

Second Quarterly Progress Report

April 1, 1996 through June 30, 1996

NIH Contract N01-DC-6-2100

Speech Processors for Auditory Prostheses

submitted by

Lorraine Delhorne
Donald K. Eddington
Nico Garcia
Victor A. Noel
William M. Rabinowitz
Joseph Tierney
Margaret E. Whearty

**Massachusetts Institute of Technology
Research Laboratory of Electronics
Cambridge Massachusetts**

1 Introduction

Work performed under the present contract is directed at the design, development, and evaluation of speech processors for use with implanted auditory prostheses in deaf humans. The major research efforts are proceeding in four areas: (1) developing and maintaining a laboratory based, software controlled, real time, speech processing facility where processor/stimulator algorithms for monaural and binaural eight-channel implants can be implemented/tested and a wide range of psychophysical measurements can be made, (2) using the laboratory facility to refine the sound processing algorithms used in the current commercial and laboratory processors, (3) using the laboratory facility to explore new sound processing algorithms for implanted subjects, and (4) designing and fabricating programmable, wearable speech processors/stimulators and using these systems to: (a) field test processor algorithms developed and tested in the laboratory, (b) evaluate the effects of learning using longitudinal evaluations of speech reception, and (c) compare asymptotic performance of different speech processors across subjects.

During this Quarter, our primary effort was concentrated on fitting Ineraid subjects with the Geneva Wearable Processor (GWP) that was designed under a collaboration between MIT, RTI, and the Geneva groups. The fabrication of the microprocessor-based processor/stimulator, as well as an NimH battery charger, was accomplished by the Geneva Engineering School. Earhooks, cables, connectors, and rechargeable NimH battery packs, were provided by MIT. At the end of this Quarter 17 subjects are now using the GWP implementing a CIS processor/stimulator algorithm, some for more than one month, and we plan to fit an additional three subjects in the next Quarter. Several hardware problems encountered this Quarter are discussed in this report, and we present speech reception scores for eleven Ineraid subjects tested after wearing the GWP for one month.

We are also continuing our study of the mapping functions used by CIS processors. This Quarter we present the results of experiments that explore how limiting the range of envelope levels mapped in a CIS processor affects consonant recognition.

2 Wearable Processors

This Quarter, 17 Ineraid subjects have been fitted with the Geneva Wearable Processor (GWP) and we anticipate fitting three additional subjects in the next Quarter. As discussed in a previous QPR[1], the new wearable sound processing system replaces all of the external components of the Ineraid hardware system. These components include the sound processor, processor/earhook cable and connectors, earhook assembly, and the earhook/pedestal cable and connector. A processor power source that consists of a four-cell battery pack using type 150AFH (roughly AA size) NimH rechargeable units, as well as a recharger, are also provided to each subject. Present subject experience demonstrates that the processors operate for a full day (18 hours) on a single battery pack without signaling low battery.

It is important to note that the power dissipated by the processor is a function of the program the processor is running (e.g., frequency of external memory access and clock rate). The

CIS processor implementation[2] we designed for the first trials of the GWP allows the processor to operate with a clock rate of 20 MHz (half of the full 40 MHz rate).

2.1 Hardware Problems with the Wearable Processor

Our experience with 17 subjects fitted with the GWP, some for over one month, confirms the general reliability of the processor hardware and software. Problems that have been encountered fall into three categories: (1) failure of the cable connecting the processor and the earhook, (2) processor shutting down while in use requiring power cycling for reactivation, and (3) failure of the NimH-battery charger.

A weak assembly point in the processor/earhook cable at the processor connector results in failures that occur with repeated and extreme flexing of the cable at that connector. Given the way many subjects wear the device on their belts, the cable flexes at the weak point whenever they bend over. Precision Interconnect (PI), the cable manufacturer, is studying and testing a proposed solution to this problem and we anticipate that a modification will be available in the near future.

The most likely causes of the processor shutting down at random times appear to be either a build up of static charge somewhere in the system or an intermittent failure of the battery contact. The static charge mechanism was suspected when one subject reported failure upon leaving his car by sliding across the seat. However, other subjects have reported this failure when walking upstairs or turning around. The failures during physical movements are consistent with battery contacts that have lost elastic strength and are not forming a pressure contact with the battery pack. Our colleagues in Geneva have reported such contact failure and are producing a modification that strengthens and preserves the contact elasticity. If this modification does not cure this failure mode, we will return to exploring the static charge hypothesis.

Because all failures of the battery-pack charger have been traced to a failed regulator IC, modifications of the charger have been introduced to protect this regulator IC. We replaced the wallplug-transformer assembly that delivers a nominal 12V (rms) voltage to power the charger circuit with one that includes rectification and filtering to reduce the likelihood of producing an inductive surge when the transformer is disconnected from the charger while connected to the wall socket. The addition of a surge suppressor element provides additional protection.

2.2 Initial Longitudinal Data from the Geneva Wearable Processor

To date, 17 Ineraid subjects have been fitted with and are using the Geneva Wearable Processor (GWP) running a Continuous Interleaved Sampling (CIS) sound processing algorithm. All 17 subjects had at least two years of experience with their Ineraid sound processor before switching to the CIS system. Sixteen subjects were switched directly from the Ineraid to the GWP. The remaining subject (I01) has made two processor transitions over the past year. In 1995, he ended more than 8 years of Ineraid use and switched to a CIS system implemented using a wearable research processor developed at the University of Innsbruck. Longitudinal speech reception scores documenting his performance over a period of 9 months were presented in an earlier report where he was identified as S1[3]. More recently, I01 switched from his Innsbruck system to the GWP unit using our standard CIS processor algorithm.

