in the bottom row of Table 2), the critical value provided by the ANOVA is 0.23. This applies to the “without *i* and *j*” averages. Thus, any two of the 9 audio material averages (“without *i* and *j*”) across systems must differ by at least 0.23 before they can be considered significantly different on statistical grounds.

The “two” systems (*i* and *j*) rate differences in the references against which subjects compared them. System are actually the same coding system. But they were treated differently in the experiment because of sampling rate differences in the references against which subjects compared them;. System *i* was always compared with 32 kHz sampling rate references, while for system *j*, the references were always sampled at 48 kHz. The ANOVA showed that the overall difference between *i* and *j* were 0.01, well below the 0.17 needed for a conclusion of significant difference.

## V. Interaction of Systems with Audio Materials

The ANOVA reveals that the resolution for the interaction of audio materials and systems in this experiment is 0.45 of a grade. This too is a very fine degree of resolution for interactions of this type. When comparing diffgrades between any two systems for any given audio material in Figure 2(a) and 2(b), Table 4 and Figure 3, a numerical difference of 0.45 or greater is required before it can be concluded that those two diffgrades are statistically different from each other rather than being due to chance ( $p < 0.05$ ).

## VI. Summary

Table 3 shows system identifications in the first column, summarizing the major outcomes using the three criteria developed and used by the ITU-R to evaluate the relative merits of audio coding systems.

First, the overall average diffgrade is shown for each system. This is presented in the second column of the table. Secondly, to summarize the interaction of audio materials by systems and to indicate the size of the variability of each system, the number of times each system fell below a diffgrade of -1.0 for the 9 materials is presented in the third column of the table. To take statistical error into account, the number of times that any system’s lower error bar fell “below -1.0” for any material in Figure 3 provided the count shown in this third column. Finally, another ITU-R criterion related to the variability or consistency of each system is shown in the fourth column. This is the number of times that a system could be considered “transparent” for an audio item. The number of times that any system’s upper error bar fell above 0.0 in the charts of Figure 3 provided the count shown in this fourth column. Table 3 also shows the systems associated with their letter codes.

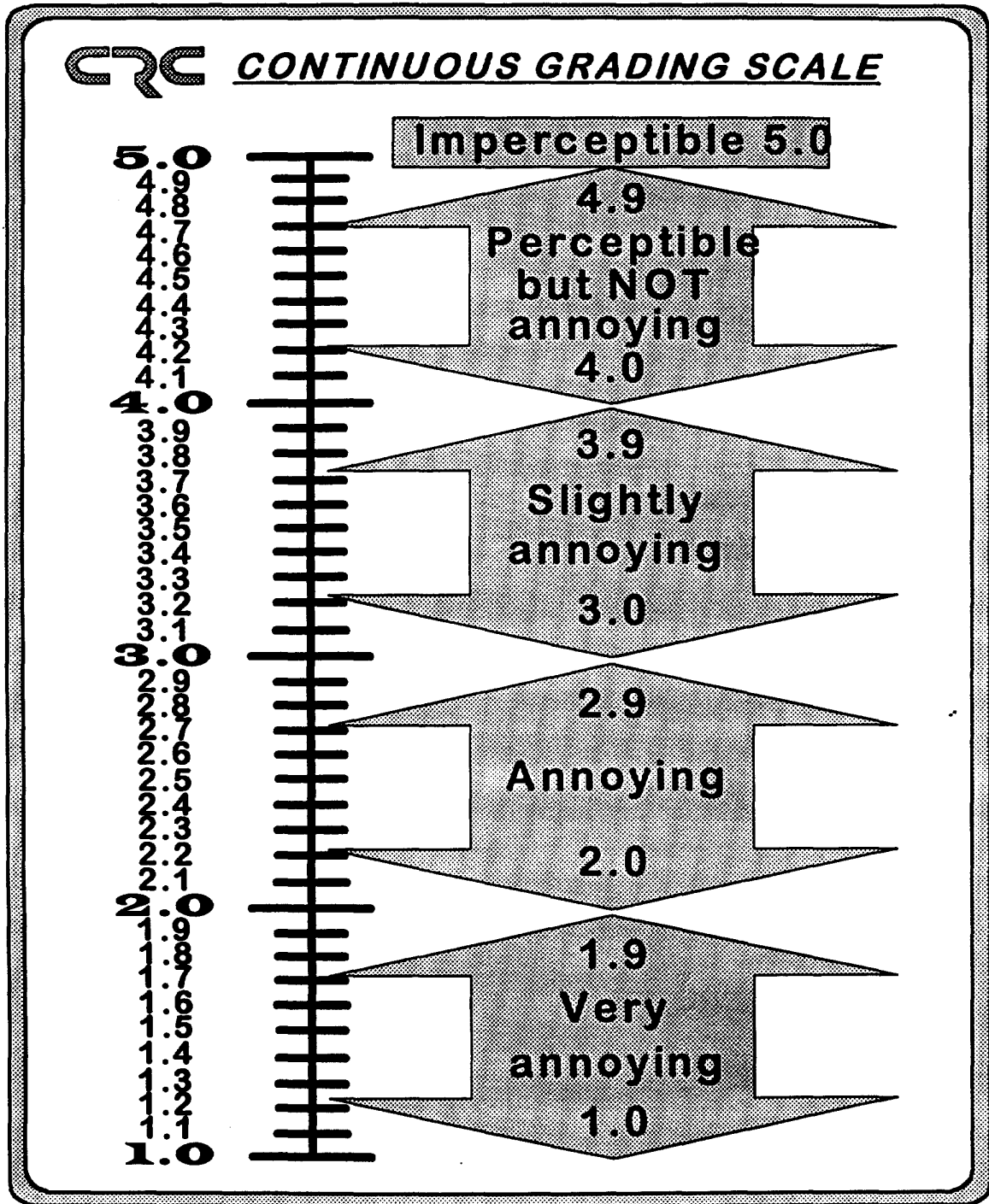


Figure 1 ITU-R continuous 5-grade impairment scale



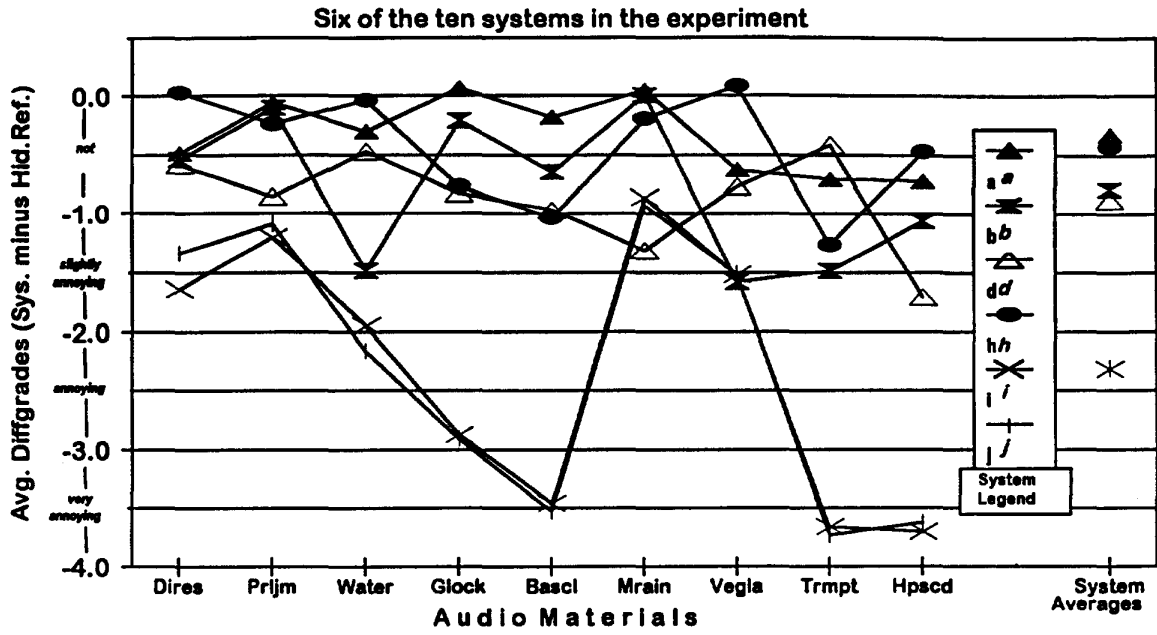


Figure 2(a) Quality test results - systems a,b,d,h,i & j

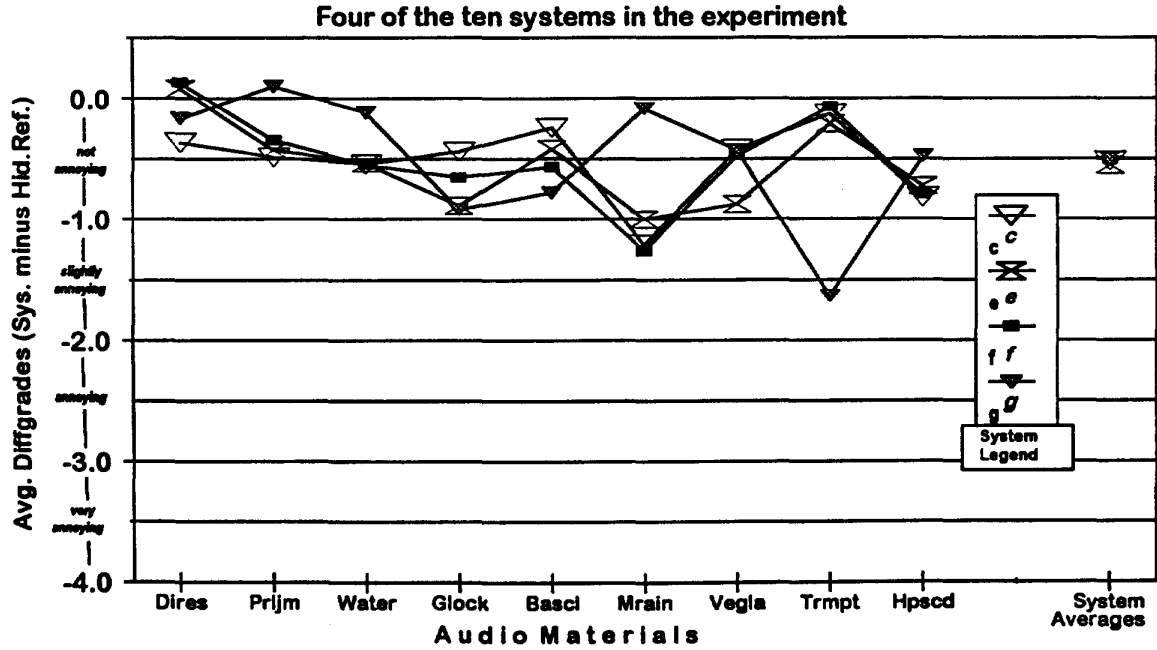
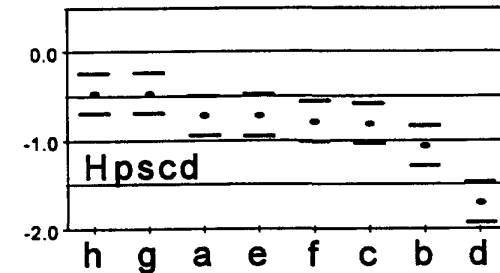
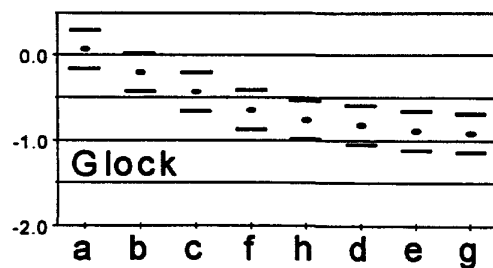
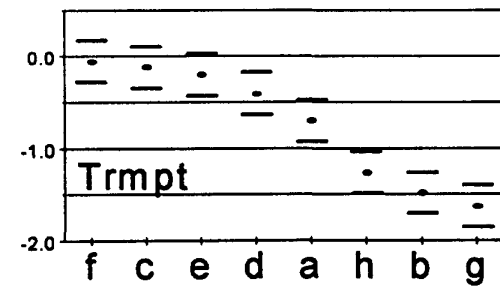
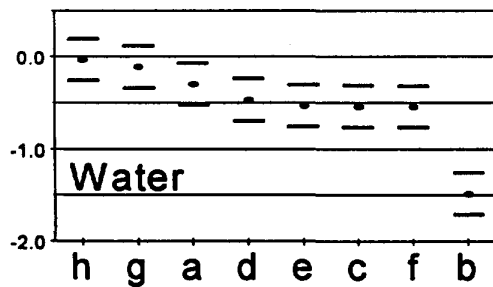
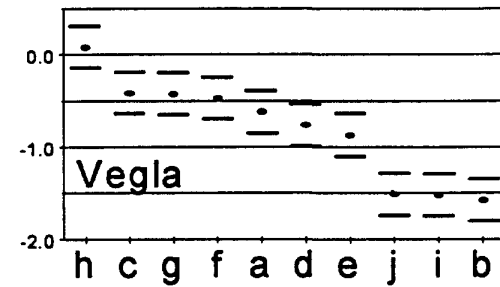
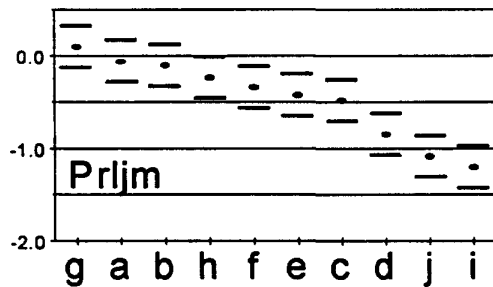
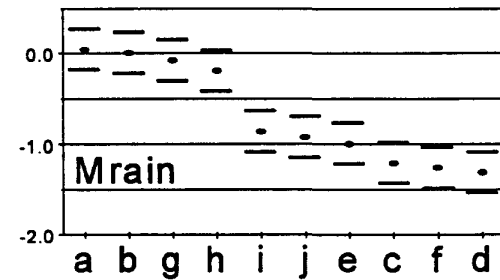
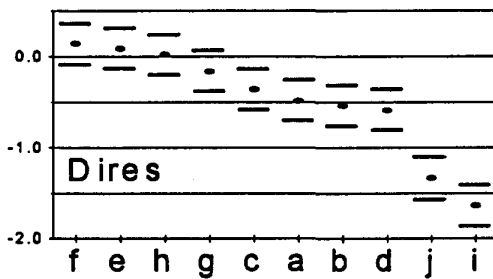
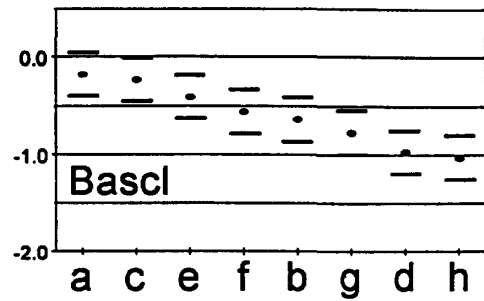


Fig. 2(b) Quality test results - systems c, e, f, and g

**Fig. 3 System Differences Within Audio Materials**  
 Upper and lower statistical boundaries are shown for the average of each system within each audio material. *Only systems with no horizontal overlaps among their boundaries are statistically different.* Within each chart, systems are ordered along the X-axis by the magnitude of their averages.  
 The vertical axes start at -2.0 rather than, as in Figs. 1a and b, at -4.0. Systems *i* and *j* are omitted from those charts where their averages fall below -2.0. At those low values, *i* and *j* are significantly different from all the other 8 systems in those audio materials without ambiguity.



Code	Description	Duration	Source
Dires	Dire Straits cut	30 s	Warner Bros. CD 7599-25264-2 (track 6)
Prjlm	Pearl Jam cut	30 s	Sony/Epic CD ZK53136 (track 3) with processing <sup>1</sup>
Water	Sounds of water	30 s	Roland Dimensional Space Processor Demo. CD
Glock	Glockenspiel	16 s	EBU SQAM CD (track 35/Index 1)
Bascl	Bass Clarinet arpeggio	30 s	EBU SQAM CD (track 17/Index 1) with processing <sup>1</sup>
Mrain	Music and rain	11 s	AT&T mix
Vegla	Susan Vega with glass	11 s	AT&T mix
Trmpt	Muted trumpet	9 s	Original DAT recording, University of Miami
Hpscd	Harpsichord arpeggio	12 s	EBU SQAM CD (track 40/Index 1)

<sup>1</sup> Processing chain used: Aphex Compellor Model 300 (set for leveling only)  
 Dolby Spectral Processor Model 740  
 Aphex Dominator II Model 720

**Table 1 List of audio test materials used in the quality tests**

*The data for a single system are shown throughout each row.*

System	Dires	Prjlm	Water	Glock	Bascl	Mrain	Vegla	Trmpt	Hpscd	Overall Averages
<i>a</i>	-0.49	-0.06	-0.30	0.07	-0.18	0.04	-0.62	-0.70	-0.72	<i>a</i> -0.33
<i>b</i>	-0.54	-0.10	-1.49	-0.21	-0.64	0.00	-1.58	-1.49	-1.07	<i>b</i> -0.79
<i>c</i>	-0.36	-0.49	-0.54	-0.44	-0.24	-1.21	-0.42	-0.12	-0.82	<i>c</i> -0.52
<i>d</i>	-0.59	-0.85	-0.47	-0.82	-0.97	-1.31	-0.77	-0.41	-1.70	<i>d</i> -0.88
<i>e</i>	0.09	-0.43	-0.53	-0.89	-0.41	-1.00	-0.88	-0.20	-0.72	<i>e</i> -0.55
<i>f</i>	0.14	-0.34	-0.55	-0.65	-0.57	-1.26	-0.47	-0.06	-0.80	<i>f</i> -0.51
<i>g</i>	-0.16	0.10	-0.11	-0.92	-0.78	-0.08	-0.43	-1.63	-0.48	<i>g</i> -0.50
<i>h</i>	0.02	-0.24	-0.04	-0.77	-1.04	-0.20	0.08	-1.27	-0.47	<i>h</i> -0.43
<i>i</i>	-1.64	-1.20	-1.95	-2.87	-3.46	-0.86	-1.52	-3.66	-3.70	<i>i</i> -2.32
<i>j</i>	-1.34	-1.09	-2.16	-2.91	-3.52	-0.93	-1.51	-3.73	-3.62	<i>j</i> -2.31
Audio Material Averages	-0.49	-0.47	-0.81	-1.04	-1.18	-0.68	-0.81	-1.33	-1.41	-0.91
Averages Without <i>i</i> and <i>j</i>	-0.24	-0.30	-0.50	-0.58	-0.60	-0.63	-0.64	-0.74	-0.85	-0.56

System *i* received a grade of -1.95 for Water. In view of the statistical error (0.45 of a grade), *i* was omitted from Water in Fig. 2.3 on the next page, along with other instances of *i* and *j* in materials where either of these two systems obtained a diffgrade lower than -2.00. (No systems other than *i* and *j* received any diffgrades below -2.00.)

**Table 2: Average Difference Grades for each of the 9 Audio Materials (columns) by each of the 10 Systems**

System Designation	Overall Average Diffgrade	Number of transparent materials	Number of materials below -1.0
A - Eureka 147, MUSICAM @ 224 kbps	-0.33	4	0
B - Eureka 147, MUSICAM @ 192 kbps	-0.79	3	4
C - AT&T/Lucent, PAC @ 160 kbps	-0.52	2	1
D - AT&T/Amati, DSB PAC @ 160 kbps	-0.88	5	0
E - AT&T/Amati, LSB PAC @ 160 kbps	-0.55	3	2
F - VOA/JPL, PAC @ 160 kbps	-0.51	2	2
G - USADR FM-2, MUSICAM @ 256 kbps	-0.50	2	4
H - USADR FM-1, MUSICAM @ 256 kbps	-0.43	2	4
I - USADR AM, MUSICAM @ 96 kbps (32 kHz reference)	-2.32	0	9
J - USADR AM, MUSICAM @ 96 kbps (48 kHz reference)	-2.31	0	9

**Table 3**  
**Summary of Audio Quality Tests**





---

**Communications  
Research Centre  
Centre de recherches  
sur les communications**

---

## **EIA/NRSC DAR Systems**

# **Subjective Tests of Audio Quality and Transmission Impairments**

## **Final Report**

**21 July 1995**

**TED GRUSEC, GILBERT SOULODRE and LOUIS THIBAUT**

Signal Processing and Psychoacoustics  
Radio Broadcast Technologies Research  
Communications Research Centre  
3701 Carling Avenue  
Ottawa, Ontario, Canada K2H 8S2

---

---

## References

- [1] *Methods for the Subjective Assessment of Small Impairments in Audio Systems Including Multichannel Sound Systems*, ITU-R Recommendation BS.1116, Radiocommunications Study Group, International Telecommunication Union, Geneva, Switzerland
- [2] Thibault L. and Grusec T., "*EIA-DAR Listening Tests - Quality and Impairment Tests Procedures*:", Signal Processing and Psychoacoustics, Radio Broadcast Technologies Research, CRC, Ottawa, Ont., Document submitted to EIA Working Group B, Updated version, 1 December 1993
- [3] Londa, Dave, et al. "*EIA-DAR/NRSC DAB objective tests of transmission impairments*".
- [4] Beaton R.J. and Wong P., *A Disk-based System for the Subjective Assessment of High Quality Audio*, Preprint 3497, AES 94th Convention, Berlin, Germany, 16-19 March 1993

---

# Table of Contents

---

List of Figures .....	iii
List of Tables .....	vii
1. Introduction.....	1
2. Subjective Tests of Audio Quality .....	2
2.1 Selection of critical materials .....	2
2.2 Test procedures.....	2
2.2.1 Listener expertise .....	2
2.2.2 Grading scale.....	4
2.3 Test results.....	6
2.3.1 Graphical and tabular presentations .....	6
2.3.2 Overall systems comparison.....	7
2.3.3 Systems comparison by audio material.....	8
2.4 Discussion and conclusions .....	9
3. Subjective Tests of Transmission Impairments .....	17
3.1 Background .....	17
3.2 Test procedures.....	17
3.2.1 Test materials .....	17
3.2.2 Impairment levels .....	18
3.2.3 Subjective test procedures .....	18
3.3 Test results.....	22
3.3.1 Data presentation and statistical analysis .....	22
3.3.2 Gaussian noise.....	24
3.3.3 Co-channel interference .....	25
3.3.4 Urban slow multipath .....	26
3.3.5 Urban fast multipath .....	27
3.3.6 Rural fast multipath .....	28
3.3.7 Obstructed fast multipath.....	28
3.3.8 DAR lower 1st adjacent to DAR channel interference.....	29
4. Subjective Retests of Audio Quality .....	79
4.1 Background .....	79
4.2 Test procedures.....	79
4.3 Results and discussion.....	81
4.4 Conclusions.....	84



---

5.	Subjective Retests of Transmission Impairments .....	85
5.1	Introduction .....	85
5.2	Test procedures.....	85
5.2.1	Test materials .....	85
5.2.2	Impairment levels .....	85
5.2.3	Subjective test procedures .....	86
5.3	Test results.....	87
5.3.1	Gaussian noise.....	87
5.3.2	Co-channel interference .....	88
6.	Facilities .....	90
6.1	Playback system .....	90
6.2	Listening room.....	94
7.	Acknowledgements.....	96
	References .....	97

---

## List of Figures

Figure 2.1	ITU-R continuous 5-grade impairment scale .....	5
Figure 2.2a	Quality test results - systems <i>a, b, d, h, i</i> and <i>j</i> .....	12
Figure 2.2b	Quality test results - systems <i>c, e, f</i> and <i>g</i> .....	12
Figure 2.3	System differences within audio materials .....	14
Figure 3.1	Failure characteristic curves of DAR systems <i>a</i> through <i>h</i> tested with the gaussian noise impairment and the Glockenspiel test material .....	30
Figure 3.2	Failure characteristic curves of DAR systems <i>a</i> through <i>h</i> tested with the gaussian noise impairment and the Soprano test material .....	31
Figure 3.3	Failure characteristic curves of DAR systems <i>a</i> through <i>h</i> tested with the gaussian noise impairment and the Clarinet test material .....	32
Figure 3.4	Failure characteristic curves for system <i>a</i> for all three test materials and with the gaussian noise impairment .....	33
Figure 3.5	Failure characteristic curves for system <i>b</i> for all three test materials and with the gaussian noise impairment .....	34
Figure 3.6	Failure characteristic curves for system <i>c</i> for all three test materials and with the gaussian noise impairment.....	35
Figure 3.7	Failure characteristic curves for system <i>d</i> for all three test materials and with the gaussian noise impairment .....	36
Figure 3.8	Failure characteristic curves for system <i>e</i> for all three test materials and with the gaussian noise impairment.....	37
Figure 3.9	Failure characteristic curves for system <i>f</i> for all three test materials and with the gaussian noise impairment.....	38
Figure 3.10	Failure characteristic curves for system <i>g</i> for all three test materials and with the gaussian noise impairment .....	39
Figure 3.11	Failure characteristic curves for system <i>h</i> for all three test materials and with the gaussian noise impairment .....	40

---

Figure 3.12	Failure characteristic curves for system <i>i</i> for all three test materials and with the gaussian noise impairment.....	41
Figure 3.13	Failure characteristic curves of DAR systems <i>a</i> through <i>h</i> tested with co-channel interference and the Glockenspiel test material .....	42
Figure 3.14	Failure characteristic curves of DAR systems <i>a</i> through <i>h</i> tested with co-channel interference and the Soprano test material .....	43
Figure 3.15	Failure characteristic curves of DAR systems <i>a</i> through <i>h</i> tested with co-channel interference and the Clarinet test material .....	44
Figure 3.16	Failure characteristic curves for system <i>a</i> for all three test materials and with the co-channel interference.....	45
Figure 3.17	Failure characteristic curves for system <i>b</i> for all three test materials and with the co-channel interference.....	46
Figure 3.18	Failure characteristic curves for system <i>c</i> for all three test materials and with the co-channel interference.....	47
Figure 3.19	Failure characteristic curves for system <i>d</i> for all three test materials and with the co-channel interference.....	48
Figure 3.20	Failure characteristic curves for system <i>e</i> for all three test materials and with the co-channel interference.....	49
Figure 3.21	Failure characteristic curves for system <i>f</i> for all three test materials and with the co-channel interference.....	50
Figure 3.22	Failure characteristic curves for system <i>g</i> for all three test materials and with the co-channel interference.....	51
Figure 3.23	Failure characteristic curves for system <i>h</i> for all three test materials and with the co-channel interference.....	52
Figure 3.24	Failure characteristic curves for system <i>i</i> for all three test materials and with the co-channel interference.....	53
Figure 3.25	Failure characteristic curves of the three DAR systems tested with the urban slow multipath impairment and the Glockenspiel test material.....	54
Figure 3.26	Failure characteristic curves of the three DAR systems tested with the urban slow multipath impairment and the Soprano test material .....	55

---

---

Figure 3.27	Failure characteristic curves of the three DAR systems tested with the urban slow multipath impairment and the Clarinet test material .....	56
Figure 3.28	Failure characteristic curves for system <i>a</i> for all three test materials and with the urban slow multipath impairment.....	57
Figure 3.29	Failure characteristic curves for system <i>b</i> for all three test materials and with the urban slow multipath impairment.....	58
Figure 3.30	Failure characteristic curves for system <i>c</i> for all three test materials and with the urban slow multipath impairment.....	59
Figure 3.31	Failure characteristic curves of the four DAR systems tested with the urban fast multipath impairment and the Glockenspiel test material .....	60
Figure 3.32	Failure characteristic curves of the four DAR systems tested with the urban fast multipath impairment and the Soprano test material .....	61
Figure 3.33	Failure characteristic curves of the four DAR systems tested with the urban fast multipath impairment and the Clarinet test material .....	62
Figure 3.34	Failure characteristic curves for system <i>a</i> for all three test materials and with the urban fast multipath impairment.....	63
Figure 3.35	Failure characteristic curves for system <i>b</i> for all three test materials and with the urban fast multipath impairment.....	64
Figure 3.36	Failure characteristic curves for system <i>c</i> for all three test materials and with the urban fast multipath impairment.....	65
Figure 3.37	Failure characteristic curves for system <i>e</i> for all three test materials and with the urban fast multipath impairment.....	66
Figure 3.38	Failure characteristic curves for the two DAR systems tested with the rural fast multipath impairment and the Glockenspiel test material.....	67
Figure 3.39	Failure characteristic curves for the two DAR systems tested with the rural fast multipath impairment and the Soprano test material.....	68
Figure 3.40	Failure characteristic curves for the two DAR systems tested with the rural fast multipath impairment and the Clarinet test material .....	69
Figure 3.41	Failure characteristic curves for system <i>e</i> for all three test materials with the rural fast multipath impairment.....	70

---

Figure 3.42	Failure characteristic curves for system <i>f</i> for all three test materials with the rural fast multipath impairment.....	71
Figure 3.43	Failure characteristic curves of the two DAR systems tested under obstructed fast multipath and the Glockenspiel test material.....	72
Figure 3.44	Failure characteristic curves of the two DAR systems tested under obstructed fast multipath and the Soprano test material.....	73
Figure 3.45	Failure characteristic curves of the three DAR systems tested under obstructed fast multipath and the Clarinet test material.....	74
Figure 3.46	Failure characteristic curves for system <i>a</i> for all three test materials with the obstructed fast multipath impairment.....	75
Figure 3.47	Failure characteristic curves for system <i>b</i> for all three test materials with the obstructed fast multipath impairment.....	76
Figure 3.48	Failure characteristic curve for system <i>c</i> for the Clarinet test material with the obstructed fast multipath impairment.....	77
Figure 3.49	Failure characteristic curves for system <i>i</i> for all three test materials with the DAR lower 1st adjacent to DAR channel interference.....	78
Figure 4.1	Average diffgrade of the 8 subjects in the quality retest for all the systems in that experiment .....	82
Figure 6.1	Equipment set-up .....	91
Figure 6.2	Video screen used by the listeners during blind rating sessions .....	93
Figure 6.3	CRC listening room layout.....	94
Figure 6.4	CRC listening room - Background noise levels.....	95
Figure 6.5	CRC listening room - Reverberation time (RT-60).....	95

# List of Tables

Table 2.1	Listeners' <i>t</i> -scores .....	3
Table 2.2	List of audio materials used in the quality tests.....	7
Table 2.3	Summary of the audio quality test results .....	11
Table 2.4a	Average difference grades for each of the 9 audio materials (columns) by each of the 10 systems - systems listed in descending order of overall average score .....	13
Table 2.4b	Average difference grades for each of the 9 audio materials (columns) by each of the 10 systems - systems listed in descending order of average score for each test material .....	13
Table 2.5	Magnitudes of the differences between each pair of the 10 systems (a to j) for each of the 9 audio materials .....	15
Table 3.1	List of audio materials used in the impairment tests.....	18
Table 3.2	Grouping of the seven impairments into five experiments .....	22
Table 3.3	$C_0/N_0$ ratios at TOA and POF and failure margin for the AWGN impairment .....	25
Table 3.4	D/U ratios at TOA and POF and failure margin for co-channel interference.....	26
Table 3.5	$C_0/N_0$ ratios at TOA and POF and failure margin for the urban slow multipath impairment.....	27
Table 3.6	$C_0/N_0$ ratios at TOA and POF and failure margin for the urban fast multipath impairment.....	28
Table 3.7	$C_0/N_0$ ratios at TOA and POF and failure margin for the rural fast multipath impairment.....	28
Table 3.8	$C_0/N_0$ ratios at TOA and POF and failure margin for the obstructed fast multipath impairment.....	29
Table 3.9	D/U ratios at TOA and POF and failure margin for the DAR lower 1st adjacent to DAR channel interference .....	29

Table 4.1	Average diffgrades for each system at each audio material and for overall system averages .....	82
Table 5.1	Comparison of Co/No ratios at TOA and POF and failure margin (FM) for AWGN impairment for systems <i>e</i> and <i>k</i> . .....	87
Table 5.2	Comparison of Co/No ratios at TOA and POF and failure margin (FM) for AWGN impairment for systems <i>h</i> and <i>l</i> . .....	88
Table 5.3	Comparison of D/U ratios at TOA and POF and failure margin (FM) for the co-channel interference for systems <i>e</i> and <i>k</i> . .....	89
Table 5.4	Comparison of D/U ratios at TOA and POF and failure margin (FM) for the co-channel interference for systems <i>h</i> and <i>l</i> . .....	89
Table 6.1	List of equipment.....	90

---

# 1. Introduction

This final report describes the procedures as well as the results of three series of subjective tests conducted at the Communications Research Centre (CRC), Ottawa, Ontario, Canada from June 1994 to March 1995. These tests were performed to assess the audio quality of Digital Audio Radio (DAR) systems submitted to the DAR subcommittee of the Electronics Industries Association (EIA) and the DAB subcommittee of the National Radio System Committee (NRSC).

A total of nine DAR systems were submitted for testing and these are labeled as *a* to *i* in this report. In the first series of tests, the subjective audio quality of the DAR systems was assessed in the absence of any transmission error. Essentially, this test evaluates the quality of the audio source coding component of each DAR system. In the second series of tests, the subjective audio quality of the DAR systems was assessed in the presence of transmission errors. Both the onset of detectability of transmission errors (namely the threshold of audibility) as well as the failure characteristic of the DAR systems (from threshold of audibility down to point of failure) were determined for the following seven types of impairment:

1. Gaussian noise
2. Co-channel interference
3. Urban slow multipath
4. Urban fast multipath
5. Rural fast multipath
6. Obstructed fast multipath
7. DAR lower first adjacent to DAR channel interference

All nine systems were tested in the presence of additive white gaussian noise and co-channel interference. A subset of the systems were tested for the other impairments. In the third and last series of tests, two modified DAR systems were retested. The audio quality (in absence of transmission error) of one of these systems was subjectively retested while the threshold of audibility and the point of failure of both DAR systems were re-evaluated for the gaussian noise and co-channel interference.

The procedures and results of the subjective tests of audio quality in the absence of transmission errors are described in chapter 2 of this report while those in the presence of transmission impairments are contained in chapter 3. Chapter 4 contains a description of the procedures and results of the subjective retests of audio quality while those of the impairment retests are contained in chapter 5. In chapter 6, the facilities used for the tests are described. Finally, acknowledgments are given in chapter 7.



---

## 2. Subjective Tests of Audio Quality

### 2.1 Selection of critical materials

Various organizations submitted potentially critical audio test materials from a wide variety of sources. Processing of these materials through the 9 DAR systems to be assessed for subjective quality took place at the NASA Lewis Research Center (LeRC) in Cleveland, USA. A panel of three expert listeners convened at the CRC for the selection of final test materials from the initial pool received from Cleveland. Two of these experts were US citizens and one was a Canadian.

The task of the panel was to find at least 2 stressful materials for each system. The panel uncovered a total of nine materials that met this requirement. The experimental results suggest that the panel's choices were excellent ones. Some evidence that supports this is presented in section 2.2.1.

One of the 9 systems was tested with two different comparison references because the sampling rate for that system was lower than for the other 8 systems. Accordingly, our report refers to "10" systems rather than 9. The identity of the systems was unknown to the CRC and, at the time of writing, still remains unrevealed. In the present report, the "10" systems are designated "a" to "j".

### 2.2 Test procedures

#### 2.2.1 Listener expertise

A total of 21 listeners went through the test process for two days to complete the 90 rating trials (10 systems x 9 materials). Thirteen of these subjects were from the US. The other listeners included one from the UK, one from France and 6 from Canada.

The equipment, listening environment and procedures were the standard ones used in subjective tests at the CRC as described in ITU-R Rec. BS.1116 [1] and in the December 1st, 1993 revision of the "*Quality and impairment tests procedures*" document submitted to the EIA-DAR Working Group B [2]. As in all similar experiments at the CRC, we applied a rigorous statistical criterion to the data to assess each individual's listening expertise.

This criterion makes use of the fact that on each of the trials in a rating experiment (90 trials in the present case), the listener supplies two ratings; one for what the subject concluded was the "hidden reference", and one for the presumed system version of the audio material heard on that trial. (section 2.2.2 below provides more details about the subject's task and the rating procedure.) Over all the 90 trials, then, we have two distributions of grades for each subject - one distribution for the true hidden references; the other for the DAR system versions. If the subject was discriminating between these two classes of events in an objectively correct way, then these two distributions will have averages which are reliably different from each other. On the other hand, if the subject's discriminations between system and hidden reference versions was

faulty, then the grades in the two distributions will resemble a "chance guessing" pattern, and the two distributions will be statistically identical.

Evaluation of these two possibilities is by way of a *t*-test, a standard statistical test applicable to comparing such distributions, as outlined in most textbooks of statistical analysis. For 90 data points, as in the present experiment, the magnitude of *t* must equal or exceed a value of approximately 2.00 in order that one may conclude ( $p < 0.05$ ) that the subject was truly discriminating beyond chance between hidden reference and system versions. If the value of *t* for a listener is smaller than 2.00, one must conclude that the subject was, overall, guessing.

As will become apparent in the presentation of results later in this report, there were two DAR systems in the experiment where detecting degradations appeared to be obviously "easy". Listeners gave those systems (*i* and *j*, Fig. 2.2a) consistently low grades and all subjects made very few errors in correctly discriminating the system from the hidden reference versions. Including these easy trials in the calculation of the *t*-scores would artificially inflate the merit scores of each subject. To avoid this, the *t*-scores which we present in the following table conservatively omit those 18 very easy trials and are calculated on only the other 72 trials. The scores for all subjects in the quality experiment are shown below in descending order of merit.

Subject	t-score
1	11.62
2	8.51
3	7.35
4	6.50
5	6.34
6	7.56
7	6.49
8	6.38
9	6.19
10	5.68
11	5.61
12	5.16
13	5.00
14	4.57
15	4.30
16	3.48
17	3.35
18	2.84
19	2.57
20	2.46
21	2.43

Table 2.1 Listeners' *t*-scores

As is clear in the above Table, no listener who took part in the experiment scored below 2.00. Therefore, they all showed that they were able to discriminate correctly between hidden reference and system versions across all the trials in the experiment. No-one's data had to be rejected for

---

insufficient listener merit. And so, the data of all the 21 participants are included in the analyses. As mentioned, this *t*-score was very conservative since the easiest two systems, *i* and *j*, were omitted in the calculation.

We must emphasize that these are high expertise scores. This fine listener performance was due, in addition to other things, to very well-chosen materials which were stressful to the DAR systems under test, while, at the same time, relatively easy for listeners to grade. This reveals that the three experts who chose the materials performed their task extremely well.

We should point out that listeners can all be, as in the present case, sufficiently expert, yet disagree with each other in the relative ratings they assign. If this were the case, then the results would be inconsistent, hence statistically unreliable. As the results presented below will show, the opposite was found. In other words, the listeners, all sufficiently expert, were highly consistent with each other. Hence the experiment provided clear, interpretable outcomes as will be seen in sections 2.3 and 2.4.

### 2.2.2 Grading scale

The test procedures used complied with those described in ITU-R Recommendation BS.1116 [1]. The 90 trials of the experiment were rated in two consecutive days of 45 trials each by the listeners. The procedure included a training phase followed by a blind rating phase. The morning of each day was devoted to the training session using those audio materials to be rated in the subsequent blind rating. Subjects worked together in groups of 2 or 3 and were allowed as much time as they required for training. Subjects trained using the same hard disk playback system used in the blind rating tests.

Blind rating took place in the afternoon of each day and was performed individually by each listener. The subjects alternated with each other in the listening room, so while one was working, the other(s) rested. Listeners used the disk-based playback system described in section 5 of this report which allowed seamless switching between the stimuli to be compared. Listeners were able to take as much time as they needed on each trial, switching as often as they liked, until satisfied with the numerical ratings they were asked to assign. They were also free to use either the loudspeakers or headphones to make a judgment, whichever they felt was the most critical transducer on any given trial.

The actual scale used by the subjects is shown in Fig. 2.1. It is a 5 grade rating scale (1.0 to 5.0) where listeners were instructed to use a single decimal point. In effect, this is a 41 point scale. The subjects were instructed to treat this as a continuous scale but, to facilitate the subjects' orientation, category labels were associated with the scale. Thus, 1.0 to 1.9 is a "very annoying" range; 2.0 to 2.9 is "annoying"; 3.0 to 3.9 is "slightly annoying"; 4.0 to 4.9 is "perceptible but not annoying". Finally, 5.0 is "imperceptible".

