

oticon
MEDICAL

Because
sound matters



SORBONNE
UNIVERSITÉ

TELECOM
Paris



IP PARIS

ircam
Centre
Pompidou



ATIAM



Cochlear®

Security relevant sounds perception in cochlear implant users

ALEXANDRE PHILIPPON

SIGNAL PROCESSING, COMPUTER SCIENCE, ACOUSTIC
MASTER ATIAM

2nd april 2024 - 13th september 2024

*University referent
teacher:*

EMMANUEL
SAINT-JAMES
saint-james@lip4.fr

Industrial mentor:

ADRIEN DANIEL
addaniel@cochlear.com

Industrial mentor:

YUE ZHANG
yuezhang@cochlear.com

Abstract

This work focuses on the analysis of how cochlear implants process relevant security sounds and aims to explain the poor performance of CI users in terms of recognition and detection. This analysis involves visualizing the electrograms and spectrograms of the processed sounds, as well as comparing the acoustic features of the input and output sounds. The output sounds are generated using relevant CI vocoders. Additionally, the work proposes a design to improve CI transduction and provides experimental protocol guidelines to test its efficiency.

Acknowledgment

I'd like to express my gratitude and affection to everyone who has supported and helped me during this internship. My first thanks go naturally to my dedicated supervisors, Dr. Adrien Daniel and Dr. Yue Zhang. It is hard for me to express how thrilling and fulfilling our meetings and discussions were. Thanks to your kindness, open-mindedness, and availability, I had the opportunity to work in an extremely pleasant environment and to learn a great deal.

I would also like to extend my gratitude to the entire RD team at Cochlear Vallauris, especially to Dr. Rafael Attili Chiea and Dr. Behnam Molae-Ardekani, who assisted me with their work and provided valuable and insightful advice.

Naturally, I also want to thank my loving family and friends, who, as always, support me in every project I undertake, even the most ambitious. Life would likely be less meaningful without the love we share, and I am truly grateful to have you all in my life.

Nomenclature

Devices and electronic

Bone anchored hearing system Bone anchored hearing system is a surgically implantable device. It sends sound waves through the bones of the skull directly into the cochlea and is target for people having conductive hearing loss

CI CI is the abbreviation of Cochlear Implant, an electronic prosthesis used to compensate the hearing loss of the wearer. It emulates the natural processing chains of hearing from the wave-detection to the electrical stimulation of neurons and is target for people suffering profound sensory-neural hearing loss.

Psycho-acoustic

Alarm detection We call in this study alarm detection the capacity of a subject to detect the presence of an alarm sound even without identifying this sound as an alarm

Alarm recognition We call in this study alarm recognition the capacity of a subject to identify a sound stimuli as an alarm and to differentiate those (police siren, fire fighter, evacuation city signal...)

Alarm audibility We call in this study alarm audibility the capacity of a subject to answer correctly at the question "Was the alarm audible" in case a mixture containing alarm sound were presented.

The current spread issue This phenomena occurs because of the tendency of electric current to spread when put in a aquatic environment, because of this, when an electrodes fires a electrical current, adjacent neurons associated with other electrodes are being stimulated causing a multiple channel interaction in term of "sound sensation".

SIR SIR stands for Speech Intelligibility Ratio, a score measuring the capacity of a subject to understand speech

SMRT SMRT stands for Spectral-temporally modulated ripple test, this test was created by David Landsberg and Justin Aronof and is design to either measure the ripple per octave threshold of the ripple rate threshold of a subject. You can find a more precise description of this test at [link to site \[2\]](#)

SRS SRS stands for Safety Related Sounds, a category of sounds providing safety related information (like the presence of a danger or a dangerous situation).

Signal processing

Beamforming Beamforming or spatial filtering is a signal processing technique used to select spatially a signal having multiple channel as input. This is achieved by combining elements in an antenna array in such a way that signals at particular angles experience constructive interference while others experience destructive interference. Beamforming can be used at both the transmitting and receiving ends in order to achieve spatial selectivity

CIS CIS stands for Continuous Interleaved Sampling. It's a technique for neural stimulation that consist of delaying the impulses of the different channels/electrodes in order to stimulates only one electrode at a time. It has been observed in various study according to [18] that the interaction between channels caused by the electrical spread decreases the performance of subject in speech understanding. This technique decreases the problem compared to non-CIS process, thought it implies a kind of temporal distortion which do not significantly impair CI user's speech recognition.

EDG EDG stands for electrodogram. An electrodogram is a 2D representation (1 dimension being the electrodes number, and the second the time) of the electrical stimulation provided at the electrodes array by the implant in response of a sound stimulus (depending on the coding strategy, the amplitude of stimulation can be encoded as an electrical tension (the higher the level of the sound, the higher the electrical tension of is) or as a duration of stimulation (the higher the level of the sounds, the longer the electrical stimulation is). Those representation are close to spectrograms as each electrodes is associated with a certain frequency band.

STFT STFT stands for Short time Fourier transform, it's a signal processing technique which consists of applying the Fourier transform on short temporal frame of a signal (eventually with an overlapping on those transforms) to obtain spectrum's of those short temporal frames. Joined in respect of the temporal axis, we obtain a Spectro/temporal representation of the signal also called spectrogram.

Acoustic

Reverberation Reverberation is a natural acoustic phenomenon occurring when a direct sound is undergoing multiple reflection to the point where the sound confused. In our study, we characterize reverberation using only the [RT60](#) coefficient

Audiofeature An audio feature is an acoustic and or perceptual characteristic who can be estimate numerically through a calculation on the raw waveform of the acoustic signal it's associate with (it associates one coefficient with an

audio waveform). Depending on the level of representation, those characteristic may be close to the physics of sounds (zero-crossing rate for example), close to a human perception (roughness coefficient) or close to both (pitch/F0)

RIR The RIR is the abbreviation for room impulse response, the Green function of a room when excited at a certain point by an acoustic impulse.

RT60 The RT60 (reverberation time to -60dB) is the time needed for a room impulse response to lose 60 dB of level compared to the impulse intensity.

Contents

1	Introduction	8
1.1	Companies presentation	8
1.2	Neurelec today	9
1.3	The R&T team	9
1.4	Context and interest of the project	9
2	Literature Review and state of the art	9
2.1	Normal functioning of hearing	9
2.2	Functioning of CIs	11
2.3	Safety related sounds	13
2.4	The Problematic Issue of Relevant Security Sounds for Cochlear Im- plant Users	16
2.5	Coping Strategies for Integrating Security Sounds	17
3	Analysis of SRS in CI Processing	20
3.1	Data Collection and Process of Analysis	20
3.2	Analysis of the SRS Processing Alone	22
3.2.1	Analysis of Electrograms	22
3.2.2	Description of the Vocoders	24
3.2.3	Analysis of the Spectral Modulation	25
3.2.4	Analysis of CI Vcoded Sounds Features vs. Actual Sounds Features	26
3.2.5	Results of the Feature Correlation Analysis	30
3.3	Analysis of SRS with Highly Masking Noise	33
4	Proposed solution for improvements of the SRS perception	37
5	Clinical Protocol for Testing the Proposed Solutions	37
5.1	Exclusion and Inclusion Criteria	37
5.2	Protocol Proposed	38
5.2.1	Experiment 1	38
5.2.2	Experiment 2	41
5.2.3	Experiment 3	41
5.3	Limitations of Experiments	43
5.3.1	Participant Profile	44
5.3.2	Experimental Setup	44
5.3.3	Material (Hardware and Software)	44
5.3.4	Balance and Representativity	45

6	Conclusion and Discussion	45
6.1	Conclusion	45
6.2	Future Work	45
7	Annex	46
7.1	Some other figures	46

1 Introduction

"Give every man thine ear, but few thy voice" [27]. The ability to hear is undoubtedly crucial for social interaction, and research in hearing aids has primarily focused on restoring the capacity to hear and understand speech. Although this is an important issue, losing the ability to hear also means losing significant environmental information carried by sounds, such as security-relevant information or general environmental context. This work aims to study, characterize, and address the issue of low recognition of environmental sounds by cochlear implant users, particularly concerning security-relevant sounds, and to propose designs for improving this perception.

We advise our readers to first consult [18] for an introduction to the subject of cochlear implants and the history of this technology, and then read [7] for more details on current coding strategies and electronic designs.

1.1 Companies presentation

Neurelec was originally founded as a laboratory called MXM in 1977, focusing primarily on developing implanted medical devices. In 2006, to enhance the quality of the company's activities, the MXM laboratory was divided into three subsidiaries, each dedicated to a specific focus area. One of these was Neurelec, specializing in the development of cochlear implants.

In April 2013, Neurelec became part of the William Demant Holding Group, a global leader in hearing healthcare established in 1904. This group includes Oticon and Oticon Medical. With 16,500 employees in over 130 countries, the group has unparalleled access to the latest technological advances and insights into hearing care.

Oticon Medical is a renowned global company specializing in implantable hearing solutions. Its activities revolve around two types of implants: the Bone Anchored Hearing System (BAHS) and the Cochlear Implant (CI). The development and commercialization of the BAHS implant are carried out in Sweden, while the cochlear implant is managed by the Vallauris site (formerly Neurelec) in France.

The cochlear implant and the bone anchored hearing system have different functionalities and indications. Bone anchored hearing systems are primarily indicated for patients with middle ear dysfunction and preserved inner ear structures, with sound transmitted as vibrations through the skull bone. In contrast, cochlear implants are designed for patients with sensorineural hearing loss (SNHL), often related to cochlear damage.

In 2021, Oticon Medical CI had over 280 employees and generated revenue exceeding 38 million euros. The company's main clients include hospitals, private clinics, various distributors, and patients benefiting from the produced hearing aids.

All activities related to the development and market introduction of the cochlear implant system are based at the Vallauris site. More specifically, the following

departments can be found there: Research and Development (R&D), Clinical Team, Logistics, Marketing, Production, Quality, and Regulatory Affairs (RA).

1.2 Neurelec today

In May 21, 2024, Neurelec has officially become a part of Cochlear limited. Based in Sydney, Cochlear was established in 1981 as a subsidiary of Nucleus, with funding from the Australian government. By 2022, the company controlled 50% of the cochlear implant market, with over 250,000 people having received a Cochlear implant by 2015 . Their Nucleus 22 implant gained FDA approval in 1985, becoming the first multi-channel cochlear implant to receive such approval . Despite my internship coinciding with a major transitional period (involving changes to the offices, training courses, etc.), my work was not disturb thanks to the thoughtfulness of my tutors and to our weekly updates to keep track and discuss the advancement and any issue I could have.

1.3 The R&T team

The goal of the **Research and Technology Team** is to discover and test new solutions that can be introduced into the production chain, where they are adapted and implemented by the R&D and embedded software teams. The R&T team is divided into two groups: the *AIR* specialized in Artificial Intelligence, Image processing and Robotics. and the signal processing and acoustic team, which I was a part of .

1.4 Context and interest of the project

The project finds its interest in the augmentation of CI users well being. It was financed by Oticon medical then Cochlear company and aimed to develop a better understanding of why **CI** users have difficulties with the recognition and perception of security relevant sounds and then to leads to the implementation of signal processing and or informatics design in order to improve this perception.

2 Literature Review and state of the art

2.1 Normal functioning of hearing

The hearing process (also called auditory transduction) can be described as follows:

1. As a sound wave occurs, it propagates through the ear canal.
2. The sound wave undergoes mechanical impedance adaptation through the middle ear, causing the eardrum and the ossicles (the small bones) to vibrate.

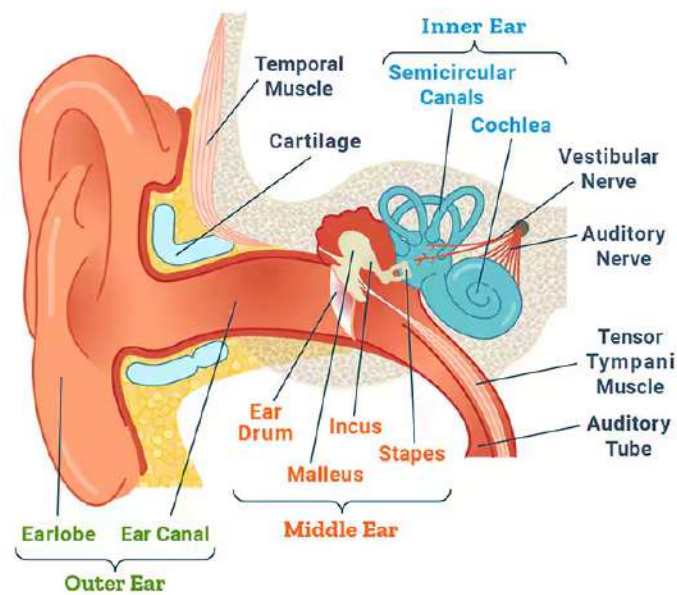


Figure 2: Anatomy of the ear (scheme from life-extension.com)

3. These vibrations cause the fluid in the scala vestibuli and scala tympani to move, making the basilar membrane within the cochlea vibrate and activating the organ of Corti.
4. Depending on the frequency of excitation, a localized part of the basilar membrane with maximal displacement will vibrate, causing the hair cells in the organ of Corti to move. These hair cells are innervated, and upon deformation, they generate neural impulses that are interpreted as sound sensations by the brain.

Within this process, localized parts of the basilar membrane are associated with specific frequencies of excitation. This property, known as tonotopy, largely explains the frequency discrimination capabilities of the ear. Neural processing in the brain leads to psychoacoustic phenomena, such as temporal and spectral masking [16].

There are three main types of hearing loss:

- **Conductive Hearing Loss:** This type occurs when there is a problem with the transmission of sound waves from the outer ear to the inner ear (e.g., when the eardrum is damaged). In such cases, surgery or the use of a **BAHS** ([Link to site](#)) can be effective solutions.
- **Sensorineural Hearing Loss:** This type originates in the inner ear, sensory organs, or the vestibulocochlear nerve and is prevalent according to the Pasteur

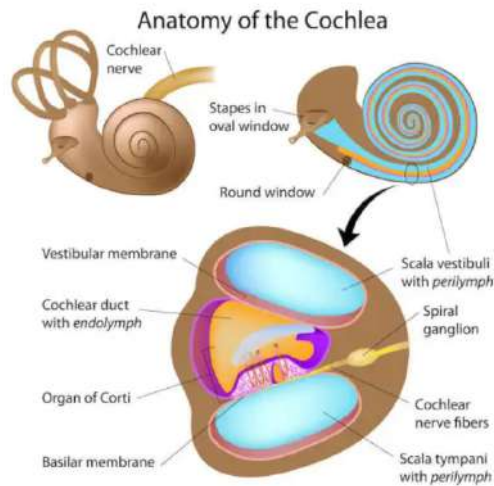


Figure 3: Focus on the inner ear (scheme from Audiocardio.com)

Institute ([Link to site](#)). The causes of sensorineural hearing loss are varied (genetic, trauma, infections, diseases, congenital malformations, etc.) and are not always treatable. Cochlear implants (CI) can address severe cases by bypassing the cochlea and directly stimulating the auditory nerves.