Three of the 17 current GWP users are not participating in our longitudinal study. While we will follow these three subjects over time, they are unable to comply with our formal testing schedule because of poor health and/or distance from Boston.

Figure 1 presents the results of four speech reception tests for the 11 subjects that have used the GWP for at least one month. Each panel is a bar graph that shows two scores for each subject. The left (lightly shaded) bar represents the subject's latest Ineraid score and the right (darkly shaded) bar depicts the subject's CIS score after one month of continuous GWP use. Also shown at the far right of each panel are the Ineraid and CIS scores averaged across subjects. All tests are conducted without speechreading and without feedback.

The bottom panel of fig. 1 shows data from a 12-consonant identification test[3]. The subject's task is to identify the initial consonant (Ci) of Ci-ah-Cf utterances without speechreading. Each of twelve initial consonants are represented by six tokens spoken by a single, female talker. The test is organized in blocks of seventy-two trials with each of the twelve initial consonants being presented six times in randomized order in each block. Two or three blocks (144-216 trials) are tested in a session. Differences between the Ineraid and CIS scores across subjects are small for this test.

The second panel from the bottom in figure 1 compares Ineraid and one-month CIS performance for identification of NU-6 Monosyllabic Words. This list of fifty words is played from a cassette tape recorded by Auditec as part of the MAC Battery of Tests[4]. The subjects have not been shown the word list but they have been tested with these materials up to twice each year since they received their Ineraid implant. For this test, the differences between the Ineraid and CIS scores for some subjects are larger than any of those for the consonant test (e.g., subjects I11, I01, I18, and I19) and the number of subjects scoring lower with CIS than with the Ineraid system are fewer (e.g., I26).

The top panel of Figure 1 presents the percentage of words identified correctly for CUNY sentences[5] presented in quiet without lipreading. These sentences are presented open-set and the subjects have much less experience with these materials than with the NU-6 words. In fact, subjects have not received a repeated list in any test session to date. This test is probably the closest to normal conversational listening, allowing context to aid in word identification. Like the results from the NU-6 test, differences between the Ineraid and CIS scores for these sentences are larger than those observed for consonant identification (e.g., I11, I16, I01, I05, I26, I27, I10, I25, and I19). Notice also that the largest decreases in CIS scores relative to the Ineraid scores occur for the CUNY Sentences (e.g., I16 and I27).

The second panel from the top of Figure 1 shows the percentage of words identified correctly for HINT[6] sentences presented in background noise (S/N: +10 dB). Three subjects demonstrate relatively large gains with CIS processing that result in scores greater than 50% (e.g., I01, I25 and I19), but in general, the differences between the Ineraid and CIS systems are modest.

Previously presented[3], longitudinal test scores for subject I01 using the Innsbruck processor are shown in Figure 2. These results show that the rate of improvement for word identification using the CUNY Sentences (top panel) is higher than that for the consonant identification task. At the time the Innsbruck system was fitted, I01's consonant identification score dropped from 57% using the Ineraid system to 35% using the Innsbruck processor and remained near 40% for more than 9 weeks. Word identification scores measured with the CUNY Sentences using the Innsbruck processor also decreased below those of the Ineraid on the day the Innsbruck processor was first used but, after two weeks, improved 30 percentage points to equal the Ineraid's performance and continued improving at a steady pace to end 30 percentage points higher than the Ineraid after 36 weeks.

This difference in the rate of improvement is consistent with the general pattern of the data

presented in Figure 1. At 4 weeks, none of the subjects show consonant scores with CIS that are substantially higher than those measured with the Ineraid system. However, several subjects do show large gains for the CUNY Sentence material. This same pattern also occurred with subject I01 when he switched from the Ineraid to the Innsbruck system and, interestingly, it appears to occur again when he switched from the Innsbruck to the GWP system. With consonants, his Innsbruck score had been 68%; after four weeks using the GWP system, it had decreased to 60%. In contrast, his CUNY score with the Innsbruck system had been about 70%, but after four weeks using the GWP system, it had increased to 94%.

By the end of next Quarter, a significant number of subjects will have reached their three-month testing date. At that time we should have a better idea whether subject I01's Ineraid-to-Innsbruck pattern will be a general trend seen across subjects. Already, however, these preliminary longitudinal data indicate that one must be cautious using acute laboratory tests to evaluate the relative benefit of sound processing systems. Initial decreases in consonant and word identification observed when moving from one processor to another may not be a reliable indicator of asymptotic performance.

3 Laboratory CIS Mapping Studies

3.1 Standard Mapping Strategy

The standard mapping strategy we use as a reference point in both our laboratory and wearable processor studies has been described in previous QPRs[1] and, for convenience, is summarized here. The purpose of a channel's mapping stage in a CIS processor (see Figure 3) is to map a relative wide range of envelope levels into a relatively narrow range of output current levels. Usually a nonlinear mapping function is employed (logarithmic in our case). An output dynamic range (DR in Figure 3) is determined for each channel based on the psychophysical threshold and most comfortable stimulus level measured using the electrode associated with that channel and stimuli that roughly correspond to those used by the processor of interest (e.g., biphasic, cathodic-first pulse train, 2 kHz repetition rate, 31.25 μ s/phase, 50 ms duration in our current five and six-channel CIS implementations). The input gain (G_{in}) determines the range of envelope levels that will fall within the 60 dB input range that is mapped. We set G_{in} such that for the TIMIT[7] database of speech materials played at a conversational level, 1% of the envelope levels will be clipped in the channel with the most energy (Channel 2). For a given G_{in} , the DR is adjusted to produce a comfortable listening level by varying a scale factor, SF, where $DR = (SF * (MCL - THR) + THR) / THR$. The operating point described by this fitting procedure is the one currently used by all of our GWP subjects.