The listener's task on a trial is to compare each of two alternative versions, labeled as "B" and "C", of an audio material with a known Reference version, labeled "A", of the same material. The subject knows that one of the alternatives ("B" or "C") is a "hidden reference", identical to

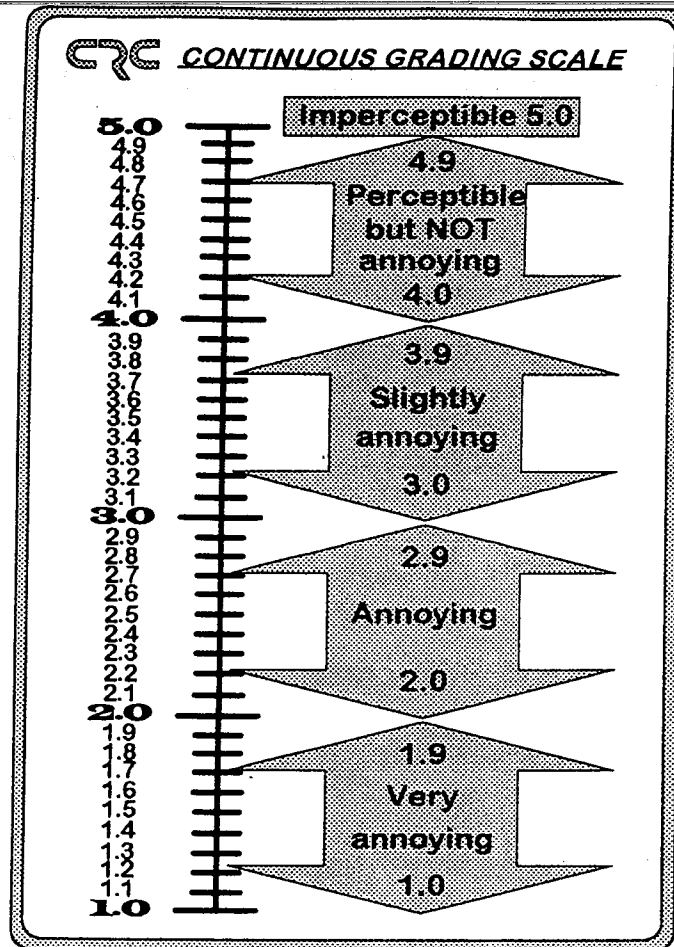


Fig. 2.1 ITU-R continuous 5-grade impairment scale

the Reference, and that the other alternative is one that has been processed through a DAR system. The subject does not know which is which, but must decide this through listening. He or she then assigns a grade to both “B” and “C” alternatives, as compared to the known Reference “A”, using the 1.0 to 5.0 scale. A constraint is that the alternative the subject has decided is the “hidden reference” must be graded 5.0. And so, *at least one* of the two grades on each trial must be a 5.0.

Thus two totally interdependent scores from the listener are recorded on each trial. This deliberate interdependence is handled by subtracting the score given to the true hidden reference from the score given the true processed version (i.e., DAR System minus Reference). The reason for subtracting in this direction rather than the opposite one (Reference minus DAR System) is only so that in a graphical plot of outcomes, the data will fall in the same geometric quadrant as they would if the actual 1.0 to 5.0 scores used by the subjects were plotted.

This subtraction means, however, that the scores are transformed so that the 1.0 to 5.0 range of the original scale becomes, instead, -4.0 to 0.0 in the analysis and presentation of results. These

difference grades or “diffgrades” represent the relative differences between the grades given to the hidden reference and the ones given to the DAR system under test. For practical purposes, it can be assumed that a direct analogy holds between the original 1.0 to 5.0 scale and the -4.0 to 0.0 diffgrades in the report. Accordingly, to facilitate interpretation for the reader who is familiar with the 1.0 to 5.0 scale used by subjects, the “annoyance” category labels are shown between the Y-axis numbers in Fig. 2.2a and 2.2b. Thus, in these figures, the words “perceptible but not annoying” are shown between 0.0 and -1.0, “slightly annoying” appears between -1.0 and -2.0, and so on.

The 90 materials to be rated were divided into 6 sessions of 15 trials each. Within each session, items were ordered in a quasi-random, unpredictable way. Moreover, the order of presentation of the sessions was changed from listener to listener to ensure that time-correlated factors (such as fatigue) would not differentially affect any level of any of the factors under test.

## 2.3 Tests results

### 2.3.1 Graphical and tabular presentations

For visual clarity, the average quality diffgrades of the experiment are divided between Fig. 2.2a and 2.2b rather than being shown within a single graph. The two highest ranking systems overall are shown in Fig. 2.2a along with the 7th through 10th ranking systems. The 3rd to 6th ones are found in Fig. 2.2b. The ranking referred to here is an ordering by the overall average diffgrade (average across all listeners and all audio materials for each system). These overall averages are plotted in the “Overall Averages” column at the right-hand side of these Figures. Also shown are the rank order of the systems in the legends in these figures.

Table 2.4a shows the overall average diffgrade for each system in the right-hand column. That Table, as well as Table 2.4b on the same page, shows all the numbers that are plotted in Fig. 2.2a and 2.2b. In Table 2.4a, the average diffgrades across all listeners for each audio material occupy a separate row for each DAR system. In Table 2.4b, the average diffgrades in each column are ordered by the magnitude of those diffgrades for each audio material. The system represented in each cell is clearly indicated in that table.

The ordering chosen for the audio materials along the X-axis in Fig. 2.2a and 2.2b reflects the average diffgrades given to each material across all systems. The averages for audio materials across systems are shown in the two bottom rows of Table 2.4a. It is the ones “without  $i$  and  $j$ ” that were used to determine the order used in Fig. 2.2a and 2.2b since gradings for those two systems were erratic. This ordering means that the most transparent materials (those receiving the highest grades) are shown at the left and, as we move toward the right, the materials are less and less transparent, i.e., more and more critical for revealing system differences.

The five-letter abbreviations used in the text and the sources for the audio materials are shown in the table below:

Code	Description	Duration	Source
Dires	Dire Straits cut	30 s	Warner Bros. CD 7599-25264-2 (track 6)
Prijm	Pearl Jam cut	30 s	Sony/Epic CD ZK53136 (track 3) with processing <sup>1</sup>
Water	Sounds of water	30 s	Roland Dimensional Space Processor Demo. CD
Glock	Glockenspiel	16 s	EBU SQAM CD (track 35/Index 1)
Bascl	Bass Clarinet arpeggio	30 s	EBU SQAM CD (track 17/Index 1) with processing <sup>1</sup>
Mrain	Music and rain	11 s	AT&T mix
Vegla	Susan Vega with glass	11 s	AT&T mix
Tmpt	Muted trumpet	9 s	Original DAT recording, University of Miami
Hpscd	Harpsichord arpeggio	12 s	EBU SQAM CD (track 40/Index 1)

<sup>1</sup> Processing chain used: Aphex Compellor Model 300 (set for leveling only)  
Dolby Spectral Processor Model 740  
Aphex Dominator II Model 720

Table 2.2 List of audio test materials used in the quality tests

### 2.3.2 Overall systems comparison

The first criterion by which the DAR systems under test can be compared and ranked is by their overall average diffgrade. Statistical analysis (Analysis of Variance, or ANOVA) reveals that the overall experimental differences among systems have a very fine resolution of 0.17 of a grade in the transformed diffgrade scale. In other words, any two systems that are numerically different by 0.17 or more in overall average diffgrades are reliably different ( $p < 0.05$ ) from each other. If they differ by less than 0.17, this difference is not considered statistically significant since it could too easily be due to chance (i.e., with  $p > 0.05$ , the generally accepted cut-off).

By this rigorous criterion, systems *a* and *h* (Fig. 2.2a, Table 2.4a) are not reliably different from each other overall, although system *a* is rated 0.10 of a grade higher than *h*. System *a* is rated significantly higher than the group consisting of systems *g*, *f*, *c* and *e* (Fig. 2.2b); however, system *h* does not differ reliably from these four systems. These four (Fig. 2.2b) form a cohesive group with little to differentiate among them in overall grades. There is only a 0.05 difference in overall average diffgrade between the highest of the four (*g*) and the lowest (*e*).

Another cohesive group of systems consists of systems *b* and *d* which differ from each other by only 0.09. These two (Fig. 2.2a) are reliably different from the four systems (*g*, *f*, *c* and *e*) of Fig 2.2b and from both *a* and *h*. Finally, systems *i* and *j* are virtually identical with only a 0.01 difference magnitude between their overall average diffgrade. By this measure, the different audio sampling rates (32 kHz for system *i* and 48 kHz for system *j*) of the two references that were used with this system made no significant difference (statistically speaking) in its overall diffgrades; nor, as will be seen in more detail in the next section, in its pattern among the 9 audio test materials. Overall, these “two” (*i* and *j*) rank lower by more than a full grade, in the 5-grade scale, from the systems that are nearest to them; thus they comprise a significantly different “group” from all the others.

To summarize, there are four distinct groups of systems, reliably different from each other on statistical grounds. Only the grouping of *h* is questionable, since it is not reliably different from either *a* or from the systems in the second ranking group :

Highest	1.	<i>a</i> , (and possibly <i>h</i> )
	2.	<i>g, f, c</i> and <i>e</i> (and possibly <i>h</i> )
	3.	<i>b</i> and <i>d</i>
Lowest	4.	<i>i</i> and <i>j</i>

It should be noted that, in overall average diffgrades, all the systems, except for the low ranking *i* and *j*, fall within the “perceptible but not annoying” category at the top of the grading scale. In fact, both *a* and *h* are in the very top half of that range, with overall diffgrades higher than -0.5 (equivalent, as explained before, to 4.5 in the 1.0 to 5.0 scale actually used by the subjects). The second ranking systems (*g, f, c, e*) are just slightly below -0.5; and the third group (*b, d*) is lower than that, but is still above -1.0 in the “perceptible but not annoying” range. Systems *i* and *j* stand apart quite clearly, falling into the “annoying” range, far below all the other systems.

### 2.3.3 Systems comparison by audio material

When comparing the audio quality of DAR systems, the kinds of overall averages discussed above need qualification by examining the pattern of interactions that occur between combinations of specific audio materials on the one hand, and specific systems on the other. These interactions may sometimes affect the interpretation of the results made from the overall picture as presented above. In the present experiment, these interactions do not greatly disturb the interpretations presented in the previous section. Rather they confirm them and provide some further illumination.

First, statistical analysis (ANOVA) reveals that the resolution for the interaction of audio materials and systems in this experiment is 0.45 of a grade. This too is a very fine degree of resolution for interactions of this type. It means that when comparing any two average diffgrades obtained for any given audio material and for any given system in Fig. 2.2a and 2.2b or Table 2.4a and 2.4b, a numerical difference of 0.45 or greater is required before it can be concluded that those two diffgrades are statistically different from each other rather than being due to chance ( $p < 0.05$ ).

In this experiment, there is little interest in comparing average diffgrades between different audio materials. However, comparing the average diffgrade of each system within a given material is of importance. In particular, it is pertinent to see for any given material, whether any particular obtained difference between systems is meaningful.

To facilitate such comparisons, a chart of 9 figures, one for each audio material, was prepared. These are shown on a single page as Fig. 2.3. It should be noted that, in order to fit all 9 figures in one page, the full range of diffgrades is not shown in Fig. 2.3, but only the range between -2.0 and 0.0. This space saving was possible because only systems *i* and *j* have any average diffgrades

below -2.0, and in all those cases,  $i$  and  $j$  were very reliably different from all the other systems (see the footnote under Table 2.4b).

Fig 2.3 shows the average diffgrades per system for each of the 9 audio test materials (corresponding to the abscissa or “columns” of Fig. 2.2a and 2.2b), arranged in descending order of magnitude from left to right. Table 2.4b contains the actual data shown in Fig. 2.3. The reader can make any comparisons desired using the 0.45 critical value mentioned above for interactions between systems and audio materials. To facilitate this comparison, “error bars” are shown in Fig. 2.3 for each system average in each audio material chart. Each error bar deviates from the average for a system by + and - one half of 0.45 (the critical magnitude). Thus, one can evaluate visually whether any two systems within each audio material should be considered to be statistically different by noting whether or not there is any overlap (along the Y-axis) between the error bars for the two systems being compared. If there is overlap, those two systems do not statistically differ on that material. If there is no overlap, it indicates that the difference between the average diffgrades of those two systems on that material is statistically reliable ( $p < 0.05$ ).

The actual numerical differences underlying the charts of Fig. 2.3 are presented in the two pages of Table 2.5. There, in a manner that parallels the charts of Fig. 2.3, differences are seen numerically rather than visually as in the Figure. Those that are statistically reliable are shown as **bold** numbers. There is redundancy in Table 2.5 in that the numbers along the right hand side above the blank diagonal are mirror images of the left hand ones below that diagonal. However, retaining this redundancy means that a reader can get all the information needed more easily. For completeness, we should also point out that if a reader is interested in evaluating overall differences among audio materials independent of systems (as are shown in the averages in the bottom row of Table 2.4a), the critical value shown by the ANOVA is 0.23. This applies to the “without  $i$  and  $j$ ” averages. Thus, any two of the 9 audio material averages (“without  $i$  and  $j$ ”) across systems must differ by at least 0.23 before they can be considered significantly different on statistical grounds.

## 2.4 Discussion and conclusions

The quality assessment results show that system  $a$  has the highest overall rank and the most consistent ratings across the audio materials tested (Fig. 2.2a). For all audio materials,  $a$  never falls below the “perceptible but not annoying” range. Even on the most revealing materials (Vegla, Trmpt, Hpscd) it only dips to the lower half of the “perceptible but not annoying” range. Taking the statistical error into consideration,  $a$  remains within the “perceptible but not annoying” range even under the most pessimistic view (i.e., by looking at the lower bound of error bars for Vegla, Trmpt, Hpscd in Fig 2.3).

System  $h$  is less consistent. Though it ranks significantly higher statistically than  $a$  in two materials, Vegla and Dires, the opposite is true for three (Glock, Bascl and Trmpt) where  $a$  is rated reliably higher than  $h$ . System  $h$  looks transparent for four stimuli (Dires, Water, Mrain and Vegla), but is rated in the “slightly annoying” range on Bascl and Trmpt. System  $a$ , like  $h$ , also appears transparent for four materials (Prilm, Glock, Bascl, and Mrain), but, as discussed



previously, *a* is never rated in the “slightly annoying” range. Other systems (i.e., *c*, *e* and *f*) perform significantly better than *a* on two materials, namely Dires and Trmpt. But all those three systems that are better than *a* at any point, fall below *a* on other materials.

The inconsistent behavior across audio materials of the second group (*g*, *f*, *c*, *e*) produces the lower overall ranking obtained by those systems. Their descent into the “slightly annoying” range is quite clear for Mrain and Trmpt (Fig. 2.2b).

Turning to the third group, generally lower performance in the “perceptible but not annoying” range, coupled with more frequent and more severe dips into the “slightly annoying” range, accounts for the lower overall ranking of systems *b* and *d*.

Finally, the performance of systems *i* and *j* across audio materials is entirely consistent with its low “annoying” overall rating. The average diffgrades of *i* are virtually identical to those of *j* on every material, showing that they are perceived as being really the same system. This is true even though the reference used for system *i* was sampled at 32 kHz and that for system *j* was a CD quality signal sampled at 48 kHz. The sampling rate variation of the reference was therefore found to be irrelevant. The highest rated materials for systems *i* and *j* (Prilm and Mrain) are close to -1.0. Three materials (Bascl, Trmpt and Hpscd) are rated well within the “very annoying” range of the scale. The remaining four materials (Dires, Water, Glock and Vegla) were rated in the “slightly annoying” or “annoying” ranges.

One way to summarize the audio materials by systems interaction is to see how many times each system fell below -1.0 for the 9 materials. A reasonable way to do this that takes the statistical error into account is to count as “below -1.0” any system whose lower error bar in Fig. 2.3 falls into that range. The following list shows those counts:

<u>Count</u>	<u>Systems</u>
0	<i>a</i>
1	no system
2	<i>c, g, h</i>
3	<i>e</i>
4	<i>b</i>
5	<i>d</i>
9	<i>i</i> and <i>j</i>

Comparing this to the previous list of “Highest” to “Lowest” ranking which was based on overall averages (section 2.3.2) shows an almost identical pattern. And so these interactions support the previous conclusions completely but place system *h* with the second best group rather than with system *a*.

The tabulation below summarizes the entire discussion and presents the systems in descending order of overall merit, with groupings to reflect the patterns of significant differences.

	System	Overall Average Diffgrade	Number of transparent materials	Number of materials below -1.0
Group 1	<i>a</i>	-0.33	4	0
Group 2	<i>h</i>	-0.43	4	2
	<i>g</i>	-0.50	4	2
	<i>f</i>	-0.51	2	2
	<i>c</i>	-0.52	1	2
	<i>e</i>	-0.55	2	3
Group 3	<i>b</i>	-0.79	3	4
	<i>d</i>	-0.88	0	5
Group 4	<i>j</i>	-2.31	0	9
	<i>i</i>	-2.32	0	9

Table 2.3 Summary of the audio quality test results

This experiment is one of the least ambiguous in outcomes among related experiments undertaken at the CRC. There are very few interaction artifacts that obscure the meaning of the overall system averages. The very fine-grained data resolution of 0.17 of a grade for overall system effects and 0.45 of a grade for the interaction of audio materials and systems are among the smallest magnitudes seen in any CRC quality assessment.

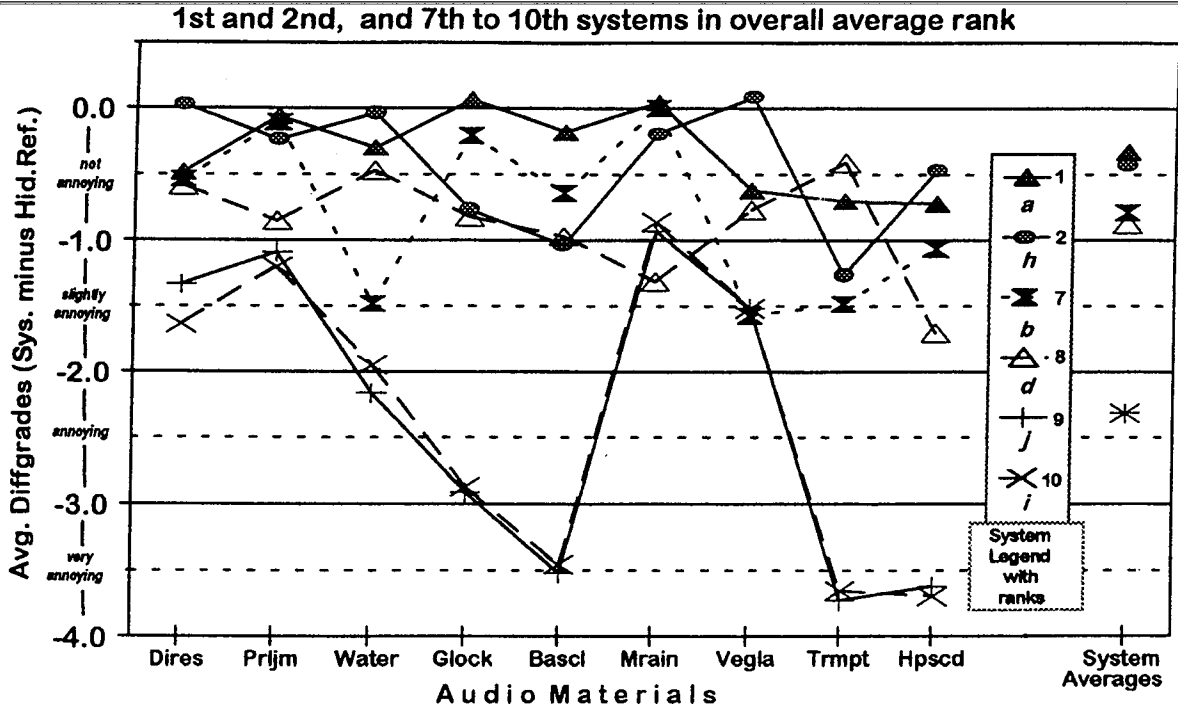


Fig. 2.2a. Quality test results - systems *a*, *b*, *d*, *h*, *i* and *j*

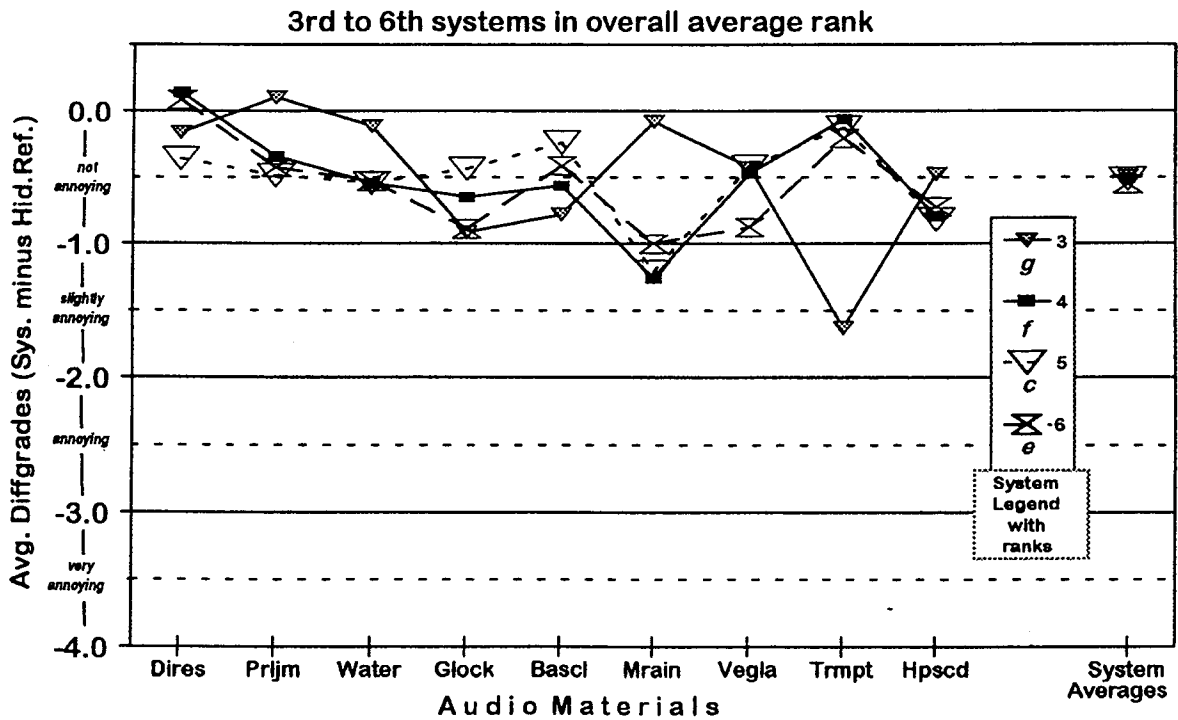


Fig. 2.2b. Quality test results - systems *c*, *e*, *f*, and *g*

**Table 2.4a** Average Difference Grades for each of the 9 Audio Materials (columns) by each of the 10 Systems

The data for a single system are shown throughout each row. Systems are arranged in descending average merit, with the highest rated system at the top, as shown in the far right-hand column (see Fig. 2.2a and 2.2b).

System	Dires	Prljm	Water	Glock	Bascl	Mrain	Vegla	Trmpt	Hpscd	Overall Averages
a	-0.49	-0.06	-0.30	0.07	-0.18	0.04	-0.62	-0.70	-0.72	a -0.33
h	0.02	-0.24	-0.04	-0.77	-1.04	-0.20	0.08	-1.27	-0.47	h -0.43
g	-0.16	0.10	-0.11	-0.92	-0.78	-0.08	-0.43	-1.63	-0.48	g -0.50
f	0.14	-0.34	-0.55	-0.65	-0.57	-1.26	-0.47	-0.06	-0.80	f -0.51
c	-0.36	-0.49	-0.54	-0.44	-0.24	-1.21	-0.42	-0.12	-0.82	c -0.52
e	0.09	-0.43	-0.53	-0.89	-0.41	-1.00	-0.88	-0.20	-0.72	e -0.55
b	-0.54	-0.10	-1.49	-0.21	-0.64	0.00	-1.58	-1.49	-1.07	b -0.79
d	-0.59	-0.85	-0.47	-0.82	-0.97	-1.31	-0.77	-0.41	-1.70	d -0.88
j	-1.34	-1.09	-2.16	-2.91	-3.52	-0.93	-1.51	-3.73	-3.62	j -2.31
i	-1.64	-1.20	-1.95 <sup>1</sup>	-2.87	-3.46	-0.86	-1.52	-3.66	-3.70	i -2.32
Audio Material Averages	-0.49	-0.47	-0.81	-1.04	-1.18	-0.68	-0.81	-1.33	-1.41	-0.91
Averages Without i and j	-0.24	-0.30	-0.50	-0.58	-0.60	-0.63	-0.64	-0.74	-0.85	-0.56

**Table 2.4b** Average Difference Grades for each of the 9 Audio Materials (columns) by each of the 10 Systems

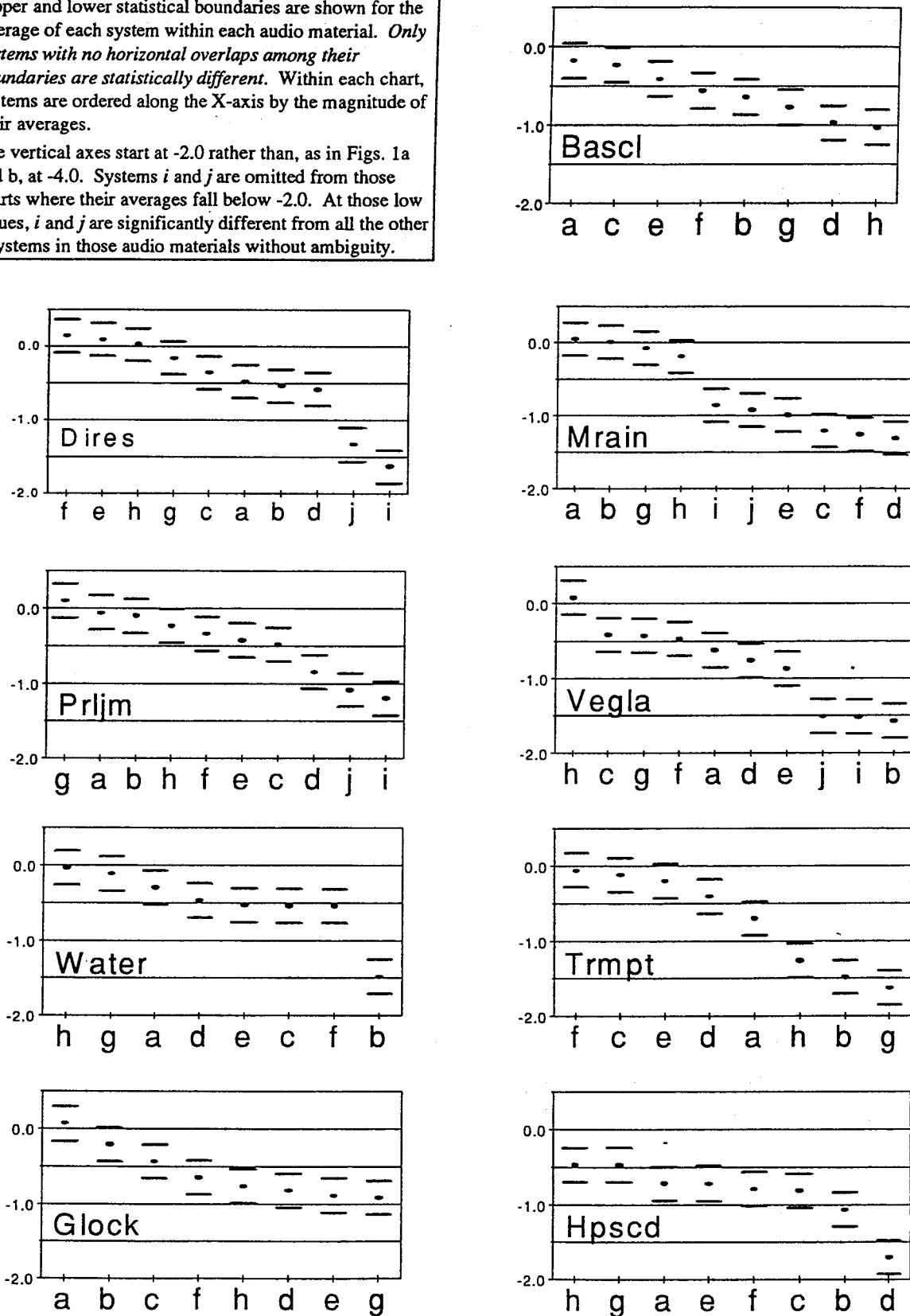
The data in the cells under each audio material are arranged in descending order of system merit, with the highest ranking system for that material at the top (see Fig. 2.2a and 2.2b).

	Dires	Prljm	Water	Glock	Bascl	Mrain	Vegla	Trmpt	Hpscd
f	0.14	g 0.10	h -0.04	a 0.07	a -0.18	a 0.04	h 0.08	f -0.06	h -0.47
e	0.09	a -0.06	g -0.11	b -0.21	c -0.24	b 0.00	c -0.42	c -0.12	g -0.48
h	0.02	b -0.10	a -0.30	c -0.44	e -0.41	g -0.08	g -0.43	e -0.20	a -0.72
g	-0.16	h -0.24	d -0.47	f -0.65	f -0.57	h -0.20	f -0.47	d -0.41	e -0.72
c	-0.36	f -0.34	e -0.53	h -0.77	b -0.64	i -0.86	a -0.62	a -0.70	f -0.80
a	-0.49	e -0.43	c -0.54	d -0.82	g -0.78	j -0.93	d -0.77	h -1.27	c -0.82
b	-0.54	c -0.49	f -0.55	e -0.89	d -0.97	e -1.00	e -0.88	b -1.49	b -1.07
d	-0.59	d -0.85	b -1.49	g -0.92	h -1.04	c -1.21	j -1.51	g -1.63	d -1.70
j	-1.34	j -1.09	i -1.95 <sup>1</sup>	i -2.87	i -3.46	f -1.26	i -1.52	i -3.66	j -3.62
i	-1.64	i -1.20	j -2.16	j -2.91	j -3.52	d -1.31	b -1.58	j -3.73	i -3.70

<sup>1</sup>System i received a grade of -1.95 for Water. In view of the statistical error, and also because i was not significantly different from the next higher rated system b on this audio material, i was omitted from Water in Fig. 2.3 on the next page, along with other instances of i and j in materials where either of these two systems obtained a diffgrade lower than -2.00. (No systems other than i and j received any diffgrades below -2.00.)

**FIG. 2.3 System Differences Within Audio Materials**  
 Upper and lower statistical boundaries are shown for the average of each system within each audio material. *Only systems with no horizontal overlaps among their boundaries are statistically different.* Within each chart, systems are ordered along the X-axis by the magnitude of their averages.

The vertical axes start at -2.0 rather than, as in Figs. 1a and b, at -4.0. Systems *i* and *j* are omitted from those charts where their averages fall below -2.0. At those low values, *i* and *j* are significantly different from all the other 8 systems in those audio materials without ambiguity.



**Table 2.5** Magnitudes of the differences between each pair of the 10 systems (a to j) for each of the 9 audio materials

This table parallels the presentation in Fig. 2. In both rows and columns, the systems are ordered by their average difference grades (see Fig. 2.2a and 2.2b). Significant differences (absolute magnitudes of 0.45 or greater) are shown in bold type

Dires-f	Dires-e	Dires-h	Dires-g	Dires-c	Dires-a	Dires-b	Dires-d	Dires-j	Dires-i
f 0.05	f 0.05	f 0.11	f 0.30	f 0.50	f 0.62	f 0.68	f 0.72	f 1.48	f 1.78
e 0.05	e 0.06	e 0.06	e 0.25	e 0.45	e 0.57	e 0.62	e 0.67	e 1.42	e 1.72
h 0.11	h 0.06	h 0.19	h 0.19	h 0.39	h 0.51	h 0.56	h 0.61	h 1.36	h 1.66
g 0.30	g 0.25	g 0.19	g 0.20	g 0.20	g 0.32	g 0.38	g 0.42	g 1.18	g 1.48
c 0.50	c 0.45	c 0.39	c 0.20	c 0.12	c 0.18	c 0.22	c 0.22	c 0.98	c 1.28
a 0.62	a 0.57	a 0.51	a 0.32	a 0.12	a 0.05	a 0.10	a 0.10	a 0.85	a 1.15
b 0.68	b 0.62	b 0.56	b 0.38	b 0.18	b 0.05	b 0.05	b 0.05	b 0.80	b 1.10
d 0.72	d 0.67	d 0.61	d 0.42	d 0.22	d 0.10	d 0.05	d 0.05	d 0.75	d 1.05
j 1.48	j 1.42	j 1.36	j 1.18	j 0.98	j 0.85	j 0.80	j 0.75	j 0.30	j 0.30
i 1.78	i 1.72	i 1.66	i 1.48	i 1.28	i 1.15	i 1.10	i 1.05	i 0.30	i 0.30
Prilm-g	Prilm-a	Prilm-b	Prilm-h	Prilm-f	Prilm-e	Prilm-c	Prilm-d	Prilm-j	Prilm-i
g 0.16	g 0.16	g 0.20	g 0.34	g 0.44	g 0.53	g 0.59	g 0.95	g 1.19	g 1.30
a 0.16	a 0.04	a 0.04	a 0.18	a 0.28	a 0.37	a 0.42	a 0.79	a 1.03	a 1.14
b 0.20	b 0.04	b 0.13	b 0.13	b 0.24	b 0.32	b 0.38	b 0.75	b 0.99	b 1.10
h 0.34	h 0.18	h 0.13	h 0.10	h 0.19	h 0.25	h 0.25	h 0.61	h 0.85	h 0.96
f 0.44	f 0.28	f 0.24	f 0.10	f 0.09	f 0.14	f 0.14	f 0.51	f 0.75	f 0.86
e 0.53	e 0.37	e 0.32	e 0.19	e 0.09	e 0.06	e 0.06	e 0.42	e 0.66	e 0.77
c 0.59	c 0.42	c 0.38	c 0.25	c 0.14	c 0.06	c 0.06	c 0.37	c 0.60	c 0.71
d 0.95	d 0.79	d 0.75	d 0.61	d 0.51	d 0.42	d 0.37	d 0.24	d 0.24	d 0.35
j 1.19	j 1.03	j 0.99	j 0.85	j 0.75	j 0.66	j 0.60	j 0.24	j 0.11	j 0.11
i 1.30	i 1.14	i 1.10	i 0.96	i 0.86	i 0.77	i 0.71	i 0.35	i 0.11	i 0.11
Water-h	Water-g	Water-a	Water-d	Water-e	Water-c	Water-f	Water-b	Water-i	Water-j
h 0.08	h 0.08	h 0.26	h 0.43	h 0.50	h 0.50	h 0.51	h 1.45	h 1.91	h 2.12
g 0.08	g 0.19	g 0.19	g 0.36	g 0.42	g 0.43	g 0.43	g 1.37	g 1.83	g 2.05
a 0.26	a 0.19	a 0.17	a 0.17	a 0.23	a 0.24	a 0.25	a 1.19	a 1.65	a 1.86
d 0.43	d 0.36	d 0.17	d 0.06	d 0.06	d 0.07	d 0.08	d 1.01	d 1.48	d 1.69
e 0.50	e 0.42	e 0.23	e 0.06	e 0.01	e 0.01	e 0.01	e 0.95	e 1.41	e 1.63
c 0.50	c 0.43	c 0.24	c 0.07	c 0.01	c 0.00	c 0.00	c 0.94	c 1.40	c 1.62
f 0.51	f 0.43	f 0.25	f 0.08	f 0.01	f 0.00	f 0.00	f 0.94	f 1.40	f 1.61
b 1.45	b 1.37	b 1.19	b 1.01	b 0.95	b 0.94	b 0.94	b 0.46	b 0.46	b 0.68
i 1.91	i 1.83	i 1.65	i 1.48	i 1.41	i 1.40	i 1.40	i 0.46	i 0.21	i 0.21
j 2.12	j 2.05	j 1.86	j 1.69	j 1.63	j 1.62	j 1.61	j 0.68	j 0.21	j 0.21
Glock-a	Glock-b	Glock-c	Glock-f	Glock-h	Glock-d	Glock-e	Glock-g	Glock-i	Glock-j
a 0.28	a 0.28	a 0.50	a 0.72	a 0.83	a 0.89	a 0.96	a 0.99	a 2.94	a 2.98
b 0.28	b 0.23	b 0.23	b 0.44	b 0.56	b 0.61	b 0.68	b 0.71	b 2.66	b 2.70
c 0.50	c 0.23	c 0.21	c 0.21	c 0.33	c 0.39	c 0.45	c 0.48	c 2.43	c 2.47
f 0.72	f 0.44	f 0.21	f 0.11	f 0.11	f 0.17	f 0.24	f 0.27	f 2.22	f 2.26
h 0.83	h 0.56	h 0.33	h 0.11	h 0.06	h 0.06	h 0.12	h 0.15	h 2.10	h 2.14
d 0.89	d 0.61	d 0.39	d 0.17	d 0.06	d 0.07	d 0.07	d 0.10	d 2.05	d 2.09
e 0.96	e 0.68	e 0.45	e 0.24	e 0.12	e 0.07	e 0.03	e 0.03	e 1.98	e 2.02
g 0.99	g 0.71	g 0.48	g 0.27	g 0.15	g 0.10	g 0.03	g 0.03	g 1.95	g 1.99
i 2.94	i 2.66	i 2.43	i 2.22	i 2.10	i 2.05	i 1.98	i 1.95	i 0.04	i 0.04
j 2.98	j 2.70	j 2.47	j 2.26	j 2.14	j 2.09	j 2.02	j 1.99	j 0.04	j 0.04

Bascl-a	Bascl-c	Bascl-e	Bascl-f	Bascl-b	Bascl-g	Bascl-d	Bascl-h	Bascl-i	Bascl-i
a	a	a	a	a	a	a	a	a	a
0.06	0.06	0.23	0.39	<b>0.46</b>	<b>0.60</b>	<b>0.79</b>	<b>0.86</b>	<b>3.28</b>	<b>3.28</b>
c	c	c	c	c	c	c	c	c	c
0.23	0.18	0.18	0.33	0.40	<b>0.54</b>	<b>0.73</b>	<b>0.80</b>	<b>3.22</b>	<b>3.22</b>
e	e	e	e	e	e	e	e	e	e
0.39	0.33	0.15	0.15	0.23	0.37	<b>0.56</b>	<b>0.62</b>	<b>3.04</b>	<b>3.11</b>
f	f	f	f	f	f	f	f	f	f
<b>0.46</b>	<b>0.40</b>	0.23	0.08	0.08	0.21	0.40	<b>0.47</b>	<b>2.89</b>	<b>2.96</b>
g	g	g	g	g	g	g	g	g	g
<b>0.60</b>	<b>0.54</b>	0.37	0.21	0.14	0.14	0.33	0.40	<b>2.81</b>	<b>2.88</b>
d	d	d	d	d	d	d	d	d	d
<b>0.79</b>	<b>0.73</b>	<b>0.56</b>	0.40	0.33	0.19	0.19	0.26	<b>2.68</b>	<b>2.74</b>
h	h	h	h	h	h	h	h	h	h
<b>0.86</b>	<b>0.80</b>	<b>0.62</b>	0.47	0.40	0.26	0.07	0.07	<b>2.49</b>	<b>2.55</b>
i	i	i	i	i	i	i	i	i	i
<b>3.28</b>	<b>3.22</b>	<b>3.04</b>	<b>2.89</b>	<b>2.81</b>	<b>2.68</b>	<b>2.49</b>	<b>2.42</b>	<b>2.42</b>	<b>2.49</b>
i	i	i	i	i	i	i	i	i	i
<b>3.34</b>	<b>3.29</b>	<b>3.11</b>	<b>2.96</b>	<b>2.88</b>	<b>2.74</b>	<b>2.55</b>	<b>2.49</b>	0.07	0.07
Mrain-a	Mrain-b	Mrain-g	Mrain-h	Mrain-i	Mrain-i	Mrain-e	Mrain-c	Mrain-f	Mrain-d
a	a	a	a	a	a	a	a	a	a
0.04	0.04	0.12	0.24	<b>0.90</b>	<b>0.97</b>	<b>1.04</b>	<b>1.26</b>	<b>1.30</b>	<b>1.36</b>
b	b	b	b	b	b	b	b	b	b
0.12	0.08	0.08	0.20	<b>0.86</b>	<b>0.93</b>	<b>1.00</b>	<b>1.21</b>	<b>1.26</b>	<b>1.31</b>
g	g	g	g	g	g	g	g	g	g
0.24	0.20	0.11	0.11	<b>0.78</b>	<b>0.85</b>	<b>0.92</b>	<b>1.13</b>	<b>1.18</b>	<b>1.23</b>
i	i	i	i	i	i	i	i	i	i
<b>0.90</b>	<b>0.86</b>	<b>0.78</b>	<b>0.67</b>	<b>0.67</b>	0.07	0.14	0.35	0.40	<b>0.45</b>
i	i	i	i	i	i	i	i	i	i
<b>0.97</b>	<b>0.93</b>	<b>0.85</b>	<b>0.73</b>	0.07	0.07	0.07	0.29	0.33	0.39
e	e	e	e	e	e	e	e	e	e
<b>1.04</b>	<b>1.00</b>	<b>0.92</b>	<b>0.80</b>	0.14	0.07	0.07	0.21	0.26	0.31
c	c	c	c	c	c	c	c	c	c
<b>1.26</b>	<b>1.21</b>	<b>1.13</b>	<b>1.02</b>	0.35	0.29	0.21	0.05	0.05	0.10
f	f	f	f	f	f	f	f	f	f
<b>1.30</b>	<b>1.26</b>	<b>1.18</b>	<b>1.07</b>	0.40	0.33	0.26	0.05	0.05	0.05
d	d	d	d	d	d	d	d	d	d
<b>1.36</b>	<b>1.31</b>	<b>1.23</b>	<b>1.12</b>	<b>0.45</b>	0.39	0.31	0.10	0.05	0.05
Vegla-h	Vegla-c	Vegla-g	Vegla-f	Vegla-a	Vegla-d	Vegla-e	Vegla-i	Vegla-i	Vegla-b
h	h	h	h	h	h	h	h	h	h
<b>0.50</b>	<b>0.50</b>	<b>0.51</b>	<b>0.55</b>	<b>0.70</b>	<b>0.85</b>	<b>0.96</b>	<b>1.60</b>	<b>1.60</b>	<b>1.66</b>
c	c	c	c	c	c	c	c	c	c
<b>0.50</b>	0.01	0.01	0.05	0.20	0.35	<b>0.46</b>	<b>1.10</b>	<b>1.10</b>	<b>1.16</b>
g	g	g	g	g	g	g	g	g	g
<b>0.51</b>	0.01	0.04	0.04	0.20	0.34	<b>0.45</b>	<b>1.09</b>	<b>1.09</b>	<b>1.15</b>
f	f	f	f	f	f	f	f	f	f
<b>0.55</b>	0.05	0.04	0.15	0.15	0.30	0.40	<b>1.04</b>	<b>1.05</b>	<b>1.10</b>
a	a	a	a	a	a	a	a	a	a
<b>0.70</b>	0.20	0.20	0.15	0.14	0.14	0.25	<b>0.89</b>	<b>0.90</b>	<b>0.90</b>
d	d	d	d	d	d	d	d	d	d
<b>0.85</b>	0.35	0.34	0.30	0.14	0.11	0.11	<b>0.75</b>	<b>0.75</b>	<b>0.75</b>
e	e	e	e	e	e	e	e	e	e
<b>0.96</b>	<b>0.46</b>	<b>0.45</b>	0.40	0.25	0.11	0.11	<b>0.64</b>	<b>0.64</b>	<b>0.70</b>
i	i	i	i	i	i	i	i	i	i
<b>1.60</b>	<b>1.10</b>	<b>1.09</b>	<b>1.04</b>	<b>0.89</b>	<b>0.75</b>	<b>0.64</b>	0.00	0.00	0.06
i	i	i	i	i	i	i	i	i	i
<b>1.60</b>	<b>1.10</b>	<b>1.09</b>	<b>1.05</b>	<b>0.90</b>	<b>0.75</b>	<b>0.64</b>	0.00	0.00	0.06
b	b	b	b	b	b	b	b	b	b
<b>1.66</b>	<b>1.16</b>	<b>1.15</b>	<b>1.10</b>	<b>0.95</b>	<b>0.81</b>	<b>0.70</b>	0.06	0.06	0.06
Trmpt-f	Trmpt-c	Trmpt-e	Trmpt-d	Trmpt-a	Trmpt-h				

---

## 3. Subjective Tests of Transmission Impairments

### 3.1 Background

This chapter describes the results obtained in a series of subjective tests performed to assess the effects of transmission errors on the audio quality of the DAR systems. Processing of the audio materials through the DAR systems took place at the NASA LeRC in Cleveland, USA. Reference and processed versions of the test materials were delivered to the CRC on DAT audio tapes and were transferred digitally, via the AES/EBU interface, onto a custom disk-based playback system which was used for the subjective tests. Details on this playback system as well as other facilities used during the tests are described in chapter 6 of this report.