- Mixed Hearing Loss: This type combines aspects of both conductive and sensorineural hearing loss.

Next, we will describe the functioning of cochlear implants.

2.2 Functioning of CIs

CIs are the first examples of neuro-prosthesis that can substitute a sensory organ [20]. To elaborate, CIs allow users to hear without the sound passing through the natural auditory pathway described in the previous section, as the implant replicates the hearing process through its own chain. However, for practical reasons and efficiency in transduction, state-of-the-art CIs tend to compress auditory information and do not fully match normal auditory transduction performance in terms of dynamic range, frequency resolution, and speech recognition [7]. The scheme in Figure 4 illustrates the general signal processing chain of CIs.

This scheme is quite generic and does not capture the full range of techniques and subtleties used in audio transduction. However, some general aspects of the process, also known as coding strategy, can be described:

- When used in acoustic mode (as opposed to streaming mode, where sound input comes from streamed content), the raw output from the microphones is preprocessed to obtain a single-channel digital signal. This preprocessing

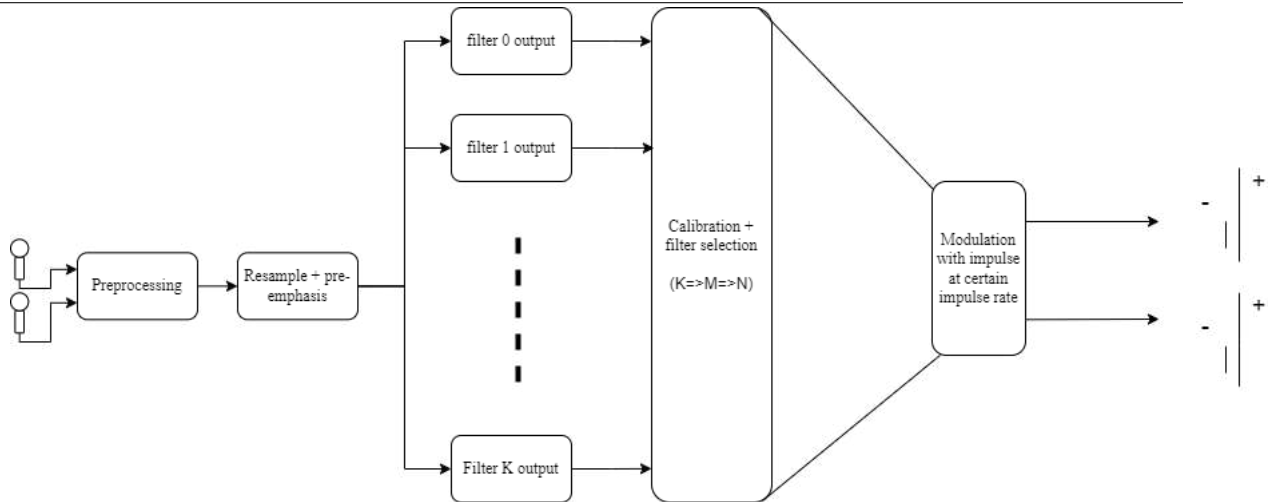


Figure 4: General processing block of CIs

may involve sampling, addition, normalization, [beamforming](#), and some noise reduction.

- The signal is then resampled and passed through a pre-emphasis filter designed to mimic the natural mechanical filtering of sounds as they pass through the ear.
- The next stage involves band-filtering the signal. Depending on the approach, this filtering may use Fourier transforms to achieve a [Short-Time Fourier Transform \(STFT\)](#), band-pass filters on short audio frames to mimic the critical bands of human hearing, or other types of band-pass filters.
- A selection algorithm is then applied to choose the "bands" of frequency to encode for each frame, which are used to stimulate the cochlea via a dedicated electrode placed at the tonotopically relevant location. Note that this description does not fully capture some coding strategies like MPEAK or F0/F1/F2 [18]. We focus on the most common schemes used in recent implementations. The number of simultaneously stimulable electrodes, referred to as *Nofm*, is a variable parameter.
- The final step involves extracting the signal envelopes and modulating these envelopes with impulses generated at a constant rate. The amplitude in a given frequency band is linked with the amplitude of stimulation in the electrodes. Other encoding schemes exist; for example, Oticon Medical CIs encode amplitude using variable-duration pulses with constant amplitude. To minimize confusion due to [electrical spread](#), pulses from different channels are subtly delayed to ensure that neurons receive the appropriate amount of stimulation. This strategy is known as [CIS](#) [18], which differs from CIs. (The frequency of

stimulation between consecutive pulses is called the pulse rate; a higher pulse rate improves amplitude envelope transmission. In our experiment, this value is fixed at 500 Hz.)

The first four parts of this process are handled by the "speech processor" of the CI. The resulting signal is then transmitted to the internal transmitter, which excites an array of electrodes placed in the cochlea at tonotopically coherent locations (or at least close to them). Note that the design of the electrode array is crucial in CI design, as it directly affects the [current spread issue](#) and the overall "electrical to neuronal" transduction.

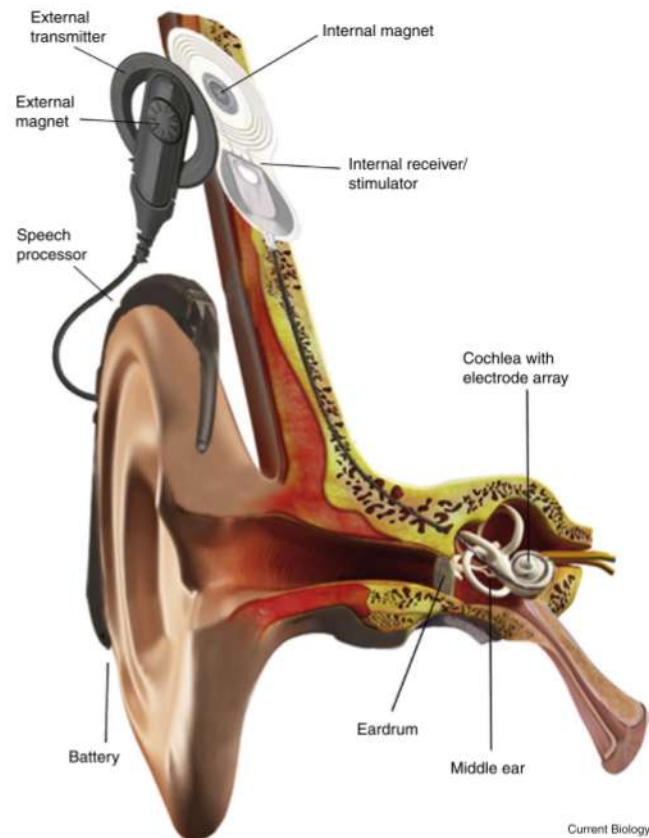


Figure 5: Scheme of an impaired patient's ear (figure from [20])

2.3 Safety related sounds

SRS (safety-related sounds) are a category of sounds carrying relevant security information for humans. Examples include alarms (especially fire alarms), human and animal vocalizations, and mechanical sounds from building sites or vehicles. We

propose a categorization based on four axes: (1) the safety rating on a scale from 1 to 5 (where 5 represents maximal danger and 1 represents no danger), (2) the need to know the source position of these sounds to ensure safety, (3) the difficulty of generating samples from a signal processing perspective, and (4) the difficulty of obtaining a database of real samples. The summarized information is shown in the following table.

	Difficulty to Obtain a Database	Difficulty to Create Samples	Need for Spatialization	Safety Rating According to [19]
Alarms	Easy	Easy	No need	4.43
Human Vocalizations	Probably Exists	Hard (Requires DL)	Need + Tracking	≈ 2.88
Animal Vocalizations	Probably Exists for Common Animals	Hard (Requires DL)	Need + Tracking	≈ 3.34
Mechanical Sounds	Uncertain	Medium to Hard	Need	≈ 4.0
Gun Sounds	Probably Exists	Easy	Need	4.62
Vehicle Sounds	Probably Exists	Medium to Hard	Need + Tracking	≈ 4.5

(\approx) is used when the category is not explicitly addressed in [19]. The score is then estimated based on the score corresponding to sounds within the category.

These sounds have specific acoustic features that make them distinguishable using signal processing methods:

- Alarms, according to ISO 7731, should have a loudness at least +15 dBA above the ambient sound environment, frequency components between 500 Hz and 2500 Hz with a main frequency between 500 and 1500 Hz, and be modulated in amplitude with a pulse rate of 0.5 to 4 Hz. Generally, they are composed of harmonic-structured sounds. Additionally, [11] notes that some acoustic features of alarms with background noise, compared to background noise alone, are significantly correlated with a "clearly audible" human assessment (see Figure 6).
- Human vocalizations (screams, cries, rattles, etc.) are quite challenging to characterize. According to [1], the distinction between different vocalizations appears to be independent of language, although their study only includes three languages. It seems reasonable to assume that human vocalizations (such as screams and laughs) are understood and recognized across cultures. According to Table 1, key features for differentiating vocalizations include: pitch, standard deviation of inter-burst intervals (time between amplitude peaks in

the vocalization), and the first quartile of spectral energy distribution (i.e., the frequency below which a quarter of the spectral energy is gathered).

- More generally, concerning environmental sound differentiation (and not only SRS), [6] Table 3 shows that among various types of audio features, spectral and cepstral features such as MFCC and brightness are most useful for classifying environmental sounds (based on the BDlib-2 library). This is corroborated by [5], which also includes temporal cues such as zero-crossing rate. Figure 7 shows the ranking of the 10 best features for separation in [5] and [6]. However, it should be emphasized that the models used in these studies are quite far in performance from state-of-the-art models. Thus, the most relevant features for these simpler models may not be the ones allowing the best differentiation in the most advanced models.

In conclusion, the recognition of SRS and environmental sounds strongly depends on spectral and cepstral (MFCC) analysis, and to a lesser extent on temporal analysis.

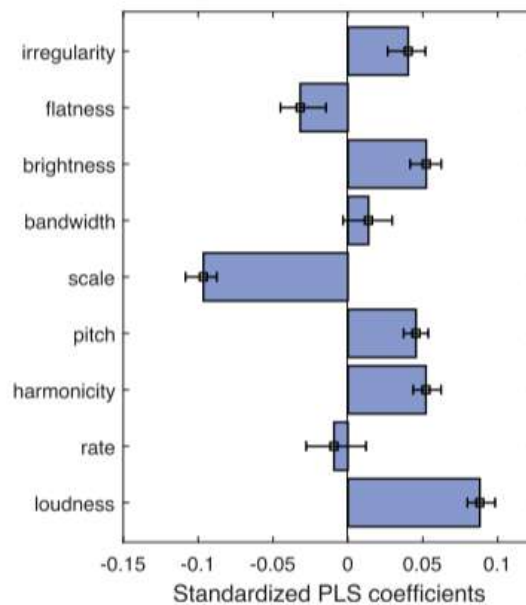


Figure 6: Coefficient of the first PLSR component. Acoustic features whose coefficients have bootstrap distributions above or below zero are considered statistically significant predictors. Figure from [11].

2.4 The Problematic Issue of Relevant Security Sounds for Cochlear Implant Users

As discussed in Section 1.2, **cochlear implants (CIs)** encode sound into electrical stimulation using one of the coding strategies and material configurations detailed in [7]. Initially optimized for processing speech [8], these implants have evolved to more generic designs [18], though they are still primarily evaluated based on speech perception.

Additionally, as mentioned in Section 1.3, **acoustic features** that are useful for differentiating environmental sounds depend on spectral (e.g., pitch, brightness, and frequential centroid) or cepstral analysis (e.g., MFCC, LPCC). Due to the constraints of CI design—such as having a limited number of electrodes or using broad band filters to decompose input sounds—CIs often evaluate spectral content imprecisely and, consequently, the cepstral content as well. This limitation likely contributes to the difficulties CI users face in recognizing **speech-relevant sounds (SRS)** and other environmental sounds.

Despite these challenges, some CI users, according to [21], continue to use their implants frequently (more than 5 hours a day) even without significant improvements in speech processing (see Figure 8). This suggests that obtaining environmental sound information may be perceived as a substantial improvement in their overall perception of the world. The following figure, extracted from [21], highlights this trend.

Each line in the figure represents a subject's daily CI usage, showing their average usage at the most recent follow-up and three months prior. It is clear that a significant majority of users maintain long-term use of their implants despite poor **speech intelligibility ratios (SIR)**, likely because the ability to perceive environmental information is considered a substantial benefit, even though it does not match normal hearing perception.

In [23], the authors demonstrate that CI users have a reduced capacity to discriminate natural soundscapes across three modalities: habitat, time of day, and seasons, compared to normal hearing individuals. This limitation could potentially impact their safety.

Furthermore, [19] shows that CI users struggle to recognize certain security-relevant sounds (e.g., only 38% recognition accuracy for explosions in a quiet context, despite high confidence in the relevance and their knowledge of the sound after brief training).

Therefore, it can be hypothesized that achieving satisfactory perception of security-relevant sounds remains a challenge for CI users, and improving this aspect could provide significant benefits.

2.5 Coping Strategies for Integrating Security Sounds

To our current knowledge, there has been limited work on the integration of relevant security sounds for cochlear implant (CI) users. Currently, these sounds are often treated as noise in relation to speech, as evidenced by systems such as Oticon's SuddenSound stabilizer [Santurette2023] and Apple's Conversation Awareness with Loud Sound Reduction [Apple2023]. Additionally, even without considering the issue of removing [speech-relevant sounds \(SRS\)](#) from audio transduction, the signal processing aspects that contribute to the salience and audibility of SRS for CI users remain unclear. One of the goals of this work is to identify these components to develop an efficient response mechanism to alert users.

Table 3. Top 20 features for the SFb method.

Rank Order	Ranking Algorithms		
	IGR	GLM	SVM
1	MFCC1	MFCC3	MFCC3
2	ENTR	MFCC1	MFCC2
3	BRI	MFCC2	ENTR
4	ROL	SFM15	MFCC1
5	MFCC2	MFCC4	BRI
6	CEN	ZCR	CEN
7	SKEW	VAR	VAR
8	SFM12	SFM18	MFCC4
9	KURT	SFM14	SMO
10	SFM10	ROU	ROL

(a) Best features for environmental sound classification depending on the model of classification from [6].