While this adjustment procedure specifies a repeatable operating point for the mapping stage, it involves several arbitrary choices. The logarithmic mapping function was selected because others have used it successfully. Whether this function is optimum is a question we are currently pursuing and we plan to report on that work in a future QPR. Our decision to represent this mapping function as a table of 1024 values is also somewhat arbitrary and may be addressed in future work. The work described in the next section begins to explore our arbitrary choice of 60 dB for the range of mapped input envelope levels.

3.2 Mapping Manipulations: Range of Envelope Levels Mapped

The work described in this section begins to explore the effect of reducing the range of input envelope levels mapped across channels. As discussed earlier, the CIS algorithm used by our GWP subjects maps a 60 dB range of envelope levels into a range of stimulation currents that varies depending on the DR of the particular electrode. In the work reported in this section, two modifications of the standard CIS processing algorithm were made. First, instead of using a constant G_{in} across channels that is based on the 1% clipping criterion in Channel 2 as described above, the G_{in} of each channel was adjusted separately so that all channels clipped 1% of their levels when the TIMIT sentences were played at a conversational level.

With the CIS processor set in the above fashion, the mapping function of each channel was manipulated to reduce the range of envelope levels mapped as shown in Figure 4. In one set of experiments the range of envelope levels mapped was reduced by moving the low-level boundary of the mapping function's input range to the right (left column of Figure 4). In a second set of experiments, the range of mapped levels was reduced by moving the high-level boundary of the range to the left (right column of Figure 4). The position of the high or low-level boundary was selected to include a specific percentage of envelope levels rather than to include a specific input mapping range as shown in Figure 4. This can be seen in Table 1 where the low-level and high-level boundaries are specified as the percentage of levels lower than the boundary level. Thus, for a low-level boundary position of 80%, the mapped range in each channel begins at a level that is higher than 79% of all the levels in each channel. Similarly, a high-level boundary position of 87% means that clipping begins at a level that is higher than 87% of all the levels in that channel.

Table 1 presents the effects of these two types of mapping manipulations upon subject I04's ability to identify 24 medial consonants spoken by a male talker[8]. In the case where the high-level boundary level was reduced, increasing the percentage of levels clipped, significant reductions in consonant scores began at the 1-80% condition. In a conversational setting, the subject could not detect a qualitative difference between the 1-99% and the 1-95% conditions, but did begin to hear increased levels of distortion at the 1-90% condition.

The left half of Table 1 shows the consonant identification scores for manipulations of the low-level boundary of the mapping function. In this case, the lowest 96% of the envelope levels must be discarded to reach a consonant score comparable to the condition of clipping the highest 20% of the envelope levels. Qualitatively, the subject preferred the 40-99% condition to the 1-99% condition in a conversational setting because of a reduction in the ambient background noise; voice quality seemed to be identical. As the low-level boundary is moved higher and a larger percentage of the envelope levels are discarded, the voice quality begins to suffer from "dropout effects." The subject began to detect this type of distortion when more than 40% of the envelope values were discarded, even though the consonant scores remain high.

To some extent, the effects of these manipulations are minimized by the controlled levels of the speech segments used in the consonant test. Nevertheless, it is clear that level variations at the high end of the range of input envelope levels are more important to represent than those at the low end of the range. Thus, given a limited input range to be mapped, one should adjust G_{in} to avoid clipping much beyond a 1% level.

3.3 Restricted Envelope Ranges Compared to Our Standard Mapping Ranges

Our standard CIS processor maps a 60 dB range of input envelope levels into an output

current range that is specified for each channel. The subjective judgment of I04 (reinforced by the 24-consonant test results) that the lower 40% of envelope levels can be discarded with no reduction in sound quality, suggests that it may be possible to reduce the mapped range of envelope levels without compromising performance. If the length of the table representing the mapping function remains the same, a reduction of this range would mean an increase in the amplitude resolution. Alternatively, maintaining the same amplitude resolution and reducing the range of mapped envelope levels would mean the mapping function could be represented in a shorter table thereby reducing the memory required by the CIS program.

The range of envelope levels that correspond to 40-99% boundary positions shown in Table 1 can be determined from the cumulative distribution of envelope levels computed for each channel before the mapping stage and shown as tables in Figure 5. For example, the cumulative level distribution for channel #3 shows that the 99% high-level boundary corresponds to -10 dB and the 40% low-level boundary corresponds to -46 dB. This means that a range of 36 dB is required to map the envelope levels within the 40-99% boundary positions in channel #3. The same process applied to the other channels gives the following ranges for channels #1 through #6: 23, 34, 36, 30, 29 and 32 dB.

In our standard CIS system, G_{in} is constant across channels and set to clip 1% of the envelope levels in channel #2. For the case shown in Figure 5, G_{in} (see Figure 3) would be set to 5 dB because the cumulative level distribution for the TIMIT sentences shows that 99% of the levels in channel #2 are 5dB below the channel's maximum (16 bit) level. This means that the input ranges for channels #1 through #6 that would position the lower-level boundary at 40% of the channel's cumulative level distribution are: 29, 34, 41, 39, 40 and 44 dB re the channel's maximum level.

If measurements in additional subjects show similar results to those of I04, they would suggest that we might be able to reduce our standard 60 dB mapping range to one of 40-45 dB or less without adversely affecting speech reception and sound quality. However, three factors may work to minimize the actual reduction realizable without a decrease in sound quality: (1) the CIS implementation used with the GWP system does not include an AGC to reduce the variations in overall sound levels encountered in the real world, (2) the medial consonant test used to measure speech reception for the various level boundary conditions does not include a wide range of overall amplitude within or across the various speech tokens and (3) the TIMIT sentences used to produce the cumulative level distributions are also well controlled in terms of overall level. We plan to continue these measurements in additional subjects and use that information to optimize the range of envelope levels mapped.

