



# **Implementation on FPGA and Evaluation Of a Prosody Modification of Speech for Impaired Persons using DWT-OLA**

L. Bendaouia<sup>1</sup>, S.M. Karabernou<sup>2</sup>, L. Kessal<sup>3</sup>, H. Salhi<sup>4</sup>

Ph.D.Student, ETIS-ENSEA (Cergy, France), SDU (Blida, Algeria), CDTA (Algiers, Algeria)<sup>1</sup>

Assistant Professor, Dept. Head of ENE, ETIS, ENSEA, CNRS UMR 8051, Cergy, France<sup>2</sup>

Assistant Professor, ETIS, ENSEA, CNRS UMR 8051, Cergy, France<sup>3</sup>

Assistant Professor, Department of Electronic Systems, SaadDahleb University, Blida, Algeria<sup>4</sup>

**ABSTRACT:** In this paper, we describe the design of a platform on Field Programmable Gate Array (FPGA) for a Prosody Modification of Speech (PMS). As it is not possible to reproduce the hearing capacities to the deficient ear, the speech signal can be shaped on the temporal or spectral domain to preserve the relevant information which could be inaccessible for the auditing system of the hearing-impaired person. So, in designing such platform, we aim to analyze the functioning of the cochlea and to perform efficient processing for improving intelligibility by shifting the frequency components of certain defective bands to different healthy ones. The design combines both Discrete Wavelet Transform (DWT) and Over-Lap and Add (OLA) technique and permits to analyze the input data at different segmental length in order to detect and manipulate the pitches. The DWT-OLA gives efficient real time results as compared to classical designs. We conducted experiments on speech data base get from Arctic corpus and used Mean Opinion Score (MOS) tests to subjectively evaluate the speech intelligibility. We obtained a gain improvement of speech intelligibility reaching 70%. Moreover, the proposed FPGA platform involves significantly fewer resources, reduced memory-size and less dynamic power consumption as compared to some previous wavelet-based implementations.

**KEYWORDS:** Hearing impairment, Noise, Intelligibility, Prosody, Synthetic Speech, Pitch.

## **I. INTRODUCTION**

Hearing aids are now used to alleviate hearing impairments. However, more than 60% of impaired persons feel uncomfortable when using their hearing-aids because of the worse intelligibility resulting from bad speech comprehensibility. We believe that hearing impairments can be alleviated by a system with characteristics closer to being body's ones. Research has been deeply involved in developing new algorithms to improve speech intelligibility.

Although, lot of researches was held to enhance speech for impaired people [1], few authors deal with the problem of power consuming. No numerical results are given, making the comparison only with the hearing aids of the market. Most closely related to our approach is the work of [2], who provide a method for power reducing based on algorithm and hardware optimizations along with the architecture uses the odd/even data lifting. Our work is improved by avoiding the access to the memory for data storage. Instead, we formulate our algorithm under which a given data of the speech signal is segmented at the input and each segment is processed individually. In order to avoid losing of data information, the segments are overlapped and an overlap technique is then used for treatment. This approach provides perfect analysis and efficient computation. Consequently, the framework devised here is made generic and requires simulations and empirical evaluation of the routing scheme in order to be applied.

Traditionally, Digital Signal Processing (DSP) algorithms are implemented using General Purpose Processors (GPP) for low rate applications. These devices showed limited capabilities for processing high volume data efficiently in real time. The trends had then been shifted to Special Purpose DSP (SPDSP) and Application Specific Integrated Circuits (ASICs) in order to meet the increased complexity and to gain in performance requirements of these algorithms but



# International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 8, August 2014

with high cost functions [3]. Today, FPGAs are highly preferred for their relatively high capacity, low cost, short design cycle and short time to market. FPGA affords the capability of constant reconfiguration to meet application performances [4]. Dealing with digital speech processing as it pertains to the hearing impaired persons especially for miniaturized system applications; FPGA allows increasing sophisticated features to be built for better sound reproduction while keeping small size and low power consumption of the devices.

Fortunately, simulation tools provide us a rapid design and basic information. Similarly, a high-level programming language is an efficient comparison tool for the final output results and system evaluation. In practice, the implementations are often subject to lot of limitations [5]. Using DWT at Multiresolution over disjoint bands remains up to now a practical necessity for perfect design [6, 7] for digital speech processing in particular and herein some references from our work, where the goal is to investigate noise reduction and hardware implementation.

This work extends previous research described in [8, 9]. In this paper we present the implementation of a multi-level one dimension DWT combined to an OLA on FPGA for a bio-inspired medical hearing aid application. The methodology aims to improve in one side better speech quality and in the other side, an efficient flexible reconfiguration and reduced cost functions. The scheme represents architecture for denoising and frequency shifting. It is realized targeting a DE2 development kit board of Altera (EP2C70F896) and results are compared to that obtained in Matlab. The system provides a generic framework allowing the use of DWT analysis / synthesis with frequency shaping of the speech signal to improved speech intelligibility. We present some simulation results under VHDL and Matlab. Hence, a comparative study is done based on the Mean Square Error (MSE) and the Signal to Noise Ratio (SNR). MOS evaluations are presented for speech intelligibility and the gain obtained by the proposed architecture.