A total of seven different types of impairment were assessed. These are:

1. Gaussian noise
2. Co-channel interference
3. Urban slow multipath
4. Urban fast multipath
5. Rural fast multipath
6. Obstructed fast multipath
7. DAR lower 1st adjacent to DAR channel interference

As will be shown in the section discussing the results, a total of nine DAR systems, labeled from  $a$  to  $i$ , were tested for additive gaussian noise and co-channel interference. A subset of four systems were tested in the urban fast multipath, three in both the urban slow and the obstructed fast multipath and two systems in the rural fast multipath. Only one system was tested for DAR lower 1st adjacent to DAR channel interference. Details on the characteristics of the simulated multipath mobile channels can be found in [3].

As described in [2], two series of subjective tests were performed, namely the threshold of audibility test, to accurately determine the point at which transmission impairment are just perceptible, and the failure characteristic test to determine the way in which the subjective quality degrades with increasing levels of transmission impairment.

### 3.2 Test procedures

#### 3.2.1 Test materials

The three critical audio materials described in Table 3.1 below were used in all the impairment tests reported in this chapter. These materials were carefully selected by the staff at the NASA LeRC laboratory because they were particularly sensitive to revealing artifacts resulting from



transmission impairments. Only one of these materials, namely Glockenspiel, was also used in the audio quality subjective tests described in chapter 2.

Code	Description	Duration	Source
Glock	Glockenspiel	16 s	EBU SQAM CD (track 35/Index 1)
Clarn	Clarinet solo	20 s	EBU SQAM CD (track 16/Index 2)
Sopra	Soprano solo	31 s	EBU SQAM CD (track 62 /Index 1)

Table 3.1 List of audio materials used in impairments tests

### 3.2.2 Impairment levels

As outlined in [2], each of the three audio test materials was recorded, for each DAR system and each type of impairment, at the following levels of impairment: CC, TOA<sub>1</sub>, TOA<sub>2</sub>, ..., TOA<sub>N</sub>, S<sub>1</sub>, S<sub>2</sub>, ..., S<sub>M</sub>, POF where:

CC = coded audio in a Clear Channel

TOA<sub>1</sub>, TOA<sub>2</sub>, ..., TOA<sub>N</sub> = N stimuli (three or more) in the close neighborhood of the approximate TOA (threshold of audibility)

S<sub>1</sub>, S<sub>2</sub>, ..., S<sub>M</sub> = M intermediate levels of impairments

POF = point of failure

For each of the above impairment levels, the C<sub>0</sub>/N<sub>0</sub> (or D/U for co-channel interference) ratio was noted. The CC level was recorded at a high value of C<sub>0</sub>/N<sub>0</sub> (or D/U) so that transmission errors could be considered as negligible, hence the label "clear channel" given to that level. The increments in the C<sub>0</sub>/N<sub>0</sub> (or D/U) ratios at which the other impairment levels were recorded varied from 0.25 dB, for those DAR systems which failed abruptly, up to 1.0 dB for those systems which failed more gracefully, with increments of 0.5 dB being the typical value used for most systems and most of the impairments.

### 3.2.3 Subjective test procedures

As described in [2], two separate experiments were performed for each of the seven types of impairment considered:

#### Experiment 1: Threshold of Audibility

The purpose of this experiment was to provide a sensitive and reliable measurement of the threshold of audibility (TOA) for the various types of channel impairments. The threshold of audibility is defined as the highest C<sub>0</sub>/N<sub>0</sub> (or D/U) level where an audible artifact due to a transmission error can be detected. Determining TOA consisted of a two part process and was performed using CC (the coded audio in clear channel) as the reference signal. Simply stated, the

TOA values were determined by agreement between the in-house expert listeners at the NASA LeRC and the CRC.

a) Determining  $TOA_{Lab}$ :

The first step in this experiment was to determine  $TOA_{Lab}$ .  $TOA_{Lab}$  is the value for TOA determined initially by the LeRC team. To determine  $TOA_{Lab}$ , the LeRC team processed the critical source materials through a given system and listened to the output in real-time. The level of the impairment was adjusted in 0.25dB increments until the expert listeners determined  $TOA_{Lab}$  for that system/critical material/impairment combination. The determination of  $TOA_{Lab}$  was based on listening to a given critical material at least 3 times at each level of impairment or until the impairment was detected. This procedure was repeated for each combination of system, critical material, and impairment. During the process of determining  $TOA_{Lab}$ , the LeRC team also established  $POF_{Lab}$ .

Having determined  $TOA_{Lab}$  and  $POF_{Lab}$ , the LeRC team processed the source materials through the systems and recorded the output onto DAT tapes. Since these recordings were to be used in the failure characteristic tests, recordings were also made for several impairment levels on either side of  $TOA_{Lab}$  and  $POF_{Lab}$  as well as for intermediate levels between  $TOA_{Lab}$  and  $POF_{Lab}$  for those systems that failed gracefully. In particular, for the conditions near the  $TOA_{Lab}$  impairment level (i.e.  $TOA_{Lab} \pm 0.25\text{dB}$  and  $\pm 0.5\text{dB}$ ), additional recordings were made. This was done to account for the statistical nature of the occurrence of impairment conditions near TOA. The DAT tapes containing the recorded materials were then sent to the CRC team while the LeRC team kept a duplicate set.

b) Determining the True TOA:

Upon receiving the DAT tapes the CRC team transferred the audio recordings to the CRC's hard disk playback system. The CRC expert listeners then auditioned the recordings in order to evaluate the values of  $TOA_{Lab}$  established by the LeRC team. Three possible scenarios resulted from these listening tests and are discussed below.

***Scenario 1: The CRC team agrees with the LeRC team that  $TOA_{Lab}$  is in fact the true TOA.***

In this case, the two teams agreed that  $TOA_{Lab}$  should be designated as the true TOA. This was by far the most common scenario.

***Scenario 2: The CRC team cannot hear any audible artifact at  $TOA_{Lab}$ .***

Due to the statistical nature of the occurrence of impairments, it was quite possible that no audible artifacts existed (at  $TOA_{Lab}$ ) in some of the recordings sent to CRC. To minimize this possibility, the LeRC team recorded several samples at  $TOA_{Lab}$ . If the CRC team could not hear any audible artifacts in any of the samples (at  $TOA_{Lab}$ ) then they notified the LeRC team. The LeRC team then listened to their duplicate tapes to determine whether or not they could detect any artifacts in the

samples at  $TOA_{Lab}$ . If the LeRC team was not able to detect any artifacts at  $TOA_{Lab}$  then the next level of impairment (typically  $TOA_{Lab} + 0.5\text{dB}$ ) was designated as the true TOA.

If the LeRC team was able to identify audible artifacts at  $TOA_{Lab}$  then they provided the CRC team with a detailed description of each artifact and where it occurred. The CRC team then listened to the recordings to confirm the LeRC team's observations. Once the two teams agreed that an audible artifact existed,  $TOA_{Lab}$  was designated as the true TOA.

### ***Scenario 3: The CRC team can hear artifacts at impairment levels lower than $TOA_{Lab}$ .***

During the course of evaluating  $TOA_{Lab}$ , the CRC team also listened to recordings with less impairment than  $TOA_{Lab}$  (i.e.  $TOA_{Lab} - 0.25$  or  $-0.5\text{dB}$ ) to determine if any audible artifacts could be detected. Due to the statistical nature of the occurrence of the impairments, it was possible for audible artifacts to exist at these impairment levels in some of the recordings sent to CRC.

In those instances where the CRC team was able to detect audible artifacts at lower impairment levels, a detailed description of each artifact and where it occurred was sent to the LeRC team. The LeRC team then listened to their duplicate recordings to confirm the CRC team's evaluation. Once the two teams agreed that an audible artifact existed, that impairment level was designated as the true TOA.

### **Experiment 2: Failure Characteristic**

The purpose of the second experiment was to determine how the subjective quality of each DAR system degraded with increasing levels of impairment. The uncoded CD original was used as the reference against which the following stimuli were compared and rated: CC, TOA (as determined in experiment 1 above),  $S_1$ ,  $S_2$ , ...,  $S_M$ , POF.

The test procedure used was essentially the same as the one used in the audio quality tests (see section 2.2.2). The informed reader may skip the rest of this section up to the paragraph preceding Table 3.2. The procedure included a training phase followed by a blind rating phase. At least half a day (the morning of the first day) was devoted to the training session. Subjects worked together in groups of 2 or 3 and were allowed as much time as they required for training. Subjects were trained using the same hard disk playback system used in the blind rating tests. For each experiment, a subset of the stimuli to be rated in the blind rating test was used for the training session. This subset consisted of stimuli which were considered to be representative of those to be rated in the subsequent blind rating test.

For the blind rating phase, the triple-stimulus A-B-C presentation was used, where "A" was the known reference (unprocessed CD signals). One of "B" or "C" was the stimulus to be rated (i.e. one of the impairment levels) and the other one was the hidden reference (i.e. a perfect replica of "A"). The assignment of stimulus and hidden reference to "B" and "C" was not known to the subjects and was arranged to be unpredictable to the listeners from trial to trial. For each trial, subjects were asked to rate the *difference* between the known reference "A" and version "B" as

---

well as the difference between "A" and "C" using the same ITU-R continuous 5-grade impairment scale as in the audio quality tests and shown in Fig. 2.1 of chapter 2.

Listeners were instructed to use scores with a single decimal point. In effect, this is a 41 point scale. The subjects were instructed to treat this as a continuous scale but, to facilitate the subjects' orientation, category labels were associated with the scale. Thus, 1.0 to 1.9 is a "very annoying" range; 2.0 to 2.9 is "annoying"; 3.0 to 3.9 is "slightly annoying"; 4.0 to 4.9 is "perceptible but not annoying". Finally, 5.0 is "imperceptible". A constraint was that the alternative the subject has decided is the "hidden reference" must be graded 5.0. And so, *at least one* of the two grades on each trial must be a 5.0.

Thus two totally interdependent scores from the listener were recorded on each trial. This deliberate interdependence was handled by subtracting the score given to the true hidden reference from the score given the true impaired version (i.e., DAR System minus Reference). The reason for subtracting in this direction rather than the opposite one (Reference minus DAR System) is only so that in a graphical plot of outcomes, the data will fall in the same geometric quadrant as they would if the actual 1.0 to 5.0 scores used by the subjects were plotted.

This subtraction means, however, that the scores are transformed so that the 1.0 to 5.0 range of the original scale becomes, instead, -4.0 to 0.0 in the analysis and presentation of results. These difference grades or "diffgrades" represent the relative differences between the grades given to the hidden reference and the ones given to the DAR system under test. For practical purposes, it can be assumed that a direct analogy holds between the original 1.0 to 5.0 scale and the -4.0 to 0.0 diffgrades. Accordingly, to facilitate comparison with the 1.0 to 5.0 scale used by subjects, the "annoyance" category labels are shown between the Y-axis numbers in the various figures shown in this chapter.

Blind rating was performed individually by each listener. Listeners used the disk-based playback system described in chapter 5 of this report which allowed seamless switching between the stimuli to be compared. Listeners were able to take as much time as they needed on each trial, switching as often as they liked, until satisfied with the numerical ratings they were asked to assign. They were also free to use either the loudspeakers or headphones to make a judgment, whichever they felt was the most critical transducer on any trial.

The total number of materials to be rated was divided into sessions of 10-15 trials each and, within each session, items were ordered in a way which was unpredictable by the listeners. Moreover, the order of presentation of the individual sessions was changed from listener to listener over the course of the test to ensure that time-correlated factors (such as fatigue) would not differentially affect any level of any of the factors under test.

The seven impairments tested were grouped into five separate experiments as shown in Table 3.2. The table also shows the number of DAR systems tested and the number of trials for each impairment. The number of days spent by each listener for each experiment is also shown in this table. The listening panel included a total of 6 expert-listeners for each of the experiments. Two

subjects took part in all of the experiments. The remaining four panel members were in general different from experiment to experiment.

Experiment	Impairment	Number of systems tested	Number of trials	Number of days per listener
1	AWGN	9	150	3
2	Co-channel interference	9	150	3
3	Urban Fast Multipath	4	72	2
4	Urban Slow Multipath	3	54	2
	Rural Fast Multipath	2	36	
5	Obstructed Fast Multipath	3	54	2
	DAR lower 1st Adj. to DAR channel interference	1	21	

Table 3.2 Grouping of the seven impairments into five experiments

### 3.3 Test results

#### 3.3.1 Data presentation and statistical analysis

The impairment test results are presented in three different forms. First, a table shows, for each impairment, each DAR system and each test material, the values of  $C_o/N_o$  (or D/U) at TOA, POF and the difference between the two values, labelled in this chapter as the *Failure Margin (FM)*. The failure margin gives an indication of how quickly the subjective quality degrades with a reduction in  $C_o/N_o$  (or D/U). **The definition of POF used in this report is any score which falls in the “very annoying” range of the rating scale, that is below -3.0.** When more than one impairment level have been rated below -3.0, the level with the **highest  $C_o/N_o$  (or D/U)** has been taken as the POF. In the results presented throughout this chapter, **the term TOA will be used to mean the true TOA level** as determined in the TOA test. Two different types of graphical representations of the results are provided. The first shows, for each impairment and each test material, the average diffgrade given to each presentation plotted against the corresponding  $C_o/N_o$  (or D/U) ratio. The resulting graph is called the failure characteristic curve. The second type of graph shows the failure characteristic curve of all three test materials for a single system and a given impairment.

In impairment studies, expert observation and commentary (EO&C) methodology is often used. In compliance with [2], a grading scale was used (the same one that was used in the audio quality assessments) in the present tests because this approach has advantages over EO&C. One advantage is that each person provides their ratings in individual blind sessions rather than in a group situation. Thus, the judgments of each individual are less influenced by those of the other judges. Also, the audio item presentation sequences are controlled by the experimenter and are varied from subject to subject to ensure that events correlated with time (fatigue, learning, etc.)

~~do not distort the average grades associated with each experimental condition. By contrast, with EO&C, all listeners hear the materials at the same time in a group situation. And so the effects of time-correlated events cannot be factored out by experimental design.~~

With EO&C, you cannot have a statistical analysis of results, since that methodology provides only a group consensus rather than sets of independent judgments, one set from each listener. However, even though we used a grading scale here, we cannot present a statistical analysis of results such as is usually associated with this methodology. The major reason for this is that, as required by the nature of the tests, gross impairments (e.g. POF) were included among the sets of stimuli, mixed in with small (e.g. clear channel) and intermediate impairments. The nature of the grading scale that was used, however, did not permit comparable discriminations along the whole range of impairments that were presented. The scale problem was at the "very annoying" end of the grading scale. Thus, the scale suffered a "floor effect" where listeners could not assign a grade worse than the minimal one allowed of 1.0, even though there might have been considerable subjective variation among those stimuli that were judged to be grossly impaired. Thus, the scale permitted rather uniform variance to occur as at the top and middle areas of the impairment scale used, but not so among the gross impairments at the bottom of the scale.

The resulting set of numerical subjective judgments, then, are "distorted". And this kind of distortion is one which precludes performing a meaningful statistical analysis since this requires that similar variance among judgments is obtained across the entire range of the scale that is used by subjects. In addition, the number of subjects was small (6 per impairment experiment) hence statistical variance was large.

Out of curiosity, ANOVA's were actually performed for each experiment, and it was found that, if these analyses had been legitimate, then the critical difference between two scores required for statistical significance was found to be half a grade, for the DAR lower 1st adjacent to DAR channel interference experiment and between 0.8 to a full grade for all the other experiments. However, because the violations of essential statistical requirements discussed above were serious ones, we cannot place any confidence in these outcomes and do not present or quote these analyses. All in all, then, we feel the results are quantitatively superior to those that would have been obtained using EO&C, but because of the inherent nature of the range of impairments, they are not, and could not be, fully equivalent to those in the quality tests.

All the conditions in the system audio quality tests which took place previous to the series of impairment tests discussed here, were "clear-channel" conditions. It might be suggested that the clear-channel baselines established in each test in the impairment series should show similar grades to those obtained by the comparable systems in the previous quality tests. However, in the quality tests, there were 9 audio materials, while there were only three in the impairment tests, and only one material, namely Glockenspiel, was in common between the two studies. Furthermore, as discussed above, the total range of impairments was quite different between the quality and the impairment tests, and introduced an unavoidable distortion in the grading scale for the impairment tests. This means that the total context of the impairment experiments was substantially different from that of the quality tests. And also, as mentioned previously, the number of subjects was

considerably smaller in the impairment experiments. Moreover, the subjects in the quality tests were in general different than those who took part in the various impairment tests.

For all of these reasons, no comparisons can legitimately be made between the quality experiments on the one hand, and the clear-channel conditions in the impairment tests. The same statement can be made for comparisons between the clear-channel conditions of the seven impairment tests.

### 3.3.2 Gaussian noise

A total of nine DAR systems were tested with the gaussian noise impairment. The complete failure characteristic curves for systems *a* through *h* are shown in Fig. 3.1, 3.2 and 3.3 for *Glockenspiel*, *Soprano* and *Clarinet* respectively. Each curve shows the mean opinion score (average diffgrade across all 6 listeners) given to each impairment level vs the  $C/N_0$  value at each of these levels. The failure characteristic curves for system *i* are given in Fig. 3.12. The score given to the coded audio in a clear channel is plotted, for each system, as a separate point identified on the X-axis by the *Clear* label. For visual clarity, the *Clear* point has not been linked to the rest of the curve. The first point to the right of the *Clear* point on each curve is the *TOA* level as determined in the *TOA* test.

It can be seen in Fig. 3.1, 3.2 and 3.3 that a few curves are not completely monotonic, that is the scores do not always decrease with a reduction of the  $C/N_0$  value. These non-monotonic fluctuations are caused by the random nature of the gaussian noise impairment. Statistically, it is possible for a momentary high level noise transient to occur at a given impairment level and generate more audible degradation than the next higher impairment level which would be free of such transients. In some instances, it is found that the score at *TOA* is worse than the next higher level of impairment. This may be explained by the fact that, for *TOA*, more time was probably taken to detect and record an impaired audio segments that contained a noise transient than for the next higher level of impairment. It is also well known to the authors of this report that listeners base their judgements on the worst impairment they heard over the entire duration of a given test material. Therefore, a single isolated and momentary artifact ("click", mute, etc...) may result in a rather poor or marginal score.

These observations also apply to the failure characteristic curves shown in section 3.3.3 for co-channel interference. The co-channel interfering signal is generally independent of the desired signal and, like the gaussian noise, exhibits random level fluctuations. In the case of multipath impairments (section 3.3.4 to 3.3.7), the impairments are the results of the combination of two random processes: the mobile channel and the gaussian noise. The mobile channel goes through a series of fades in an unpredictable and random fashion. It is consequently possible for a particularly bad fade pattern to happen at a given impairment (i.e. noise) level and for a less severe one to occur during the recording of the next higher impairment level. It must be remembered that injected noise was the parameter varied between each impairment level in the various multipath tests. The interaction of the random fades with the random noise is probably the reason why the failure characteristic curves for the multipath impairments exhibit, in general, more non-monotonicity than the curves obtained in gaussian noise alone and co-channel interference.

As shown in Table 3.3, the  $C_o/N_o$  values at TOA cover a range of approximately 23 dB from system *f* to system *g* for all three test materials. In this table, the systems have been arranged in ascending order of their  $C_o/N_o$  value at TOA. System *f* is at one extreme of the range with a  $C_o/N_o$  value at TOA of 3.26 dB for all three test materials. At the other extreme lies system *g* which required a  $C_o/N_o$  value at TOA of around 25-26 dB for the three materials tested. Systems *a* and *b* are practically identical while the values for systems *e* and *h* are very close to each other, being separated by 1 dB or less for all three test materials.

DAR System	Glockenspiel			Soprano			Clarinet		
	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)
<i>f</i>	3.26	2.76	0.5	3.26	2.76	0.5	3.26	3.01	0.25
<i>b</i>	8.46	7.46	1.0	8.71	6.71	2.0	8.96	6.46	2.5
<i>a</i>	8.48	6.48	2.0	8.23	6.73	1.5	8.98	7.98	1.0
<i>h</i>	10.51	9.01	1.5	9.51	8.51	1.0	10.01	9.51	0.5
<i>e</i>	10.76	9.76	1.0	10.51	9.51	1.0	10.76	9.76	1.0
<i>c</i>	11.36	10.86	0.5	11.11	10.61	0.5	11.11	10.61	0.5
<i>d</i>	18.85	17.85	1.0	18.1	16.6	1.5	18.6	17.1	1.5
<i>g</i>	25.1	23.1	2.0	25.1	22.1	3.0	26.35	22.85	3.5
<i>i</i>	16.82			16.82			17.07		

Table 3.3  $C_o/N_o$  ratios at TOA and POF and failure margin for the gaussian noise impairment

Systems *f* and *c* showed an average (across all three test materials) failure margin around 0.5 dB while this value was 1.0 dB for both systems *e* and *h*. Systems *d* and *a* are next with 1.33 dB and 1.5 dB respectively. Systems *b* and *g* follow with average failure margin values of 1.8 and 2.8 dB respectively.

The failure characteristic curves for each individual system and all three test materials are plotted in Fig. 3.4 to 3.12 for systems *a* to *i* respectively. In general, the results obtained with the three test materials are within 1 dB or so of each other, with *Soprano* being the (slightly) less critical of the three materials.

### 3.3.3 Co-channel interference

As for the gaussian noise, a total of nine DAR systems were tested with the co-channel interference. The co-channel failure characteristic curves for systems *a* through *h* are shown in Fig. 3.13, 3.14 and 3.15 respectively for *Glockenspiel*, *Soprano* and *Clarinet*. The failure characteristic curves for system *i* are given in Fig. 3.24.



As shown in Table 3.4, the D/U values at TOA cover a range of approximately 36 to 38 dB from system *f* to system *g* for all three test materials. In this table, the systems have been arranged in ascending order of their D/U values at TOA. System *f* is at one extreme of the range with a D/U value at TOA of 5.5 dB for all three test materials. At the other extreme lies system *g* which required a D/U value at TOA of around 42-43 dB for the three materials tested. As for the gaussian noise impairment, systems *a* and *b* are practically identical while the D/U values at TOA for systems *e* and *h* are very close to each other, being separated by 1 dB or less for all three test materials. Systems *e* and *c* are also separated, in general, by less than 1 dB from TOA to POF.

DAR System	Glockenspiel			Soprano			Clarinet		
	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)
<i>f</i>	5.5	4.5	1.0	5.25	4.75	0.5	5.5	5.0	0.5
<i>b</i>	8.5	6.0	2.5	8.25	6.25	2.0	8.5	7.0	1.5
<i>a</i>	8.6	5.85	2.75	7.85	5.85	2.0	8.35	6.85	1.5
<i>h</i>	11.37	8.87	2.5	10.87	7.87	3.0	10.87	9.87	1.0
<i>e</i>	11.12	9.62	1.5	10.87	9.87	1.0	10.87	9.87	1.0
<i>c</i>	11.64	10.64	1.0	11.14	10.64	0.5	11.64	10.64	1.0
<i>d</i>	17.4	15.4	2.0	17.15	15.65	1.5	17.4	15.9	1.5
<i>g</i>	42.6	39.6	3.0	41.6	39.1	2.5	43.1	40.1	3.0
<i>i</i>	26.98	23.23	3.75	26.23	24.23	2.0	26.73	25.73	1.0

Table 3.4 D/U ratios at TOA and POF and failure margin for the co-channel interference impairment

Systems *f* and *c* showed an average (across all three test materials) failure margin of around 0.7-0.8 dB while this value was 1.2 dB for system *e* and 1.7 for system *d*. Systems *b*, *a* and *h* are next with an average failure margin of around 2.0 dB. Finally, systems *g* has the slowest degradation rate with an average margin of 2.8 dB.

The failure characteristic curves for each individual system and all three test materials are plotted in Fig. 3.16 to 3.24 for systems *a* to *i* respectively. In general, the results for the three test materials were within 1 dB or so of each other, with *Clarinet* being the (slightly) more critical of the three materials if one considers the overall failure curves. By considering the TOA point only (Table 3.4), all three materials are within 0.5 dB of each other for most of the cases.

### 3.3.4 Urban slow multipath

Only systems *a*, *b* and *c* were tested subjectively with the urban slow multipath impairment. The failure characteristic curves of these systems are shown in Fig. 3.25, 3.26 and 3.27 respectively for *Glockenspiel*, *Soprano* and *Clarinet*. These figures (and Table 3.5) show that the performance of systems *a* and *b* are very close to each other, with system *b* requiring a  $C_0/N_0$  ratio 0.5 to 1 dB higher than that of system *a* at TOA across the three test materials. The  $C_0/N_0$  ratios at TOA for system *c* are around 32 dB for all three test materials, roughly 10 dB higher than for systems *a* and *b*. The failure margins for systems *a* and *b* are between 6 and 7 dB (except for system *b* on *Clarinet* where a failure margin of 5 dB was obtained (by applying the strict definition of POF of section 3.3.1) while those of system *c* are between 7 and 8 dB.

Examination of Fig. 3.28, 3.29 and 3.30, which show respectively the failure characteristics of systems *a*, *b* and *c* for all three test materials, reveals that, overall, *Clarinet* is slightly more critical than *Glockenspiel* and *Soprano*. This is particularly true for systems *a* and *b*.

DAR System	Glockenspiel			Soprano			Clarinet		
	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)
<i>a</i>	21.79	15.04	6.75	21.29	15.29	6.0	22.04	15.04	7.0
<i>b</i>	22.74	16.24	6.5	21.74	15.74	6.0	22.74	17.74	5.0
<i>c</i>	32.14	24.14 <sup>1</sup>	8.0	31.64	24.64	7.0	32.14	24.64	8.0

Note 1: This value has been extrapolated from the lowest measured score of 2.9 at a  $C_0/N_0$  value of 24.64 dB (see Figures 3.25, 3.30). It has been assumed that if the  $C_0/N_0$  ratio had been reduced to 24.14 (an additional increment of 0.5 dB), the subjective quality would have been rated below -3.0.

Table 3.5  $C_0/N_0$  ratios at TOA and POF and failure margin for the urban slow multipath impairment

### 3.3.5 Urban fast multipath

Only the four systems, *a*, *b*, *c*, and *e*, were tested subjectively under the urban fast multipath impairment condition. The failure characteristic curves for these systems are given in Fig. 3.31, 3.32, and 3.33 respectively for *Glockenspiel*, *Soprano* and *Clarinet*. Table 3.6 provides the  $C_0/N_0$  ratios at TOA and POF, as well as the failure margins (FM) for the four systems with the three test materials.

The failure characteristic curves for each system with all three test materials are plotted in Fig. 3.34, 3.35, 3.36, and 3.37 for systems *a*, *b*, *c*, and *e* respectively.

DAR System	Glockenspiel			Soprano			Clarinet		
	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)
<i>a</i>	17.79	14.79	3.0	17.54	14.54	3.0	19.29	14.79	4.5
<i>b</i>	20.24	15.49	4.75	20.99	12.79	8.2	19.24	16.99	2.25
<i>e</i>	26.12	22.62	3.5	24.87	19.12	5.75	25.62	21.12	4.5
<i>c</i>	26.64	22.64	4.0	27.14	24.64	2.5	26.64	20.64	6.0

Table 3.6  $C_0/N_0$  ratios at TOA and POF and failure margin for the urban fast multipath impairment

### 3.3.6 Rural fast multipath

Only the two systems, *e* and *f*, were tested subjectively under the rural fast multipath impairment condition. The failure characteristic curves for these systems are given in Fig. 3.38, 3.39, and 3.40 respectively for *Glockenspiel*, *Soprano* and *Clarinet*. Table 3.7 provides the  $C_0/N_0$  ratios at TOA and POF, as well as the failure margins (FM) for the two systems with the three test materials.

The failure characteristic curves for systems *e* and *f* for all three test materials are plotted in Fig. 3.41 and 3.42 respectively.

DAR System	Glockenspiel			Soprano			Clarinet		
	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)
<i>f</i>	5.96	2.46	3.5	6.46	2.46	4.0	6.96	2.46	4.5
<i>e</i>	31.62	25.12	6.5	30.62	24.62	6.0	31.62	24.62	7.0

Table 3.7  $C_0/N_0$  ratios at TOA and POF and failure margin for the rural fast multipath impairment

### 3.3.7 Obstructed fast multipath

Three systems, *a*, *b*, and *c*, were tested subjectively under the obstructed fast multipath impairment condition. System *c* was only tested with the *Clarinet* test material. The failure

characteristic curves for these systems are given in Fig. 3.43, 3.44, and 3.45 respectively for *Glockenspiel*, *Soprano* and *Clarinet*. Table 3.8 provides the  $C_o/N_o$  ratios at TOA and POF, as well as the failure margins (FM) for the three systems with the three test materials. Note that the entries under *Glockenspiel* and *Soprano* are intentionally left blank for system *c* since this system was not tested subjectively with these materials.

DAR System	Glockenspiel			Soprano			Clarinet		
	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)
<i>a</i>	19.54	15.54	4.0	19.04	15.54	3.5	19.54	15.54	4.0
<i>b</i>	20.49	15.99	4.5	20.49	15.99	4.5	20.49	16.49	4.0
<i>c</i>							25.14	22.64	2.5

Table 3.8  $C_o/N_o$  ratios at TOA and POF and failure margin for the obstructed fast multipath impairment

The failure characteristic curves for each system with all three test materials are plotted in Fig. 3.46, 3.47, and 3.48 for systems *a*, *b*, and *c* respectively. Again, note that for system *c* (i.e. Fig. 3.48) a failure characteristic curve is only available for the *Clarinet* test material.

### 3.3.8 DAR lower 1st adjacent to DAR channel interference

Only system *i* was tested subjectively under the DAR lower 1st adjacent to DAR channel interference condition. The failure characteristic curves for this system are given in Fig. 3.49. All three source materials, *Glockenspiel*, *Soprano* and *Clarinet*, are plotted together in this one figure. The  $C_o/N_o$  ratios at TOA and POF, as well as the failure margins (FM) are given below in Table 3.9 for the three test materials.

DAR System	Glockenspiel			Soprano			Clarinet		
	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)
<i>i</i>	31.63	30.13	1.5	31.13	30.63	0.5	31.13	31.13	0.0

Table 3.9 D/U values at TOA and POF and failure margin for the DAR lower 1st adjacent to DAR channel interference

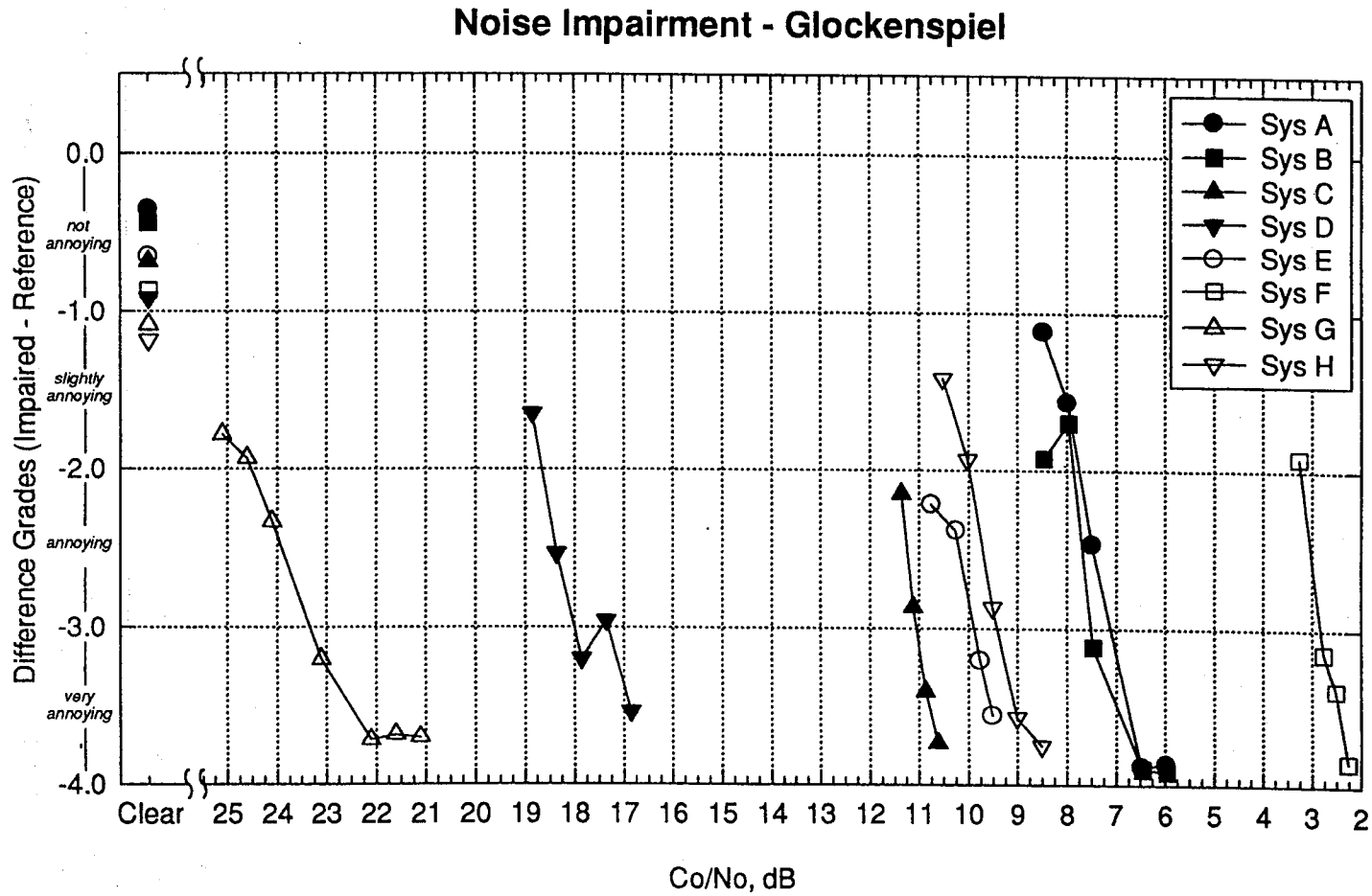


Figure 3.1. Failure characteristic curves of DAR systems a through h tested with the gaussian noise impairment and the Glockenspiel test material. (DAR laboratory was unable to certify digital power levels for system g)

### Noise Impairment - Soprano

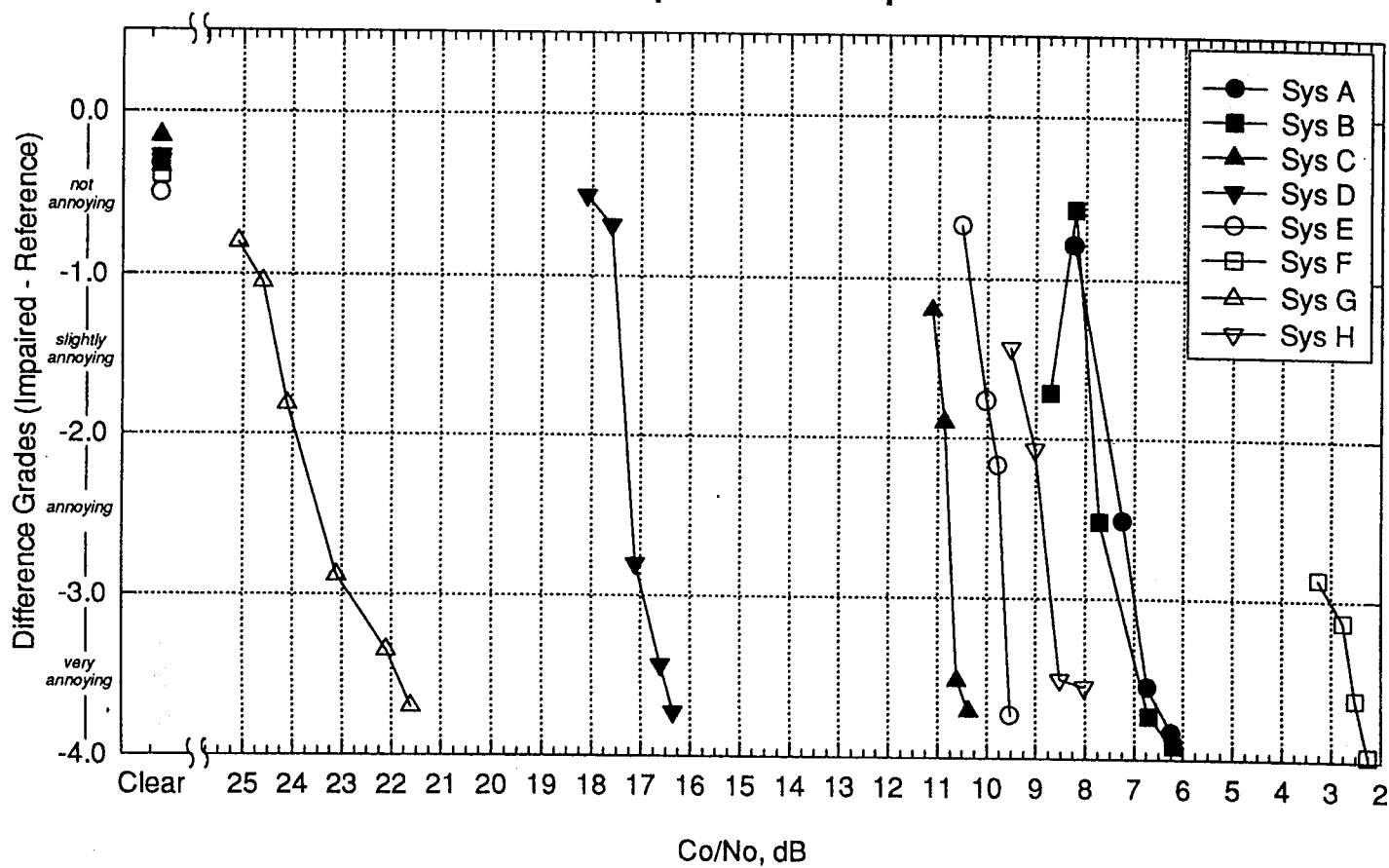


Figure 3.2. Failure characteristic curves of DAR systems a through h tested with the gaussian noise impairment and the Soprano test material. (DAR laboratory was unable to certify digital power levels for system g)