3.4 Choices for Channel Gain Equalization

We have mentioned two strategies for setting G_{in} (see Figure 3) using a standard set of speech materials (the TIMIT sentences in our case). In one, G_{in} is the same across all channels and is set such that 1% of the levels will be clipped in the channel with the greatest energy. Another strategy is to set G_{in} separately for each channel. In this case the G_{in} of each channel is set such that 1% of the levels are clipped in each channel. The first method maintains the spectral shape of the speech because the gain across channels is the same. In the second method, the gains vary across channels and the spectral shape is distorted. When the TIMIT sentences are used to set these individual gains, a high frequency emphasis is produced because the G_{in} needed at each

channel to move the 99% point of the cumulative level distribution is, in general, greater for the high-frequency channels (e.g., 11, 5, 10, 14, 16, and 17 dB for channels #1 - #6 respectively as can be seen from Figure 5).

To determine whether the equalization on a channel basis enhances speech reception even though the overall speech spectrum suffers a constant distortion, we tested four subjects for medial consonant reception when using our six-channel processor. In one case each channel's gain was adjusted so that all six mapping inputs saturated 1% of the time when the TIMIT speech materials were used as input. In the second case an equal gain was used for all six channels such that the strongest channel's output (channel#2) saturated 1% of the time. For our CIS processor implementation, which uses a 60 dB mapping band level, the results for four subjects are shown in Table 2. In each of the four subjects, average percent-correct scores for the 16 consonants spoken by a male talker were higher for the case of equal gain. However the limited amount of data shows significance for only three of the four subjects as indicated.

Although the evidence is not strong, it certainly does not indicate a clear advantage for the equalized gain condition. Intuition would suggest that if the input band levels will accommodate the range of envelope values required for good speech reception (as discussed earlier), then there is no advantage to the introduction of a spectrum-distorting equalization. However, if the band levels are highly reduced, then the distortion may buy some reception advantage.

4 Future Work

Next Quarter we will continue longitudinal trials of the GWP. By then we expect to be following 20 of our research subjects and should be able to report three-month test results for approximately 11 subjects.

We also plan to conduct experiments associated with amplitude mapping. These will include continuing the measures that explore the relationship between speech reception and the range of envelope levels mapped in CIS sound processing and will also include evaluation of mapping functions designed to restore normal loudness growth.

References

- [1] D.K. Eddington et al. (1996): "Speech Processors for Auditory Prostheses, First Quarterly Progress Report (January 1 through March 31, 1996)." *NIH Contract N01-DC-6-2100*.
- [2] D.K. Eddington et al. (1994): "Speech Processors for Auditory Prostheses, Ninth Quarterly Progress Report (October 1 through December 31, 1994)." *NIH Contract N01-DC-2-2402*.
- [3] D.K. Eddington et al. (1996): "Speech Processors for Auditory Prostheses, Final Report (September 30, 1992 through December 31, 1995)." *NIH Contract N01-DC-2-2402*.
- [4] E. Owens et al. (1985): "Analysis and Revision of the Minimal Auditory Capabilities (MAC) battery." *Ear Hear.* 6:280-287.
- [5] A. Boothroyd, L. Hanin, et al. "A sentence test of speech perception: Reliability, set equivalence, and short term learning." *Speech and Hearing Sci. Rept. RC110* City Univ. of N.Y. (New York, NY) 1985.

- [6] M. Nilsson, S.D. Soli, J.A. Sullivan, "Development of the Hearing-In-Noise-Test for the measurement of speech reception thresholds in quiet and in noise." *Journal of the Acoustical Soc. of Amer.*, Vol 92: pp1085-1099, 1994.
- [7] W.M. Fisher et al. (1986): "The DARPA Speech Recognition Research Database: Specifications and Status." *Proc. DARPA Workshop on Speech Recognition, pp93-99, Palo Alto California.*
- [8] R.S. Tyler et al. (1987): "The Iowa Audiovisual Speech Perception Laser Videodisc. Laser Videodisc and Laboratory Report." *Dept. of Otolaryngology, University of Iowa Hospital Clinic. (Iowa City, IA).*

Table 1: Percent correct scores for subject I04 measured with the Iowa, 24 medial consonant test (male talker) for CIS mapping functions that vary in their low-level and high-level input boundaries as shown in Figure 4. The table's left side shows results for moving the low-level boundary to higher levels, while the right side shows results obtained when moving the high-level boundary to lower levels.

Low-level Boundary Variation			High-level Boundary Variation		
Boundary Position		24 Consonant Score	Boundary Position		24 Consonant Score
Low	High		Low	High	
1	99	99%(1.86)	1	99	99%(1.86)
40	99	98%(3.73)	1	95	98%(3.73)
80	99	99%(1.86)	1	90	94%(3.73)
90	99	92%(1.86)	1	80	87%(4.56)
97	99	87%(5.43)	1	70	77%(2.28)
98.5	99	75%(6.59)	1	60	58%(6.59)
98.7	99	69%(2.28)			

Table 2: A comparison of speech reception scores for 16 medial consonants with two conditions of channel gain. In the first case, all channel gains are independently adjusted to provide input envelope values that produce saturation (i.e. are at the maximum input level) 1% of the time, when the input speech is a representative sampling of the TIMIT database. In the second case, all channel gains are set at the same value, the gain needed to produce 1% saturation for the channel with the highest energy of envelopes (channel #2).