Table 2. Ranking of the salient features

Rank Order	Ranking Algorithms	
	InfoGainAttributeEval	OneRAttributeEval
1	Brightness	ZCR
2	ZCR	SFM_12
3	SFM_12	MFCC_2
4	MFCC_1	Brightness
5	MFCC_2	LPCC_12
6	SFM_13	MFCC_1
7	Spectral Smoothness	SFM_10
8	LPCC_12	LPCC_11
9	MFCC_3	LPCC_10
10	SFM_11	Spectral Smoothness

(b) Best features for environmental sound classification depending on the method of relevance evaluation from [5].

Figure 7: The best features for environmental sound classification on BDlib.

3 Analysis of SRS in CI Processing

To address the issue of SRS recognition, we analyze the effects of CI audio transduction by comparing [electrodograms \(EDG\)](#) and [short-time Fourier transform \(STFT\)](#), and by examining important [audio features](#). Initially, we focus on isolated SRS (due to time constraints, we worked with alarm signals) and then on SRS combined with environmental noises (in this case, we used alarms signals with specific morphed white noise, which will be described in the relevant section). Our analysis supports two claims about CI transduction that may explain the poor alarm detection and recognition experienced by CI users:

- Some spectral modulations of alarms are above the detection thresholds of CI users, making it difficult for them to recognize these alarms.
- Features important for the salience and audibility of alarms, according to [11], are degraded throughout the processing chain.

3.1 Data Collection and Process of Analysis

To begin the study, the first task was to collect relevant data. In our case, these data are sound files of [SRS](#) and environmental sounds (such as speech and vocalizations from humans and animals). We chose to use sounds from peer-reviewed databases such as:

- The Acted Emotional Speech Dynamic Database [24] [29]. This database contains speech (Greek in this case) expressed with 5 different emotions (anger, disgust, fear, happiness, and sadness). We consider speech expressing anger, disgust, and fear as SRS, and the other emotions as environmental sounds.
- The Montreal Affective Voices [4]. This database contains human vocalizations pronounced with the vowel "A," where subjects try to reproduce vocalizations associated with various emotions (such as anger, happiness, and fear). We use the stimuli associated with negative emotions as SRS and the others as environmental sounds.
- The audio used in the Luzum sound recognition study [19]. These are "every-day sounds" (vehicles, animals, objects), some of which are classified as SRS and others as environmental sounds.
- The BDlib-2 [5] [6]. This database includes 20 categories of environmental sounds such as rain, thunder, airplane, and applause. Sounds indicating danger (such as thunder, alarms, and screams) are categorized as SRS, while the rest are classified as environmental sounds.

- The VoxCeleb dataset [9], which contains speech pronounced by celebrities. This entire dataset is used as environmental sounds to construct babble noise, for example.
- The ESC-50 [26]. This database consists of a large variety (50 categories) of environmental sounds, such as keyboard typing, fire crackling, and brushing teeth.

All these sounds were sorted according to their relevance to safety. We classified sounds carrying security information (such as alarms, human vocalizations with negative emotions, and vehicle sounds) as signals, while all other sounds were classified as noise.

As mentioned earlier, we only worked with alarm sounds as SRS in this study, leaving the investigation of other SRS for future work.

To create the test sounds, denoted as stimuli used in our investigation, we followed the processing scheme shown in Figure 9.

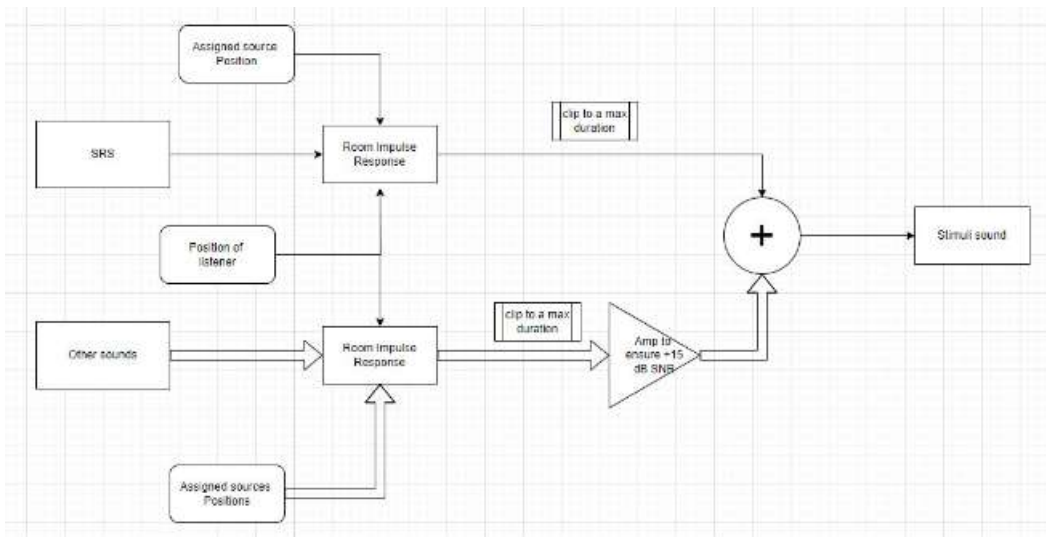


Figure 9: Process to generate the stimuli sounds

Thick arrows indicate potential multiple sources/sounds (i.e., in this process, a stimulus sound contains at most one SRS and multiple noise/environmental sounds). The general process involves applying a room effect to the SRS and the noise based on the algorithm described in [22] (an implementation of the image/source model of room impulse response). The amplitude of the noise is then controlled to obtain a stimulus sound with a specified SNR (calculated over the total duration of the signal).

Since the [22] algorithm uses the localization of the source/microphone, room dimensions, and acoustic reflections (same for all surfaces) to calculate the RIR, and uses the T60 as an input, we fixed the room dimensions and used the Eyring

formula to calculate the reflection coefficient corresponding to the desired T60. The same positions and room dimensions were used for all signals.

$$1 - C_r = \exp\left(\frac{0.163 \cdot V}{-T_{60} \sum_i S_i}\right)$$

The general nomenclature of these sound files was:

$$\textit{Signal} - \textit{XX.wav} - \textit{with} - \textit{Ys} - \textit{T60} - \textit{SNR} - \textit{Z}$$

where XX is the name of the original wave file, Y is the T60 value applied to the file, and Z is the SNR. In this first stage of analysis, only alarms without noise were studied (i.e., with infinite SNR).

3.2 Analysis of the SRS Processing Alone

Having correctly created those sounds, we proceed to our first analysis. We use Oticon Medical's CI process chain to obtain the [EDGs](#) of these signals and compare them with the spectrogram of the input. The goal is to gain an initial understanding of the influence of various parameters and the CI on the acoustic information.

3.2.1 Analysis of Electrodiagrams

The EDGs were created by varying three parameters:

- The T60 of the signal (i.e., the T60 used to create the [RIR](#) filter that processes the input signal). The values tested were 0s (simulating a free field situation), 2s (simulating a typical room environment), and 15s (simulating an over-reverberated environment such as a church).
- The Input Level, a parameter from the CI process chain that normalizes the sound to simulate different dB SPL levels. The tested values were 50 dB SPL (quiet environment level), 70 dB SPL (conversation level), and 85 dB SPL.
- The NofM, another parameter of the chain, which is the maximum number of channels (i.e., electrodes) that can be stimulated simultaneously (due to [CIS](#), "simultaneously" refers to a certain frame of stimulation rather than absolute simultaneity).

An example of these [EDGs](#) is shown in [Figure 10](#). In these plots, the spectrogram is shown in the background (with darker colors representing higher amplitudes). The [EDG](#) is superimposed in orange, with each line representing a channel (i.e., an electrode) and aligned with the center frequency of its channel. The amplitude is

represented by the thickness of the line, and you may notice that these plots are often superimposed.

The first observation is that the cochlear implant drastically reduces spectral resolution compared to normal hearing listeners, who have up to 144 channels between 100 and 8000 Hz (considering 6 octaves and 24 steps per octave). In practice, most CI users have fewer than 20 independent channels due to their neurological condition and [current spread issues](#).

A second observation is that the [EDGs](#) have a minimum threshold to trigger an amplitude (notably visible in [Figure 10\(a\)](#)) and a maximum amplitude, above which higher levels are clipped (as seen in [Figure 10\(c\)](#) compared to [Figure 10\(b\)](#)). These maximum and minimum levels are known as C and T levels (C for comfort, the maximum stimulation level, and T for threshold, the minimum stimulation level). These levels depend on the patient and the specific channel and can change over time. In summary, the dynamic range of CI users is significantly reduced by the process chain to match the user's capabilities.

Lastly, the [EDGs](#) exhibit good temporal resolution, with a temporal frame typically consisting of 128 samples at a 16 kHz sample rate (i.e., a 2 ms frame).

Regarding parameter influence, reverberation has the most noticeable impact, blurring frequency patterns over time, particularly in high frequencies (above 2290 Hz), which can decrease alarm recognition due to the loss of harmonicity. Other parameters have more limited effects. Input levels increase EDG stimulation amplitude until clipping occurs and can trigger additional channels if the stimulation level exceeds the T levels. The NofM parameter modifies the number of available channels in a given frame but does not necessarily cause different [EDGs](#) when it varies.

At this point, we aimed to study the most correlated feature with audibility as discussed in [11]. Specifically, we wanted to compare the features of input sounds with their "encoded version" in the EDG. However, we encountered a methodological challenge. While the [EDG](#) provides an objective representation of sound and probably close to the perception for CI users, it is not a unique representation of sound, in particular concerning spectral feature evaluation. Therefore, we decided to use vocoders to reconstruct sounds from the EDG and compare the features of these vocoder outputs with those of the input sounds.

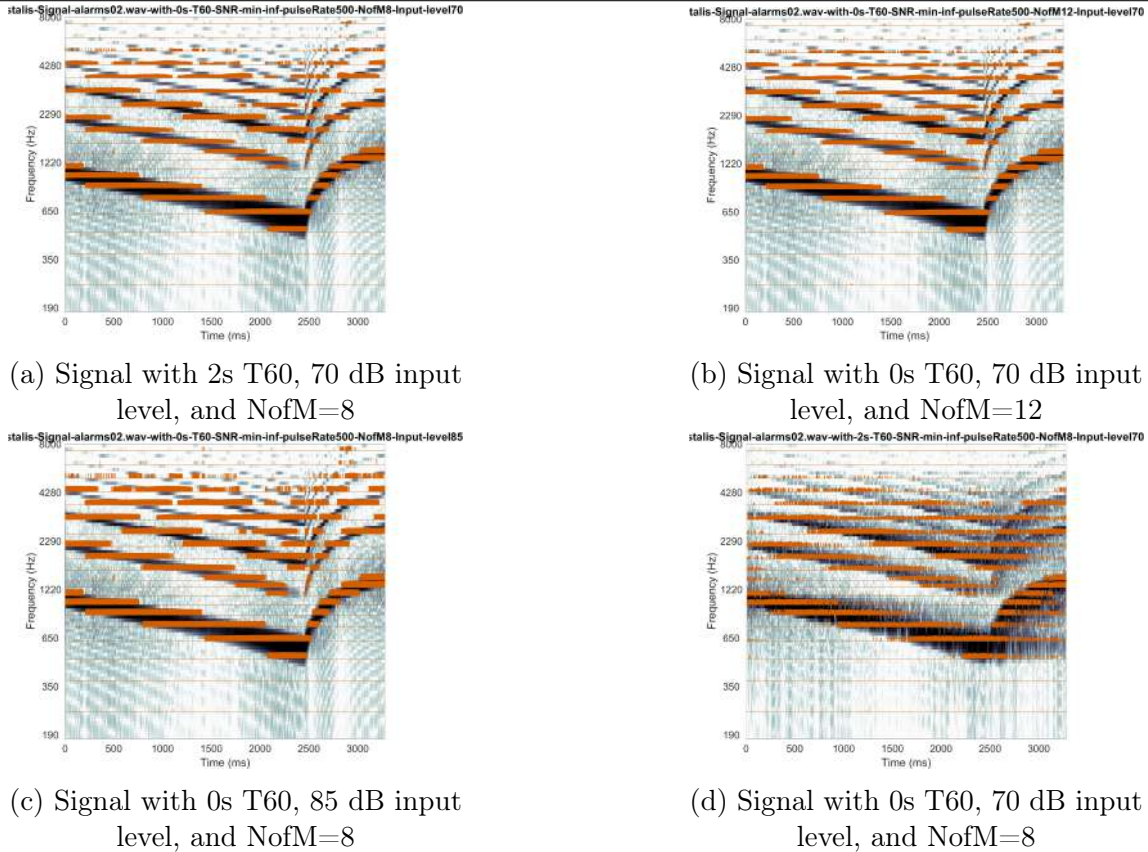


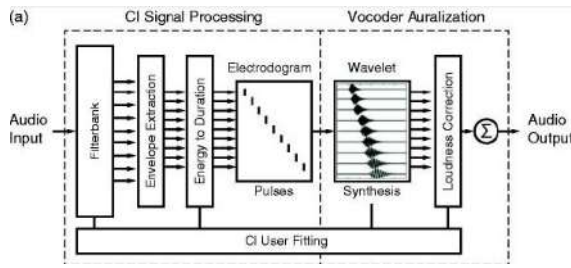
Figure 10: Examples of EDGs superimposed with the spectrogram

3.2.2 Description of the Vocoders

We used two different vocoders during our analysis: the sine vocoder and the spiral vocoder.

These vocoders take as input the EDG, which can be seen as an N-channel envelopes array in [Figure 11 \(b\)](#), and compute the audio auralization from it. The process consists of extracting the amplitude envelopes for each channel and multiplying them with a sine wave at the frequency equal to the center frequency of the corresponding channel. Note that this is done under the assumption that when a channel is triggered, it is only by a sine wave with the center frequency of that channel, which is not always the case in our situation. The sine waves are then summed across the different channels.

The spiral vocoder differs in that it adds an envelope mixing stage where neighboring channel envelopes are combined with a decremting coefficient as they become further neighbors. This additional stage aims to mimic the [current spread phenomenon](#), providing a closer representation of what CI users may hear. Consequently, we expect the correlation score of spiral/input features to be lower than the correlation score of sine/input features. We may also observe some effects of



(a) Scheme of the sine vocoder from [3]

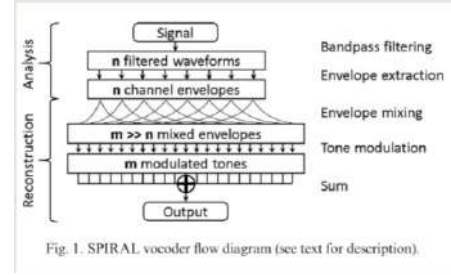


Fig. 1. SPIRAL vocoder flow diagram (see text for description).