### Noise Impairment - Clarinet

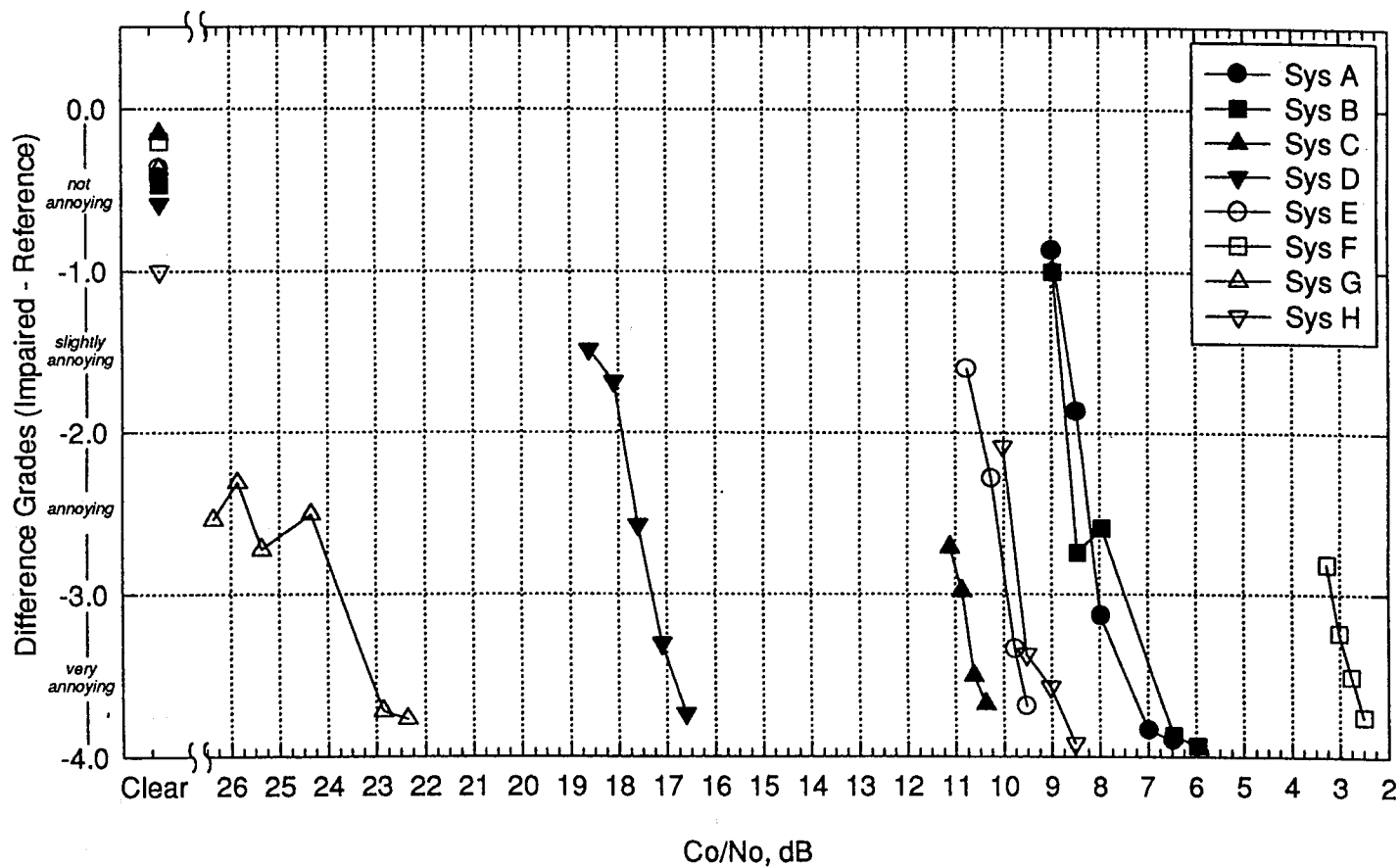


Figure 3.3. Failure characteristic curves of DAR systems a through h tested with the gaussian noise impairment and the Clarinet test material. (DAR laboratory was unable to certify digital power levels for system g)

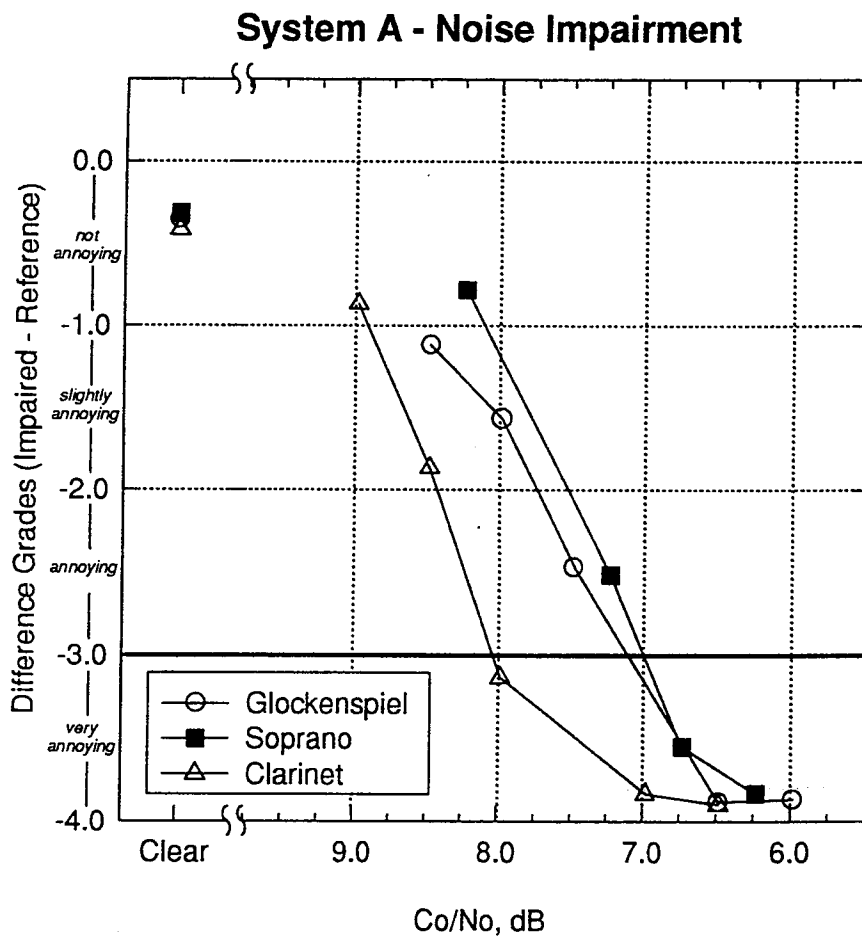


Figure 3.4. Failure characteristic curves for system a for all three test materials with the gaussian noise impairment.



### System B - Noise Impairment

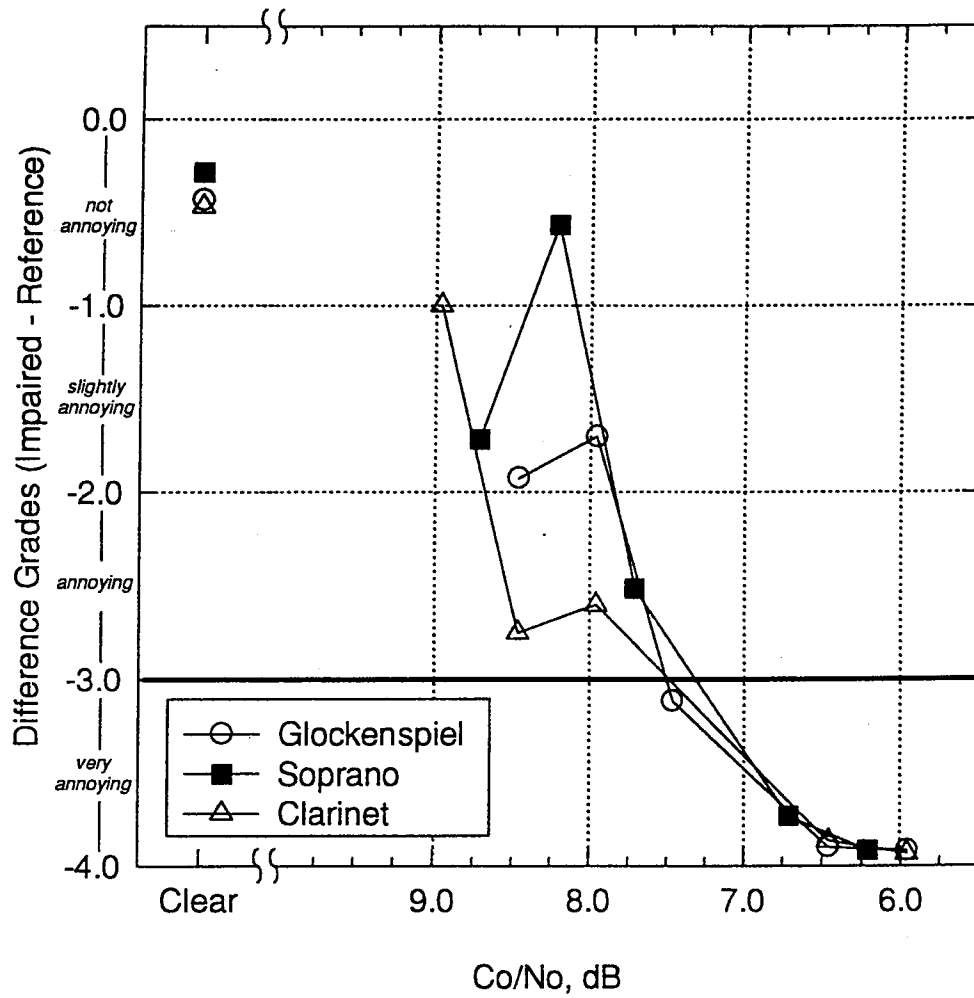


Figure 3.5. Failure characteristic curves for system *b* for all three test materials with the gaussian noise impairment.

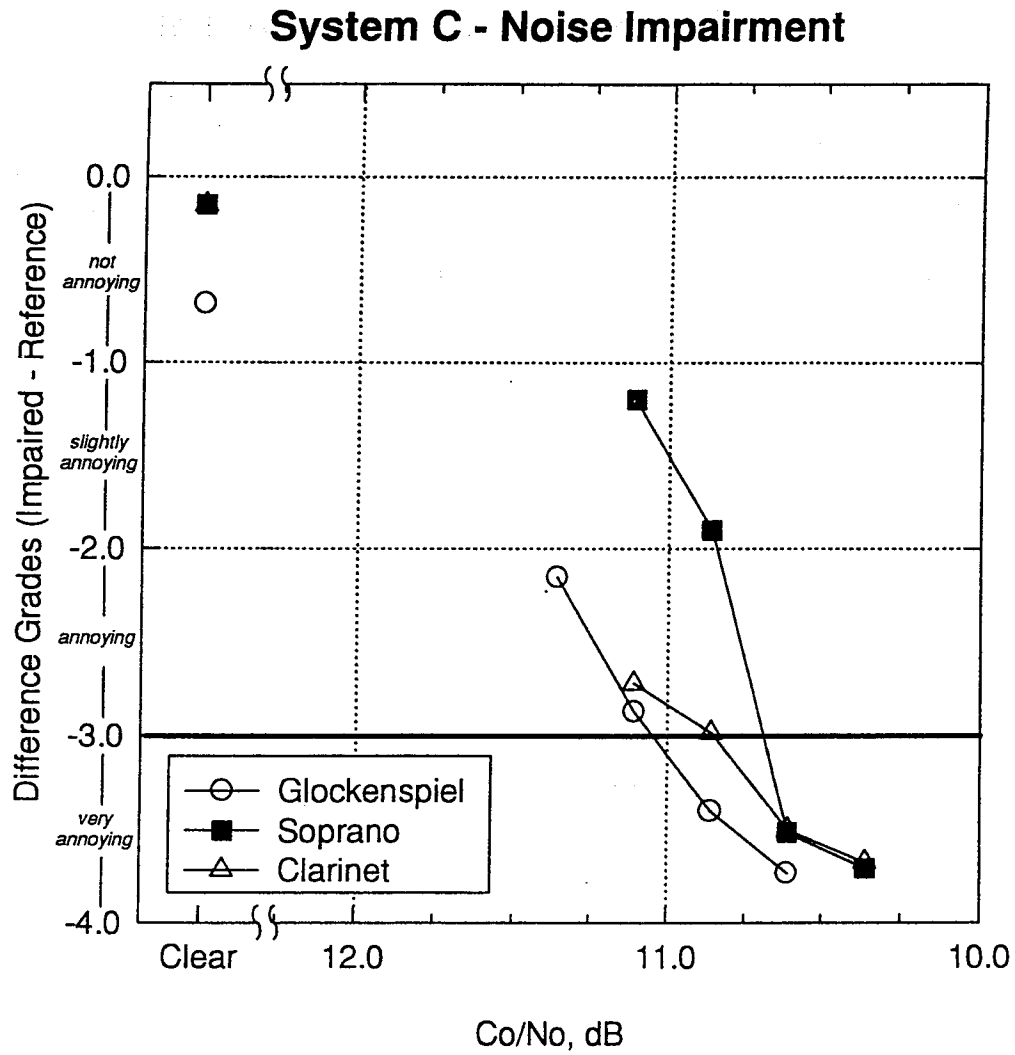


Figure 3.6. Failure characteristic curves for system c for all three test materials with the gaussian noise impairment.

### System D - Noise Impairment

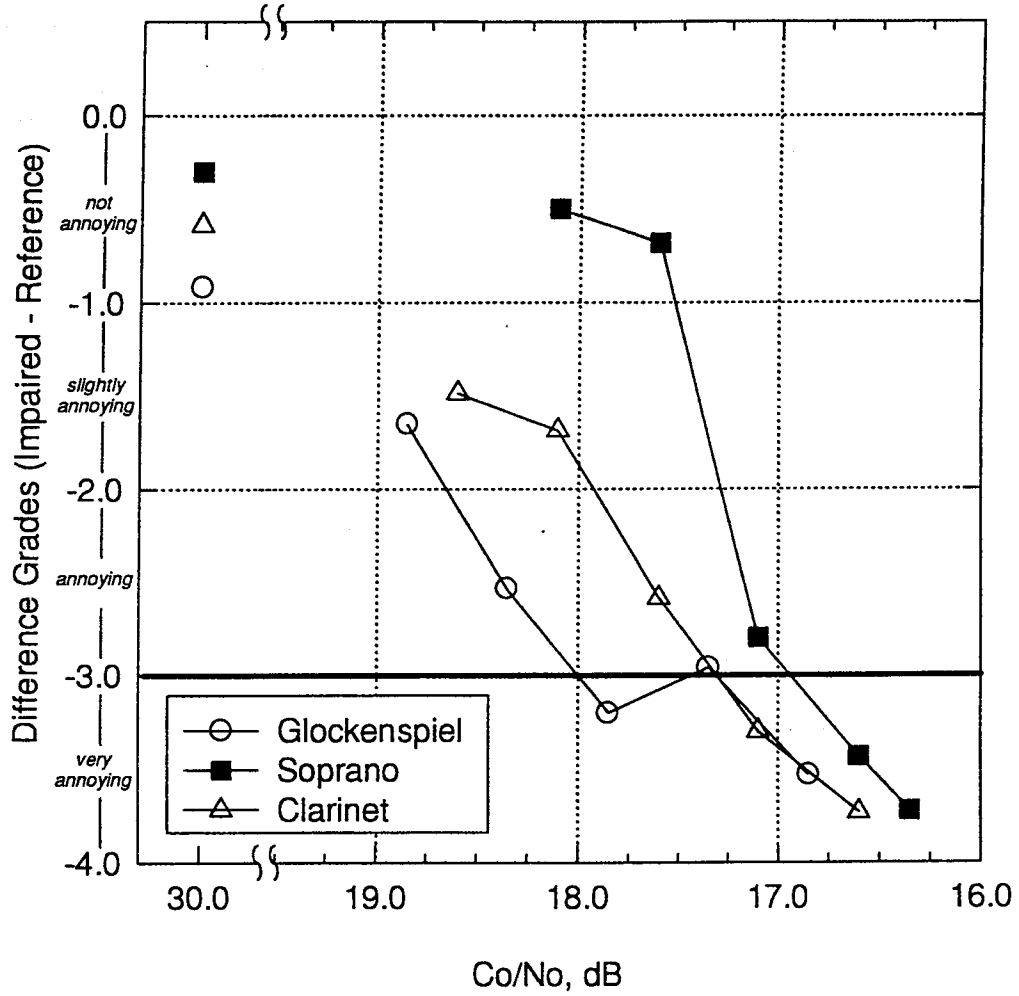


Figure 3.7. Failure characteristic curves for system *d* for all three test materials with the gaussian noise impairment.

### System E - Noise Impairment

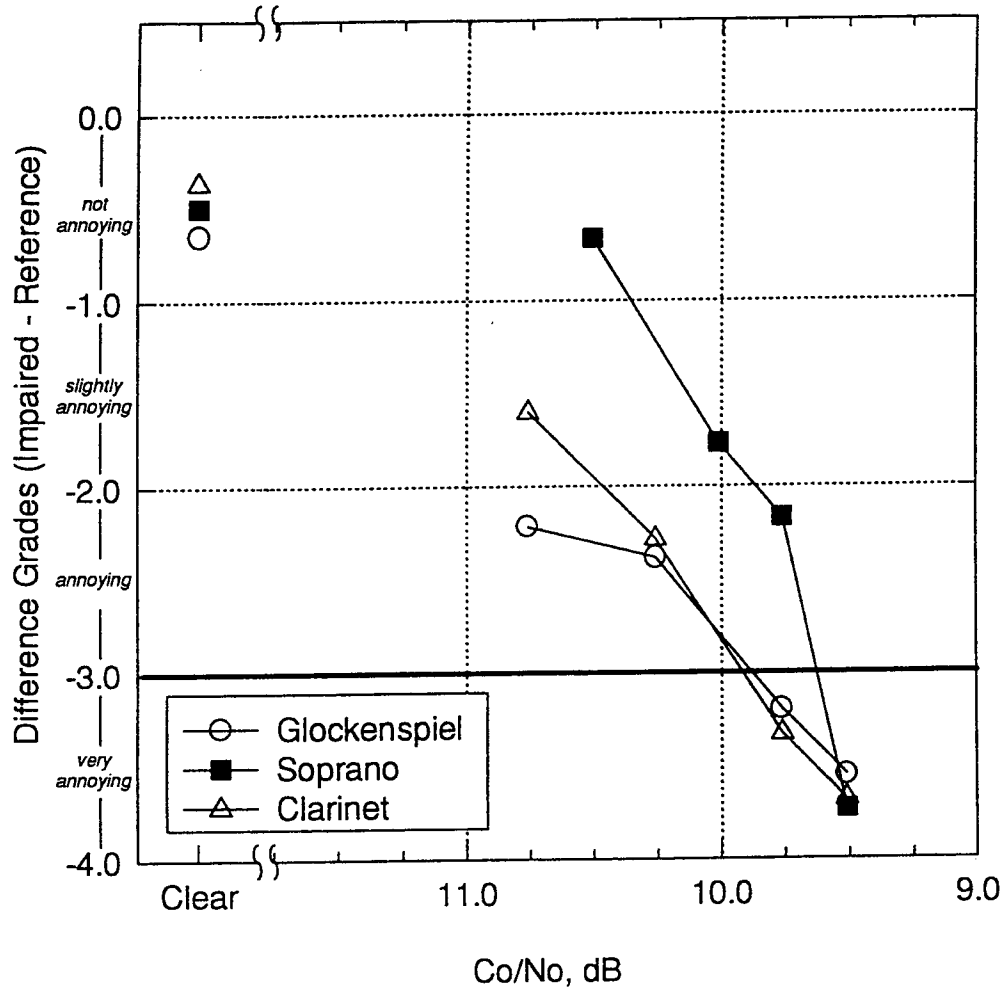


Figure 3.8. Failure characteristic curves for system e for all three test materials with the gaussian noise impairment.

### System F - Noise Impairment

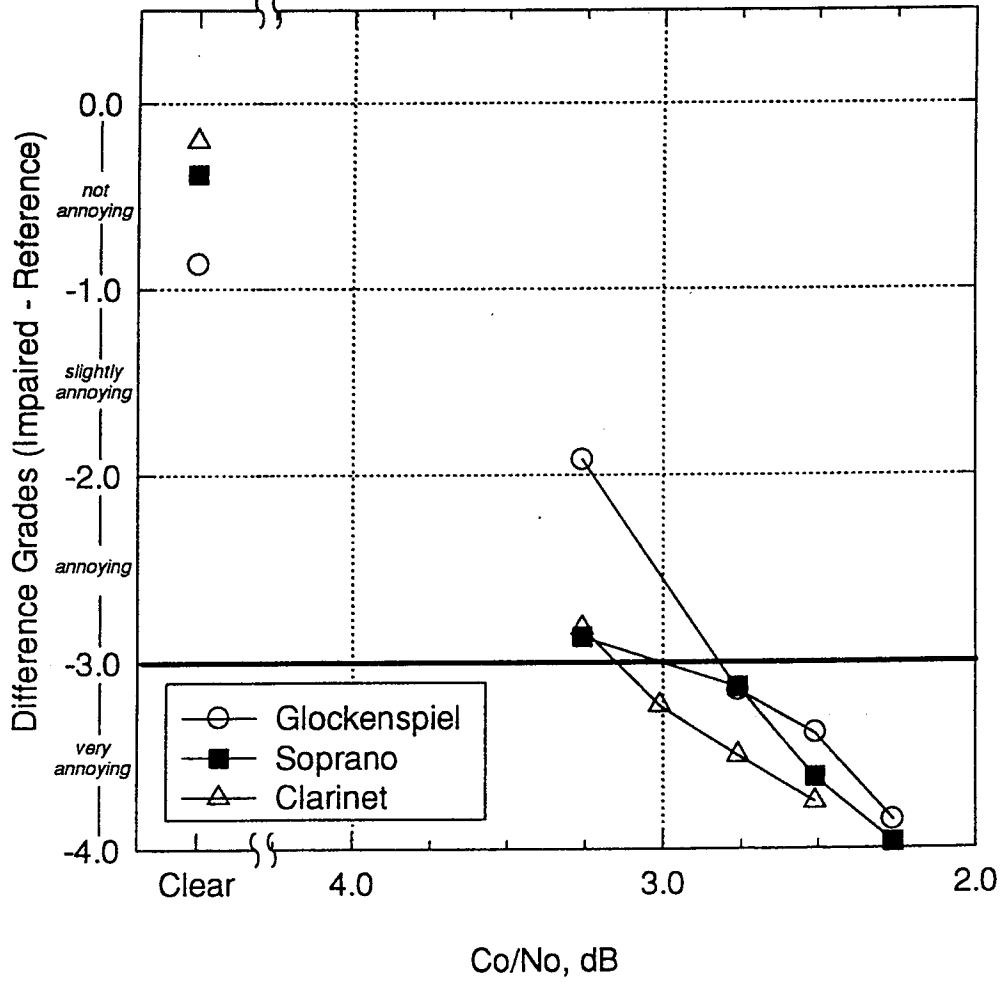


Figure 3.9. Failure characteristic curves for system *f* for all three test materials with the gaussian noise impairment.

### System G - Noise Impairment

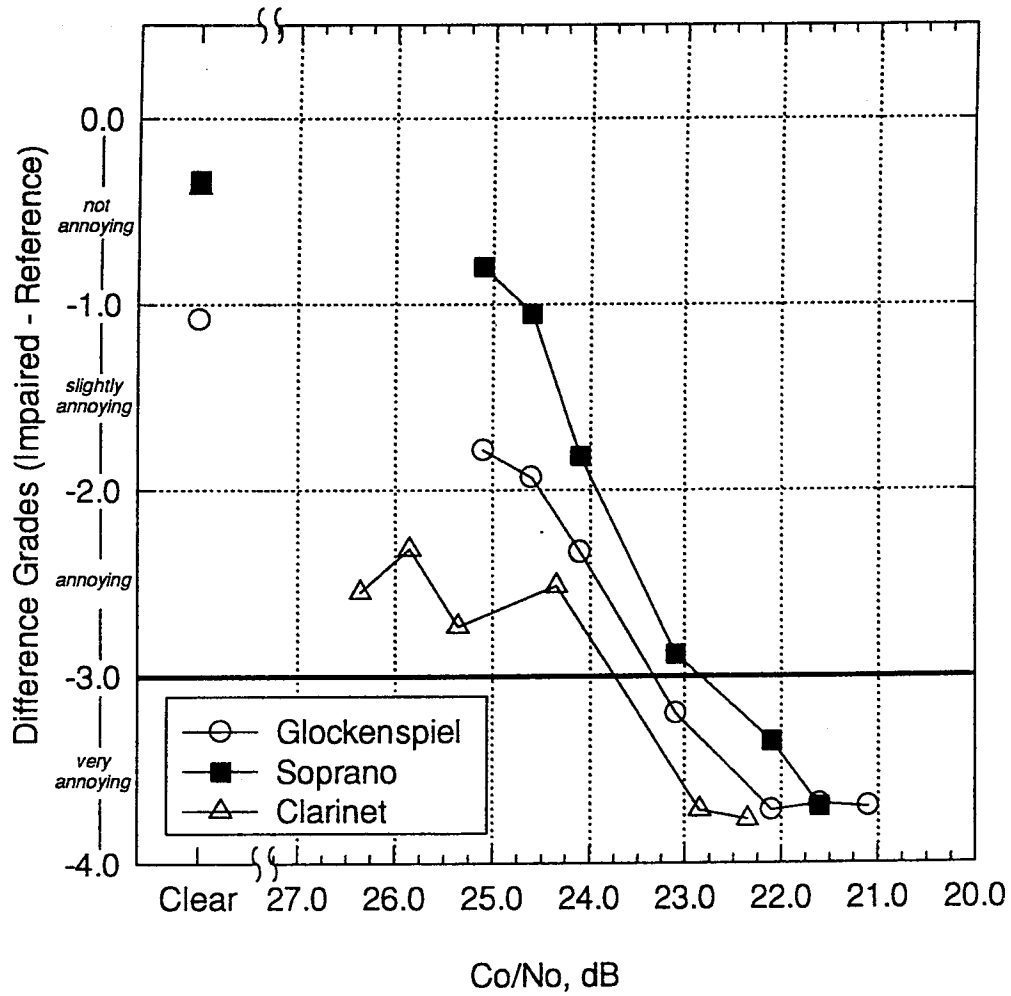


Figure 3.10. Failure characteristic curves for system *g* for all three test materials with the gaussian noise impairment. (DAR laboratory was unable to certify digital power levels for this system)

### System H - Noise Impairment

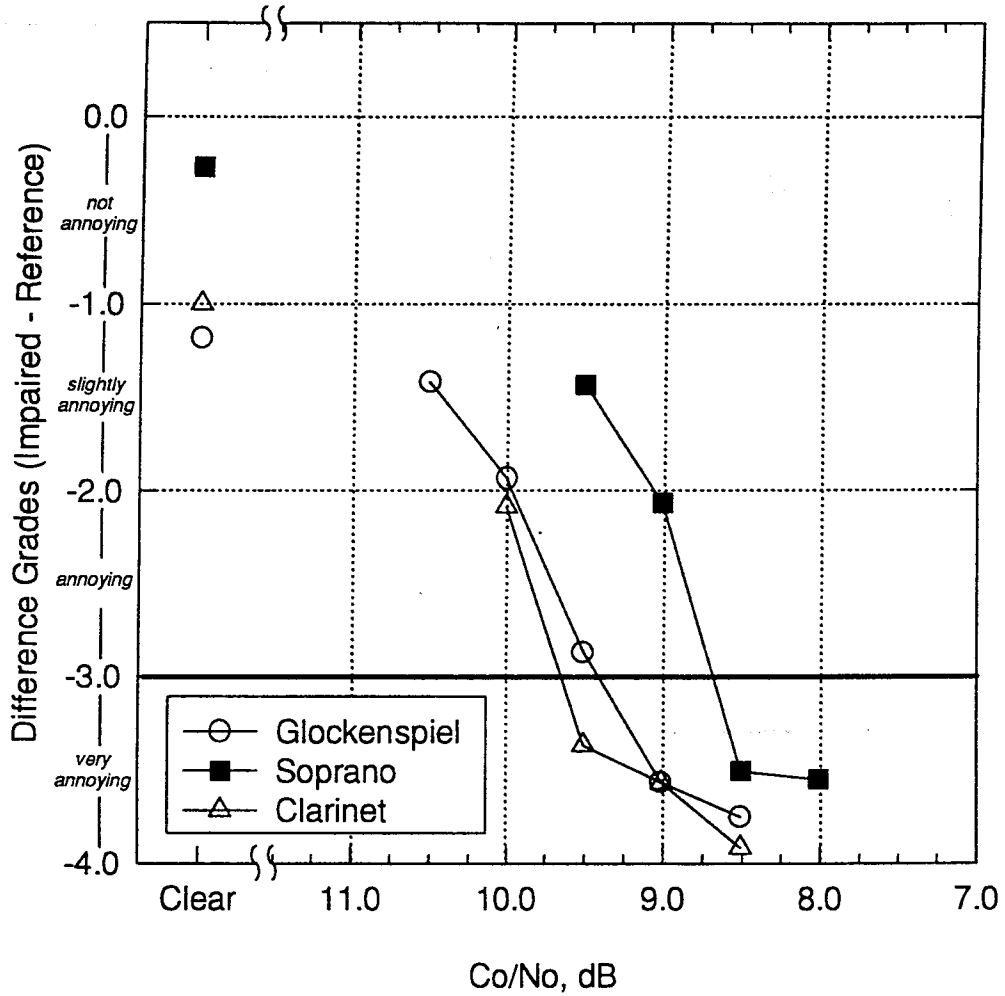


Figure 3.11. Failure characteristic curves for system *h* for all three test materials with the gaussian noise impairment.

### System I - Noise Impairment

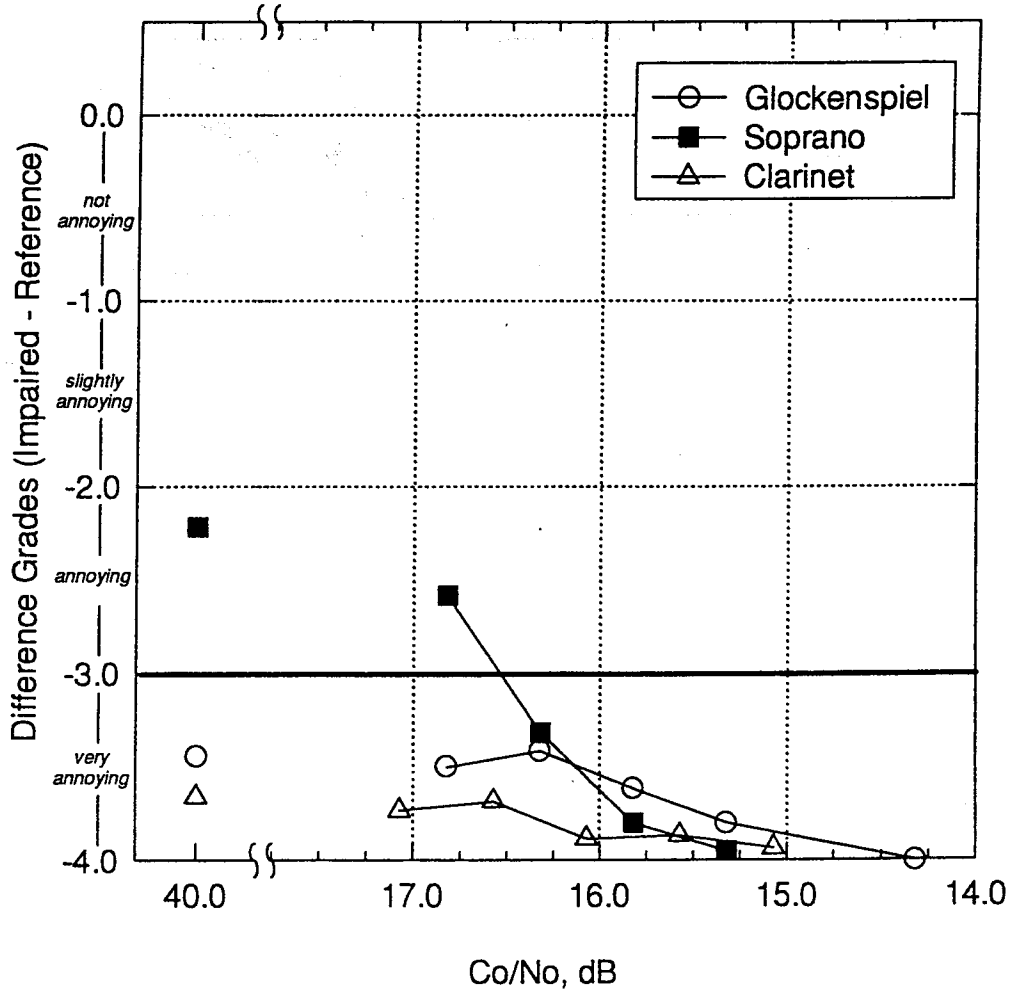


Figure 3.12. Failure characteristic curves for system *i* for all three test materials with the gaussian noise impairment. (DAR laboratory was unable to certify digital power levels for this system)



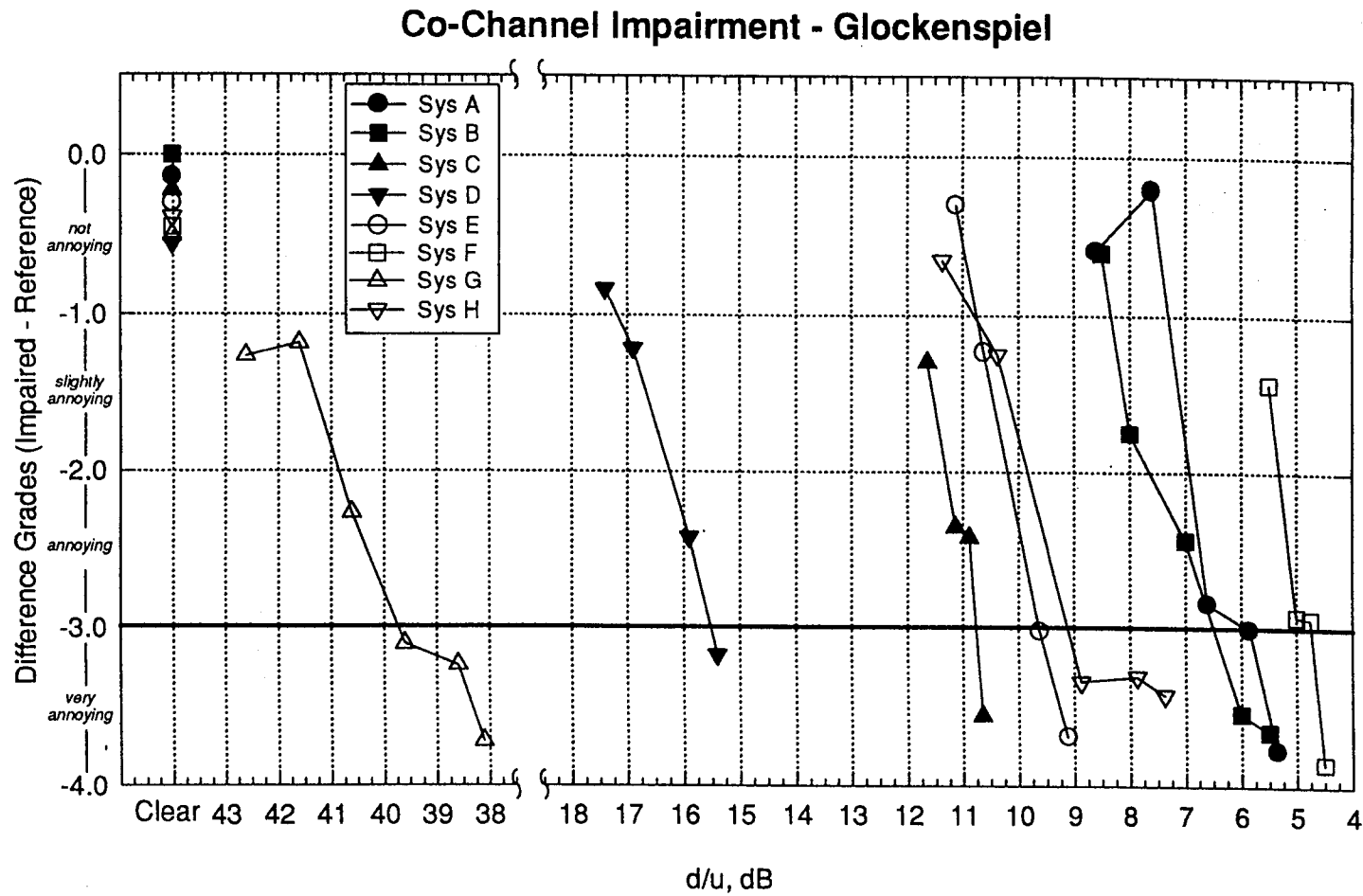


Figure 3.13. Failure characteristic curves of DAR systems a through h tested with co-channel interference and the Glockenspiel test material. (DAR laboratory was unable to certify digital power levels for system g)

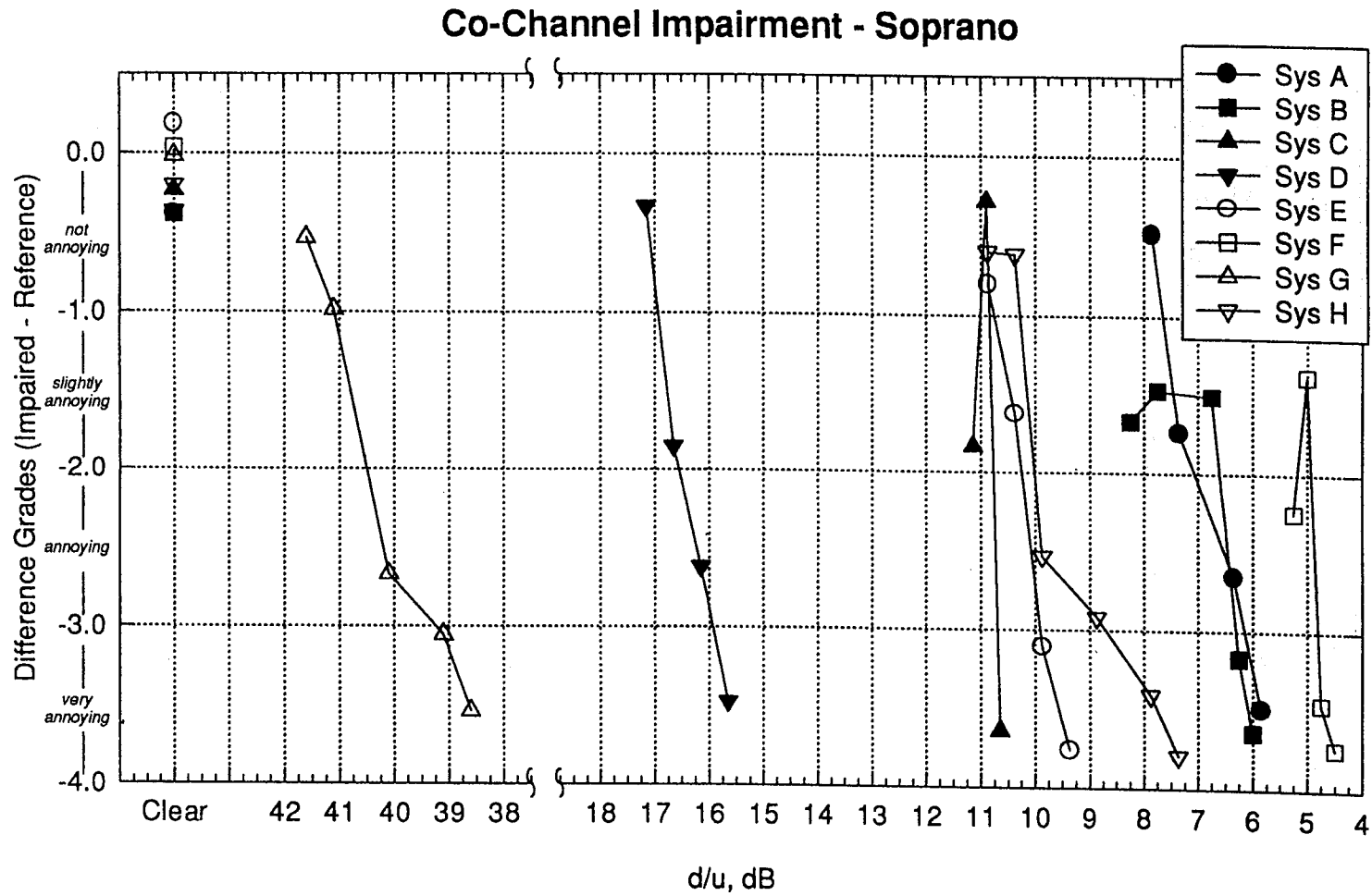


Figure 3.14. Failure characteristic curves of DAR systems *a* through *h* tested with co-channel interference and the Soprano test material. (DAR laboratory was unable to certify digital power levels for system *g*)