Subject	All at 1%	Only channel #2	comments
S02	50%(12.4)	56%(8.5)	significant at 0.1
S11	45.2%(4.5)	51.6%(4.6)	significant at 0.25
S01	73.8%(3.8)	79.5%(8.4)	significant at 0.05
S27	29.4%(14.2)	37.6%(13.9)	not significant at 0.1

Figure 1: Data for eleven subjects fitted with the Geneva Wearable Processor (GWP) implementing a CIS algorithm. Each panel is a bar graph that shows two scores for each subject. The left (lightly shaded) bar represents the subject's latest Ineraid score and the right (darkly shaded) bar depicts the subject's CIS score after one month of continuous GWP use.

Figure 2: Longitudinal test scores for Ineraid subject I01. The gray-filled symbols show test scores measured using his Ineraid processor and the dark symbols are results using the Innsbruck system.

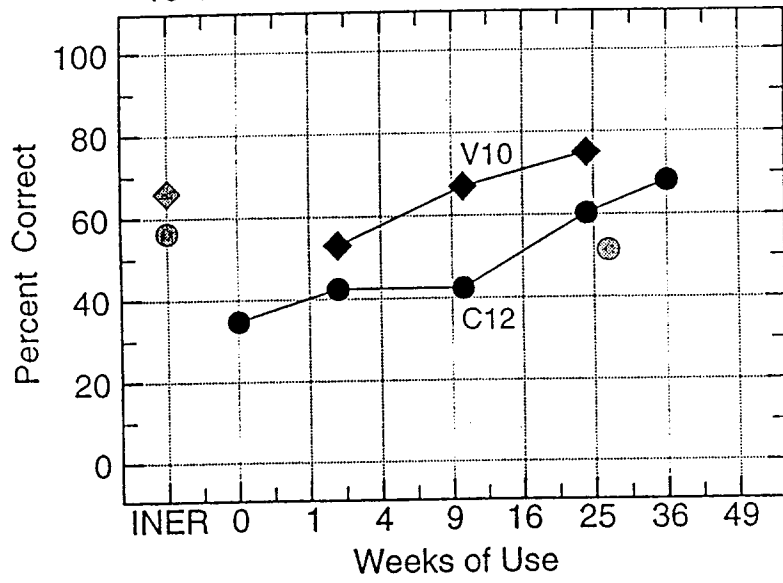
Figure 3: Diagram of the non-linear mapping function used in a CIS processing scheme. The total gain before the mapper is characterized by a single input gain ("Gin") that determines where the distribution of envelope levels will be positioned at the input to the mapping function. The mapper output amplitude modulates the pulse train that is converted to a current stimulus by the voltage-to-current source converter ("V/I"). The amplitude of a current pulse generated by a specific input level is determined by the mapping function and "Gout." From QPR1 figure #2 NIH Contract #N01-DC-6-2100.

Figure 4: A set of mapping functions demonstrating how the low and high-level boundaries of the mapping function's inputs were varied to study the effects of input range variations. The boundaries are given for channel #2 in terms of its cumulative level distribution (see Figure 5).

Figure 5: Level histograms and cumulative level distributions for envelope signals in each channel of our standard CIS processor using the standard TIMIT database as input. These statistics were computed on the envelope signals before Gin of the mapping stage (see Figure 3).

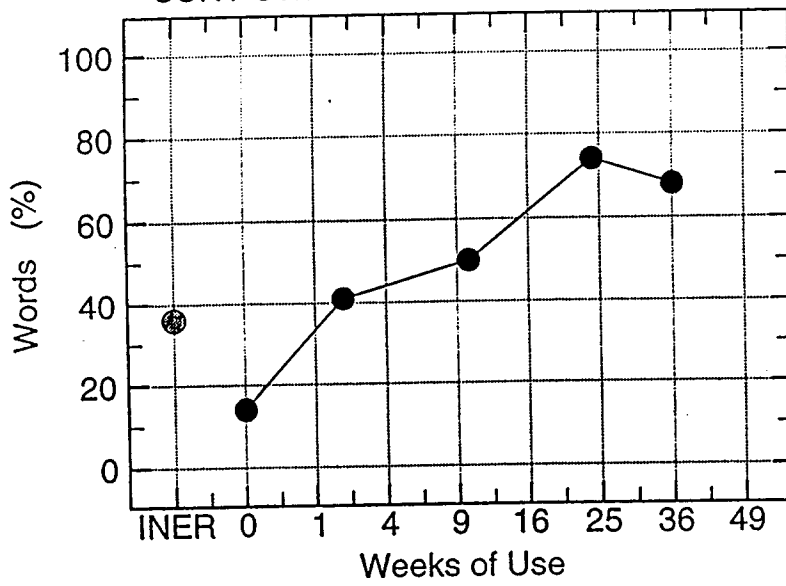
12-Consonant and
10-Vowel Identification