(b) Scheme of the spiral vocoder from [13]

Figure 11: Concept schemes of the vocoders

spreading on the degradation of features.

3.2.3 Analysis of the Spectral Modulation

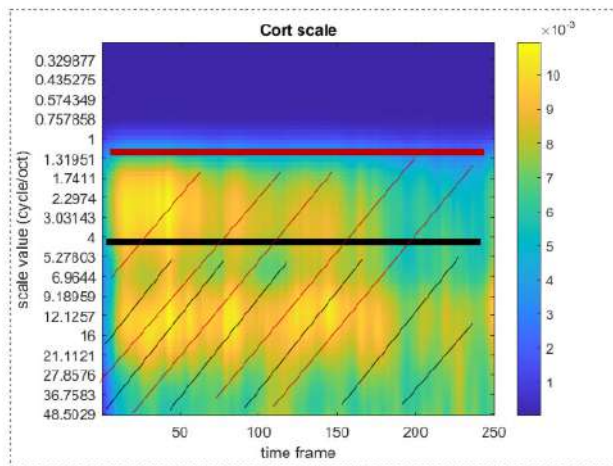
Among the features we studied (to be described in the next subsection), the scale feature was particularly important. The perception of spectral modulation has already been studied concerning CI users, notably through the use of the SMRT. We aimed to investigate whether part of the alarm's frequency modulation might be imperceptible to CI users.

To verify this, we calculated the frequency modulation spectrogram for all our signals (using the NSL Matlab toolbox). Knowing the detection threshold from previous studies, such as [15] (around 9 ripples per octave (RPO) for adults with normal hearing and around 4 RPO for CI users), we calculated the proportion of modulation components above the threshold divided by the total number of components for all time frames and signals. Figure 12 shows an example for one signal. We then obtained the rate of imperceptible modulation for CI users compared to normal hearing subjects.

The calculation follows the formula:

$$\text{Not perceived energy ratio}(t) = \frac{\sum_{n=\text{threshold}}^{\max \text{ RPO}} \text{spec modulation}(t,n)^2}{\sum_{n=\min \text{ RPO}}^{\max \text{ RPO}} \text{spec modulation}(t,n)^2}$$

We obtained a mean value of 7.97% for normal hearing subjects and 23% for CI users, with a standard deviation of 3%. This score even rises to 45.29% when calculated from internal clinical values of Oticon Medical (mean RPO at 2.145), validating our initial claim. See [link to claim](#).



In red: the example threshold of a CI USER
 In black: the example threshold of a normal hearer.

13

Figure 12: Example of frequency modulation spectrogram. The red hashed area represents the part of the modulation spectrum not perceived by CI users, while the black area represents that perceived by normal hearing individuals. Threshold values are arbitrary, just to illustrate the discussion.

3.2.4 Analysis of CI Vcoded Sounds Features vs. Actual Sounds Features

Let's first describe more precisely the set of features we will study:

- Loudness: Corresponding roughly to an input level defined on a time frame, extracted using the Zwicker model. Calculated from the MATLAB function 'acoustic-loudness'.
- Harmonicity: A value describing the proximity of a spectrum to a template spectrum representing a totally harmonic sound (see [14] for more details). Due to inconsistent results with our implementation following Pr. Ehlilali's original code, we instead used the harmonic-ratio measure from MATLAB, estimating this ratio using the autocorrelation method.
- Pitch: In our case, this refers to the F0 value. Due to inconsistent results with our implementation following Pr. Ehlilali's original code, we used the autocorrelation estimation method from MATLAB.
- Scale: The centroid along the frequency modulation axis of a cortical representation of frequency modulations, calculated using the [NSL MATLAB toolbox](#).
- Brightness: The frequency centroid of a spectral profile, calculated using the formula:

$$BR(t) = \frac{\sum_f f \cdot spec(t,f)^2}{\sum_f spec(t,f)^2}$$

- Flatness: Calculated as the geometric mean of the spectrum divided by the arithmetic mean:

$$FL(t) = \frac{N \sqrt{\prod_{f=0}^{N-1} spec(t,f)}}{\sum_{f=0}^{N-1} spec(t,f)/N}$$

- Irregularity: Extracted as a measure of the difference in strength between adjacent frequency channels, defined as:

$$IR(t) = \frac{\sum_f (spec(t,f+1) - spec(t,f))^2}{\sum_f spec(t,f)^2}$$

These features were chosen because they are, according to [11] Figure 6, significant predictors of audibility attributes for normal hearing individuals. Thus, checking their degradation can help explain the lack of detection and recognition of alarms by CI users. They were calculated on 256-sample frames with a 128-sample overlap, aiming to use analysis parameters close to the sound processor's capacity to facilitate potential future feature estimation on this processor and to provide an initial idea of the precision to be expected.

We chose the Pearson correlation score as the metric to evaluate degradation, defined as:

$$Corr\ score(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

We hypothesize that there is a monotonic relationship between the correlation score of our input and vocoded features and a perceptual alarms audibility score. Hence, the correlation score of features provides an indication of the perceptual loss experienced by CI users (an experiment to verify this hypothesis is proposed in the Clinical Protocol section).

Additionally, we included two new parameters in our analysis:

- Number of Deactivated Electrodes (values: 1, 4, 8, 12, 14, 17). Until our hypothesis about the existence of a mapping is validated (or not), we wanted a measure of feature perception degradation by the CI. Since the effect of electrode deactivation has been more studied, we expected to establish anchors for each feature and make comparisons such as: "The mean degradation (correlation score) of the distribution of feature N is comparable to the mean

degradation of this distribution for M deactivated electrodes." Unfortunately, as we'll see, the CI behaves differently. Electrodes are deactivated "sparsely," meaning they are spaced out consistently (e.g., for 1, the 10th electrode is deactivated; for 2, the 7th and 14th electrodes, etc.).

- Compression (values: "auto", "loud", "medium", "quiet"). This parameter sets the dB HL of compression knee points for each selected channel, as shown in Figure 13. When set to "auto", the knee points are calculated automatically based on RMS. We expected that having the largest dynamic range possible would improve feature conservation (i.e., that the "auto" setting would give the best results concerning feature correlation scores).

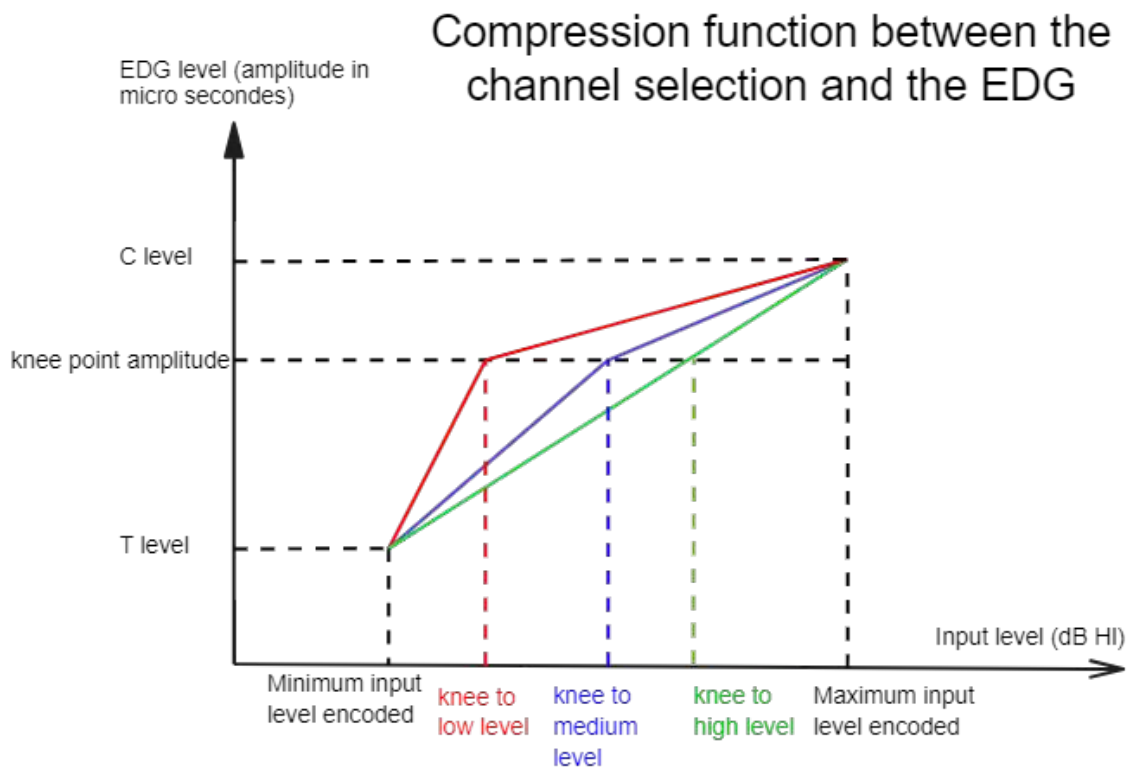


Figure 13: Example of compression function for one channel

We analyzed a total of 20 different alarm sounds, resulting in 432 configurations (20 × 6 × 4 × 2 × 3 × 3 = 8640 different sound files processed). The following table summarizes the different configurations tested:

Sound File Name	Input Level (dB)	T60 (s)	Nofm	Number of Electrodes	Compression
XX.wav	50/70/85	0/2/15	8/12	1/4/8/12/14/17	auto/high/medium/low

Figure 14 describes the process of calculating the correlation score starting from an audio file. Different colors emphasize the type of data: red for audio signal, green for EDG, blue for audio features, and purple for correlation scores.

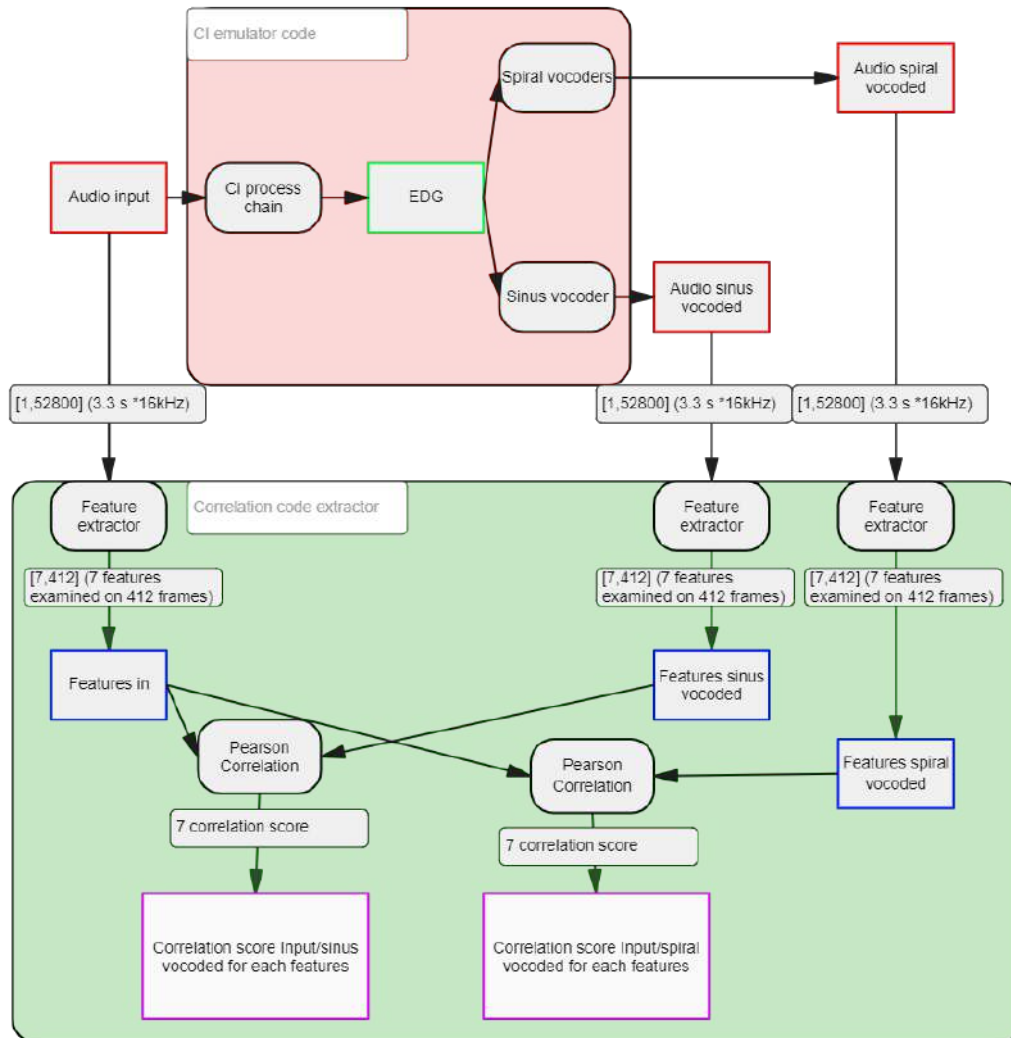


Figure 14: Process of correlation score calculation

Note that, in addition to the 7 features mentioned, we also used the STOI [28] measured between input and vocoded sounds (interpreted as a global degradation measure of the sound and of acoustic features relevant to speech perception by the vocoder).

3.2.5 Results of the Feature Correlation Analysis

We aimed to observe the global effect of the CI on acoustic features. For the 8640 files, we calculated the correlation scores of the inputs and outputs, resulting in the data shown in Figure 15.

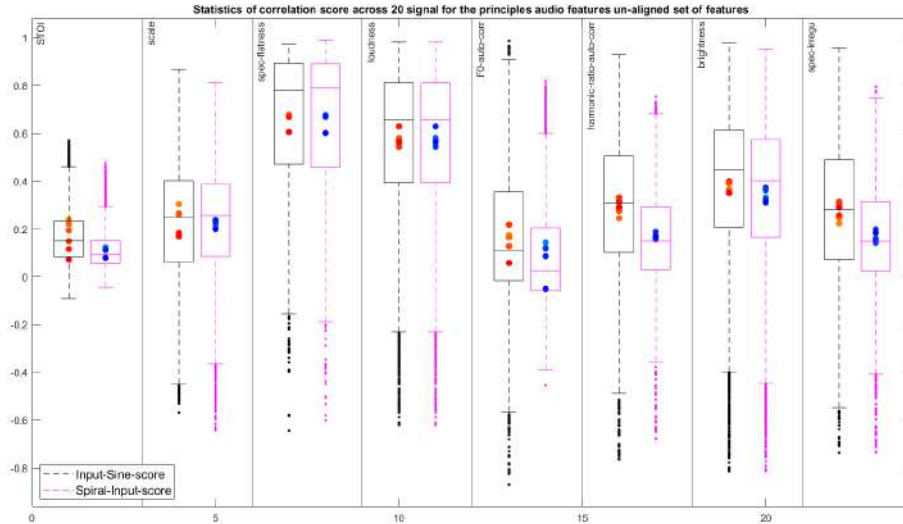


Figure 15: Statistics of correlation scores of features according to the vocoder. The Y-axis represents the correlation score, and the X-axis distinguishes the distributions based on vocoder type and feature. The features are referenced in the top left part of different columns.