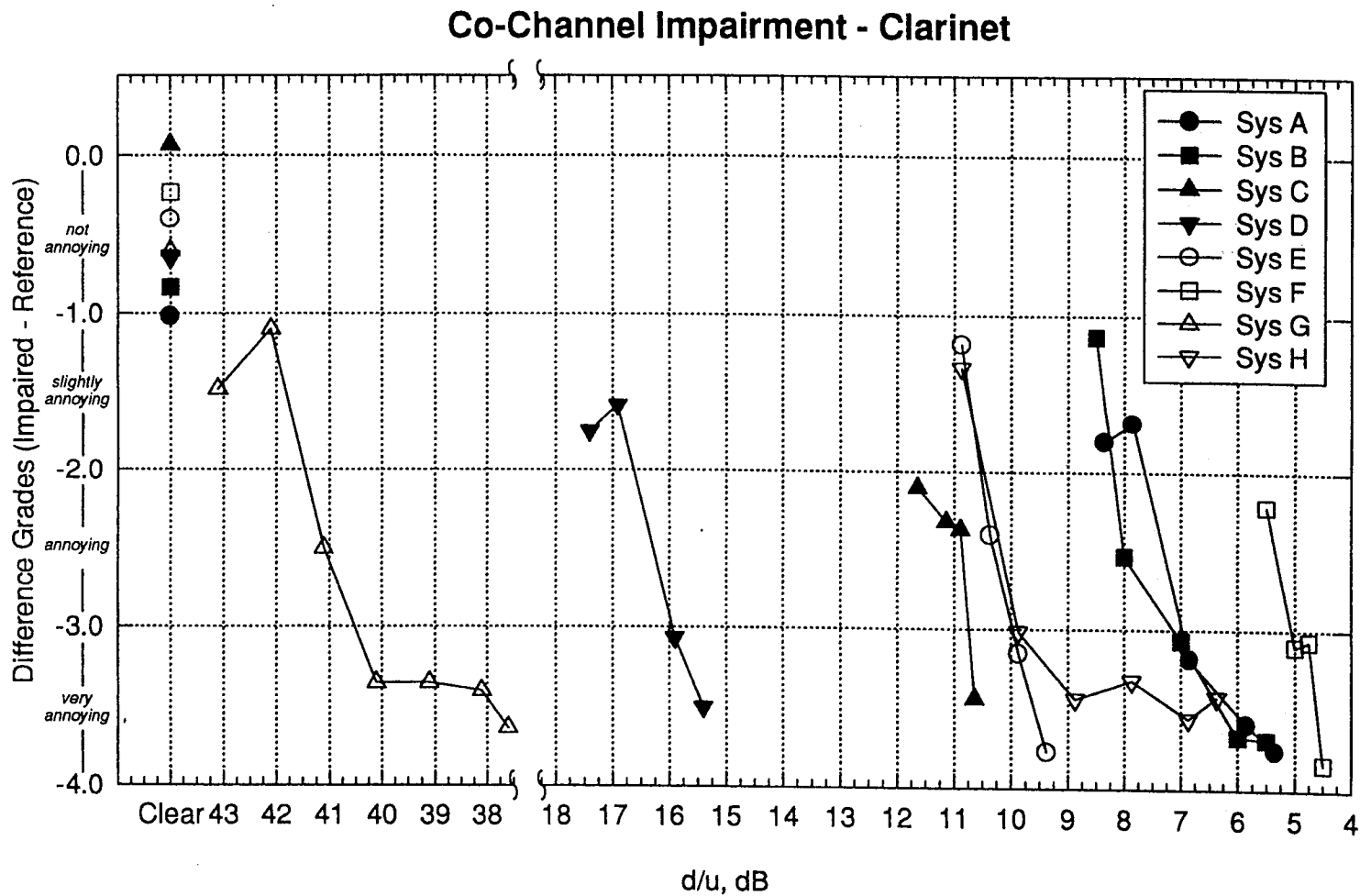


Figure 3.15. Failure characteristic curves of DAR systems *a* through *h* tested with co-channel interference and the Clarinet test material. (DAR laboratory was unable to certify digital power levels for system *g*)

### System A - Co-Channel Impairment

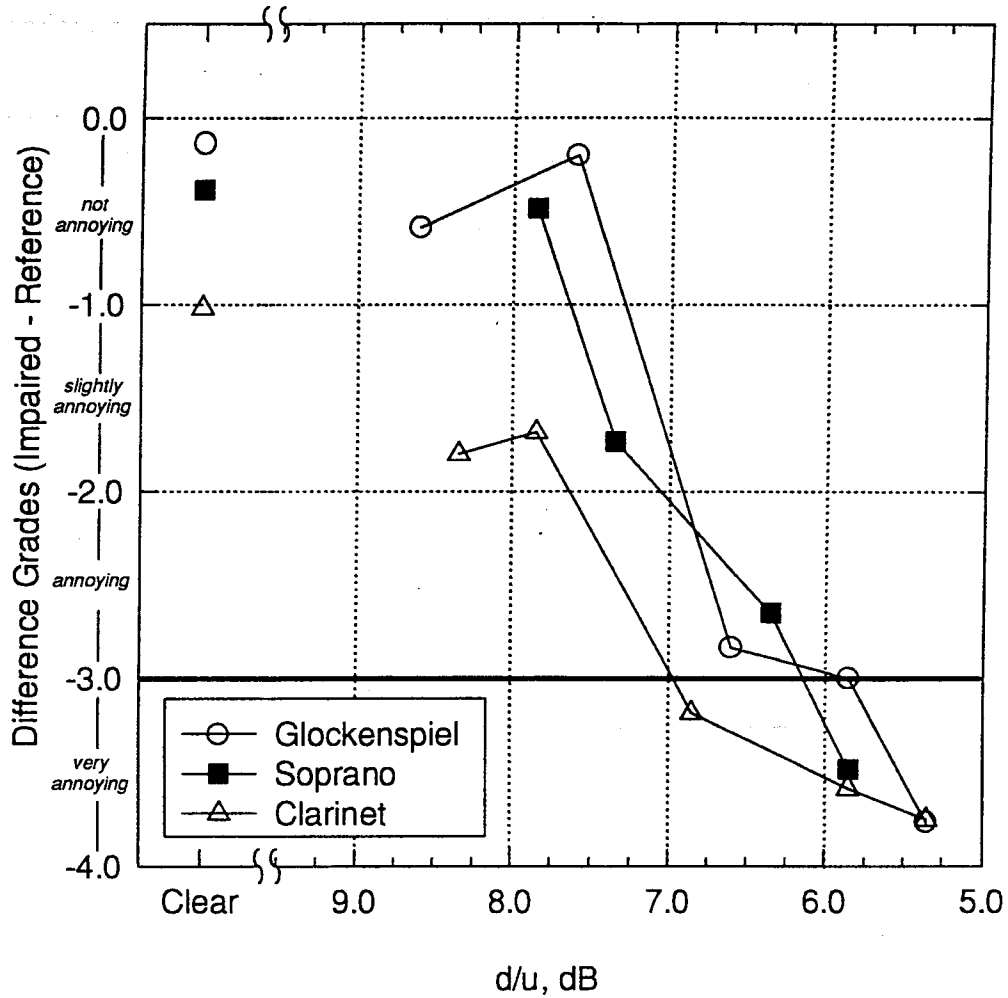


Figure 3.16. Failure characteristic curves for system a for all three test materials with the co-channel interference.

### System B - Co-Channel Impairment

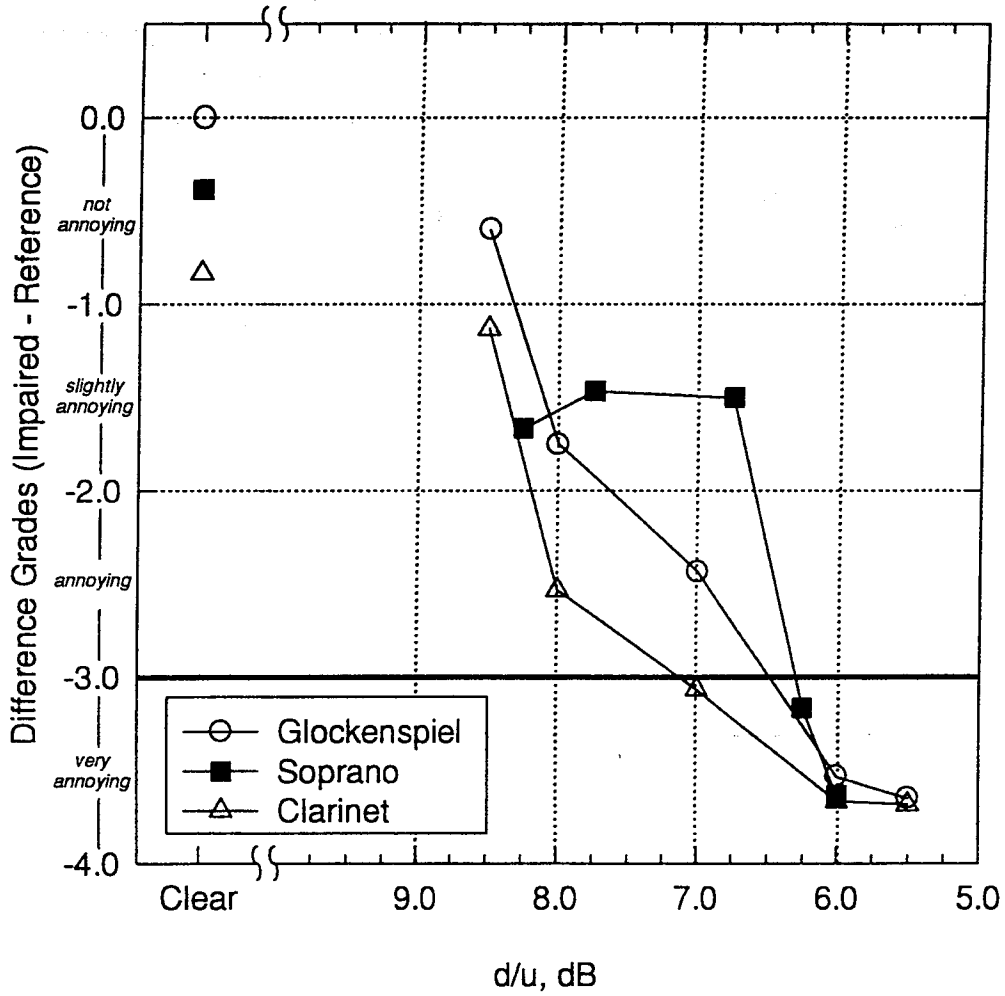


Figure 3.17. Failure characteristic curves for system *b* for all three test materials with the co-channel interference.

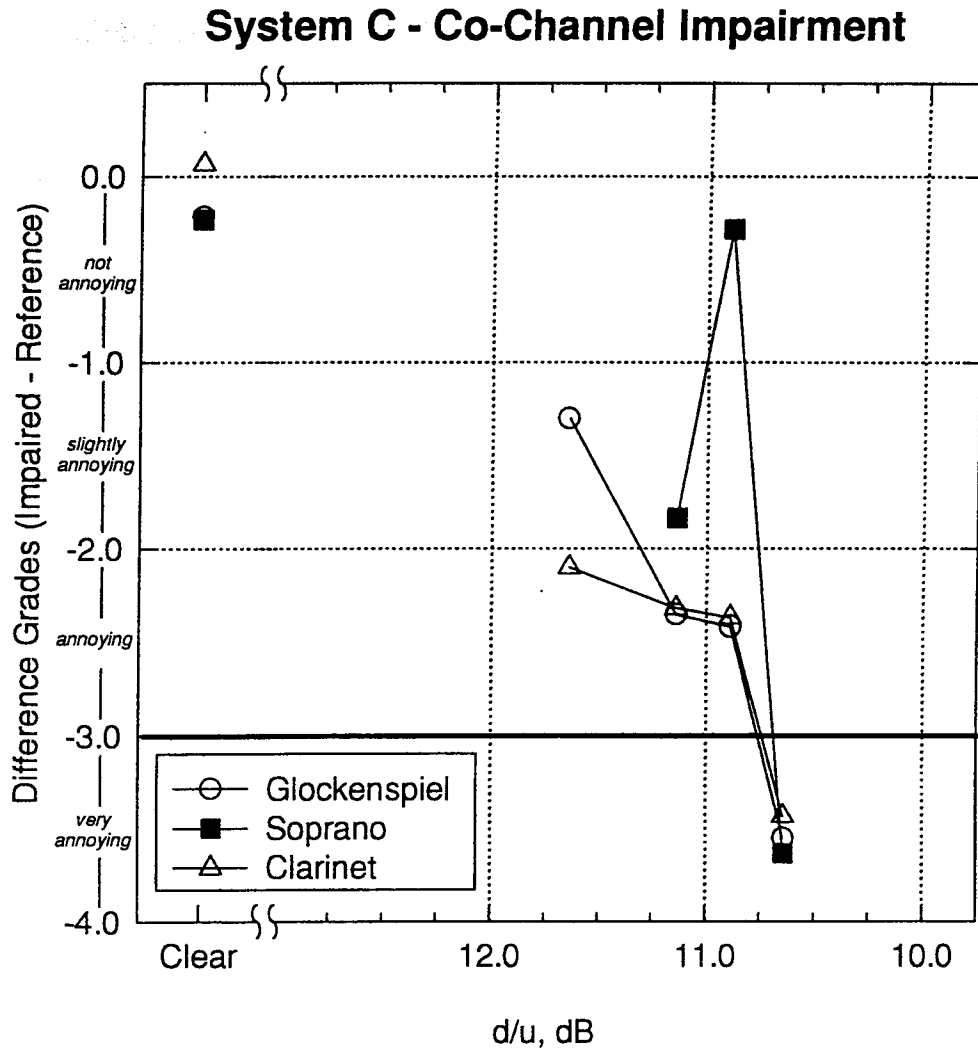


Figure 3.18. Failure characteristic curves for system c for all three test materials with the co-channel interference.

### System D - Co-Channel Impairment

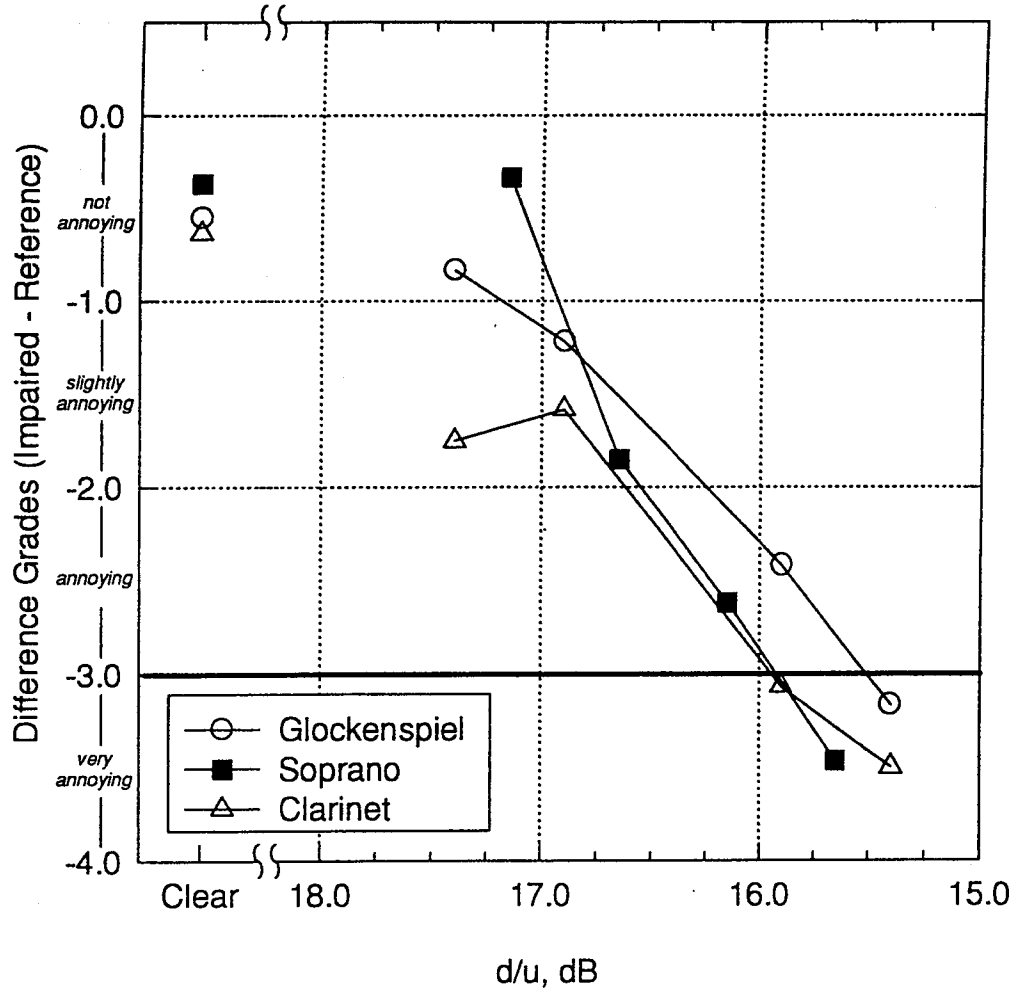


Figure 3.19. Failure characteristic curves for system *d* for all three test materials with the co-channel interference.

### System E - Co-Channel Impairment

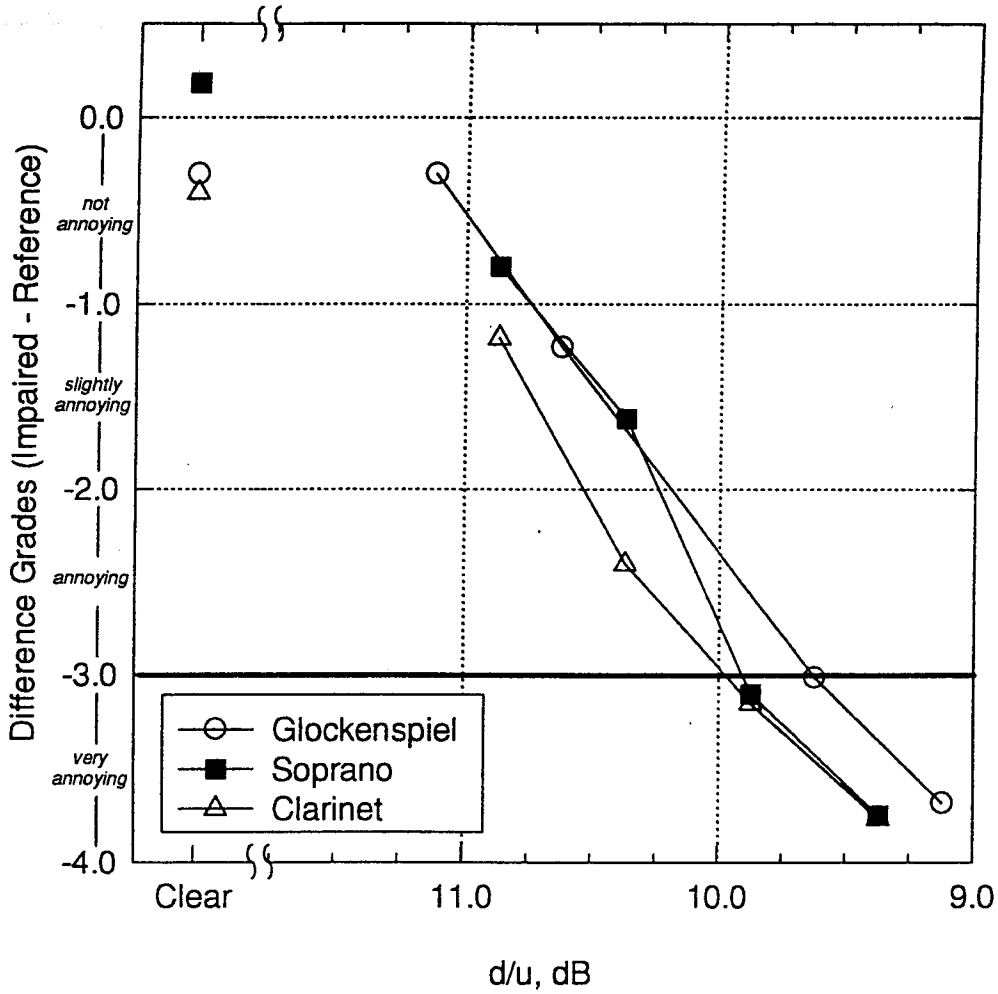


Figure 3.20. Failure characteristic curves for system e for all three test materials with the co-channel interference.



### System F - Co-Channel Impairment

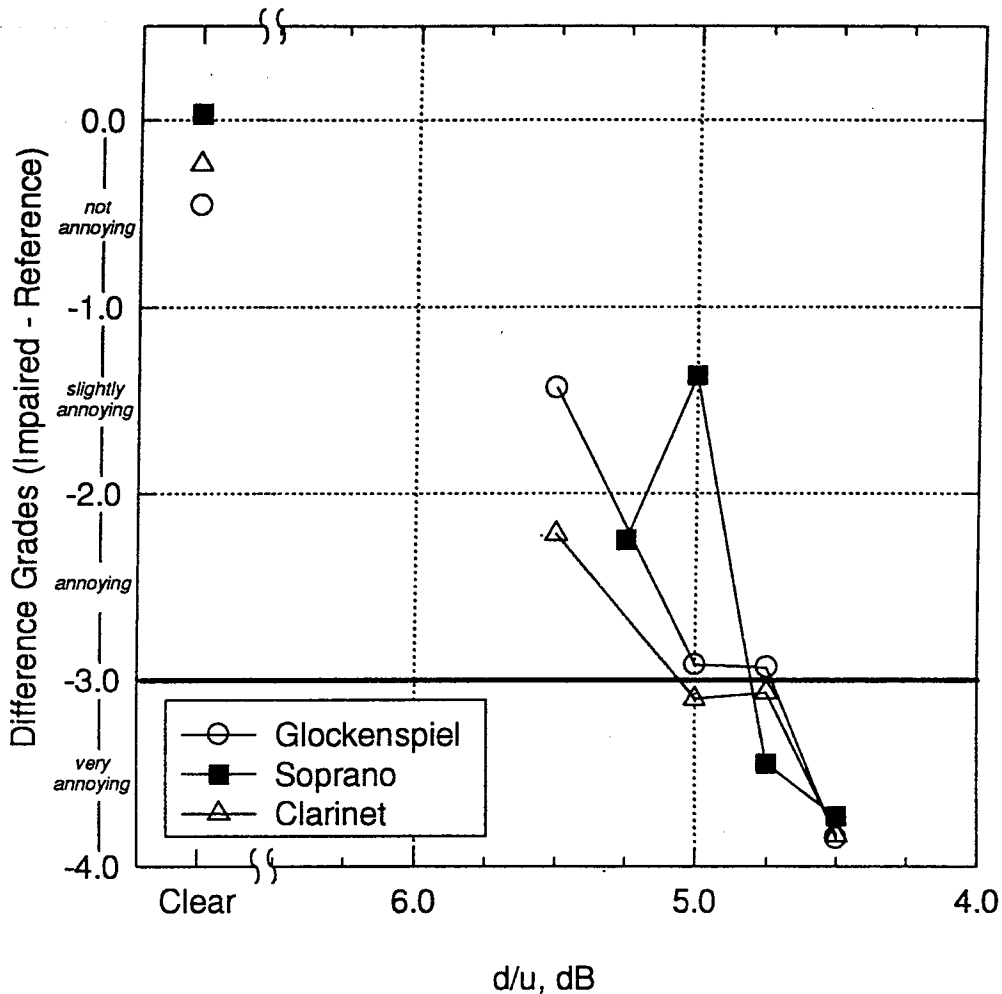


Figure 3.21. Failure characteristic curves for system *f* for all three test materials with the co-channel interference.

### System G - Co-Channel Impairment

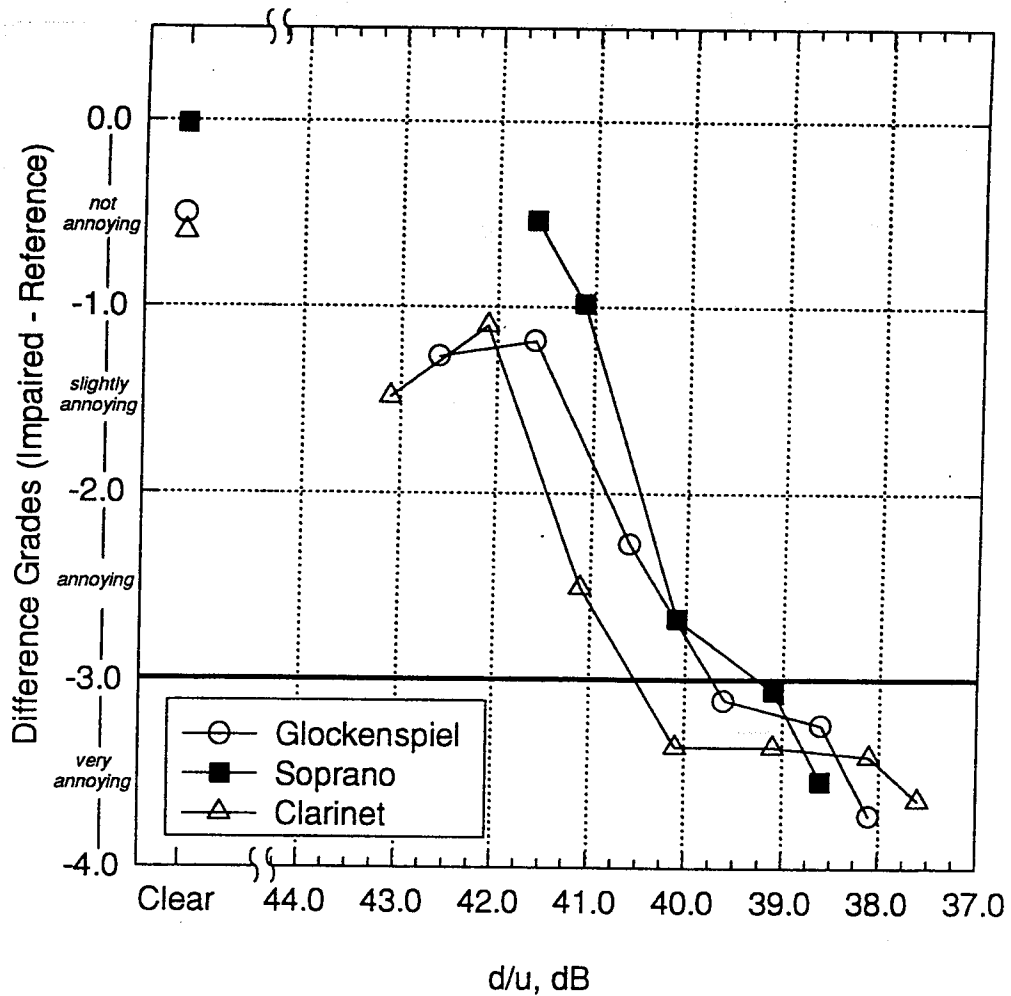


Figure 3.22. Failure characteristic curves for system *g* for all three test materials with the co-channel interference. (DAR laboratory was unable to certify digital power levels for this system)

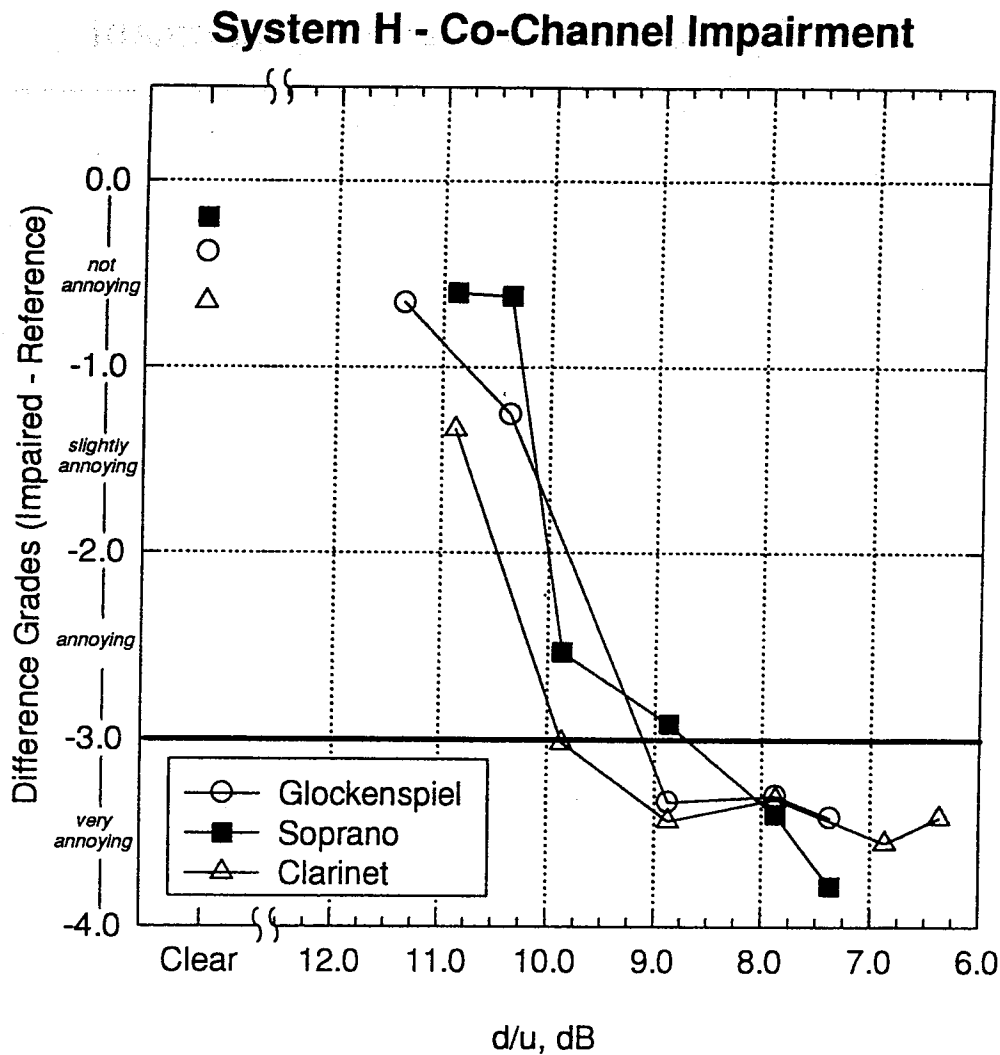


Figure 3.23. Failure characteristic curves for system *h* for all three test materials with the co-channel interference.

### System I - Co-Channel Impairment

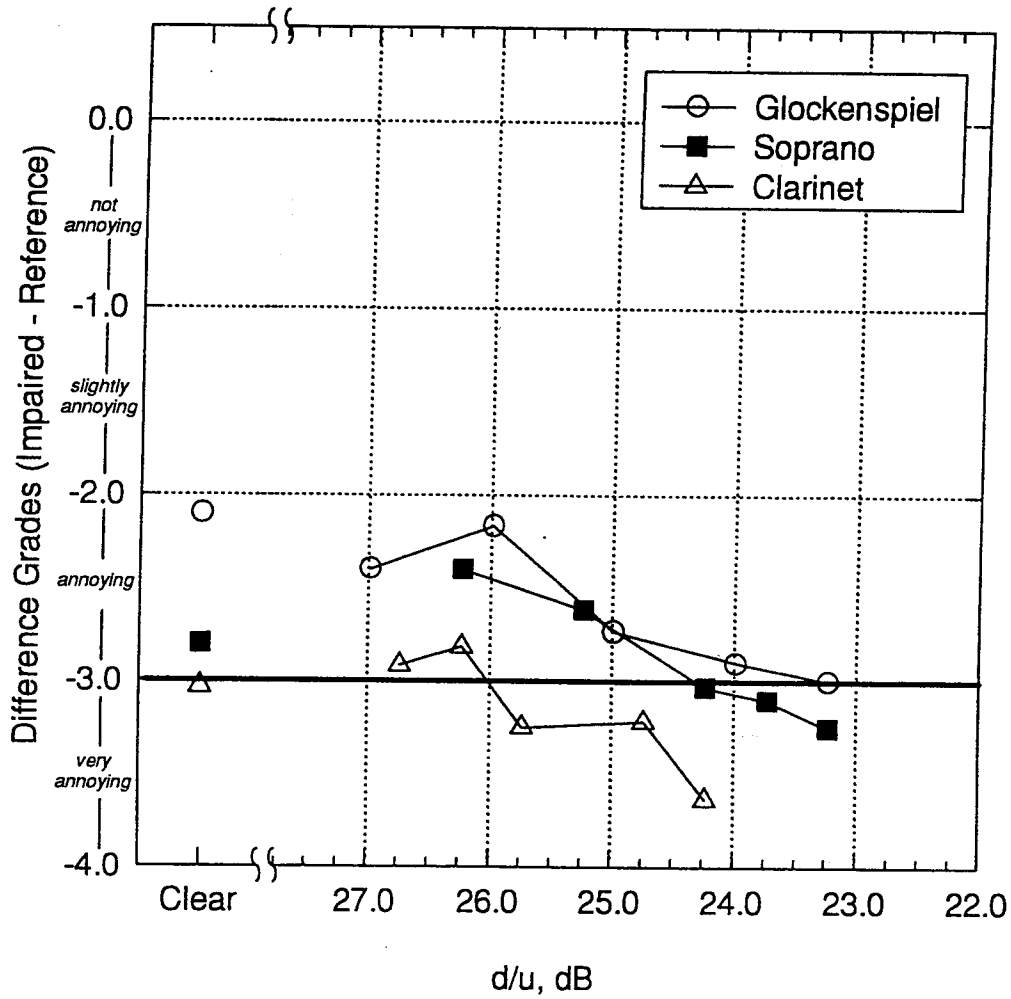


Figure 3.24. Failure characteristic curves for system *i* for all three test materials with the co-channel interference. (DAR laboratory was unable to certify digital power levels for this system)

### Urban Slow Impairment - Glockenspiel

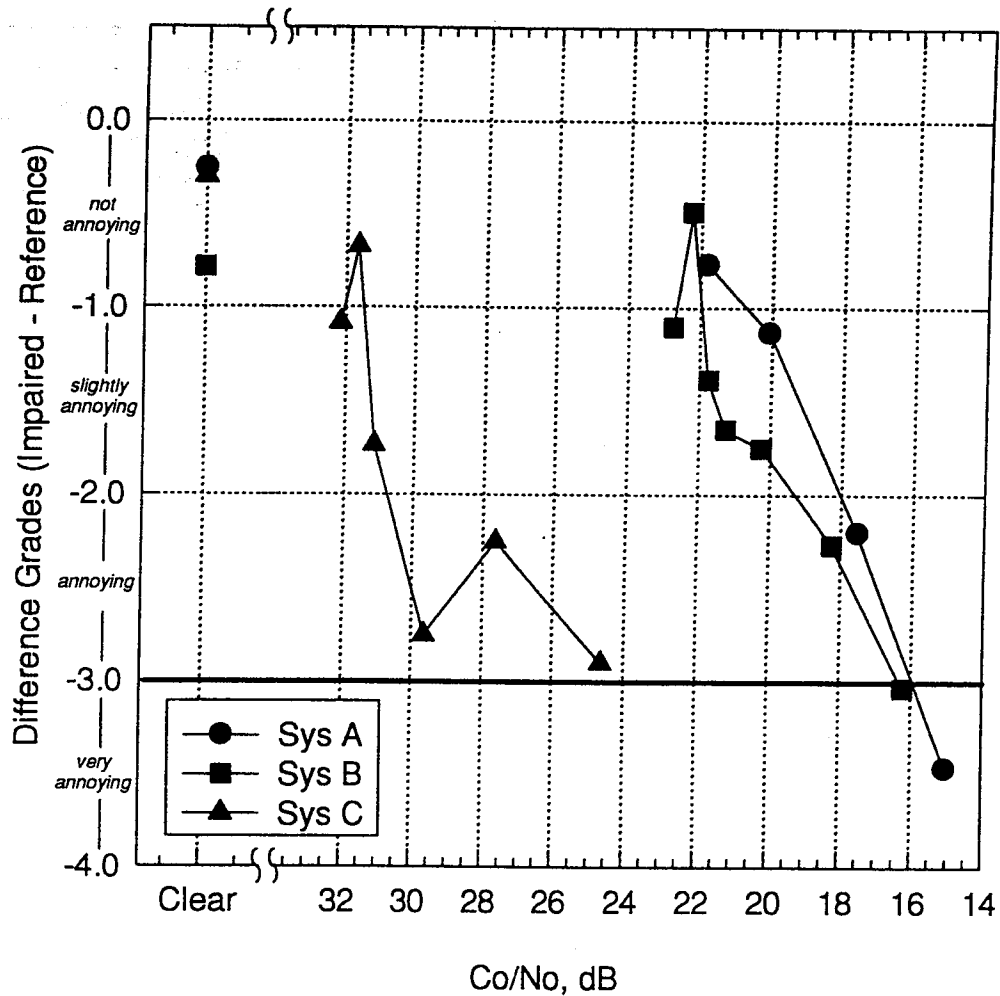


Figure 3.25. Failure characteristic curves of the three DAR systems tested with the Urban Slow multipath impairment and the Glockenspiel test material.

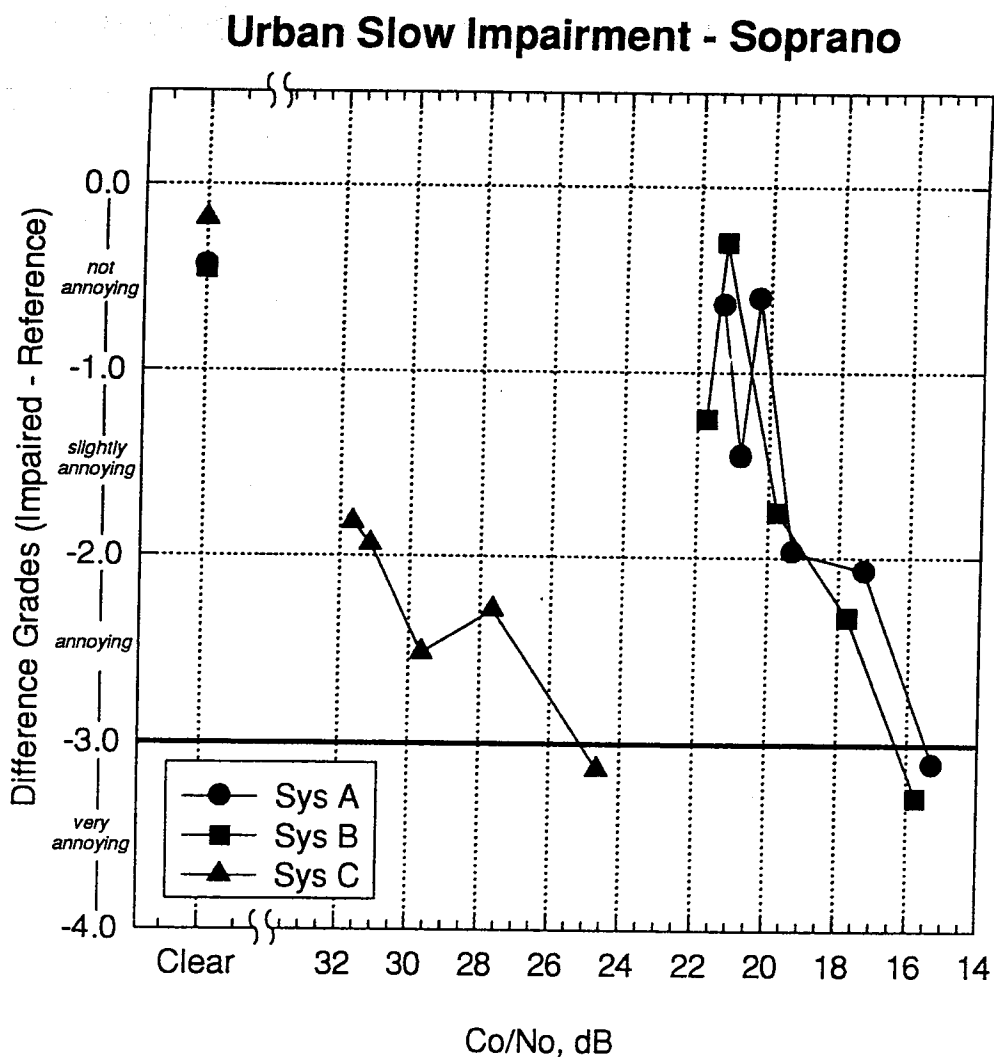


Figure 3.26. Failure characteristic curves of the three DAR systems tested with the Urban Slow multipath impairment and the Soprano test material.