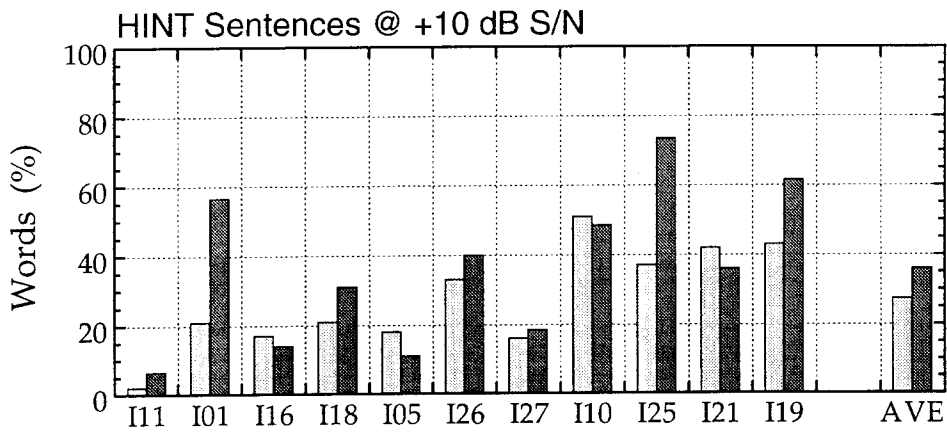
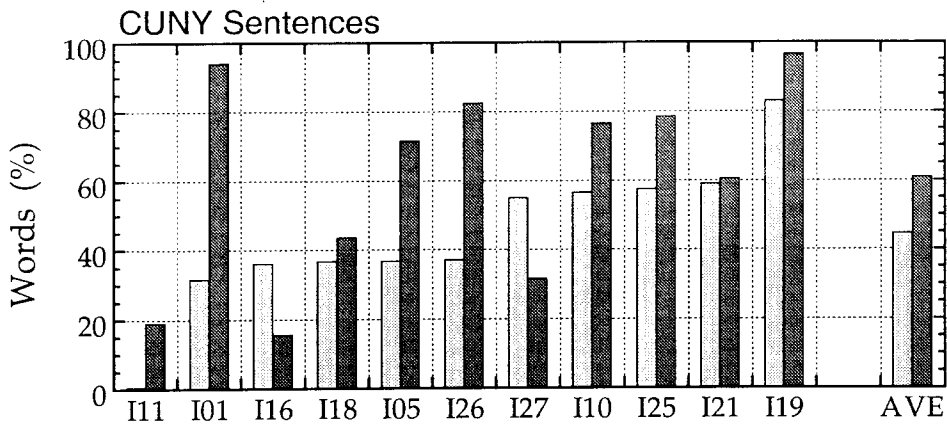
S1



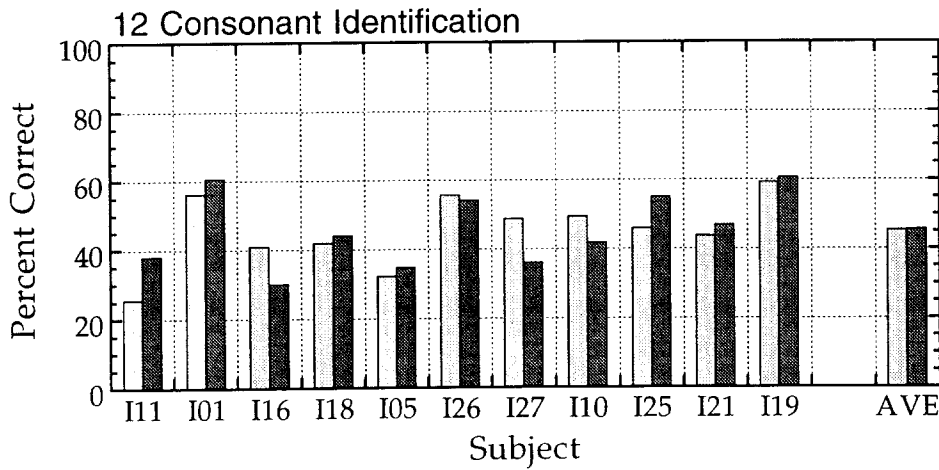
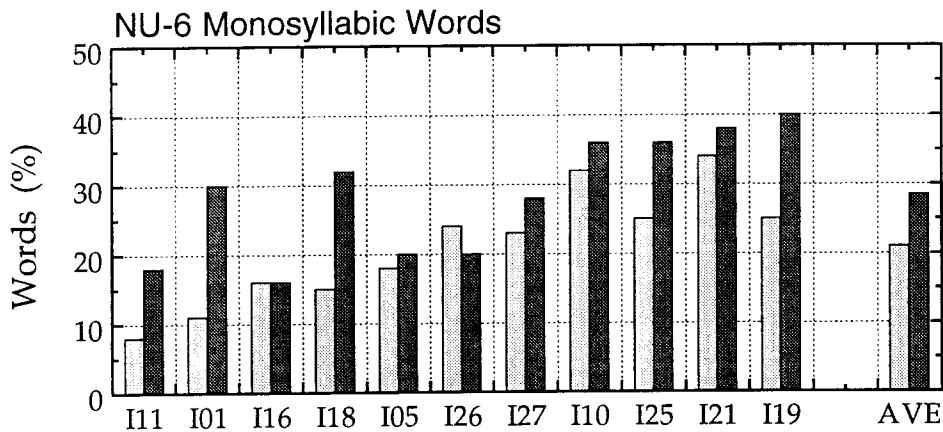
CUNY Sentences in Quiet

S1





Ineraid
 CIS @ 1 mo.