As expected, the spiral-input feature correlation scores are generally lower than the input-sine correlation scores. Additionally, the data for the loudness feature are identical because the loudness of the output cannot be computed from the vocoder's output, primarily due to rescaling during the inversion of the EDG to N-channel components. The "CI loudness" is estimated by summing the pulse amplitude of the EDG over the time frame processed and the channels. The CI loudness is calculated from the EDG using the following formula:

$$loudness(t) = \sum_{channels} \sum_{n=t_0(t)}^{t_{end}(t)} scaled_edg_charge(chan, n)$$

where $t_0(t)$ and $t_{end}(t)$ are the first and last sample numbers of the time frame t . The scaled EDG charge is obtained by subtracting the maximum accepted value (C level) for each channel. We use a standard C and T levels set of parameters for all configurations. This calculation is consistent across both vocoders.

Furthermore, we observed some configurations with negative correlation scores. Initially, we thought this issue could be resolved by artificially synchronizing the

input and vocoded features, but this did not yield a significant difference. We assume this is part of the CI's effect on acoustic features.

Apart from spectral flatness and loudness, the correlation scores of other features are generally close to 0. We suspect that the lack of frequency channels is the reason, as these features depend on spectral content. The higher scores for flatness and brightness can be attributed to the fact that they are mean values, making them less sensitive to blurring compared to F0, which theoretically depends on a unique frequency component.

At this point, we have supported our second claim (assuming a link between the loss of correlation score between features and poor audibility of CI users for alarms) [link to claim](#). The most intriguing aspect is the behavior of the dots.

The red and blue dots represent the means of the correlation scores of features for a fixed number of deactivated electrodes. From light to dark blue/red, there are 6 points representing means with 1, 4, 8, 12, 14, and 17 electrodes deactivated. The idea was to use these points as anchors to interpret correlation scores in terms of perception. Indeed, the effect of deactivating electrodes on perception has been studied for CI users [12] [10] (although most of these studies focus on speech perception). We expected to see higher scores with fewer deactivated electrodes (as adding electrodes increases information). This behavior is observed for the STOI and scale features. For spiral-input correlation scores, the order of these anchors is more uncertain, as eliminating channels results in a loss of information but also affects the interaction between channels.

Since the input-sine anchors for other features did not behave this way, we suspected interactions between parameters and the number of electrodes. For example, interactions between NofM and the number of deactivated electrodes are possible; above 12 deactivated electrodes, NofM makes no difference because there are fewer than 8 available channels.

We then conducted a repeated 5-way ANOVA for each feature, with the dependent variable being the correlation score and the parameters including NofM, Input Level, T-60, Compression, and the number of deactivated electrodes.

The ANOVAs (see [Figure 16](#) and [Figure 17](#)) revealed that some parameters indeed interacted, and the nature of these interactions differed depending on the vocoding type (green indicates $p < 0.05$, i.e., the null hypothesis is rejected, while red indicates $p \geq 0.05$, i.e., the null hypothesis failed to be rejected).

The standard error bars of these plots were calculated according to [25]. Due to the number and complexity of the plots, general conclusions about interactions across features are challenging. To verify and observe the specific interactions, we plotted the standard errors with mean plots for the relevant subgroups. For clarity and conciseness, we present and describe four examples below (all these distributions were found by our ANOVA to have at least one significantly different mean distribution, with $p < 0.05$).

We clearly see that input/sine STOI scores are more organized and follow a

	STOI	scale	loudness	F0 auto corr	harmonic ratio auto corr	brightness	spec flatness	spec irregu
Compression	0.349238-02	1.147958-08	0.852918-14	0.0178829-13	1.05498E-05	1.14988E-05	2.89949E-14	0.0072848-09
input_level	0.10201206	0.04071037	1.023496-08	0.002618779	0.4267867-02	0.011203204	0.522120E-11	0.2318762-06
T_60	1.15121E-06	0.00681965	1.10949E-08	0.011967046	0.007112348	0.005013650	0.00050306	0.0427268-01
NofM	0.50818E-05	0.04281187	0.00400121	0.0174384-08	0.4668762-05	0.4304811-05	0.4449718-07	0.0585802-05
deactivated_electrodes	1.09704E-08	0.00341E-11	0.002028181	0.0413065-07	0.002842319	0.0750594-02	1.77147E-12	0.0492970-07
Compression:input_level	0.04082401	0.20948879	0.00846E-13	0.112814628	0.030137897	0.009491418	1.00213E-14	0.2314302-06
Compression:T_60	0.000174208	0.01184281	0.0000579174	0.008090444	0.34487E-14	1.1400E-10	0.1434817-03	1.10249E-22
Input_level:T_60	0.202138383	0.00416957	0.018748013	0.061100101	1.44746E-11	0.10082059	0.34476E-21	0.0044741
Compression:NofM	0.419661051	0.00614119	0.00000088	0.266227019	0.200513182	0.110003176	0.7150169-06	0.0006048
Input_level:NofM	0.60575E-07	0.05481201	0.000143447	0.0870255-08	0.074141311	0.483874001	0.0218899-01	0.01713064
T_60:NofM	0.31803390	0.13018009	7.70739E-03	0.002027964	0.073481413	0.1077989-03	0.57708E-18	0.01760087
Compression:deactivated_electrodes	0.32051E-07	1.20051E-14	0.4091882	0.001707795	0.12871E-17	0.000000001	0.000000000	0.013861024
input_level:deactivated_electrodes	0.001010396	0.4151E-13	4.46174E-28	0.04621089	0.012123134	1.51444E-00	0.101502009	0.01191E-08
T_60:deactivated_electrodes	0.14985E-04	7.70746E-04	0.004038657	0.81133E-05	0.000514385	4.40465E-00	2.38010E-21	1.70988E-08
NofM:deactivated_electrodes	1.0583E-15	0.00221094	7.64873E-07	0.441911901	0.64488280	0.486237597	0.4660997-07	0.045381006
Compression:input_level:T_60	0.11867711	0.11770669	0.00042009	0.004171401	0.000102615	0.000404899	1.84179E-20	0.001671738
Compression:input_level:NofM	0.00564837	0.04884824	0.011415203	0.0439015747	0.07642648	0.7308857-03	0.89018805	0.001431282
Compression:T_60:NofM	0.449029463	0.020412013	0.040312065	0.010612088	0.012776218	0.413891041	0.38001730	0.210801878
input_level:T_60:NofM	0.002012154	0.04210387	0.01020E-05	0.28049146	0.42007208	0.103291002	0.770254448	0.00820066
Compression:input_level:deactivated_electrodes	1.4710E-10	1.15054E-08	0.7501E-03	0.015108204	0.005406667	0.000116964	0.040770513	1.4691E-08
Compression:T_60:deactivated_electrodes	0.47201E-06	0.00449677	1.08617E-05	0.01004753	0.02204981	0.4029E-06	1.0204E-07	0.010049006
input_level:T_60:deactivated_electrodes	0.084032261	0.01701284	0.20779E-09	0.11984176	4.11983E-10	0.1020E-17	0.001010145	0.18699E-22
Compression:NofM:deactivated_electrodes	0.00001479	0.00004088	0.1905E-11	0.180704018	0.100000000	0.110000000	0.370139868	0.011601776
input_level:NofM:deactivated_electrodes	0.10408E-04	0.03177099	1.0233E-08	0.070703207	0.074000750	0.000504204	0.017054937	0.030768993
T_60:NofM:deactivated_electrodes	0.00002748	0.04810266	2.6949E-09	0.228794807	0.00208488	0.17274404	0.06020988	0.414564119
Compression:input_level:T_60:NofM	0.475884108	0.00687781	0.07002104	0.000011644	0.790000000	0.000000000	0.000000000	0.000000000
Compression:input_level:T_60:deactivated_electrodes	0.154402404	0.01242069	1.70818E-08	0.41978970	0.00240180	0.7844E-06	0.4010E-07	7.7091E-07
Compression:input_level:NofM:deactivated_electrodes	1.4703E-12	0.01813683	0.5221E-03	0.260431903	0.06019079	0.962200848	0.70020011	0.98870889
Compression:T_60:NofM:deactivated_electrodes	0.07707820	0.00381944	0.08800017	0.00011644	0.30511138	0.07076042	0.07010014	0.001018429
input_level:T_60:NofM:deactivated_electrodes	0.109101008	0.00204973	0.40018E-09	0.44007800	0.781048218	0.07111882	0.00013110	0.017801261
Compression:input_level:T_60:NofM:deactivated_electrodes	0.06066482	0.71	0.0375214	0.118071802	0.068016471	0.003118872	0.068825306	0.000013058

Figure 16: P-value table of the 5-way repeated measures ANOVA of the correlation input/sine scores

	STOI	scale	loudness	F0 auto corr	harmonic ratio auto corr	brightness	spec flatness	spec irregu
Compression	2.344E-17	0.000447	5.80E-14	0.00332750	0.012063058	1.18E-05	0.000445061	3.0044E-14
input_level	0.71E-07	0.077432	1.03E-06	0.19724654	0.008130613	0.011737	0.612879488	2.102E-09
T_60	0.0018189	0.202954	1.19E-09	0.3429E-06	0.000980177	0.001375	0.01096020	0.2852E-20
NofM	0.4923377	0.2018113	0.00400	0.014249591	0.004424211	0.209215	0.005248503	0.207194411
deactivated_electrodes	5.75E-20	0.1113782	0.002028	0.07802E-12	0.122524329	0.020574	0.018330597	3.0070E-13
Compression:input_level	0.016607	0.024819	2.1E-13	0.017872296	0.023215406	0.009339	0.090403404	0.40177-08
Compression:T_60	0.361594	1.002296	0.003058	0.009758887	0.5314E-00	1.71E-12	5.56801E-14	0.272861386
input_level:T_60	0.072465	0.1341668	0.007168	0.1151084	0.19948E-11	3.23E-05	0.000208094	0.23112E-20
Compression:NofM	0.3111225	0.267359	0.0003	0.071187723	0.076841933	0.129727	0.211107294	0.046119568
Input_level:NofM	0.076282	0.1818166	0.000743	0.709889127	0.96653192	0.515361	0.181899117	0.10448309
T_60:NofM	0.390738	0.202819	7.77E-05	0.220578893	0.07806658	0.153582	0.251708294	0.138571039
Compression:deactivated_electrodes	1.14E-09	0.048338	0.457189	0.000198072	0.000196485	0.62E-08	0.010021405	0.308771E-06
input_level:deactivated_electrodes	0.274229	0.009426	4.48E-18	0.042013114	0.006826242	0.48E-10	0.000497143	0.000144478
T_60:deactivated_electrodes	0.008016	0.005872	0.004036	0.027899102	0.1302E-06	2.11E-10	0.000147091	1.37096E-22
NofM:deactivated_electrodes	0.403512	0.517373	7.80E-07	0.017887888	0.469407129	0.437411	0.011433736	0.141672005
Compression:input_level:T_60	0.0037	0.103419	0.002427	0.00259375	0.000017743	0.016784	0.003795403	2.33030E-28
Compression:T_60:NofM	0.717197	0.031798	0.011413	0.007805773	0.003172948	0.239165	0.140240297	0.109487487
Compression:input_level:T_60:NofM	0.500797	0.700145	0.984214	0.002914492	0.014424477	0.101333	0.140022014	0.21502054
Input_level:T_60:NofM	0.244226	0.765701	0.81E-06	0.01518256	0.000089032	0.264563	0.001226023	0.10801721
Compression:input_level:deactivated_electrodes	0.001258	0.027573	0.77E-13	0.51182E-09	0.027081893	2.33E-00	0.40296E-16	0.309357022
Compression:T_60:deactivated_electrodes	0.433345	0.784472	1.29E-05	0.048694048	0.00060278	0.000273	0.000195743	0.01149E-00
input_level:T_60:deactivated_electrodes	0.562088	0.232956	0.24E-09	0.07069277	0.142028019	0.01E-31	0.00261E-00	1.57010E-05
Compression:NofM:deactivated_electrodes	0.17761	0.721027	0.10E-11	0.241149708	0.005029581	0.112401	0.547469592	0.486411402
input_level:NofM:deactivated_electrodes	0.520946	0.185047	1.32E-00	0.017048504	0.00641272	0.430372	0.002154289	0.007438072
T_60:NofM:deactivated_electrodes	0.111009	0.650014	2.56E-00	0.56528292	0.0110205	0.000792	0.11879143	0.074009079
Compression:input_level:T_60:NofM	0.630422	0.073157	0.074923	0.083482129	0.02938734	0.239957	0.394111507	0.156660104
Compression:input_level:T_60:deactivated_electrodes	0.051189	0.384879	1.71E-09	0.009618178	0.002600079	1.2E-13	0.780291-05	0.200959E-09
Compression:input_level:NofM:deactivated_electrodes	0.00209	0.000058	0.00E-06	0.15076017	0.4020E-04	0.348967	0.044118211	0.000103808
Compression:T_60:NofM:deactivated_electrodes	0.003204	0.092273	0.099916	0.041479211	0.005314773	0.040462	0.197201413	0.072772116
input_level:T_60:NofM:deactivated_electrodes	0.466537	0.00942	3.48E-05	0.105042763	0.003344742	0.117025	0.426428117	0.487480504
Compression:input_level:T_60:NofM:deactivated_electrodes	0.029773	0.088838	0.037702	0.009753245	0.036048553	0.320653	0.443955313	0.000617657

Figure 17: P-value table of the 5-way repeated measures ANOVA of the correlation input/spiral scores

simpler pattern compared to input/spiral STOI scores. For example, Figure 18 (b) aligns with our expectations: higher NofM, less reverberation, and fewer deactivated electrodes result in better scores. We also note that after 12 electrodes are deactivated, NofM makes no further difference. Additionally, the amplitude of change for NofM is less than that for T60, and T60's effect is less significant than the deactivation of electrodes.