### Urban Slow Impairment - Clarinet

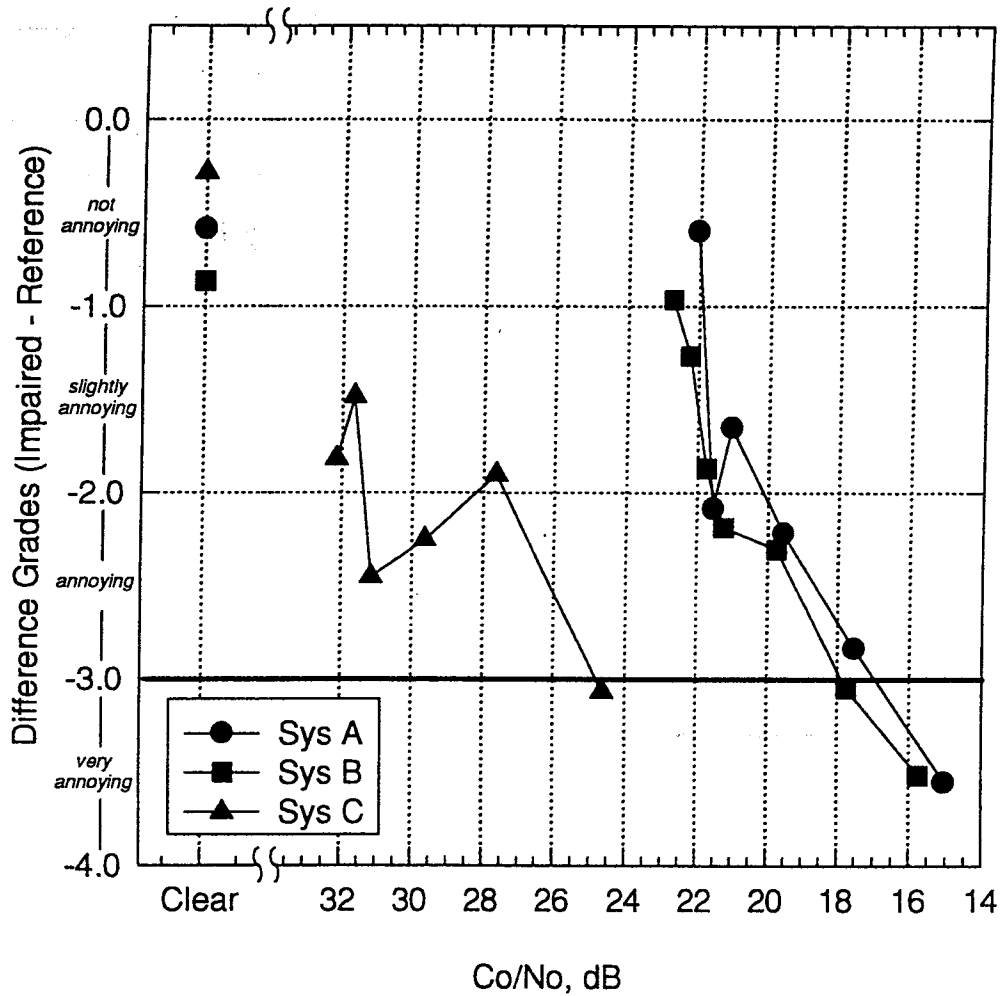


Figure 3.27. Failure characteristic curves of the three DAR systems tested with the Urban Slow multipath impairment and the Clarinet test material.

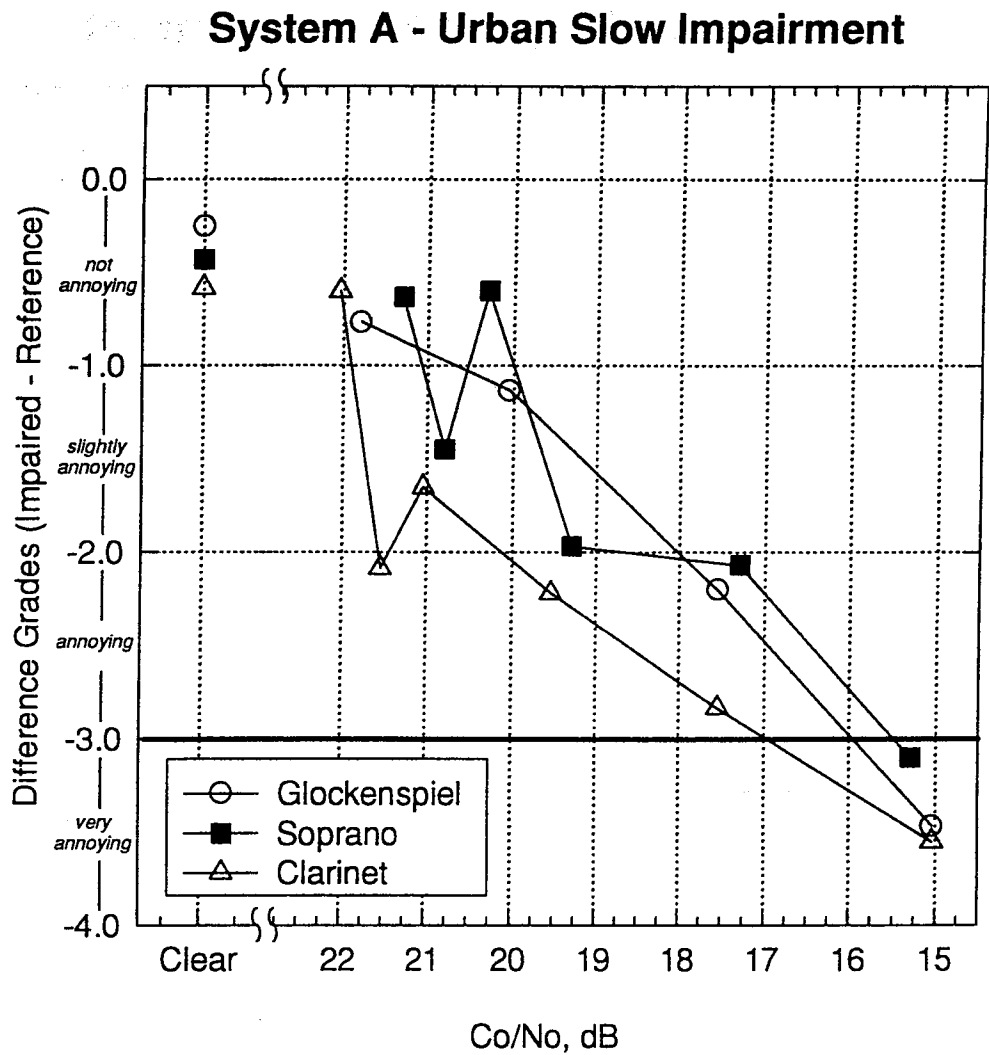


Figure 3.28. Failure characteristic curves for system a for all three test materials with the Urban Slow Multipath impairment.



### System B - Urban Slow Impairment

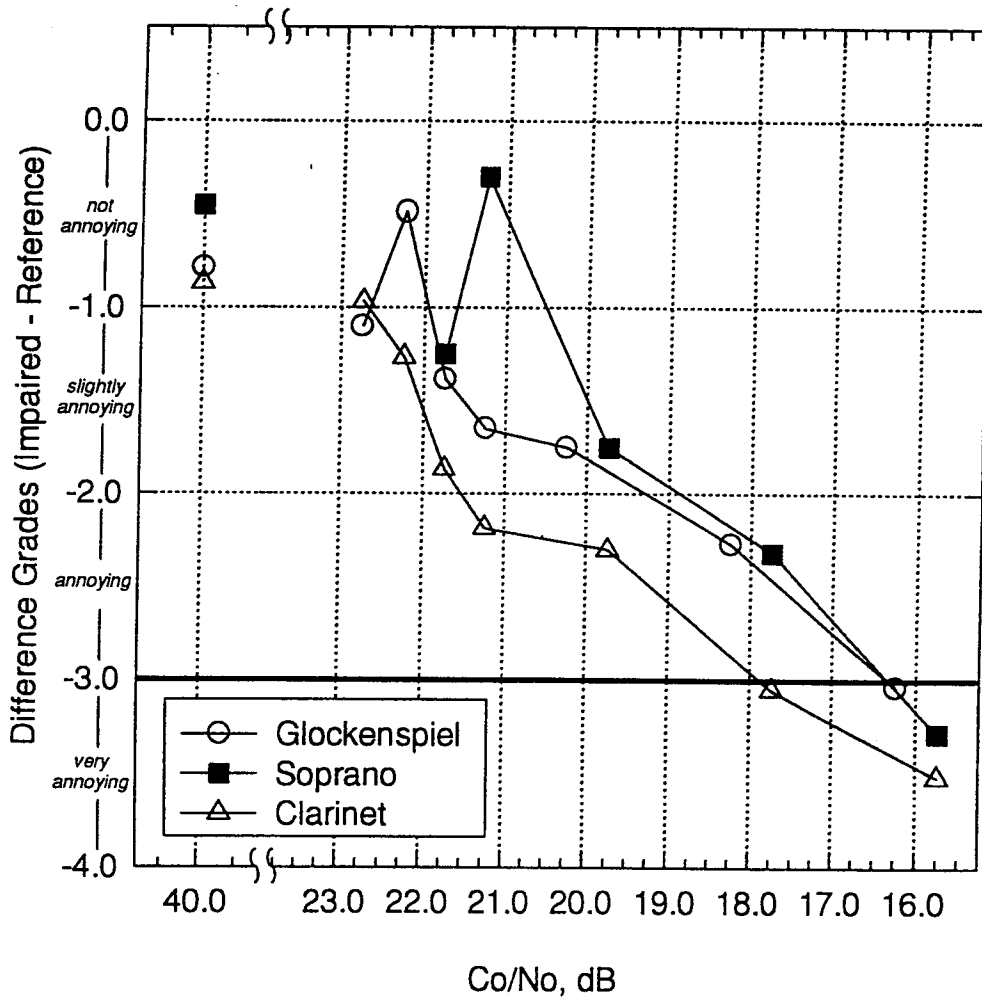


Figure 3.29. Failure characteristic curves for system *b* for all three test materials with the Urban Slow Multipath impairment.

### System C - Urban Slow Impairment

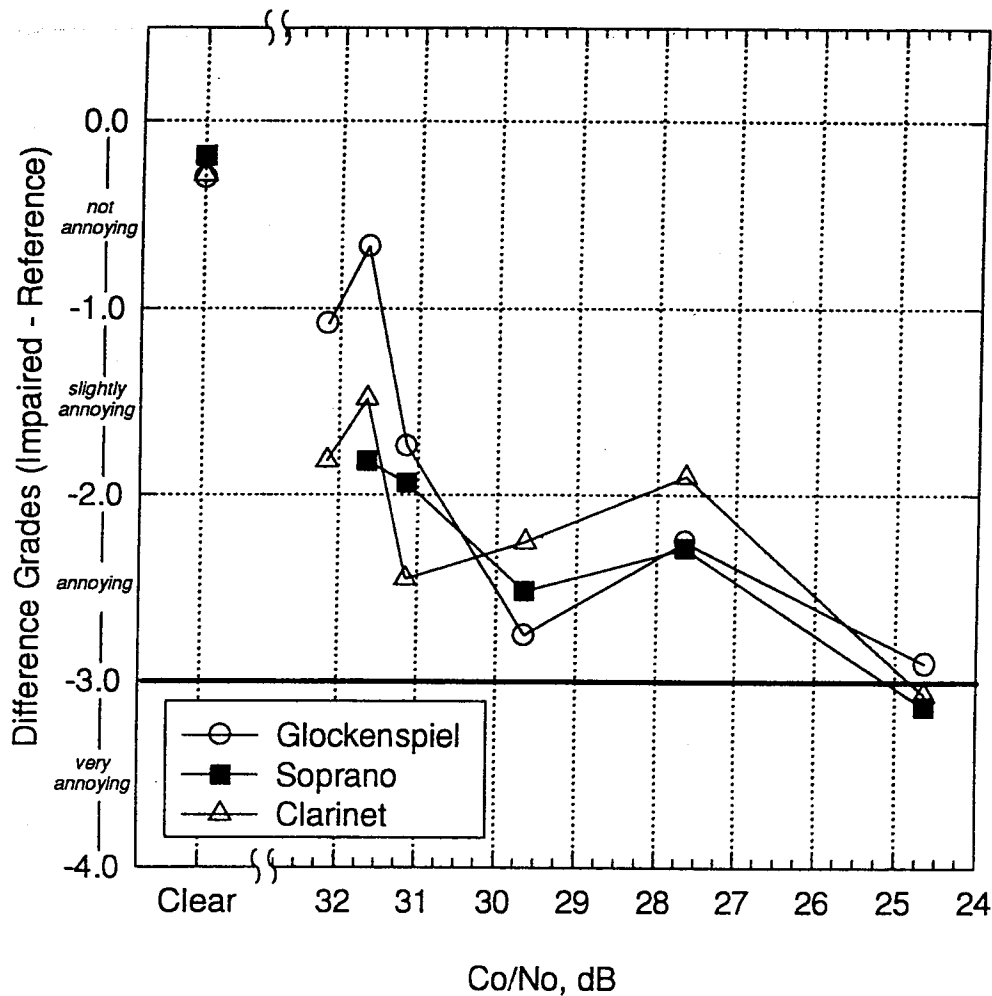


Figure 3.30. Failure characteristic curves for system c for all three test materials with the Urban Slow Multipath impairment.

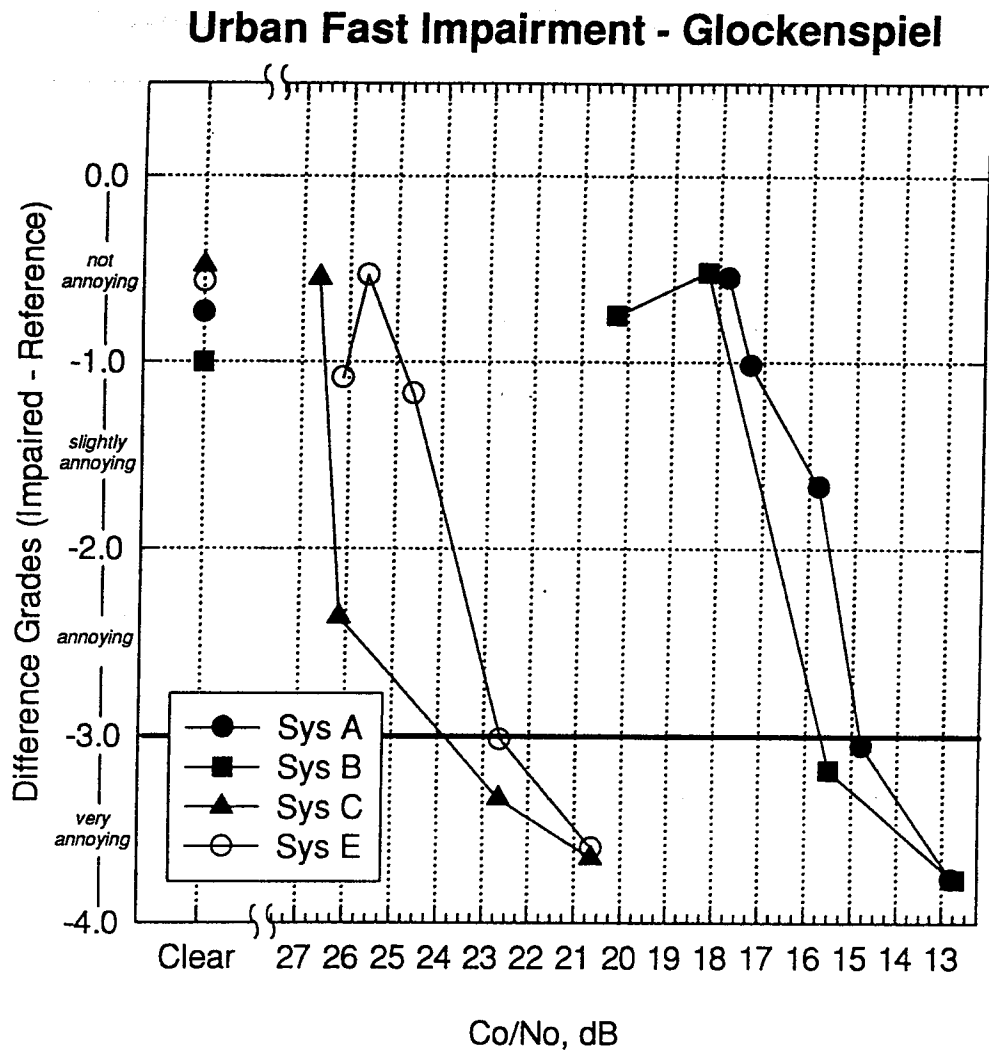


Figure 3.31. Failure characteristic curves of the four DAR systems tested with urban fast multipath impairment and the Glockenspiel test material.

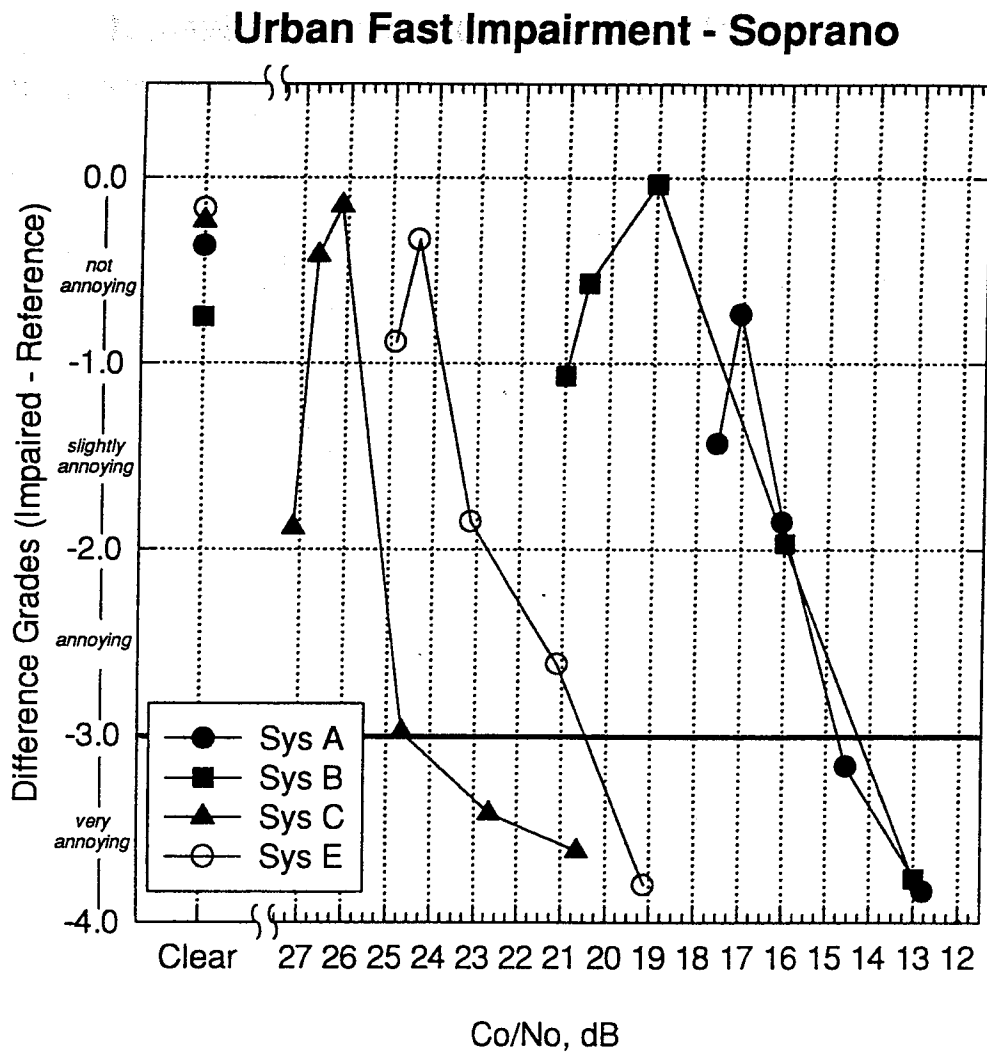


Figure 3.32. Failure characteristic curves of the four DAR systems tested with urban fast multipath impairment and the Soprano test material.

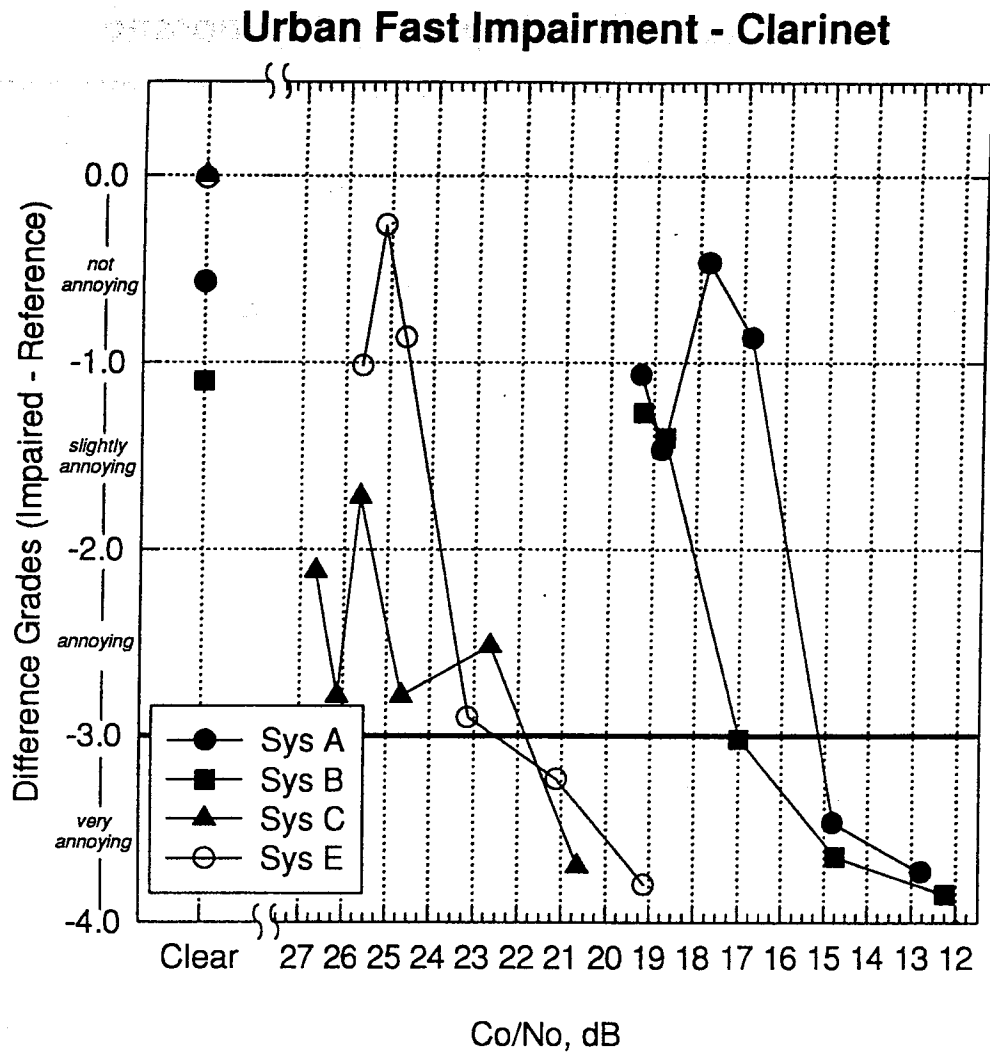


Figure 3.33. Failure characteristic curves of the four DAR systems tested with urban fast multipath impairment and the Clarinet test material.

### System A - Urban Fast Multipath Impairment

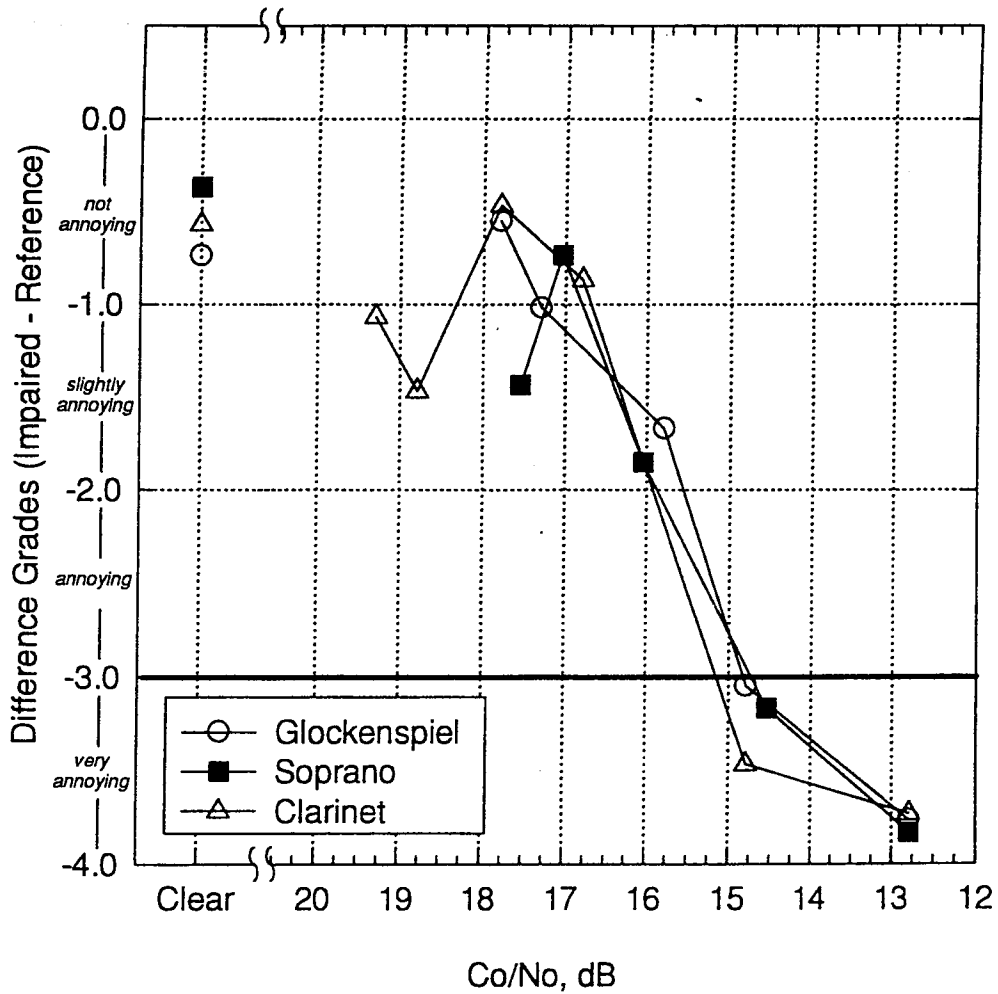


Figure 3.34. Failure characteristic curves for system a for all three test materials with the urban fast multipath impairment.

### System B - Urban Fast Multipath Impairment

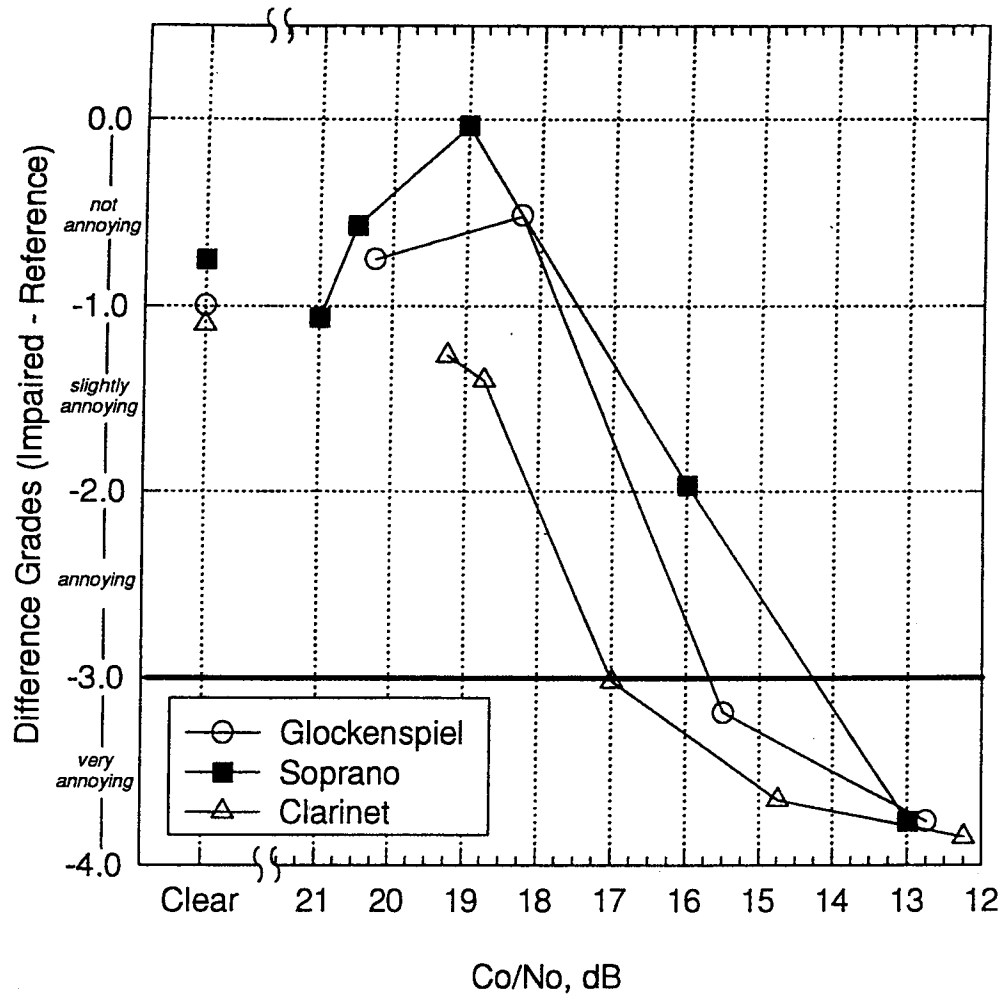


Figure 3.35. Failure characteristic curves for system *b* for all three test materials with the urban fast multipath impairment.

### System C - Urban Fast Multipath Impairment

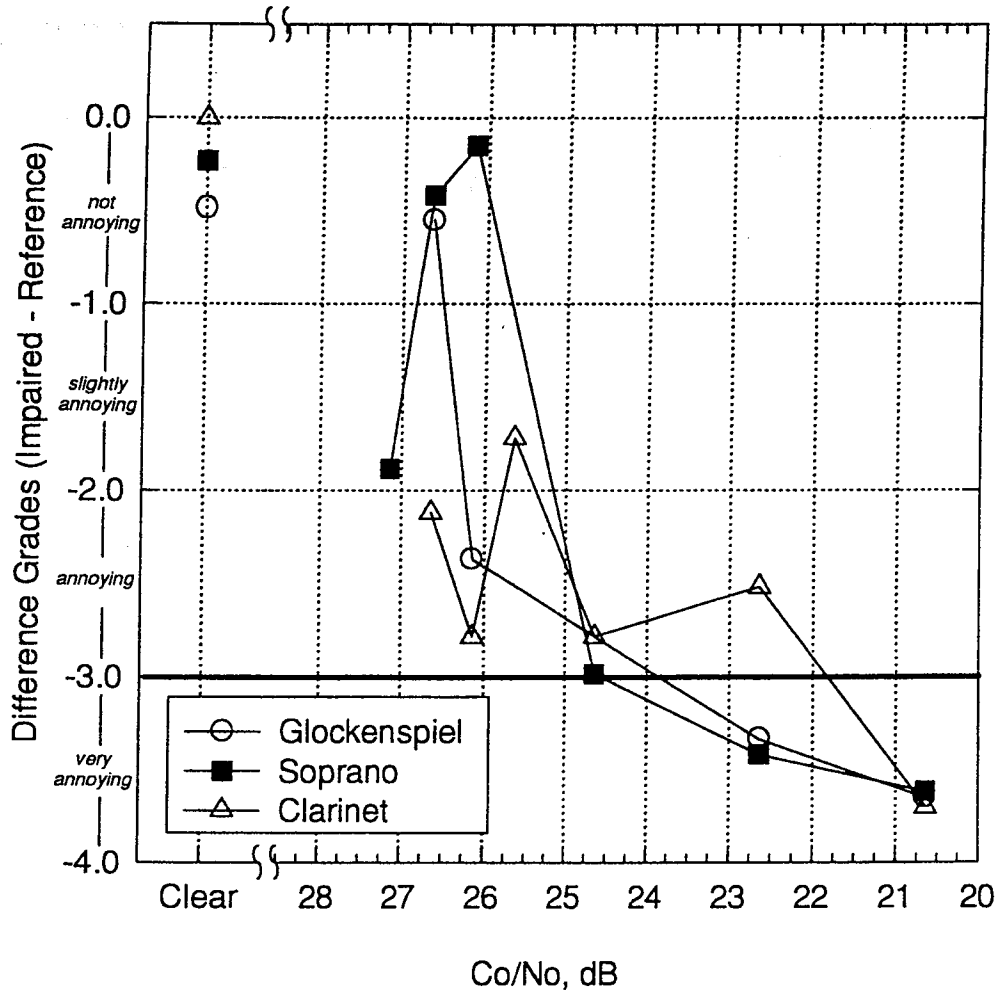


Figure 3.36. Failure characteristic curves for system c for all three test materials with the urban fast multipath impairment.



### System E - Urban Fast Multipath Impairment

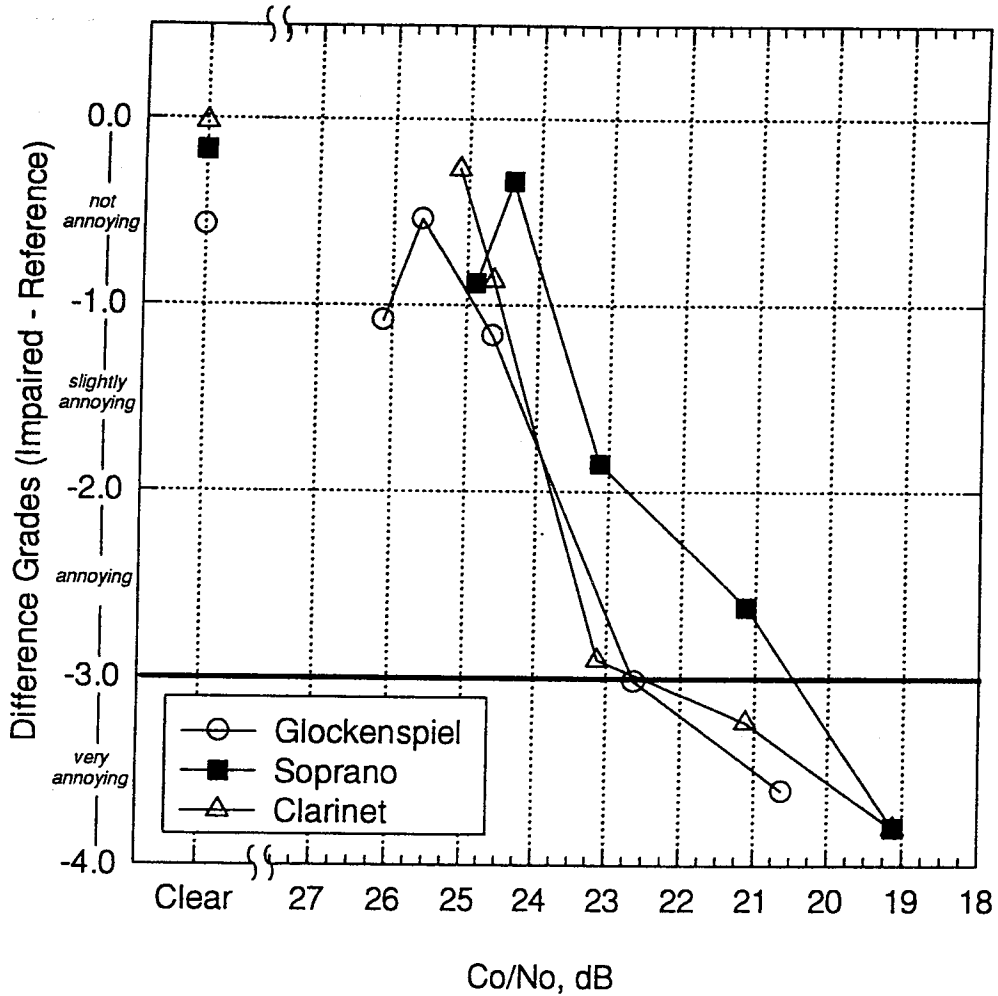


Figure 3.37. Failure characteristic curves for system e for all three test materials with the urban fast multipath impairment.

### Rural Fast Impairment - Glockenspiel

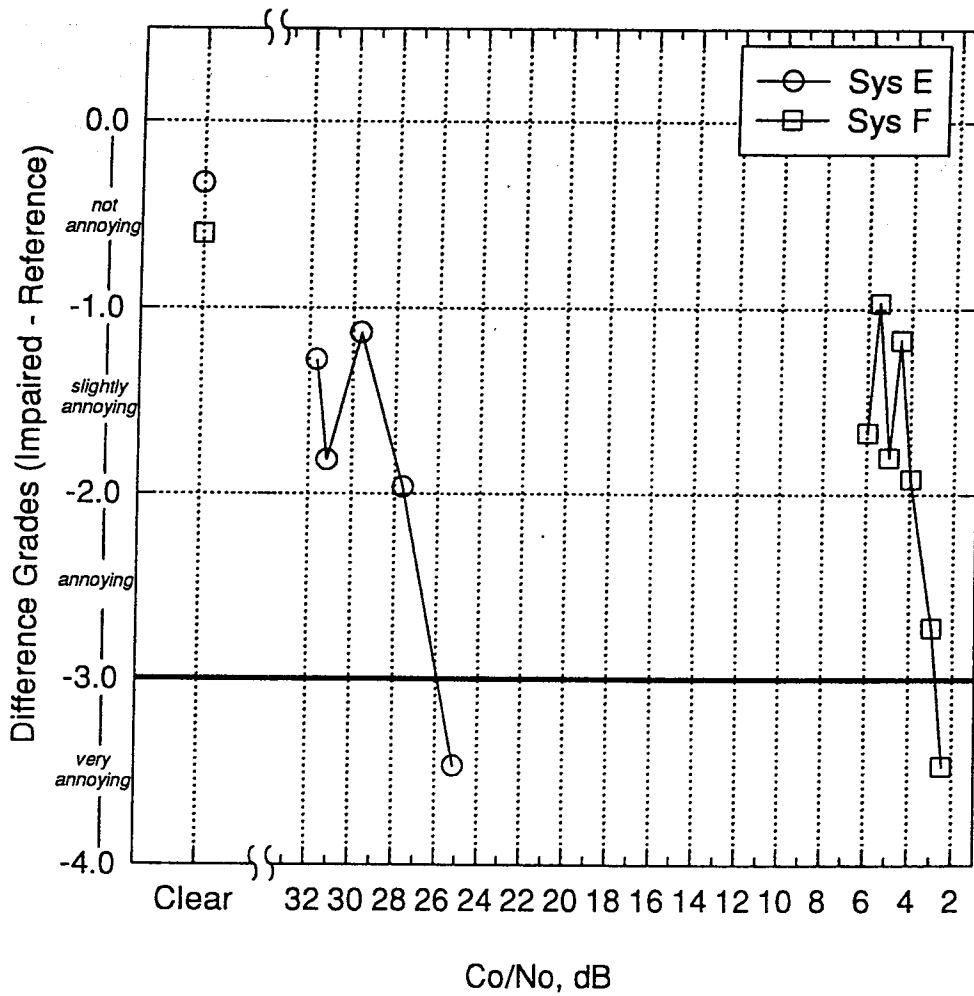


Figure 3.38. Failure characteristic curves for the two DAR systems tested with the rural fast multipath impairment and the Glockenspiel test material.