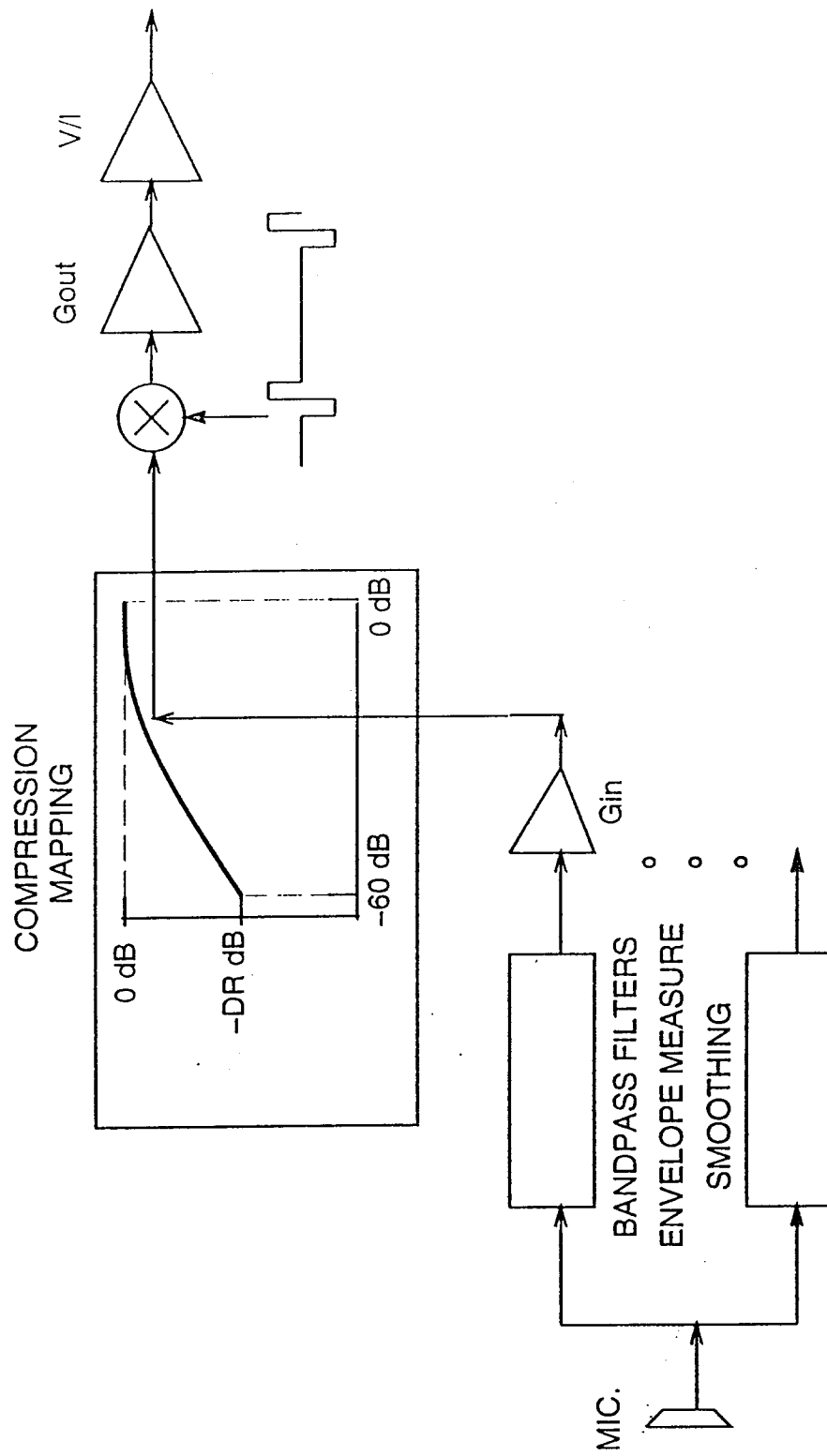
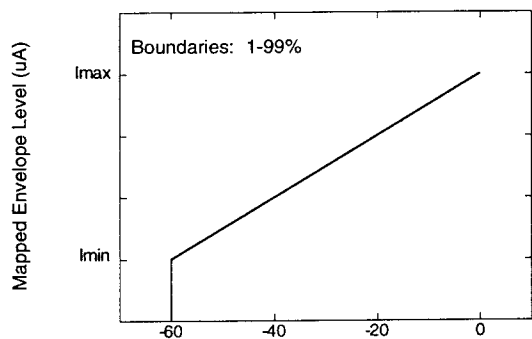


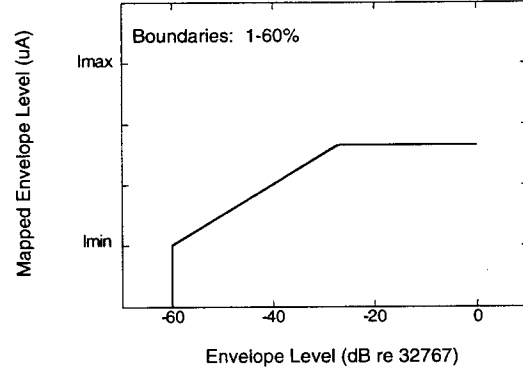
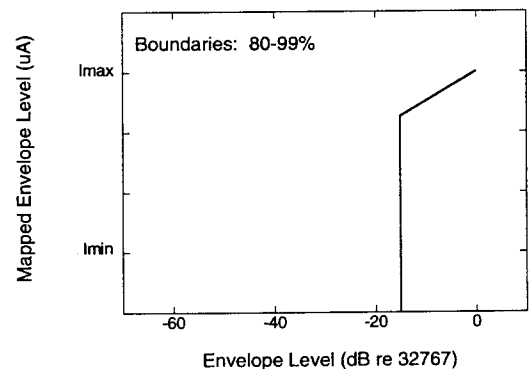
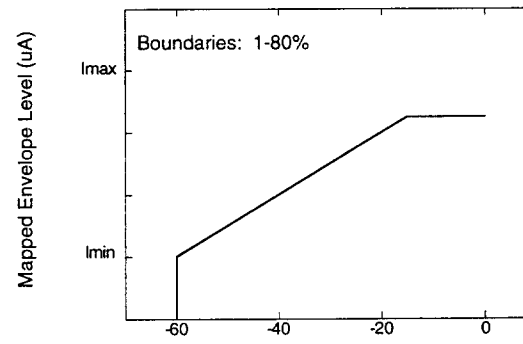
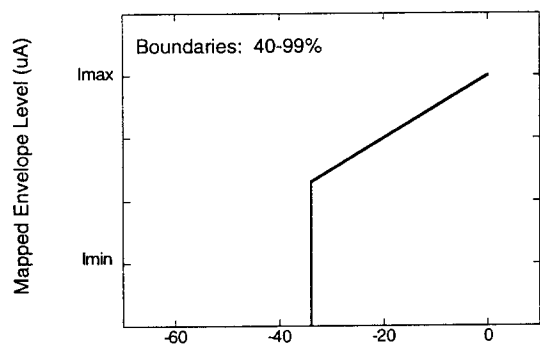
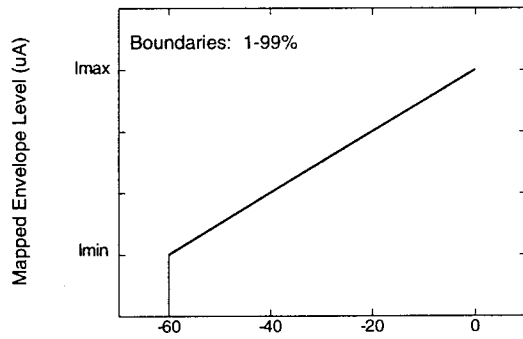
FIGURE 3