## II. LITERATURE REVIEW

Voice is an important tool for communicating and transmitting information. When deafness occurs due in most of the cases to the destruction of the Outer Hearing Cells (OHC), hearing-impaired people feel great difficulties in understanding speech in noisy and reverberant environments. In such case, the hearing thresholds at some frequencies are quantified in terms of losses in dB on a certain region of the cochlea and lead to the loss of linearity of the sound frequencies. This results to the loss of compression and amplification of the active mechanisms which will be treated in section II and, the extension of the hearing filters which will produce two main consequences. First, the phenomenon of masking is strengthened and the environmental ambient noise will be more annoying because the filters will become less selective [10]. And, the fine spectral and temporal information engendered by the effects of filtering and diffraction of the external and the middle ear will be lost and the sensation of sound will be profoundly altered.

Digital hearing aids concern light or moderate deafness and permit the controls and adjustments of number of electroacoustic parameters, among them frequency response, SPL saturation, compression aspects, fine tuning characteristics, noise reduction and acoustic feedback cancellation. Although they offer many advantages and methods of signal processing capabilities, no improvement will occur on the analyzed sound signal and is badly detected by the auditing system. It is a consequent that these prosthesis do not mainly employ their advanced digital technologies for temporal and frequency modifications concurrently. These prosthesis present lacks of frequency selectivity and several studies had suggested approaches for increasing the signal duration or for shifting the non-selected frequencies to active places of the cochlea. Speech is language dependent and generally described over several levels. At the psycho-acoustical level, voice comprehension is based on its basic parameters namely frequency, amplitude and duration [11]. The purpose of the prosody modification is to make one, two or all of these parameters change over a speech segment without affecting the timber. The signal intensity can be easily modified by a multiplication; meaning that by simply amplifying consonant energy will improve their identification. But, the changes of the fundamental frequency or pitch ( $F_0$ ) and the duration or speed are not so obvious [12]. Clear speech has better intelligibility that conversational one where significant differences in phonetic, phonological and prosodic features are observed. If the duration-rate decreases, the speech intelligibility increases. However, applying phoneme duration from conversational to clear speech did not improve the intelligibility. Making use of pitch ( $F_0$ ) in identifying initial voiced/unvoiced consonants or inserting pauses at the phrase boundaries will improve the intelligibility [13]. Some approaches use speech denoising techniques and others try to model the speech signal by parametric techniques. One of the recommended solutions is the source / filter decomposition of the vocal signal based on the knowledge of the speech production system. This technique belongs to a family of reference methods used for speech synthesis. This type of methods is effective to



# International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 8, August 2014

reduce the background noise. They have structures which are perfectly adapted to the implementation of the hearing rehabilitation process. They operate directly on the signal waveform to incorporate the prosody information. As an example, there are techniques operating in the time domain such as the Synchronous Overlap-And-Add procedure (SOLA) proposed by Roucoux and Wilgus, the Pitch Synchronous Over-Lap and Add (PSOLA) and the modified version of PSOLA using Waveform similarities (WSOLA) proposed by Verhelst and Roelands [14] and, for the Frequency Domain (FDPSOLA). The improvement of this technique by using temporal methods of decomposition in waveforms and bands based on the knowledge of the acoustical psychology is also possible. To utilize the advantages of wavelet processing for speech enhancement, lot of researches has been carried out leading to many contributions for algorithm developments and architecture designs with less complexity and fast processing frequency.

### III. BASILAR MEMBRANE MODELING AND SYSTEM DESIGN

The auditing system analyses the sound signal by means of a series of filters. These hearing filters overlap continuously over the whole range of the audible frequencies (20-20000 Hz). The phenomenon of frequency distribution on the basilar membrane was first brought by the masking experiments of the perception intensity. The detection of the sound signal is also likened to the output of a hearing filter whose central frequency is closer to that of the signal. The noise within the bandwidth determines if the signal is detectable or not. Several experiences showed that the hearing system uses the hearing filter in which the signal to noise ratio is the highest, known as off-place listening. One of these experiences was held by Fletcher in 1940. G. Von Békésy discovered that the basilar membrane positions itself selectively to specific frequencies of speech sounds. The displacement of the basilar membrane to the stimulus of various sound pressure levels was measured by B. M. Johnstone and al. They made clear the quality factor of resonance in the basilar membrane which varies depending on the pressure of an inputsound signal. Mathematical model from Békésy's data to approximate basilar membrane displacement was derived by J. L. Flanagan. The Basilar Membrane Model (BMM) is constructed based on Flanagan's mathematical models taking into consideration Johnstone's experimental data [15]. The feature extract function of this model has been examined in order to apply the hearing function to engineering models.