### Rural Fast Impairment - Soprano

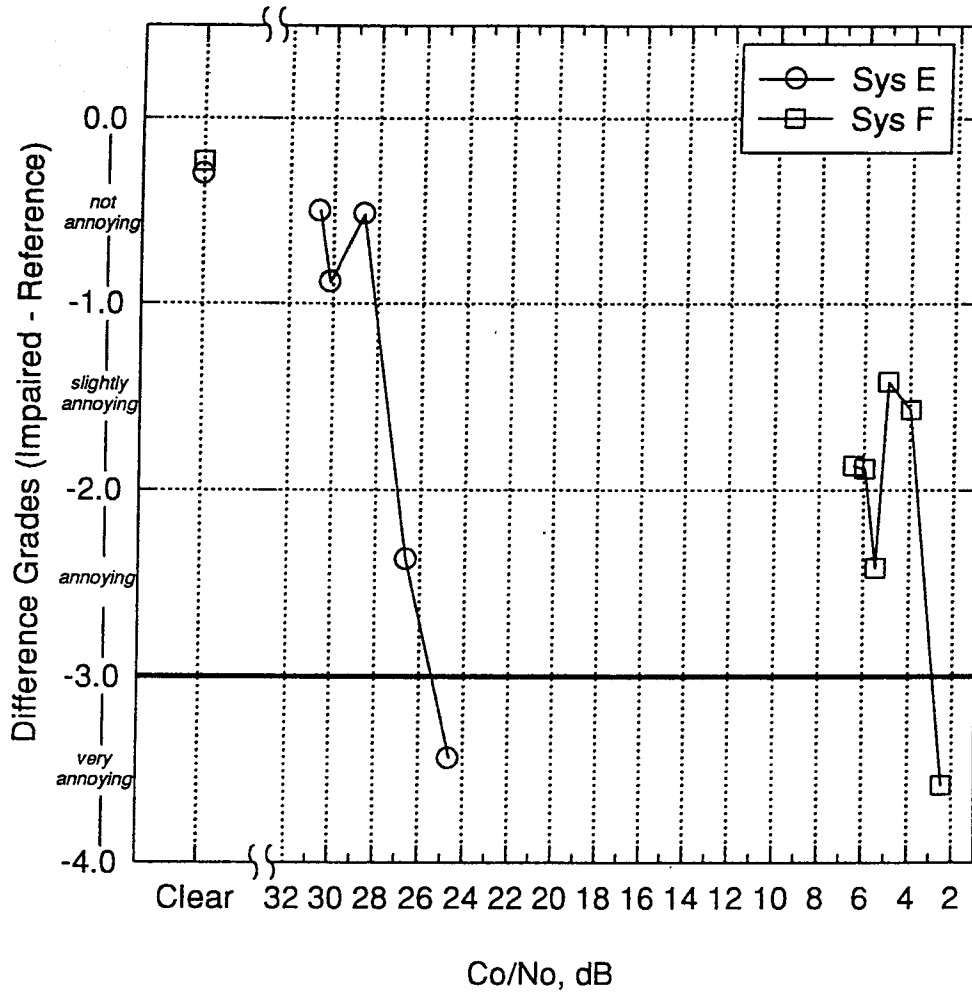


Figure 3.39. Failure characteristic curves for the two DAR systems tested with the rural fast multipath impairment and the Soprano test material.

### Rural Fast Impairment - Clarinet

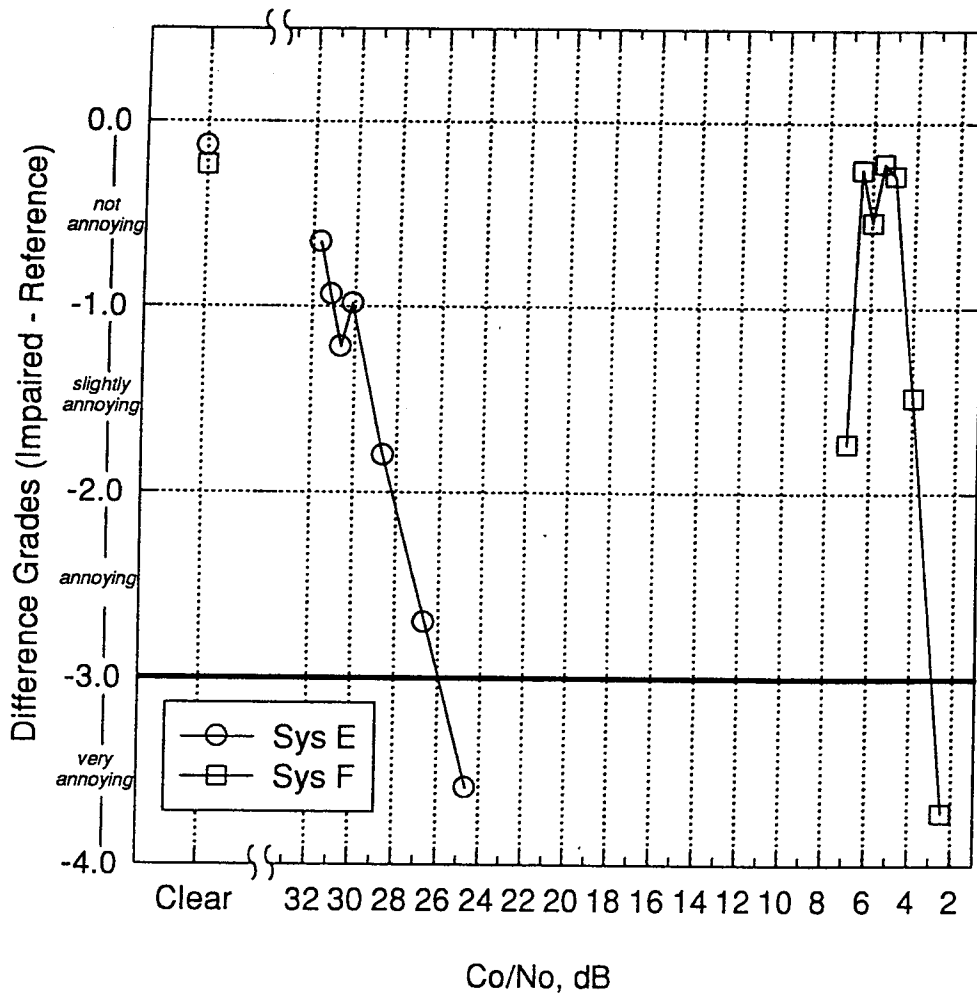


Figure 3.40. Failure characteristic curves for the two DAR systems tested with the rural fast multipath impairment and the Clarinet test material.

### System E - Rural Fast Multipath Impairment

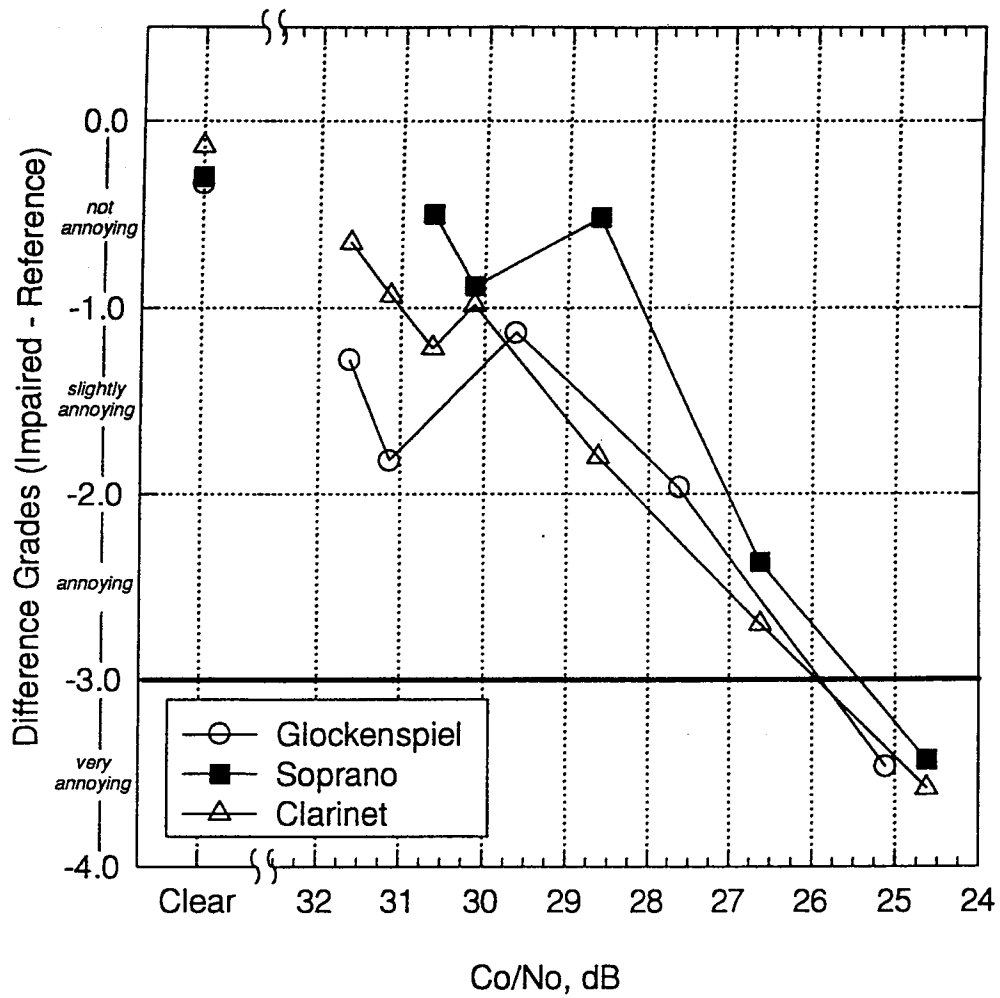


Figure 3.41. Failure characteristic curves for system e for all three test materials with the rural fast multipath impairment.

### System F - Rural Fast Multipath Impairment

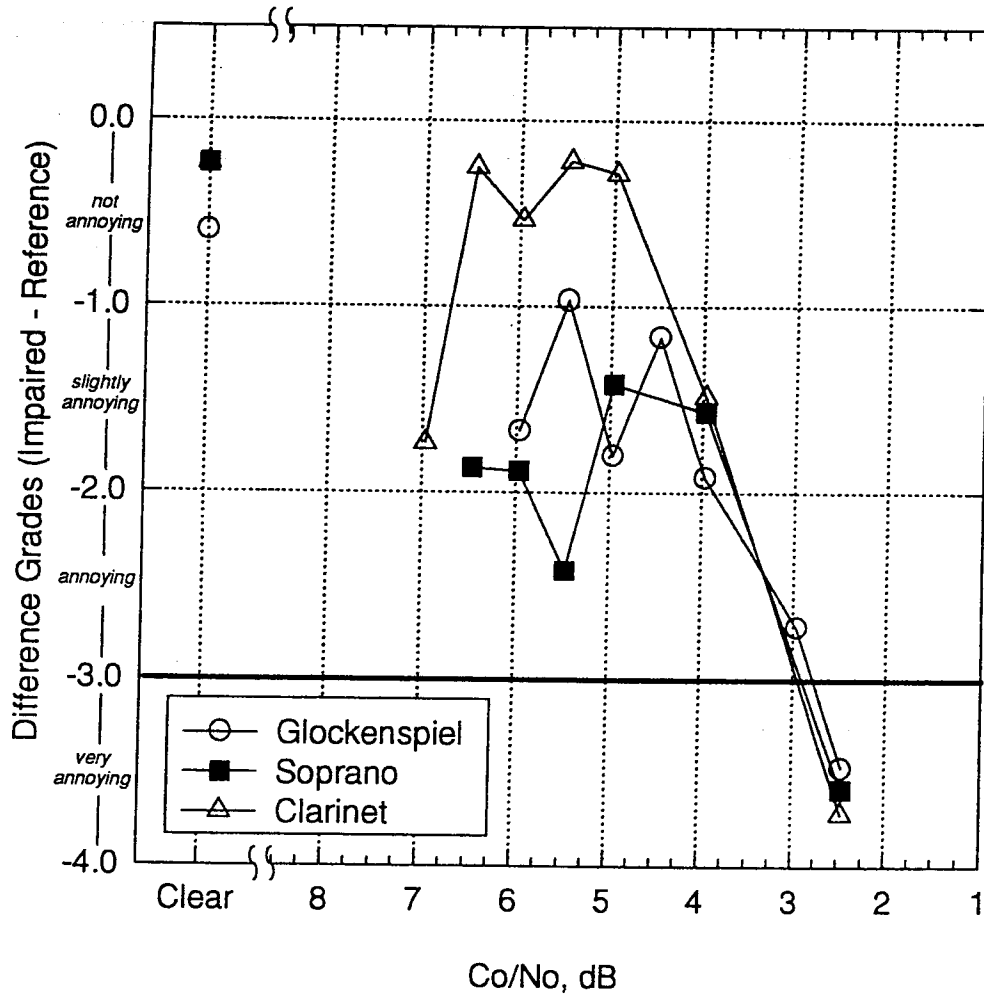


Figure 3.42. Failure characteristic curves for system *f* for all three test materials with the rural fast multipath impairment.

### Obstructed Fast Impairment - Glockenspiel

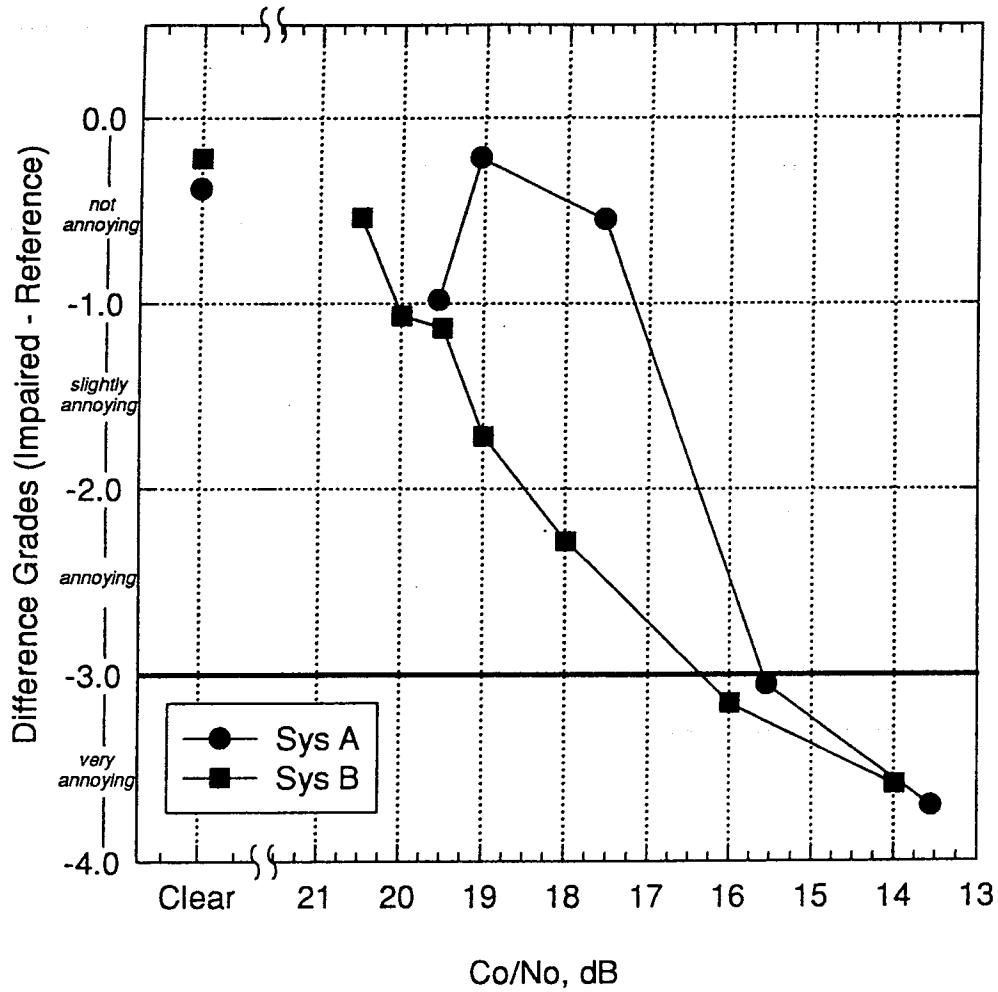


Figure 3.43. Failure characteristic curves of the two DAR systems tested under obstructed fast multipath impairment and the Glockenspiel test material.

### Obstructed Fast Impairment - Soprano

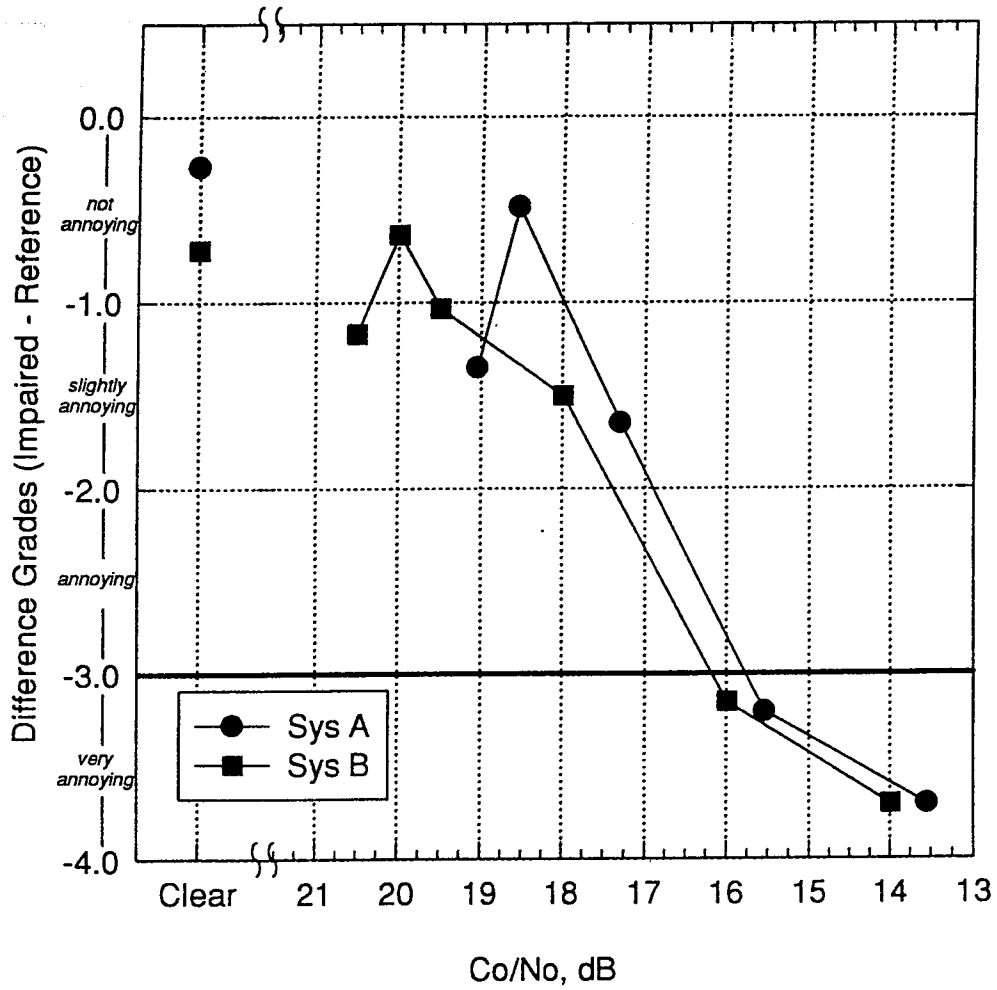


Figure 3.44. Failure characteristic curves of the two DAR systems tested under obstructed fast multipath impairment and the Soprano test material.



### Obstructed Fast Impairment - Clarinet

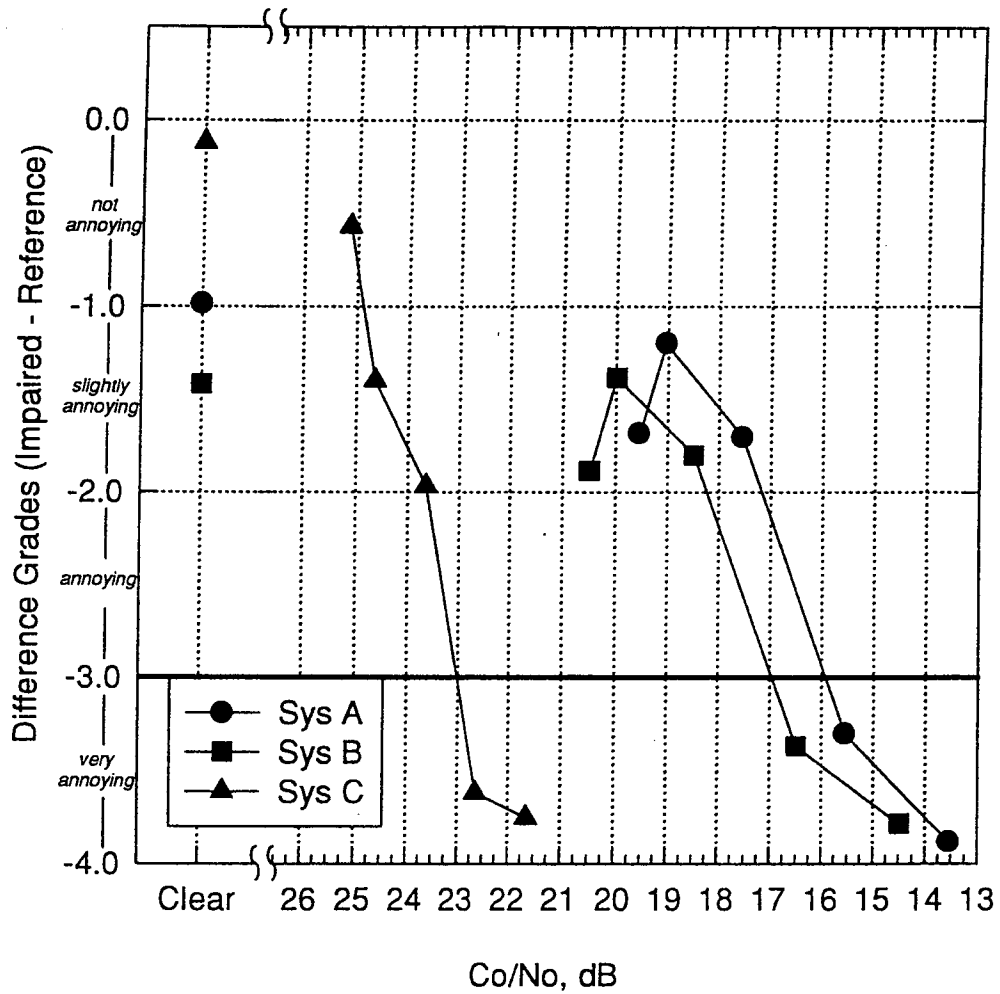


Figure 3.45. Failure characteristic curves of the three DAR systems tested under obstructed fast multipath impairment and the Clarinet test material.

### System A - Obstructed Fast Impairment

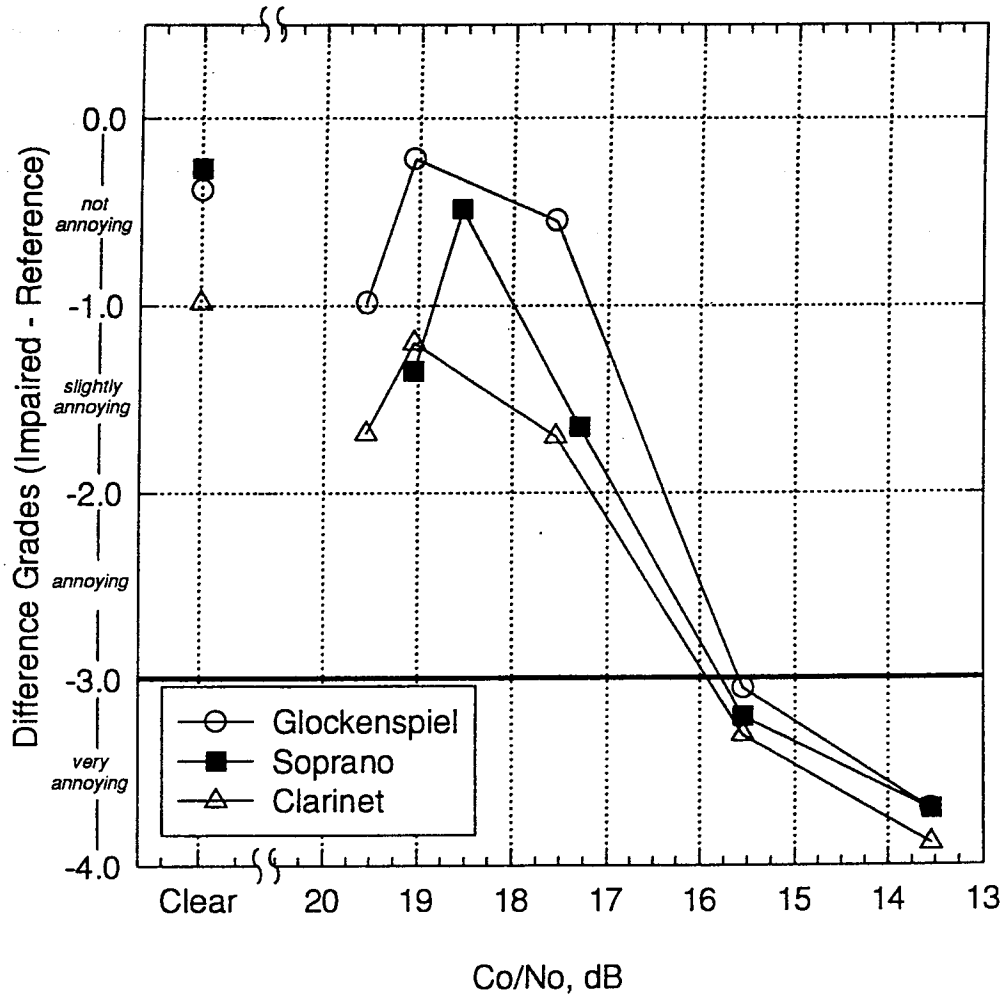


Figure 3.46. Failure characteristic curves for system a for all three test materials with the obstructed fast multipath impairment.

### System B - Obstructed Fast Impairment

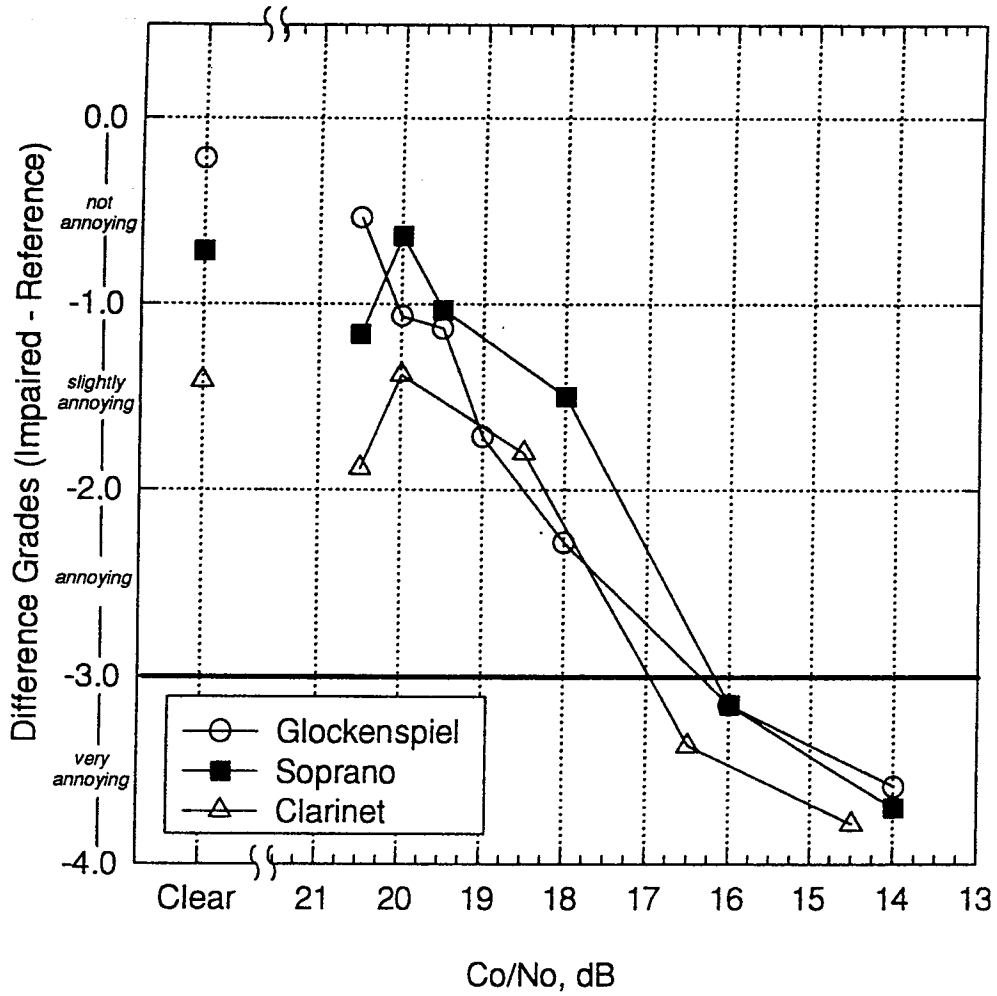


Figure 3.47. Failure characteristic curves for system *b* for all three test materials with the obstructed fast multipath impairment.

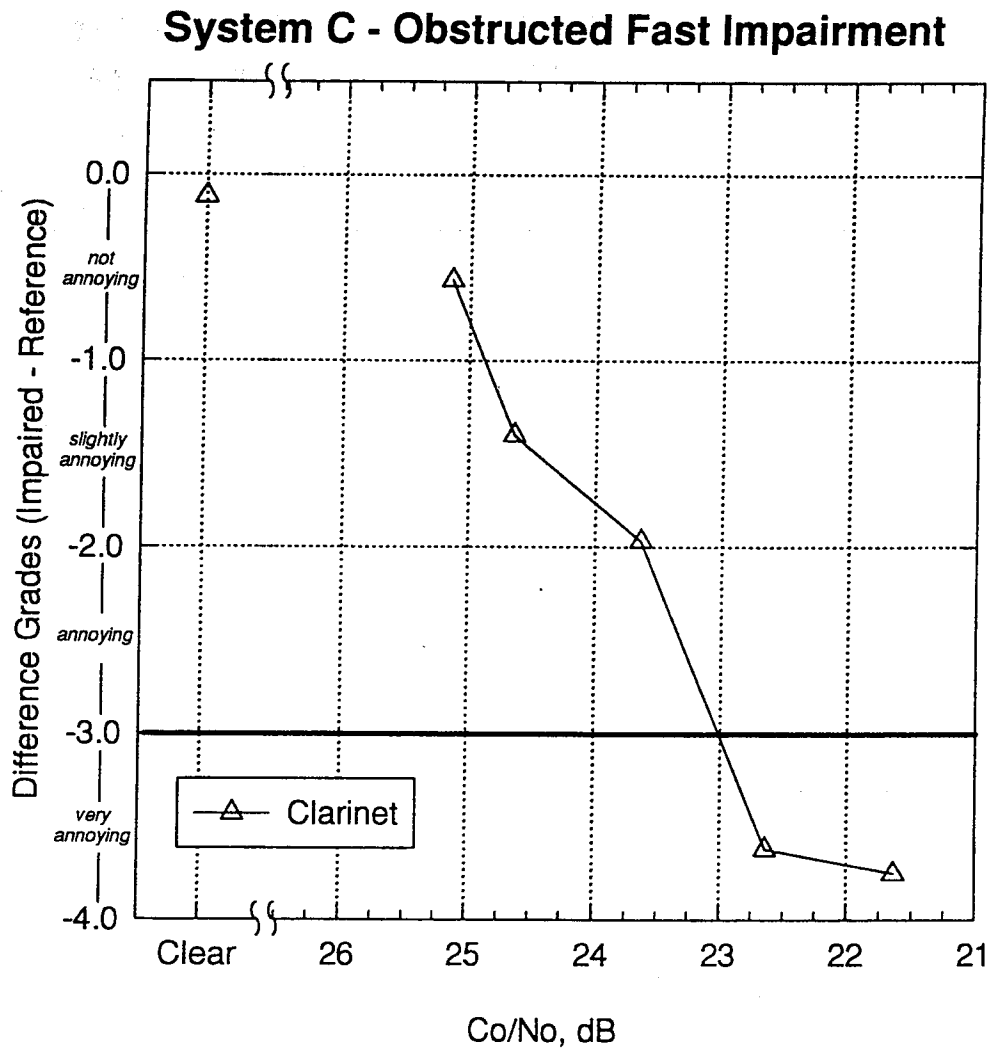


Figure 3.48. Failure characteristic curve for system *c* for the Clarinet test material with the obstructed fast multipath impairment.

### System I - DAR Lower 1st Adjacent to DAR Channel Interference

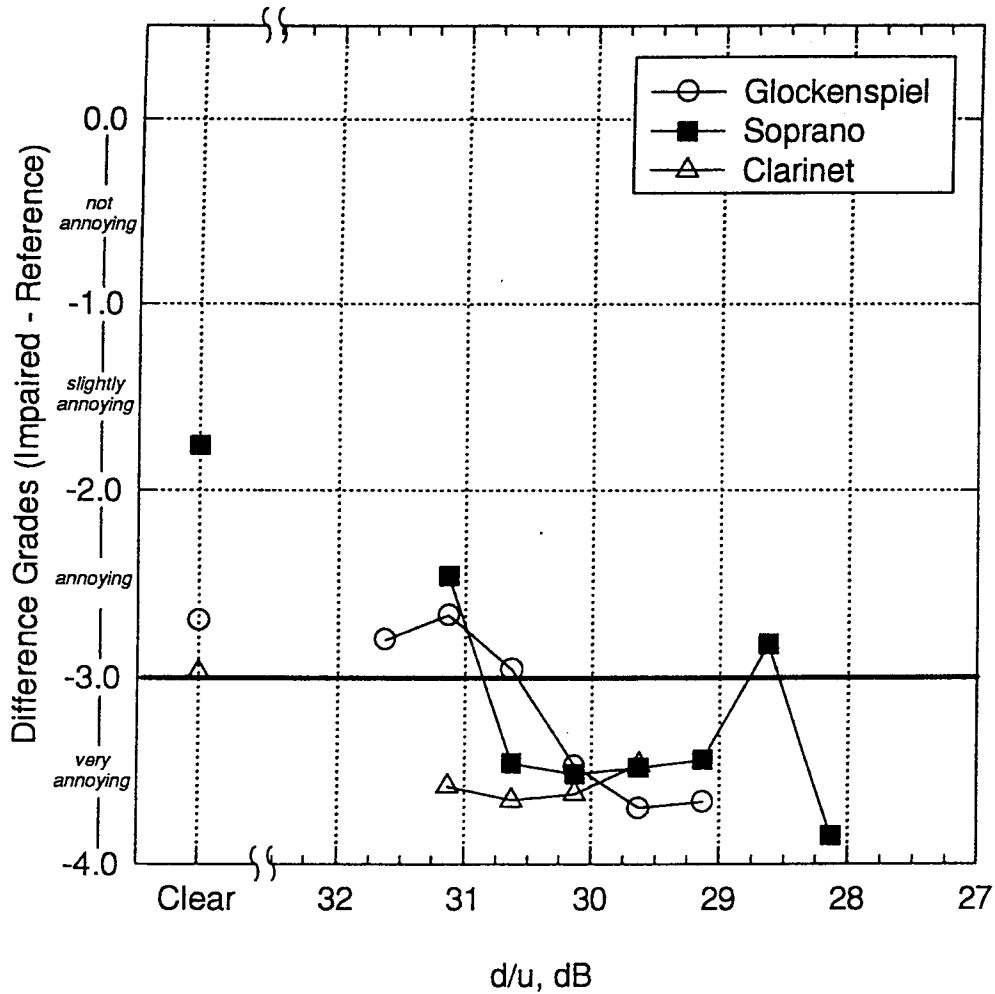


Figure 3.49. Failure characteristic curves for system *i* for all three test materials with the DAR lower 1st adjacent to DAR channel interference. (DAR laboratory was unable to certify digital power levels for this system)

## 4. Subjective Retests of Audio Quality

### 4.1 Background

After the original quality test had been completed at the CRC on 10 proposed DAR systems, system *h*, which was one of those 10 systems, was modified. The present retest was to determine whether or not the modification had altered the perceptual quality of that system.

There are many different ways to approach the question of identity in quality between the modified version (now called system *l*) and the original one (*h*). The fundamental question is whether or not *l* would have performed in the same way in the original quality test as *h* had done. The most direct way to answer the question would be to replicate the original quality experiment, but substituting system *l* for *h*. Such an approach, however, would bring many difficulties. For example, recruiting essentially equivalent subjects, might prove to be time-consuming. And the time needed for full replication of these tests would further delay the testing process and would escalate costs.

At the other extreme of cost and complexity is running an experiment in which subjects would compare *h* and *l*, exclusively and directly. The biggest problem with this approach is that one could never be sure, no matter what comparison methods were used, whether the results could be directly extrapolated to the outcomes of the original experiment. The many necessary differences in methodologies between the original experiment with 10 systems, and one with only 2 systems, would make such a projection questionable, regardless of how *h* and *l* were found to compare.

To minimize both problems - those of cost, and of veracity - a compromise approach was taken. The experiment reported here did follow the basic design of the original experiment, but version *l* was compared with only 4 of the other systems (including system *h*) rather than with all 10 systems. Other departures from the original quality experiment are outlined and discussed in the following "Test Procedures" section.

### 4.2 Test Procedures

No consequential changes in the room and equipment had occurred during the time elapsing between the original quality experiment and this retest. Accordingly, in using the same situation for the retest, it can be assumed that the basic acoustic, and other conditions were the same for the two experiments.

For reasons of efficiency and cost, 8 subjects were used rather than 21 as in the original experiment. To achieve stable results with this smaller number of listeners, exceptionally sensitive listeners were used exclusively, chosen for this characteristic from a pool of subjects who had established their expertise in prior research at the CRC labs and who were available for the present test. Of these 8 subjects, 5 had been in the original EIA quality test.

The same 9 audio materials were used as in the original experiment. A total of 5 systems were examined in the present test, rather than 10 as in the previous test. The 5 systems used were *a*, *d* and *j* as well as systems *h* and *l*. The primary systems of interest were, of course, *h* and *l*. The other 3 were chosen to represent the range of system qualities present in the original study. Thus *a* was from the highest rated Group 1 of the original test, *d* was from Group 3, and *j* came from Group 4, the lowest rated one. System *h*, itself, was from Group 2. The presence of the other systems besides *h* and *l* was an attempt to have a total evaluation context as similar to that of the original experiment as possible, short of a full replication. No matter what the outcome regarding systems *h* and *l*, the validity of that result required that the ratings of the other systems (*a*, *d* and *j*) were commensurate here with those in the original study.

Since the major question was establishing whether or not systems *h* and *l* were equivalent in subjective quality while using only a small group of subjects, an additional means of reducing variance and enhancing sensitivity (besides using exceptionally expert subjects) was implemented. This additional means was to include both of systems *h* and *l* three times, rather than only once, in all conditions. In effect, this was like having 9 systems in the experiment, with 9 observations per system, one per audio material. But of these 81 observations, systems *h* and *l* had 27 each, as though each of them were 3 separate systems, while the other 3 systems (*a*, *d*, and *j*) each had 9.

The subjects, of course, knew nothing about how many systems were in the test. They were simply evaluating items of audio material, with 9 observations for each of the 9 materials. Training was carried out in the same way as for the original test, and took up most of the morning of each subject's first day of the two day experiment.

One important departure from the original experiment was that not all 81 items that were to be rated in the later blind sessions were made available to the subjects during training. Specifically, the *h* and *l* versions of each material were included only once in the training materials, rather than three times as in the grading sessions. And so, there were 45 items in training, rather than the full 81 of the blind grading sessions. This was done so that subjects would not discover the fact of repeated presentations of identical items (the 3 inclusions of both *h* and *l* for each audio material) during training. Such discovery might have led a subject to adopt a grading strategy during the blind trials which anticipated repetitions. A strategy of this sort might detract a subject from listening to, and grading, each audio item as a stand-alone entity. This in turn, would defeat the purpose of the duplications for systems *h* and *l* which was specifically for providing independent repeated ratings.