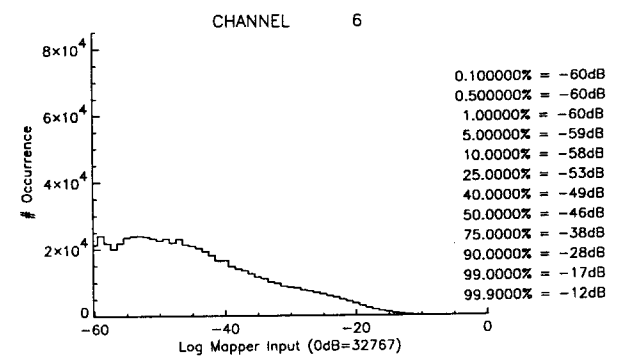
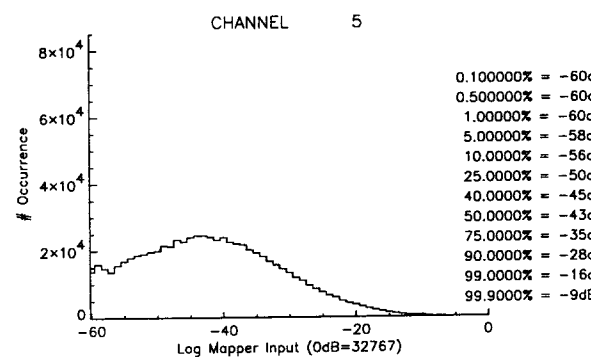
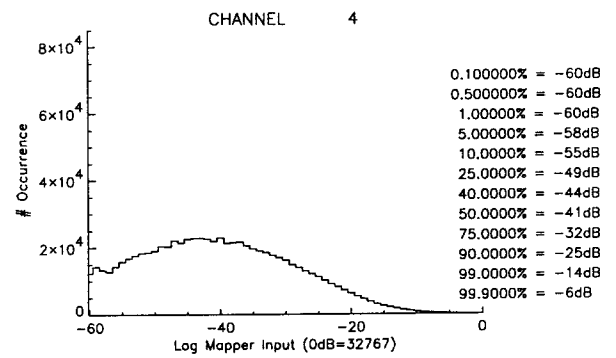
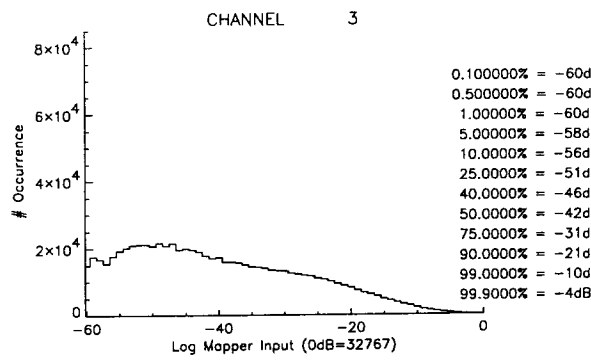
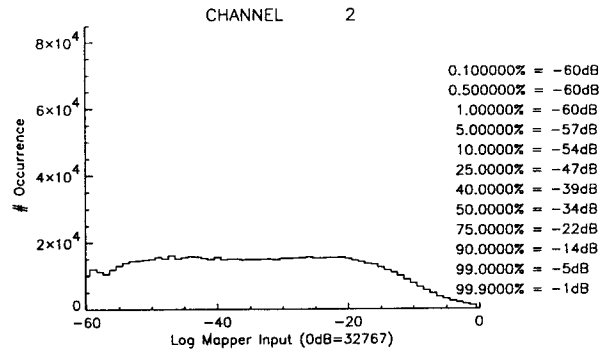
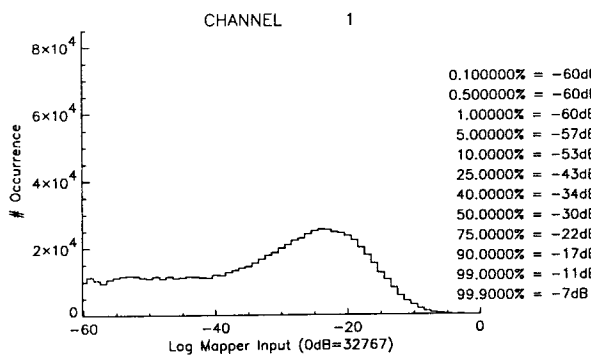
Mapping Range Variations

Variation of Low-Level Boundary



Variation of High-Level Boundary





dr1-8.pcm.hist

binsize=1dB