In comparison, the input/spiral Figure 18 (a) appears flatter regarding the number of electrodes, which could be due to interactions between channels. For input/spiral scores, no significant interactions were found within electrodes/NofM/T60, so we decided not to make comparisons.

Regarding F0 input-sine scores, Figure 19 (b), we do not observe a clear degradation of the correlation score until 17 electrodes are deactivated. It even seems that deactivating 12 electrodes might increase the score. We expected to see a clear interaction between input level and compression, with low input levels and quiet compression providing higher scores compared to higher input levels and louder

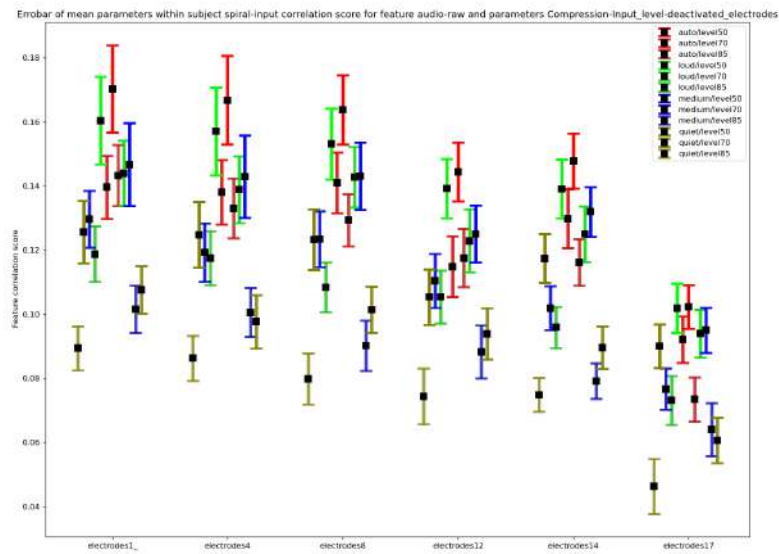
compression, and that the "auto" setting would yield the best results for feature correlation scores. These results might be explained by the fact that F0 is generally independent of input level or compression. Additionally, the alarm signals used had variable amplitude and thus variable input levels, which could introduce variability into the results.

Finally, for [Figure 19 \(a\)](#), we observe a significant decrease in score with 17 deactivated electrodes and a flattening of scores for other numbers of electrodes. The best correlation values changed from "auto" to "medium" and "loud" values depending on the number of electrodes, which is surprising but might result from data variability (note that we conducted a repeated measures ANOVA on 20 sounds across 432 configurations, so more varied sounds might yield stronger results).

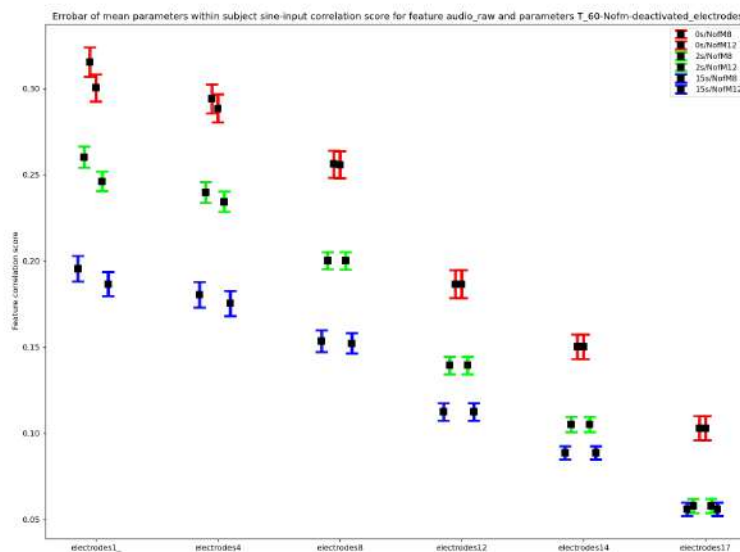
3.3 Analysis of SRS with Highly Masking Noise

The purpose of this part of the work was to analyze the behavior of CI audio transduction when noise is added. Due to time constraints, we focused our study on alarm sounds with added masking noise shaped according to the alarm spectrogram. This scenario represents a "worst-case" situation concerning environmental noise, as the noise is designed to effectively mask the alarm signal.

The shaping of the noise involves estimating the average spectrum of the input sounds over 30 ms frames, resulting in a list of spectra, each corresponding to a specific time frame. The masked noise is then constructed by filtering white noise with the estimated spectrum.

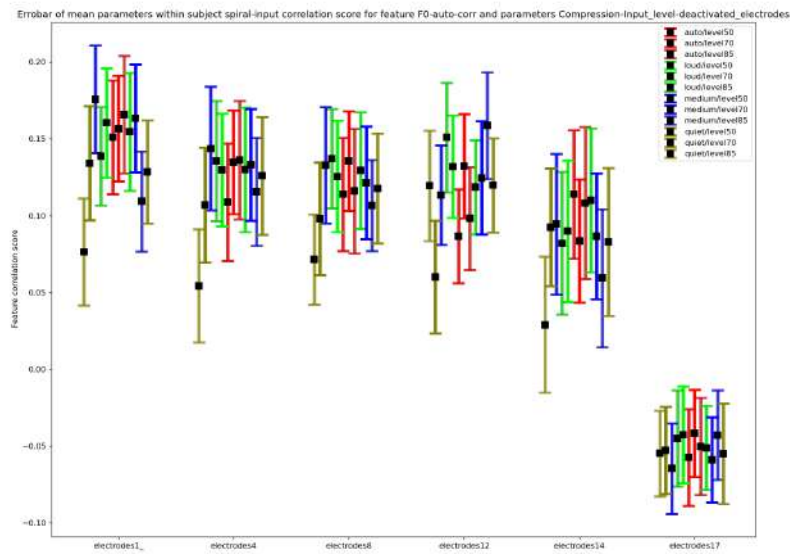


(a) Error bars of the mean for parameters compression:input-level-electrodes. The dependent variable is the spiral/input STOI score.

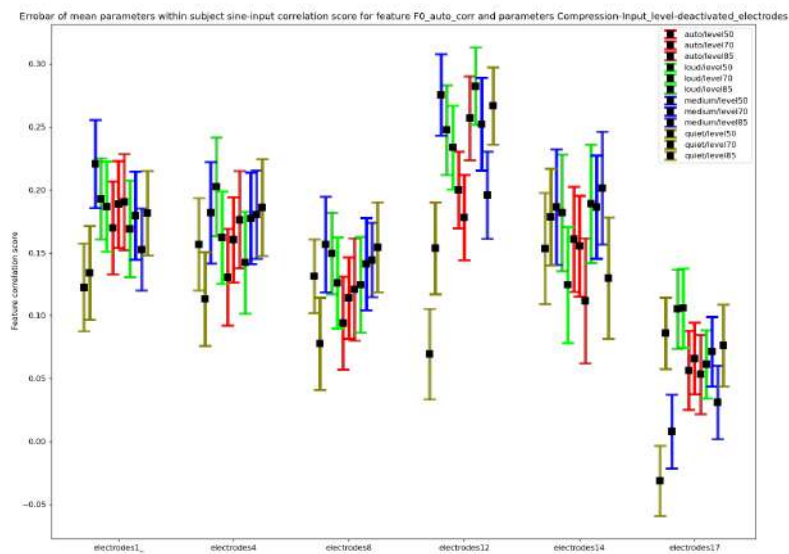


(b) Error bars of the mean for parameters T60:NofM:electrodes. The dependent variable is the sine/input STOI score.

Figure 18: Error bars of STOI scores for input-sine and input-spiral vocoders

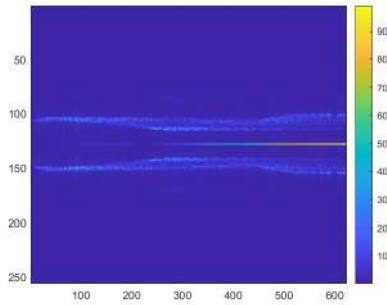


(a) Error bars of the mean for parameters compression:input-level-electrodes. The dependent variable is the spiral/input F0 correlation score.

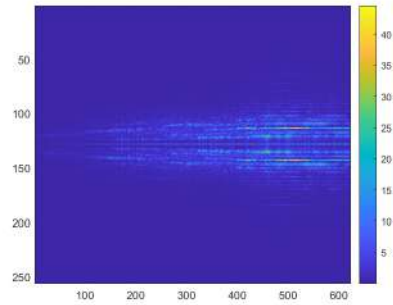


(b) Error bars of the mean for parameters compression:input-level-electrodes. The dependent variable is the sine/input F0 correlation score.

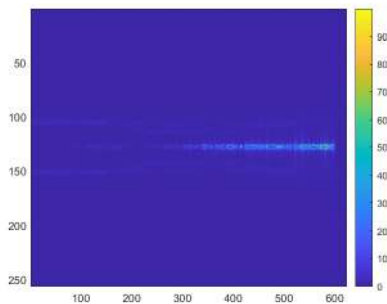
Figure 19: Error bars of F0 correlation scores for the input-sine and input-spiral



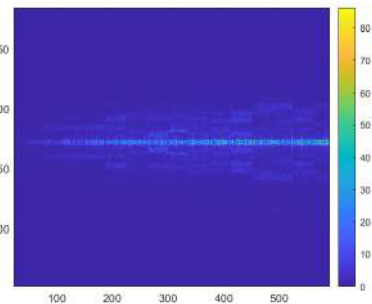
(a) First signal spectrogram example



(b) Second signal spectrogram example



(c) Masking noise spectrogram of signal 1



(d) Masking noise spectrogram of signal 2

Figure 20: Examples of signals and correspondings masking noise spectrogram, Yaxis=frequency bins,Xaxis=time frame

4 Proposed solution for improvements of the SRS perception

This part of the report cannot be described because the solution I proposed at the moment of writing is discussed to be part of a patent. (Note here that if the patent lawyer decides to drop the patent, then the description of this process will be add to this report as an appendix).

– \ _(" /)_ /–

5 Clinical Protocol for Testing the Proposed Solutions

The solution mentioned earlier is intended to be implemented within a coding strategy. We will refer to the classical Oticon Medical coding strategy as the "Crystalis" strategy and to this strategy with the proposed solution implemented as the "Crystalis Enhanced" strategy.

The experiments have three main objectives:

- Estimating the SNR [alarm detection](#) threshold to evaluate the effect of the Crystalis Enhanced solution.
- Estimating the [recognition of alarms](#) to evaluate the effect of the Crystalis Enhanced solution.
- Verifying our model, i.e., whether there is a monotonic mapping between input sound features and the corresponding CI vocoded features correlation scores, and a [audibility](#) perceptual score for alarm sounds for CI users, or at least for normal hearing individuals.

To achieve these objectives, we propose three experiments.

5.1 Exclusion and Inclusion Criteria

Two groups of participants should be formed:

- A control group consisting of individuals with normal hearing (NH), with hearing levels up to 25 dB HL for frequencies between 100 and 8000 Hz.
- A second group consisting of cochlear implant users (with unilateral or bilateral impairment) who have the classical Crystalis coding strategy implemented (implying that they use an Oticon Medical device).

Exclusion criteria for the normal hearing group: psycho/sensory impairments (e.g., neuro-sensory loss, cognitive impairment).

Exclusion criteria for the CI user group: any psycho/sensory impairments (e.g., neuro-sensory loss, cognitive impairment) other than hearing loss, or if the CI is not from Oticon Medical.

Additionally, subjects should have at least 1 year of implantation with a minimum of 6 months of cumulative use of their device (to ensure they are already adapted to their devices).

The following parameters concerning subjects must be recorded and will be used for analysis: duration of deafness, pre- or post-lingual hearing loss, age of the subject, unilateral or bilateral hearing loss/implantation, duration of implantation, duration of use, and type of hearing loss (etiology).

5.2 Protocol Proposed

5.2.1 Experiment 1

The purpose of this experiment is to measure the Signal-to-Noise Ratio (SNR) threshold for detecting alarm sounds. This experiment will be conducted in several blocks within two different groups (refer to [Figure 22](#)):

- A control group consisting of normal hearing subjects.
- A CI user group with Oticon Medical implants.

This test is a 3-Alternative Forced Choice (3AFC) test with a one-up/one-down adaptive procedure, where the variable parameter is the SNR, coupled with a 2×2 crossover design. The procedures for the NH group and the CI group during a session will be slightly different.

CI Test Session:

- The target signal is an alarm sound with an environmental composition of sounds (e.g., babble noise, street noises) at a controlled SNR (signal being the alarm sound).
- The reference signal is a target alarm sound with a modification (e.g., music instrument morphing, spectral symmetric inverting according to a constant frequency axis) and the same environmental composition of sounds at the same controlled SNR.

Other target and reference signals (referred to as dummy target and reference) may also be used in the testing session to avoid the subjects focusing specifically on the alarm. However, subjects' responses to these will not be considered during the adaptive one-up/one-down procedure (see [Figure 21](#)).

At each presentation, two reference signals and one target signal are presented in random order to the subject, who must differentiate the target from the references. A session will contain a maximum of 10 trials (see Figure 1). A trial ends after either 30 presentations or 5 consecutive reversals between two values. These trials should include at least 10 presentations of dummy target/reference signals presented in random order.

The SNR threshold for the trial is the mean of the last 5 reversal values.

NH Test Session:

The target and reference signals are defined similarly to the CI test session. However, after construction, the reference and target pairs will be vocoded using either the Crystalis strategy or the Crystalis Enhanced strategy (depending on the session and subgroups) and processed with the sinusoidal and spiral vocoders (for all testing sessions and subgroups) before being presented to the subjects (see Figure 22).

The same procedure will be applied to dummy target/reference pairs.

During a session, each trial will focus on one type of sound production (sine vocoded, spiral vocoded, or unprocessed) presented in random order. The SNR threshold estimation will then correspond to the SNR threshold for normal hearing subjects using vocoder X or no vocoder with coding strategy Y. Thus, each session will estimate 3 SNR thresholds: the SNR for the no vocoder coding strategy, sine vocoder, and spiral vocoder.