In all other essential respects, the training and blind testing procedures were the same as in the original experiment. After a morning training session on the first day (done in a group when there was more than one subject scheduled for a two-day experimental time slot), each subject had three blind grading sessions in the afternoon of the first day. The second day consisted of six grading sessions during the morning and afternoon. Each grading session consisted of 9 trials. The 81 presentations were distributed over the sessions so that each of the 9 audio materials always occurred once per session. Within that limitation, the "9" systems (5 systems with 3 repetitions for both *h* and *l*) were distributed equally among the nine sessions to the greatest extent possible. Within sessions the systems were presented as unpredictably ("randomly") as

possible as to their trial-by-trial presentation order. To control for any interaction between events in time (such as fatigue, learning) and the content of sessions, the actual temporal order of the nine programmed sessions was different for each subject.

### 4.3 Results and Discussion

Subjects had been told before training that only a subset of the materials to be graded in blind sessions were to be presented in training. One comment made by almost all of the subjects after the experiment was completed was that the grading sessions were much more difficult than the training sessions had led them to expect. It is reasonable to assume that they found this to be true because  $h$  (as seen in the original experiment) and  $l$  (the modified version of  $h$ ) were very high quality codecs. The three repetitions of each of these codecs in grading sessions would, then, have provided many more materials that were difficult to distinguish from transparency than was the case in the original quality experiment. Since these repetitions were not presented in training (for reasons noted previously), that part of the experiment would seem relatively easy compared to the grading sessions.

Five of the 8 subjects had been in the original quality test, and these 5 listeners also commented that this test was more difficult in an absolute sense than the previous one had been. Some objective evidence for this is that the expertise scores of these subjects (explained in the report of the original quality test) was consistently lower in this retest than in the original study. On average, their expertise scores in the retest were 78% of the magnitudes of those scores in the previous test. From their comments, none of these subjects seemed to suspect that this was a retest which included some of the same systems as before.

Harking back to the relative ease during training compared to blind rating sessions, we might note that the subjects did not themselves volunteer reasons for the unexpected difficulties during blind grading sessions, even though this did seem to surprise them. This might be interpreted to indicate that the subjects did not discover the fact of repetitions of items during the grading sessions. If this is true (and it is likely that if not true, subjects would have mentioned finding repetitions), then, indeed, as intended, the ratings obtained were truly independent from item to item despite repetitions.

The average diffgrades (coded minus hidden reference grades) for each of the 5 systems at each of the nine audio materials are presented in Fig. 4.1. The actual numerical diffgrades plotted in the figure are shown in Table 4.1 below the figure. Since the differences between the averages for the various audio materials are of little interest here, the bottom row of this table shows the difference between  $h$  and  $l$  at each material instead of (as would be more usual in a table of this type) the average diffgrade per material.

It can be seen (in the figure and table) that  $h$  and  $l$  both received highly similar grades across all audio materials. Statistical analysis of the total data set (ANOVA, followed by  $t$ -tests, all  $p$ 's <



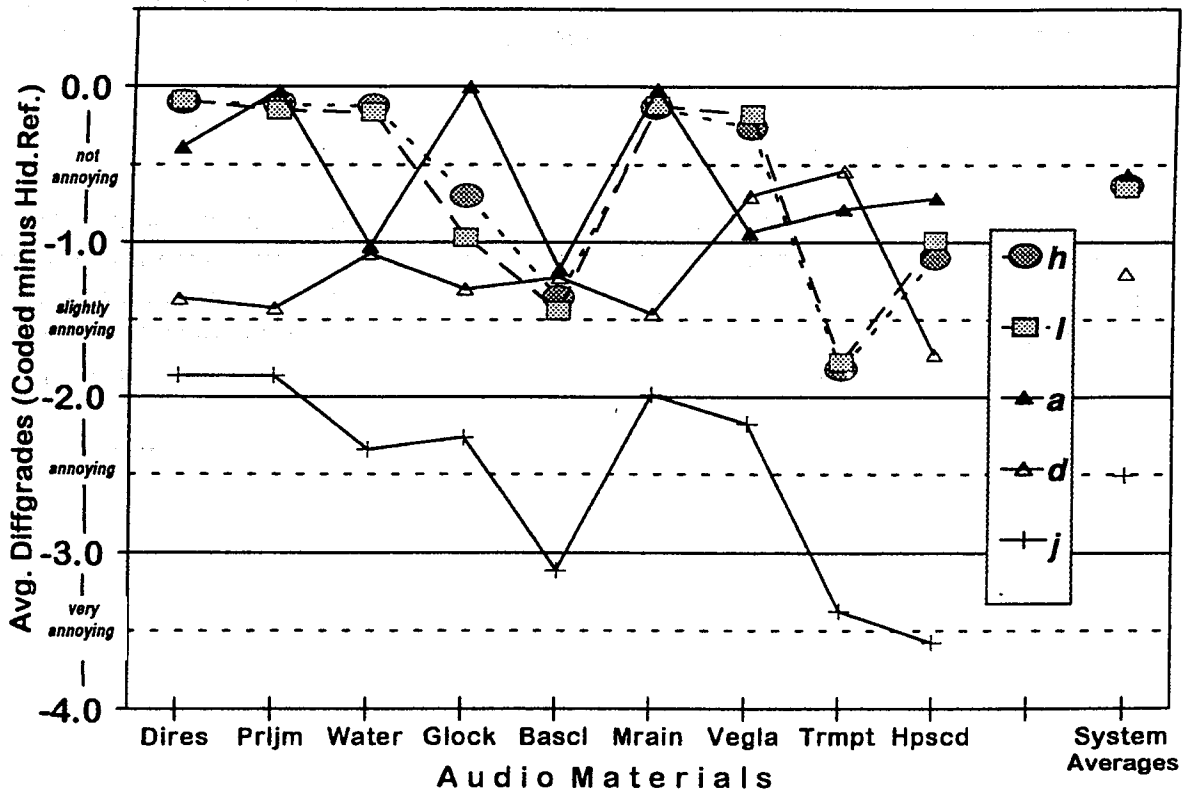


Fig. 4.1 Average diffgrades of the 8 subjects in the quality retest for all the systems in that experiment. Each audio material was presented 3 times for systems *h* and *l*, and once for each of the other 3 systems (*a*, *d*, and *j*) in the 81 trials of the retest.

	Dires	Prijm	Water	Glock	Bascl	Mrain	Vegla	Trmpt	Hpsc		SysAvg
<i>h</i>	-0.10	-0.11	-0.13	-0.70	-1.36	-0.13	-0.26	-1.82	-1.10		-0.64
<i>l</i>	-0.09	-0.15	-0.16	-0.97	-1.44	-0.12	-0.18	-1.78	-0.99		-0.65
<i>a</i>	-0.39	-0.03	-1.03	0.00	-1.18	-0.01	-0.94	-0.79	-0.71		-0.56
<i>d</i>	-1.36	-1.43	-1.08	-1.30	-1.23	-1.46	-0.70	-0.54	-1.73		-1.20
<i>j</i>	-1.86	-1.86	-2.34	-2.26	-3.11	-1.99	-2.18	-3.38	-3.58		-2.51
<i>h-l</i>	-0.01	0.04	0.04	0.27	0.08	-0.01	-0.09	-0.04	-0.11		0.01

Table 4.1 Average diffgrades for each system at each audio material and for overall system averages. The bottom row shows the algebraic difference between the diffgrades for systems *h* and *l* in the top two data rows.

.05) shows that, within any material, differences between any two systems must exceed 0.68 of a grade in absolute magnitude in order for those systems to be considered statistically different for that material. The bottom row of the table shows that the differences between  $h$  and  $l$  are less than 1/10th of a grade except for the Glockenspiel (0.27) and for the Harpsichord (0.11). All of these differences are considerably smaller than the 0.68 needed for the verdict of a reliable difference.

The relatively large 0.68 needed for significance (compared to 0.45) in the original study) is due to the smaller number of subjects used in this retest (8 versus 21 previously). The exceedingly small magnitudes of differences between  $h$  and  $l$  across materials here (mostly less than 1/10th of a grade) must be attributed to the exceptional expertise of these subjects as well as to the three repeated observations for those two systems.

It must be concluded that  $h$  and  $l$  cannot be considered to be different in quality across all 9 audio materials under the conditions of the experiment.

The results for the other systems ( $a$ ,  $d$ , and  $j$ ) are important too for several reasons. First, the results must show that the experiment was not insensitive to revealing differences. If it were, then no valid final judgment could be made about  $h$  and  $l$  since their apparent lack of difference might be explained as the outcome of a poorly executed experiment.

Examining the table and figure shows that many significant differences did emerge in the experiment. For example, looking at the "water" and "trmp" materials, the differences between  $h$  (and  $l$ ) on the one hand, and both  $a$  and  $d$  on the other, obviously exceed the criterion magnitude of 0.68. And, of course, almost all comparisons between system  $j$  and any of the other systems (with only 3 exceptions involving system  $d$ ) are larger than 0.68. Hence it cannot be said that the experiment was too crude for reliable differences to emerge. And so the lack of difference between  $h$  and  $l$  is not due to experimental insensitivity.

Second, and quite important, the results for these other three systems must relate in some sensible way to the outcomes for those same systems in the original quality experiment. If such a relationship was obscure, then it might be argued that despite the apparent identity of  $h$  and  $l$  here, the conditions of this study were so different from the original one that one cannot assume that  $l$  would have performed identically to  $h$  if it had been in the original study instead of  $h$ .

A visual comparison of Fig. 4.1 here, with Fig. 2.2a (of chapter 2) of the original quality test results shows that the general pattern of results for each of the comparable systems are strikingly similar. This is most immediately obvious with system  $j$  because it occupies a distinct region towards the bottom of both figures. But it is also true of the other systems. A point-by-point comparison shows that the up-and-down pattern across the audio materials is generally the same for each system in both figures.

An even stronger case can be made for the identity of outcomes between the two experiments. Recalling that 5 of the subjects in the present retest were also in the original quality assessment, the results for those subjects alone can be isolated and compared. These comparison will not be

presented in detail here. The reader can be assured, however, that they show virtually identical numerical results with no statistical differences between them.

Since the present results were obtained with exceptionally expert subjects, as compared with sufficiently expert ones in the original test, the implication is that whatever small differences are seen between the two experiments is due to the larger range of individual subject expertise in the first test. In fact, consistent with findings about the grades of subjects with exceptional expertise in many previous studies at the CRC, the outcomes here show somewhat harsher judgments (lower diffgrades) than in the first quality test. In other words, we are confident that had listeners with fully comparable expertise been used in this retest as in the first test (and if the number of subjects had been as large), the results would have been quite identical. This retest, then, provides excellent confirmation of the original study.

#### 4.4 Conclusions

We conclude that  $h$  and  $l$  can be considered to be identical in quality. This finding is not due to insensitivity of the retest to finding differences. The outcomes of this retest strongly confirm the quality outcomes of the original experiment.

We conclude that if  $l$  had been in the original quality experiment instead of  $h$ , the outcomes of that first experiment would have been no different than they actually were.

## 5. Subjective Retests of Transmission Impairments

### 5.1 Introduction

After the original transmission impairment tests had been completed at the CRC on the proposed DAR systems, two systems, *e* and *h* were modified. The modified version of system *e* was named system *k* while the modified version of system *h* was named system *l*. As such, the present retests were conducted to determine the failure characteristics for the modified systems. These retests were only conducted for the gaussian noise and co-channel transmission impairments.

In order to minimize costs and time requirements it was decided that the retests would be limited to determining the threshold of audibility and the point of failure only. Furthermore, these points (TOA and POF) would be determined by means of expert observations and commentary, EO&C. The procedures for determining TOA and POF are described in greater detail in the following section.

### 5.2 Test Procedures

#### 5.2.1 Test Materials

The three critical audio materials used in the impairment retests were the same as the ones used in the original impairment tests and were carefully selected by the staff at the NASA Lewis Research Laboratory because they were particularly sensitive to revealing artifacts resulting from transmission impairments. These materials are described in Table 3.1 in section 3.2.1 of this document.

#### 5.2.2 Impairment Levels

As mentioned earlier, only the TOA and POF points were to be determined in the transmission impairment retests. Nonetheless, a complete set of recordings were provided by the staff at the NASA Lewis Research Laboratory for the two systems *k* and *l* as described below.

Each of the three audio test materials was recorded, for the two DAR systems (*k* and *l*) and the two types of impairment (gaussian noise and co-channel interference), at the following levels of impairment: CC, TOA<sub>1</sub>, TOA<sub>2</sub>, ..., TOA<sub>N</sub>, S<sub>1</sub>, S<sub>2</sub>, ..., S<sub>M</sub>, POF where:

CC = coded audio in a Clear Channel

TOA<sub>1</sub>, TOA<sub>2</sub>, ..., TOA<sub>N</sub> = N stimuli (three or more) in the close neighborhood of the approximate TOA (threshold of audibility)

S<sub>1</sub>, S<sub>2</sub>, ..., S<sub>M</sub> = M intermediate levels of impairments

POF = point of failure

For each of the above impairment levels, the  $C_o/N_o$  ratio for the gaussian noise impairment, or D/U for co-channel interference, was noted. The CC level was recorded at a high value of  $C_o/N_o$  (or D/U) so that transmission errors could be considered as negligible, hence the label "clear channel" given to that level. The increments in the  $C_o/N_o$  (or D/U) ratios at which the other impairment levels were recorded varied from 0.25 dB, for those situations where the DAR system failed abruptly, up to 1.0 dB for those situations where the system failed more gracefully.

### 5.2.3 Subjective test procedures

Two separate experiments were performed for each of the two systems. The purpose of the first test was to determine the TOA point for each system/impairment-type/audio-material combination. The procedure for determining the TOA points was the same as the one used in the original impairment tests. A detailed description of this procedure is given in section 3.2.3 of this document.

The second test was designed to provide an estimate of the POF points for the two systems. The procedure used for this test is described below.

#### Point of Failure

The purpose of the second experiment was to derive an estimate of the POF points for the two systems  $k$  and  $l$  as opposed to determining full failure characteristic curves for each system under various impairment conditions as was done in chapter 3. However, it was desirable to try to find some way of relating the results of the retests to the results of the original tests. To do this, a series of "matching experiments" were conducted wherein expert listeners identified the impairment levels required for the modified systems  $k$  and  $l$  to be perceptually equivalent to the POF's found for the original systems  $e$  and  $h$  in the original tests. The expert listeners consisted of three CRC staff members and were the same listeners who took part in the TOA EO&C tests.

In the figures of chapter 3, it can be seen that for many of the failure characteristic curves there are several points which fall in the POF range. (note that a diffgrade of -3 or less is considered to be POF). Therefore, in the retests, a perceptually equivalent impairment level for the modified systems was found for each of the points (for the old systems) which fell in the POF range.

As stated above, to determine the POF points for the two modified systems, a series of matching experiments were conducted at the CRC. In these tests the CRC staff members listened to the POF points for the original systems ( $e$  and  $h$ ) and found the impairment levels for the modified systems ( $k$  and  $l$ ) which were perceptually equivalent. This process was repeated for the two systems, the three test materials, and the two types of transmission impairments (gaussian noise and co-channel interference). Furthermore, this process was done for all points for the original systems which fell in the POF range.

In some instances, it was found that the perceptually equivalent impairment level fell between the impairment levels recorded at the NASA Lewis Research Center and provided to the CRC. In

these instances, a value of  $C_0/N_0$  (or D/U) which was half-way between the levels provided was used even though this impairment level was not actually auditioned by the CRC listening team.

## 5.3 Test Results

### 5.3.1 Gaussian noise

The results of the retest for the gaussian noise impairment are given in Table 5.1 for systems  $e$  and  $k$ . The table provides the  $C_0/N_0$  ratios at TOA and POF, as well as the failure margin (FM). This is done for the three critical audio materials, Glockenspiel, Soprano, and Clarinet. It should be noted that the values for system  $e$  are taken from the original transmission impairment tests described in chapter 3 of this document. The bottom row of the table shows the *differences* between the original system and the modified system (i.e.  $e-k$ ).

The results given in Table 5.1 show that the differences between the original system and the modified system are very small. None of the differences are greater than 0.72 dB. The largest differences tend to occur at TOA while the differences in the failure margins are extremely small (no more than 0.25 dB).

As stated earlier, the matching experiments were conducted for all points which fell in the POF range in the original transmission impairment tests. The POF points given in the table below represent the highest  $C_0/N_0$  ratio which fell into the POF range. The other POF points are omitted for clarity purposes since they do not alter the basic conclusions.

DAR System	Glockenspiel			Soprano			Clarinet		
	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)
$e^*$	10.76	9.76	1.0	10.51	9.51	1.0	10.76	9.76	1.0
$k$	10.29	9.54	0.75	10.04	9.04	1.0	10.04	9.29	0.75
<i>difference</i>	<b>0.47</b>	<b>0.22</b>	<b>0.25</b>	<b>0.47</b>	<b>0.47</b>	<b>0.0</b>	<b>0.72</b>	<b>0.47</b>	<b>0.25</b>

\* These values were determined in the original transmission impairment tests described in chapter 3.

Table 5.1 Comparison of the  $C_0/N_0$  ratios at TOA and POF and failure margin (FM) for the gaussian noise impairment for systems  $e$  and  $k$ .

The results of the retest for the gaussian noise impairment are given in Table 5.2 for systems  $h$  and  $l$ . The form of the table is identical to that of Table 5.1. It should be noted that the values for system  $h$  are taken from the original transmission impairment tests described in chapter 3 of this document.

The results given in Table 5.2 show that the differences between system *h* and system *l* are small. They are, however, slightly larger than the differences shown in Table 5.1. None of the differences is greater than 1.57 dB. Again, the largest differences tend to occur at TOA. The difference in the failure margin is constant at -0.5 dB. The negative value for the failure margins indicate that they are 0.5 dB larger for the original system than for the modified system.

DAR System	Glockenspiel			Soprano			Clarinet		
	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)
<i>h</i> *	10.51	9.01	1.5	9.51	8.51	1.0	10.01	9.51	0.5
<i>l</i>	11.33	9.33	2.0	10.83	9.33	1.5	11.58	10.58	1.0
<i>difference</i>	-0.82	-0.32	-0.5	-1.32	-0.82	-0.5	-1.57	-1.07	-0.5

\* These values were determined in the original transmission impairment tests described in chapter 3.

Table 5.2 Comparison of the  $C/N_0$  ratios at TOA and POF and failure margin (FM) for the gaussian noise impairment for systems *h* and *l*.

As stated earlier, the matching experiments were conducted for all points which fell in the POF range in the original transmission impairment tests. The POF points given in the table represent the highest  $C/N_0$  ratio which fell into the POF range. The other POF points are omitted for clarity purposes since they do not alter the conclusions.

### 5.3.2 Co-channel interference

The results of the retest for the co-channel interference are given in Table 5.3 for systems *e* and *k*. The table provides the D/U ratios at TOA and POF, as well as the failure margin (FM). This is done for the three critical audio materials; Glockenspiel, Soprano, and Clarinet. It should be noted that the values for system *e* are taken from the original transmission impairment tests described in chapter 3 of this document. The bottom row of the table shows the *differences* between the original system and the modified system (i.e. *e-k*).

The results given in Table 5.3 show that the differences between the original system and the modified system are quite small. None of the differences are greater than 1.11 dB. As was the case for the gaussian noise impairment condition, the largest differences tend to occur at TOA. The differences in the failure margins are no greater than 0.5 dB for the three audio source materials.

As stated earlier, the matching experiments were conducted for all points which fell in the POF range in the original transmission impairment tests. The POF points given in the table below represent the highest D/U ratio which fell into the POF range. The other POF points are omitted for clarity purposes since they do not alter the basic conclusions.

DAR System	Glockenspiel			Soprano			Clarinet		
	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)
<i>e*</i>	11.12	9.62	1.5	10.87	9.87	1.0	10.87	9.87	1.0
<i>k</i>	10.26	9.26	1.0	9.76	9.14	0.62	10.01	9.01	1.0
<i>difference</i>	<b>0.86</b>	<b>0.36</b>	<b>0.5</b>	<b>1.11</b>	<b>0.73</b>	<b>0.38</b>	<b>0.86</b>	<b>0.86</b>	<b>0.0</b>

\* These values were determined in the original transmission impairment tests described in chapter 3.

Table 5.3 Comparison of the D/U ratios at TOA and POF and failure margin (FM) for the co-channel interference for systems *e* and *k*.

The results of the retest for the co-channel interference are given in Table 5.4 for systems *h* and *l*. The format of the table is identical to Table 5.3. It should again be noted that the values for system *h* are taken from the original transmission impairment tests described in chapter 3 of this document.

The results given in Table 5.4 show that the differences between the original system and the modified system are very small. None of the differences is greater than 0.83 dB. The largest difference in the failure margin is -0.5 dB which occurs for the Glockenspiel audio source material.

As stated earlier, the matching experiments were conducted for all points which fell in the POF range in the original transmission impairment tests. The POF points given in the table represent the highest D/U ratio which fell into the POF range. The other POF points are omitted for clarity purposes since they do not alter the conclusions.

DAR System	Glockenspiel			Soprano			Clarinet		
	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)	TOA (dB)	POF (dB)	FM (dB)
<i>h*</i>	11.37	8.87	2.5	10.87	7.87	3.0	10.87	9.87	1.0
<i>l</i>	11.04	8.04	3.0	10.54	7.54	3.0	11.04	10.04	1.0
<i>difference</i>	<b>0.33</b>	<b>0.83</b>	<b>-0.5</b>	<b>0.33</b>	<b>0.33</b>	<b>0.0</b>	<b>-0.17</b>	<b>-0.17</b>	<b>0.0</b>

\* These values were determined in the original transmission impairment tests described in chapter 3.

Table 5.4 Comparison of the D/U ratios at TOA and POF and failure margin (FM) for the co-channel interference for systems *h* and *l*.



## 6. Facilities

### 6.1 Playback system

The playback system used during the subjective tests is shown in Fig. 6.1. The particular brand name and model of equipment used in this set-up is described in Table 6.1 below.

The heart of the playback system is a custom audio workstation. It is implemented by means of a 486 personal computer (PC) equipped with a dual DSP56001 processor board and custom software. On two daughter cards, the DSP board also hosts an AES/EBU interface and a SCSI bus interface. Both the DSP board and software are developed by MPR Teltech Ltd, Burnaby, BC, Canada [4]. The workstation is also equipped with a large SCSI disk drive (1.75 Gbyte), a color VGA monitor and serial mouse. It operates in the Windows 3.1 environment.

Qty	Description
1	486-66 PC compatible with 1.75 Gbytes SCSI hard-disk and MPT Teltech Dual DSP56001 Processor Card (DSPC)
1	SVGA color video monitor
1	Spectral model ADDA 2218 D/A converter
1	Klark Teknik model DN410 parametric equalizer
1	Bryston model BP-5 professional stereo preamplifier
1	Bryston model 3B PRO professional stereo power amplifier
1	Bryston model 4B PRO professional stereo power amplifier
1	State-of-the-Art Elektronik model AAX2-2-750 active crossover system
2	State-of-the-Art Elektronik model CF 750 monitor loudspeakers
1	Stax model SRM-1/MK-2 professional headphone driver
1	Stax model ED-1 diffused field equalizer
3	Stax model Lambda Pro headphones
2	Sony model PCM-2500 DAT recorder/player
1	Panasonic SV-3700 DAT recorder/player

Table 6.1 List of equipment

The software driving the workstation consists in four different windows applications. Each of them corresponds to one of the four operations required for the preparation and the presentation of the audio materials for the listening tests. These operations are:

1. Recording of the audio materials on the audio workstation (*Record* application)
2. Synchronization of the audio files (*TimeSync* application)
3. Building of the session files (*Notepad* application)
4. Presentation of the audio materials (*ABC* application)

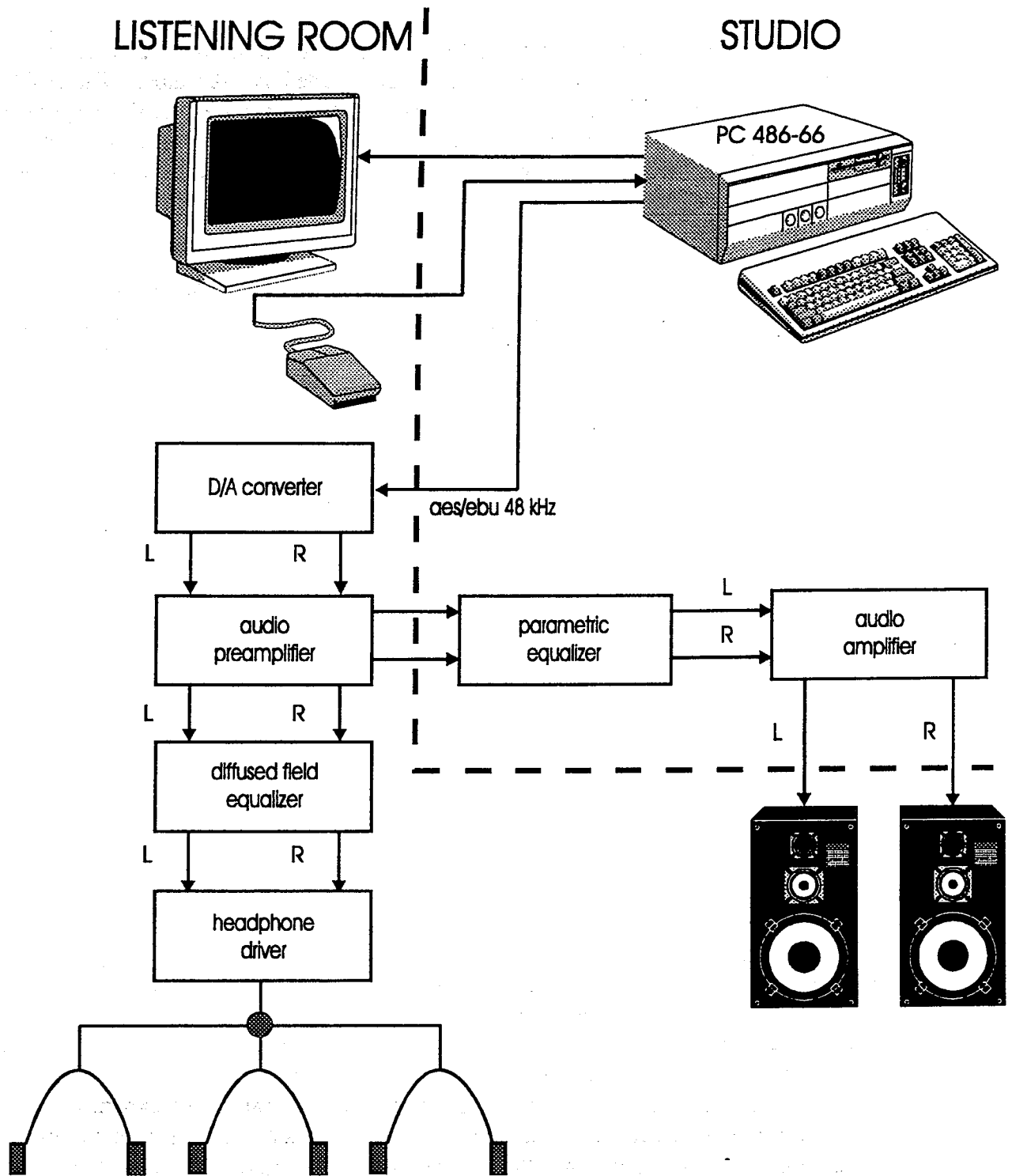


Figure 6.1 Equipment set-up

---

- **Record application**

Audio test materials were processed through the DAR systems under test at the NASA Lewis Research Centre, Cleveland, Ohio, USA and the output was recorded on DAT tapes. The DAT tapes were then shipped to the CRC where each of the audio materials required for the subjective tests (audio quality and transmission impairments) was individually transferred digitally to hard disk via the AES/EBU serial interface and stored as separate audio files with individual filenames.

- **TimeSync application**

Once on hard disk, the audio files containing the reference and all processed versions of each audio material are then precisely time aligned with *TimeSync*, a special software system developed at the CRC. This time synchronization process is required in order to ensure seamless switching between reference and processed versions of audio materials during the presentation.

- **Notepad application**

The next step is to build "session files". A trial consists of an A-B-C presentation of a given audio material and a session consists of one or more consecutive trials. A session file contains the list of audio files to be assigned to A, B and C (in the triple-stimulus A-B-C presentation) for each trial. The session file is built with any text editor that is capable of generating ASCII files, such as *Notepad* in the Windows environment.

The experimenter can subsequently re-order the audio files within a trial, or trials within a session, simply by editing a session file. Modifications are done in a few minutes compared to the many hours that an equivalent re-ordering would take to do on a DAT based playback system. This is a powerful feature of this disk-based playback system.

- **ABC application**

The screen used by the listener during the blind testing phase is shown in Fig. 6.2. During the training phase, there are additional buttons not seen in this figure, namely a *Session* button that is used to select a particular session file, and a *Trials* button that allows the listener to select a particular trial within the selected training session. These buttons are removed during the blind rating phase. This is to prevent the subject from accidentally changing the programmed session with the *Session* button, since this must remain under control of the experimental design. Removal of the *Trial* button prevents the listener from moving backwards through the trials during blind rating since trial order, as is true for session, must also remain under experimental design control. During blind rating, the listener can only move forwards (after spending as much time as he or she wants on a given trial) and this is done by using the *Next* button seen in Fig. 6.2. A pause button is provided, and by using the loopback button seen in Fig. 6.2, the listener can control whether playback of the trial stops at the end of each material, or whether looping is continuous. A small comment identifying the audio material being presented is also shown. This comment is added, as explained previously, in the session file for each trial.



Figure 6.2 Video screen used by the listeners during blind rating sessions

Switching among A, B or C is done using anyone of the following two methods: (a) by clicking with the mouse pointer on the “A”, “B” or “C” button displayed on the video screen of Fig. 6.2 or, (b) by pressing respectively the left, centre or right button of the three-button mouse (this method requires the mouse pointer to be first dragged outside any button area)

A “zoom” tool allows the listeners to listen to a smaller subsegment of the audio material they are asked to compare and assess in a given trial. The start and end time of the subsegment can be set anywhere within the audio material using the two horizontal scroll bars shown at the bottom right of Fig. 6.2. The top bar changes the start time of the subsegment while the bottom bar adjusts the end time. The continuous looping, if activated, is performed on the subsegment of the material. With this feature, listeners can focus more closely on a specific section of any material that appears to reveal suspected distortions.

Since the files are precisely time-aligned and since cross-fading is used, the switching is truly seamless in that there are no audible cues generated, and very exact continuity between the audio

materials triggered by the buttons is maintained. When ready for rating, the listener clicks on the *Next* button. This action invokes a new scoring screen (not shown). Scores for B and C are entered by sliding a cursor to the desired impairment level in the corresponding scroll bar. After the *OK* button is clicked, another small window (not shown) appears requesting the listener to confirm his/her selected ratings.

One additional distinctive feature of the disk-based system is that any selected segment (A, B or C) is output digitally via the same single AES/EBU interface, and is fed through the same D/A convertor as any other selected segment. And so, all versions of the audio materials to be compared are presented through the very same pieces of hardware. This eliminates any possible contribution by hardware to differences detected by listeners among the different versions.

## 6.2 Listening Room

The specifications of the CRC listening room are summarized in the following three figures which show respectively the room layout (Fig. 6.3), the background noise levels (Fig. 6.4) and the reverberation time curve (Fig. 6.5). Both the background noise levels and the reverberation time curve of the room comply with the requirements of ITU-R Recommendation BS.1116 [1].

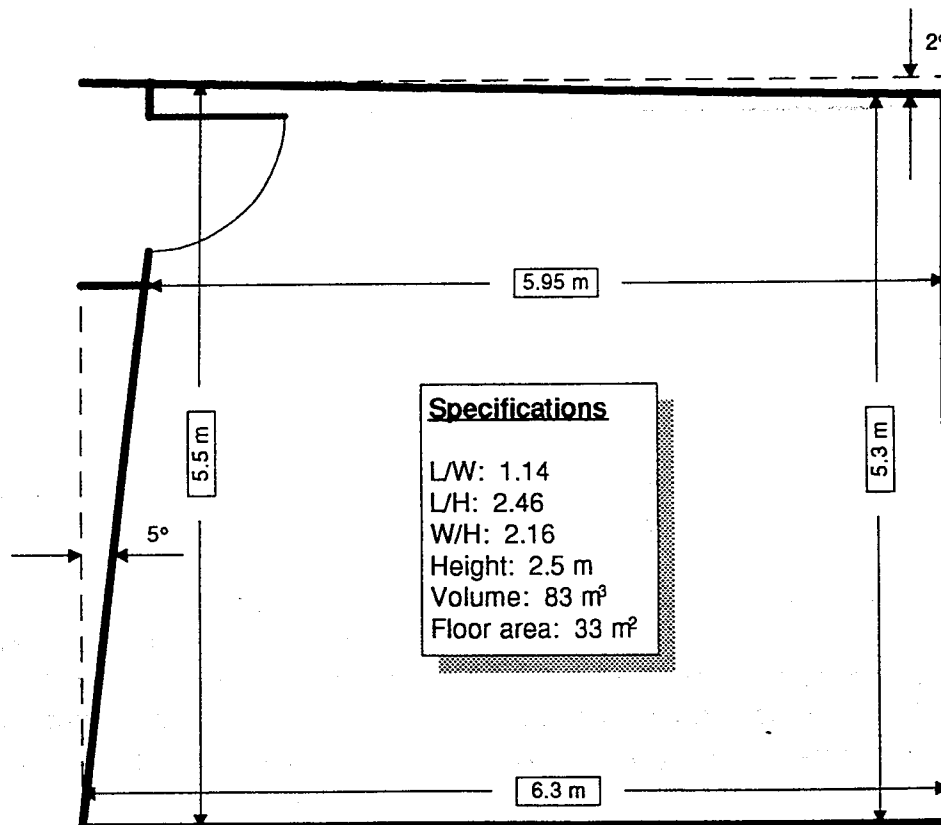


Fig. 6.3 CRC listening room layout

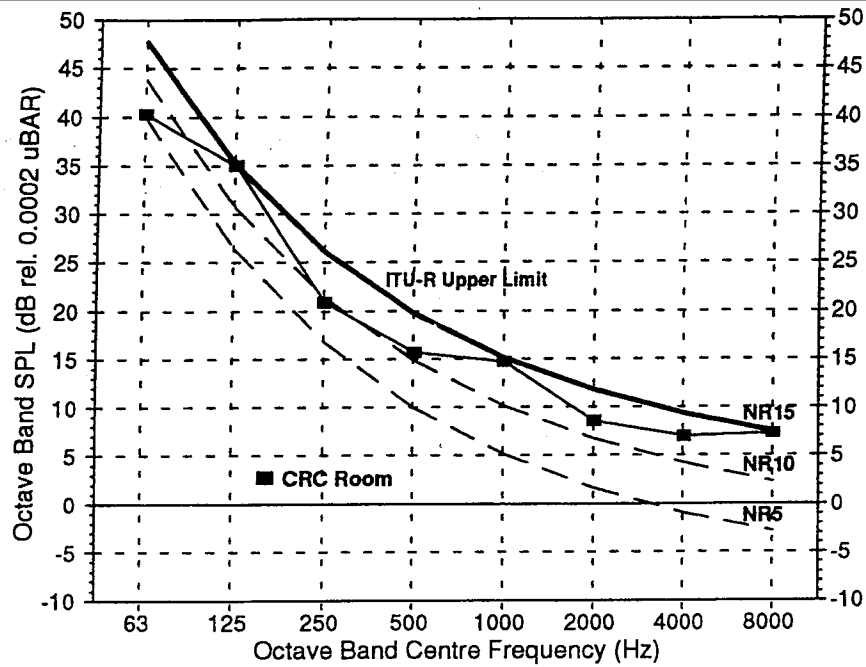


Figure 6.4 CRC listening room - Background noise levels

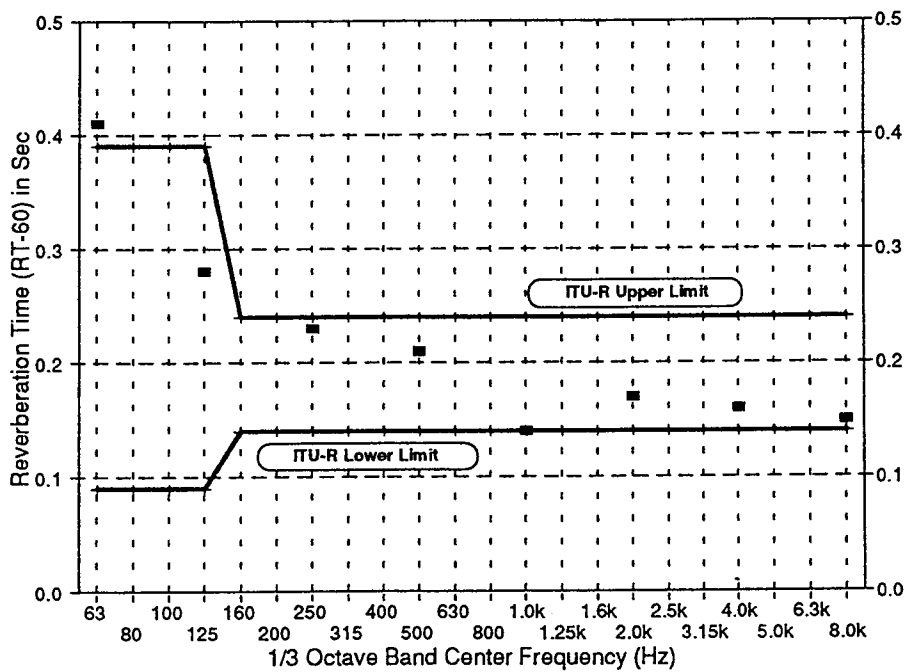


Figure 6.5 CRC listening room - Reverberation time (RT-60)

---

## 7. Acknowledgements

The authors of this report would like to gratefully acknowledge the work, support and collaboration of the following persons throughout the execution of the test described in this report:

- Mr. Ralph Justus of the EIA in Washington, USA
- Messrs. Tom Keller and Dave Londa at the NASA LeRC in Cleveland, USA
- Mr. Emil L. Torick for his participation in the selection of critical materials

The authors also wish to acknowledge the dedicated and invaluable work of Messrs. Gerald Hatton and Douglas McKinnie of the CRC who played key roles in the preparation and execution of the tests. Special thanks to Mr. Lewis Scott of the CRC for his assistance in the preparation of the various legal documents for this contract.

---

## References

- [1] *Methods for the Subjective Assessment of Small Impairments in Audio Systems Including Multichannel Sound Systems*, ITU-R Recommendation BS.1116, Radiocommunications Study Group, International Telecommunication Union, Geneva, Switzerland
- [2] Thibault L. and Grusec T., "*EIA-DAR Listening Tests - Quality and Impairment Tests Procedures.*", Signal Processing and Psychoacoustics, Radio Broadcast Technologies Research, CRC, Ottawa, Ont., Document submitted to EIA Working Group B, Updated version, 1 December 1993
- [3] Londa, Dave, et al. "*EIA-DAR/NRSC DAB objective tests of transmission impairments*".
- [4] Beaton R.J. and Wong P., *A Disk-based System for the Subjective Assessment of High Quality Audio*, Preprint 3497, AES 94th Convention, Berlin, Germany, 16-19 March 1993



NRSC-R55

NRSC Document Improvement Proposal

If in the review or use of this document a potential change appears needed for safety, health or technical reasons, please fill in the appropriate information below and email, mail or fax to:

National Radio Systems Committee  
c/o Consumer Electronics Association  
Technology & Standards Department  
1919 S. Eads St.  
Arlington, VA 22202  
FAX: 703-907-4190  
Email: [standards@ce.org](mailto:standards@ce.org)

DOCUMENT NO.	DOCUMENT TITLE:	
SUBMITTER'S NAME:	TEL:	
COMPANY:	FAX:	
	EMAIL:	
ADDRESS:		
URGENCY OF CHANGE: _____ Immediate _____ At next revision		
PROBLEM AREA (ATTACH ADDITIONAL SHEETS IF NECESSARY): a. Clause Number and/or Drawing: b. Recommended Changes: c. Reason/Rationale for Recommendation:		
ADDITIONAL REMARKS:		
SIGNATURE:		DATE:
<b>FOR NRSC USE ONLY</b>		
Date forwarded to NAB S&T:	_____	
Responsible Committee:	_____	
Co-chairmen:	_____	
Date forwarded to co-chairmen:	_____	



**CEA**<sup>®</sup>  
Consumer Electronics Association