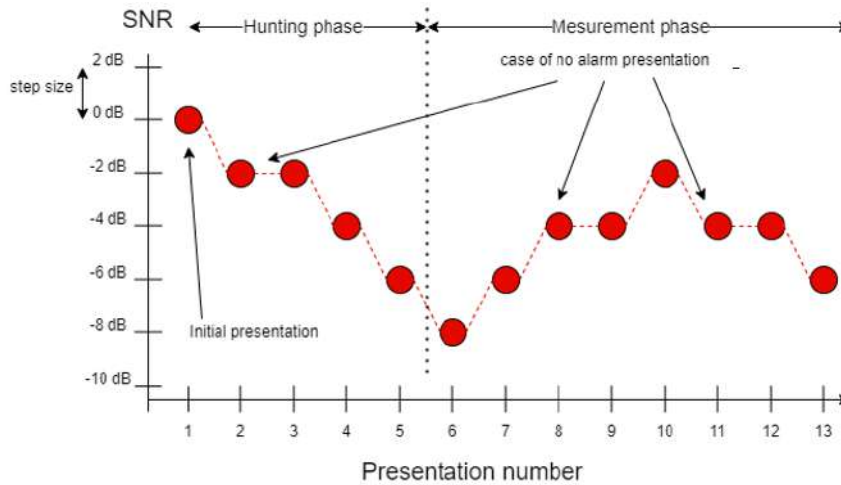


Figure 21: Example of a one-up/one-down trial

Note that the training phases for CI users will involve using their CI with the specified coding strategy for the required time before participating in the sessions.

- **Red:** Sessions testing unprocessed, Crystalis sinus vocoded, and Crystalis spiral vocoded signals.

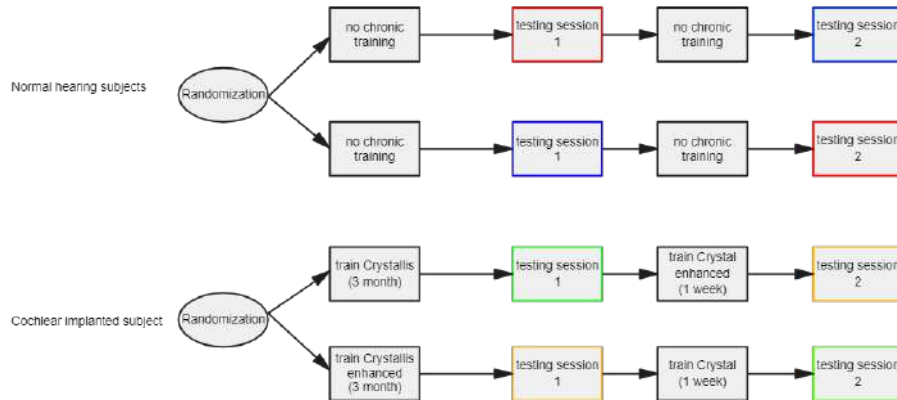


Figure 22: Session planning of Experiment 1

- **Blue:** Sessions testing unprocessed, Crystallis Enhanced sinus vocoded, and Crystallis Enhanced spiral vocoded signals.
- **Green:** Sessions for CI users under the Crystallis strategy.
- **Yellow/Orange:** Sessions for CI users under the Crystallis Enhanced strategy.

Analysis of Results and Hypotheses:

The SNR threshold for a trial will be estimated using the mean of the last 5 reversal SNR values. The mean of the estimated thresholds across trials and subjects will represent the SNR threshold for the session and subject. The result is validated only if the standard deviation (STD) estimated from at least 4 trials is less than or equal to the step size (2 dB). Some biases can be mitigated by the design, as discussed in [17].

After averaging the scores across trials and subjects, we obtain 7 SNR threshold estimates:

- NH threshold
- NH Crystallis sine vocoded threshold (NHC-sin-threshold)
- NH Crystallis Enhanced sine vocoded threshold (NHCE-sin-threshold)
- NH Crystallis spiral vocoded threshold (NHC-spi-threshold)
- NH Crystallis Enhanced spiral vocoded threshold (NHCE-spi-threshold)
- CI Crystallis threshold (CIC-threshold)
- CI Crystallis Enhanced threshold (CICE-threshold)

We expect the SNR thresholds to be in the following order:

$NH\ threshold < NHC\text{-}sin\text{-}threshold < NHCe\text{-}sin\text{-}threshold < NHC\text{-}spi\text{-}threshold < NHCe\text{-}spi\text{-}threshold < CICE\text{-}threshold < CIC\text{-}threshold$

We naturally expect that NH subjects will perform better with unprocessed sound compared to degraded sounds and will outperform CI subjects. Additionally, we anticipate that the Crystalis Enhanced design might disrupt normal hearing listeners (although the specific reasons for this will be revealed through the solution), but not to the extent that it reverses the performance benefit of the sinus vocoder compared to the spiral vocoder.

We will also compare the SNR thresholds from the green sessions and the yellow sessions to assess the impact of short adaptation times for both coding strategies. We expect the results from CI testing sessions under Crystalis Enhanced to be lower than those from sessions under the Crystalis strategy.

Finally, we anticipate that the enhanced version of the coding strategy will improve detection capabilities for CI users compared to those with the classical coding strategy.

5.2.2 Experiment 2

The second experiment, involving the same groups, is a recognition test aimed at estimating alarm recognition. Specifically, it is a 42-item computerized forced-choice, closed-set test. This test is designed to assess recognition scores based on the hearing condition of the participants.

The test should be conducted at least 2 months after the first experiment to avoid carryover effects. Additionally, the CI group should be randomly split into two subgroups: one with the Crystalis strategy and another with the Crystalis Enhanced strategy.

Participants will be asked to rate their familiarity with a set of sounds and assess the security relevance of these sounds, using a scale from 1 to 5 (where 5 represents the highest confidence or greatest security relevance, and 1 represents the lowest confidence or security relevance). They will also receive training on recognizing the sounds used in the Luzum 2023 protocol (refer to the Materials and Methods section of [19]).

We will specifically examine the recognition rate of alarm signals and compare the results for Crystalis Enhanced, Crystalis, and NH (Normal Hearing) scores. We expect the recognition score for Crystalis Enhanced to be higher than that for Crystalis but lower than the Normal Hearing score.

5.2.3 Experiment 3

The purpose of this experiment is to establish a mapping between the perceptual alarm audibility score of CI users and the correlation score of certain acoustic features of the original sounds with a vocoded version of the same sound.

The features considered are brightness, pitch, harmonicity, scale, loudness, flatness, and irregularity (defined according to Elhilali’s procedure [14]).

The processing of alarm sounds will involve modifying one of these features while attempting to keep the others as constant as possible. This will result in alarm signals with controlled correlation scores for the feature in question.

Ideally, we aim to create a table for each sound and each feature as follows:

Correlation Score of Feature F with Its Modified Version	Alarm Sounds Created
1	Original Alarm
0.8	Alarm with Feature F Score 0.8
0.6	Alarm with Feature F Score 0.6
0.4	Alarm with Feature F Score 0.4
0.2	Alarm with Feature F Score 0.2

Table 1: Example of table necessary for the experiment

Note that the method for obtaining sounds with modified features will not be explicitly detailed here. Remarks about the creation of these sounds will be provided in the limitations section.

Method

The same groups are retained (i.e., a normal hearing control group and a CI group with the Crystallis strategy). The CI group is divided and tested using the same 2×2 crossover design as shown in Figure 22.

The test is a 2-Alternative Forced Choice (2AFC) test, where the target signal includes an environmental background sound and an alarm sound, while the reference signal contains only the environmental background sound.

Before the test, subjects will be trained to recognize the unprocessed alarms that will be used. During the testing session, target and reference signals will be presented in a random order, and subjects will be asked to indicate whether the alarm was clearly audible. The target signals will include both unprocessed and processed alarms, presented in random order.

Analysis of Results

We will compute the mean of the rate of correct positive responses across subjects and alarms for each modality (i.e., feature correlation score). We expect to obtain scatter plots similar to Figure 23.

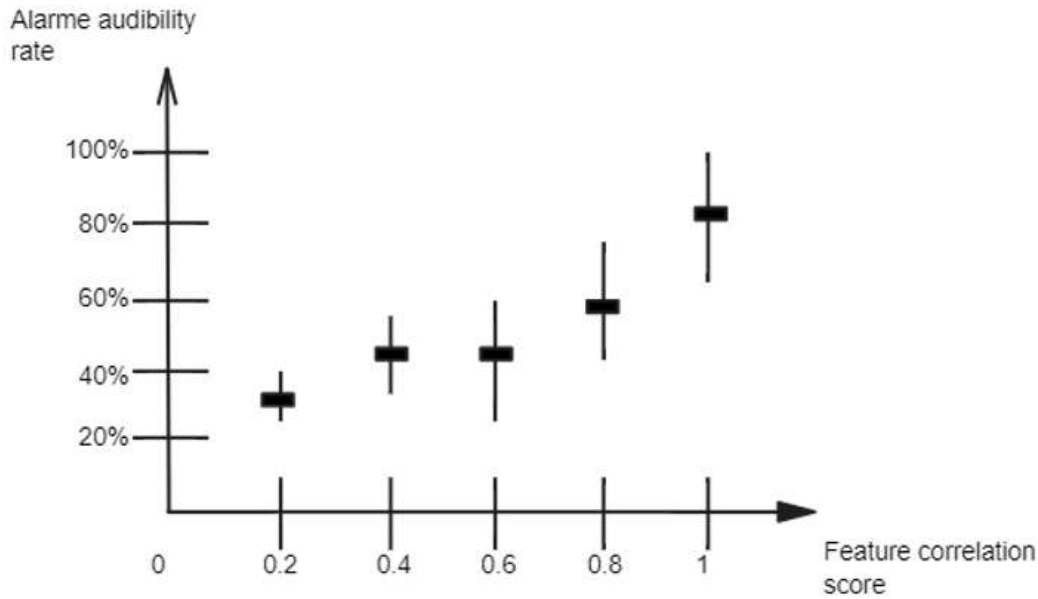


Figure 23: Example of expected mapping between the feature correlation score and the alarm audibility rates

We will use the matching procedure outlined in [28]. Other mapping functions, such as polynomial, exponential, or sigmoid functions, are also considered. Possible expressions for modeling the perceptual score based on the feature correlation include:

$$\begin{aligned} \text{perceptual score}(\text{corr score}) &= \\ a \cdot (\text{corr score})^3 + b \cdot (\text{corr score})^2 + c \cdot (\text{corr score}) + d \\ \text{perceptual score}(\text{corr score}) &= \frac{1}{1 + \exp[-(a \cdot \text{corr score} + b)]} \\ \text{perceptual score}(\text{corr score}) &= a \cdot \text{corr score} + b \end{aligned}$$

where a , b , c , and d are learnable parameters, optimized using a method such as gradient descent.

5.3 Limitations of Experiments

We can classify the limitations of these experiments into the following categories: participant profile, experimental setup, material (hardware and software), and balance and representativity.

5.3.1 Participant Profile

Participant variability may influence their adaptation capabilities, notably during the adaptation period of the first experiment (in terms of time required and performance achieved). To preserve participant focus, the experiment should be limited in duration. CI users unfortunately have heterogeneous characteristics, such as varying durations of deafness and hearing acuity, which we cannot control and which will likely affect the variability of our results.

5.3.2 Experimental Setup

- Reproducing the same hearing conditions across experiments can be very challenging. It is particularly critical that the sound stimuli produced in the experiment have consistent characteristics (e.g., loudness, pitch), which cannot be guaranteed without a controlled experimental environment. Using a crossover design for Experiment 1 may introduce undesirable effects, such as carryover or sequence effects [17].
- In Experiment 3, the CI might affect the acoustic features of their processed sounds. Therefore, we cannot precisely assess the effect of modifying alarm sound features on CI users' perception (the best approximation would be to use the vocoders again).
- For Experiment 3, there is not a single method to modify a signal x into another signal y such that $\text{corr-score}(x, y) = a$ for a determined value. It is reasonable to assume that among all y candidates where $\text{corr-score}(x, y) = \text{constant}$, some will be perceived differently by CI users. Thus, the specific way CIs modify features in relation to correlation requires further study.

5.3.3 Material (Hardware and Software)

- The equipment must be as controlled as possible to ensure the repeatability of the experiments and the randomization of stimulus presentation to avoid sequential bias.
- Modifying only one **acoustic feature** for alarms may not be feasible for spectral-related features without affecting potentially important spectral cues. For example, modifying the pitch will also impact the brightness (spectral centroid) because changing the pitch shifts the entire spectrum.
- The modifications applied to the features will depend on the input alarm sounds, especially if the flattening of feature variations is insufficient to achieve the targeted correlation score.

5.3.4 Balance and Representativity

A small number of subjects poses a problem for representativity, making it challenging to obtain a balanced sample of every characteristic of interest (such as age, duration of CI use, and pre- or post-lingual hearing loss). A preliminary sample analysis based on the expected effect size must be conducted.

6 Conclusion and Discussion

6.1 Conclusion

In this work, we created a database of SRS and environmental sounds and developed a framework for processing these sounds, including adding reverberation and controlling the SNR. These tools have not been used to their full potential due to time and computational constraints but should be capable of providing controlled signals for future experiments. We support two claims about the reasons for the poor alarm detection capabilities of CI users. We proposed a design to improve both the detection and recognition capacities of CI users, although we cannot discuss it in detail. Finally, we provided clinical guidelines to test the efficacy of the solution in both recognition and detection tasks and to validate our hypothesis regarding the link between feature correlation scores and alarm audibility rates.

6.2 Future Work

Our approach has several areas for improvement in future experiments:

- We focused only on alarm signals, which, although important, do not necessarily represent all SRS signals. For example, alarms are harmonic signals with periodic patterns, unlike broadband noises such as vehicle noise. Additionally, for other SRS signals, the pitch feature may be difficult or even impossible to define, making it challenging to generalize our analysis to all SRS sounds.
- We used STOI as a measure of "global degradation," relying on the assumption that STOI correlates with alarm recognition rates. While STOI is highly correlated with speech intelligibility, its effectiveness as a predictor for SRS recognition or detection rates is not well established.
- We used a measure of "CI loudness" that is contentious due to questionable assumptions. For instance, integrating contributions from different channels into the loudness calculation is problematic, and adding channel contributions alone is insufficient. This measure also does not account for channel interactions caused by the [current spread issue](#).

- Our analysis focused solely on correlation scores (i.e., the variation of feature values compared to the output). We observed that features depending on signals can vary significantly over short durations. Therefore, averaging values or calculating standard deviations could be misleading, as significant variations might be important for understanding CI effects but would be obscured by such statistical tools. However, mean value differences can be important for assessing audibility, as shown in [11], which compares the mean value of background + alarm with background only during the alarm period. We believe that correlation may still be a good predictor of audibility, as humans can recognize sequences of sounds based on temporal-frequency patterns, such as recognizing a melody even if transposed or played by a different instrument at a different tempo.
- Another concern is that our analysis was conducted on a limited number of different signals. This limitation is primarily due to software constraints, as calculating scale features and processing chains were demanding. We prioritized the number of parameters over the number of files, as adding more alarms would have introduced less variability than changes in parameters.
- Additionally, using the vocoder involves making assumptions about the EDG, specifically regarding frequency perception. Determining how well these models represent CI users' sound perception is a complex question with no universally correct answer.

7 Annex

7.1 Some other figures

We studied other features since the beginning of this work, notably because in other paper less related to our topics [5], [6] cite those features are best features to make environmental sound classification. [11] also mention roughness as a feature valuable of investigation.

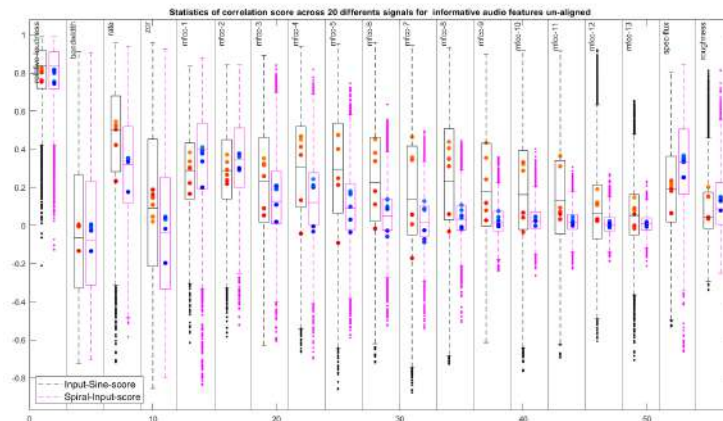


Figure 24: Statistic of correlation score for other acoustic features. On Y-axis you have the correlation score, on X-axis the different categories tested (feature/vocoder type)

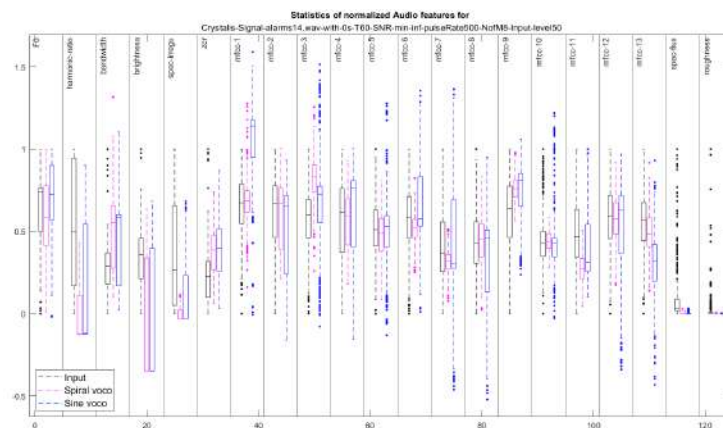


Figure 25: Example of plotting the statistic across the time of feature values for the input and the vocoders on one modality, the values are normalized to ensure that the f0 is between 0 and 1. On Yaxis, you see the means features value of a signal under a specific modality normalize for all feature to ensure the f0 values are between 0 and 1 (those values are no more interpretable by themselves). The X axis separate the different outputs (unprocessed, sine or spiral vocoded). A way to interpret this graph it is that if the distribution of the vocoded values are very different from the output, then the feature is degraded.

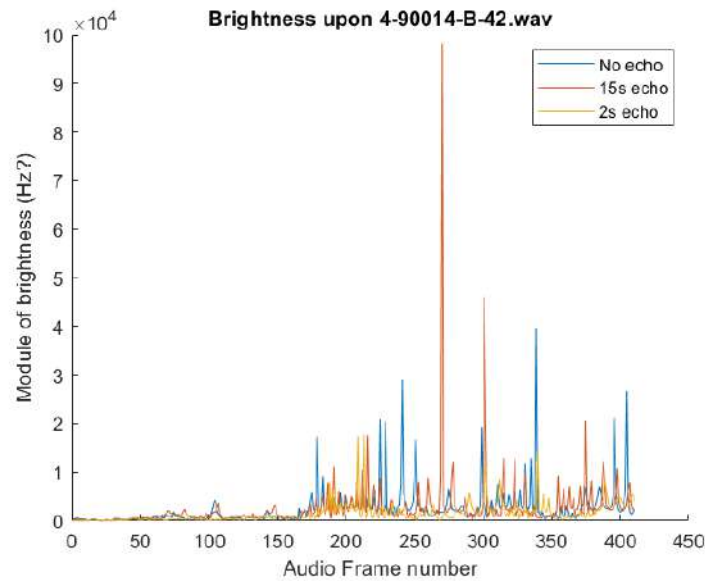


Figure 26: Brightness estimation of signal 490014-B-42.wav, we studied those kind of figure to ensure that the parameters had effects on feature and were correctly set up, note here that brightness is supposed to be a frequential centroid value for a signal at 16kHz sample rate, so the value superior to 8kHz are a bit odd



Figure 27: Example of former EDG used to our analysis, (the electrodes being numerated from apex to base of the cochlea, first number correspond to low frequency and last to high ones, so it's upside done compared to spectrogram classic plots, thanks to Dr.Rafael Attili Chiea code, we could superposed EDG and spectrogram.

References

- [1] Andrey Anikin, Rasmus Bååth, and Tomas Persson. “Human Non-linguistic Vocal Repertoire: Call Types and Their Meaning.” In: *Journal of nonverbal behavior* 42 (1 2018), pp. 53–80. ISSN: 0191-5886. DOI: [10.1007/s10919-017-0267-y](https://doi.org/10.1007/s10919-017-0267-y).
- [2] Justin M. Aronoff and David M. Landsberger. “The development of a modified spectral ripple test”. In: *The Journal of the Acoustical Society of America* 134 (2 Aug. 2013), EL217–EL222. ISSN: 0001-4966. DOI: [10.1121/1.4813802](https://doi.org/10.1121/1.4813802).
- [3] Sebastian A. Ausili et al. “Sound Localization in Real-Time Vocoder Cochlear-Implant Simulations With Normal-Hearing Listeners”. In: *Trends in Hearing* 23 (2019). ISSN: 23312165. DOI: [10.1177/2331216519847332](https://doi.org/10.1177/2331216519847332).
- [4] Pascal Belin, Sarah Fillion-Bilodeau, and Frédéric Gosselin. “The Montreal Affective Voices: A validated set of nonverbal affect bursts for research on auditory affective processing”. In: *Behavior Research Methods* 40 (2008), pp. 531–539. URL: <https://api.semanticscholar.org/CorpusID:1463309>.
- [5] Vasileios Bountourakis, Lazaros Vrysis, and George Papanikolaou. “Machine Learning Algorithms for Environmental Sound Recognition: Towards Soundscape Semantics”. In: *Proceedings of the Audio Mostly 2015 on Interaction With Sound*. AM ’15. Thessaloniki, Greece: Association for Computing Machinery, 2015. ISBN: 9781450338967. DOI: [10.1145/2814895.2814905](https://doi.org/10.1145/2814895.2814905). URL: <https://doi.org/10.1145/2814895.2814905>.
- [6] Vasileios Bountourakis et al. “An Enhanced Temporal Feature Integration Method for Environmental Sound Recognition”. In: *Acoustics* 1.2 (2019), pp. 410–422. ISSN: 2624-599X. DOI: [10.3390/acoustics1020023](https://doi.org/10.3390/acoustics1020023). URL: <https://www.mdpi.com/2624-599X/1/2/23>.
- [7] Robert P. Carlyon and Tobias Goehring. *Cochlear Implant Research and Development in the Twenty-first Century: A Critical Update*. Oct. 2021. DOI: [10.1007/s10162-021-00811-5](https://doi.org/10.1007/s10162-021-00811-5).
- [8] Lorenzi Christian et al. “Human Auditory Ecology: Extending Hearing Research to the Perception of Natural Soundscapes by Humans in Rapidly Changing Environments”. In: *Trends in Hearing* 27 (Jan. 2023). ISSN: 23312165. DOI: [10.1177/23312165231212032](https://doi.org/10.1177/23312165231212032).
- [9] J. S. Chung, A. Nagrani, and A. Zisserman. “VoxCeleb2: Deep Speaker Recognition”. In: *INTERSPEECH*. 2018.
- [10] Mishaela DiNino et al. “Vowel and consonant confusions from spectrally manipulated stimuli designed to simulate poor cochlear implant electrode-neuron interfaces”. In: *The Journal of the Acoustical Society of America* 140 (6 Dec. 2016), pp. 4404–4418. ISSN: 0001-4966. DOI: [10.1121/1.4971420](https://doi.org/10.1121/1.4971420).

- [11] F. Effa et al. “Evaluating and predicting the audibility of acoustic alarms in the workplace using experimental methods and deep learning”. In: *Applied Acoustics* 219 (Mar. 2024). ISSN: 1872910X. DOI: [10.1016/j.apacoust.2024.109955](https://doi.org/10.1016/j.apacoust.2024.109955).
- [12] Tobias Goehring et al. “The effect of increased channel interaction on speech perception with cochlear implants”. In: *Scientific Reports* 11 (1 Dec. 2021). ISSN: 20452322. DOI: [10.1038/s41598-021-89932-8](https://doi.org/10.1038/s41598-021-89932-8).
- [13] Jacques A. Grange et al. “Cochlear implant simulator with independent representation of the full spiral ganglion”. In: *The Journal of the Acoustical Society of America* 142 (5 Nov. 2017), EL484–EL489. ISSN: 0001-4966. DOI: [10.1121/1.5009602](https://doi.org/10.1121/1.5009602).
- [14] Nicholas Huang and Mounya Elhilali. “Auditory salience using natural soundscapes”. In: *The Journal of the Acoustical Society of America* 141.3 (Mar. 2017), pp. 2163–2176. ISSN: 0001-4966. DOI: [10.1121/1.4979055](https://doi.org/10.1121/1.4979055). eprint: https://pubs.aip.org/asa/jasa/article-pdf/141/3/2163/15325788/2163_1_online.pdf. URL: <https://doi.org/10.1121/1.4979055>.
- [15] David M. Landsberger et al. “Spectral-Temporal Modulated Ripple Discrimination by Children With Cochlear Implants”. In: *Ear and hearing* 39 (1 Jan. 2018), pp. 60–68. ISSN: 15384667. DOI: [10.1097/AUD.0000000000000463](https://doi.org/10.1097/AUD.0000000000000463).
- [16] Birgitta Larsby and Stig Arlinger. “Auditory Temporal and Spectral Resolution in Normal and Impaired Hearing”. In: *Journal of the American Academy of Audiology* 10 (04 Apr. 1999), pp. 198–210. ISSN: 1050-0545. DOI: [10.1055/s-0042-1748481](https://doi.org/10.1055/s-0042-1748481).
- [17] In Junyong Lim Chi-Yeon. “Considerations for crossover design in clinical study”. In: *Korean J Anesthesiol* 74.4 (2021), pp. 293–299. DOI: [10.4097/kja.21165](https://doi.org/10.4097/kja.21165). eprint: <http://ekja.org/journal/view.php?number=8753>. URL: <http://ekja.org/journal/view.php?number=8753>.
- [18] Philipos C. Loizou. “Mimicking the human ear”. In: *IEEE Signal Processing Magazine* (Sept. 1998), pp. 101–130. DOI: [10.1109/79.708543](https://doi.org/10.1109/79.708543).
- [19] Nathan R. Luzum et al. “Identification Accuracy of Safety-Relevant Environmental Sounds in Adult Cochlear Implant Users”. In: *Laryngoscope* 133 (9 Sept. 2023), pp. 2388–2393. ISSN: 15314995. DOI: [10.1002/lary.30475](https://doi.org/10.1002/lary.30475).
- [20] Olivier Macherey and Robert P. Carlyon. *Cochlear implants*. Sept. 2014. DOI: [10.1016/j.cub.2014.06.053](https://doi.org/10.1016/j.cub.2014.06.053).
- [21] Loureiro Manuel et al. “Datalogging Findings in Adult Cochlear Implant Recipients Who Never Developed Intelligible Speech”. In: *J Int Adv Otol* 20 (2024), pp. 113–118. DOI: [10.5152/iao.2024.231193](https://doi.org/10.5152/iao.2024.231193). URL: www.advancedotology.org.

- [22] Stephen G. McGovern. “Fast image method for impulse response calculations of box-shaped rooms”. In: *Applied Acoustics* 70 (1 Jan. 2009), pp. 182–189. ISSN: 0003682X. DOI: [10.1016/j.apacoust.2008.02.003](https://doi.org/10.1016/j.apacoust.2008.02.003).
- [23] Nicole Miller-Viacava et al. “Sensorineural hearing loss alters auditory discrimination of natural soundscapes”. In: *International Journal of Audiology* (2023). ISSN: 17088186. DOI: [10.1080/14992027.2023.2272559](https://doi.org/10.1080/14992027.2023.2272559).
- [24] vryzas nikolaos et al. “speech emotion recognition for performance interaction”. In: *journal of the audio engineering society* 66 (6 June 2018), pp. 457–467. DOI: [10.17743/jaes.2018.0036](https://doi.org/10.17743/jaes.2018.0036).
- [25] Fearghal O’Brien and Denis Cousineau. “Representing Error bars in within-subject designs in typical software packages”. In: *The Quantitative Methods for Psychology* 10.1 (2014), pp. 56–67. DOI: [10.20982/tqmp.10.1.p056](https://doi.org/10.20982/tqmp.10.1.p056). URL: <http://www.tqmp.org/RegularArticles/vol10-1/p056/p056.pdf>.
- [26] Karol J. Piczak. “ESC: Dataset for Environmental Sound Classification”. In: *Proceedings of the 23rd Annual ACM Conference on Multimedia*. Brisbane, Australia: ACM Press, Oct. 13, 2015, pp. 1015–1018. ISBN: 978-1-4503-3459-4. DOI: [10.1145/2733373.2806390](https://doi.org/10.1145/2733373.2806390). URL: <http://dl.acm.org/citation.cfm?doid=2733373.2806390>.
- [27] William Shakespear. *Hamlet*. 1623.
- [28] Cees H. Taal et al. “An algorithm for intelligibility prediction of time-frequency weighted noisy speech”. In: *IEEE Transactions on Audio, Speech and Language Processing* 19 (7 2011), pp. 2125–2136. ISSN: 15587916. DOI: [10.1109/TASL.2011.2114881](https://doi.org/10.1109/TASL.2011.2114881).
- [29] Nikolaos Vryzas et al. “Subjective Evaluation of a Speech Emotion Recognition Interaction Framework”. In: Sept. 2018, pp. 1–7. DOI: [10.1145/3243274.3243294](https://doi.org/10.1145/3243274.3243294).