Hearing filters are arranged on the BM in a precise nonlinear manner as shown in figure 1. They are dependent of the stimuli level. The filter width is called Critical Band (CB) referring to the measures made by Zwicker. The bandwidth can be calculated according to the central frequency ( $f$ ) using Bark scaled formula 1 [16]. In a finer way, Moore and Glasberg proposed a method of measures introduced by Patterson, formula 2.

$$\Delta f_{Bark} = 25 + 75 \left[ 1 + 1,4 \left( \frac{f}{1000} \right)^2 \right]^{0,69} \quad (1)$$

$$\Delta f_{ERB} = 24,7 \left( \frac{4,37 f}{1000} + 1 \right) \quad (2)$$

**Detection of dead zones:** Since the sound energy in the cochlea travels from the base to the apex, it is not surprising that more damage to hearing occurs at high frequencies, near the base, where all the sound energy passes, than low frequencies, near the apex which is reached only by the low frequency components of the signal. The screening test of the inert zones of the cochlea appears to be indispensable in the audio-prosthetic care of the hearing-impaired people in order to avoid the over correction of the frequency ranges which can make disturbances rather than improve their understanding. The Threshold Equalizing Noise (TEN) tests use a narrow band masking noise of 132 Hz centered at the 1000 Hz frequency. TEN levels are expressed in dB / ERB (Decibel over Equivalent Rectangular Bandwidth). TEN levels must be upper than 10 dB not masked by the better frequency. According to the conclusion made by Brian J.C. Moore, an inert zone revealed in the absolute threshold mask on the frequency of 10 dB SPL is strictly greater than the absolute threshold level and the nominal TEN values. Furthermore, according to Moore, in the case of acquired deafness, hearing losses superior to 90 dB HL on high frequencies and 80 dB HL on low frequencies are often associated to dead regions. To allow an effective detection of these frequencies, several techniques were developed. Among these, the techniques of pitch detection and transformation are the most commonly applied.

**Sound classification by DWT:** Discrete Wavelet Transform approach [17] had proven its importance for the analysis of a transient signal since the connection made between the wavelet transform and multi-rate filter bank trees by Mallat in 1989. The formulation of the DWT as a set of FIR filters establishes the foundation for modeling complex algorithm

# International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 8, August 2014

as hardware architecture. Wavelet Transform has the advantage of using variable time size windows for different frequency bands, Figure 1(a). It is useful for speech denoising, speech classification and pitch marking. The high scale low frequency components are the approximation coefficients and noted ( $A_x$ ). Whereas, the low scale high frequency components are the detail coefficients and noted ( $D_x$ ). In the classification process, the speech signal is windowed using Hamming window given by formula 3. Each window is fragmented into ( $m$ ) overlapping segments (frames) of fixed length  $L$  with ( $S_a$ ) samples in each one as shown by Figure 1(b) and the DWT-OLA is applied to the segments within the window.

$$w(i) = 0,54 - 0,46 \times \cos\left(2 \cdot \pi \cdot i / (n - 1)\right) \quad (3)$$

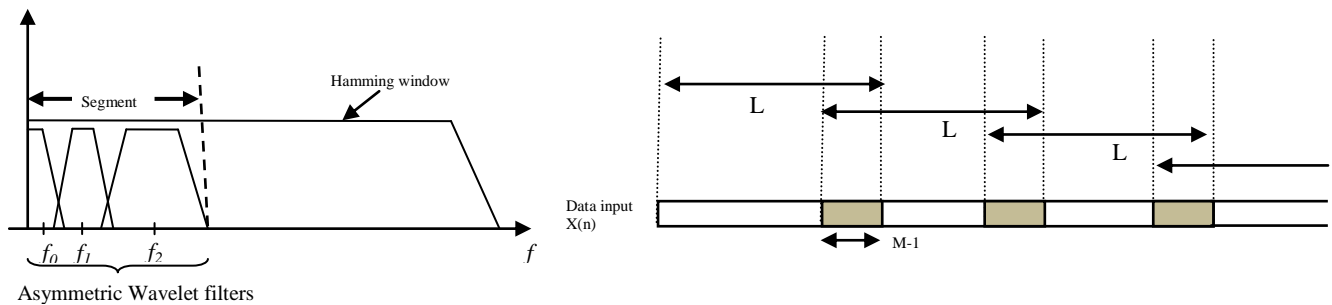


Figure 1. Three level DWT and the Over Lap and Add of the data input.

Sound is classified in a process to find the boundaries between words, syllables or phonemes. In order to perform classification, we must consider the acoustic characteristics of the spoken language. A great deal of techniques uses the segmentation of speech as the basic methodology. The approach is to determine the voiced/unvoiced or silent sections. For real time applications, energy and zero crossing rate or autocorrelation are used. The rate of zero crossing is estimated by formula 4.

$$Z(n) = \frac{1}{2} \sum_{m=-N/2}^{N/2} |\text{sgn}x(n+m) - \text{sgn}x(n+m-1)| W(m) \quad (4)$$

With  $W(m) = \begin{cases} 1/N & -N/2 \leq n \leq N/2 \\ 0 & \text{otherwise} \end{cases}$

Whereas, the autocorrelation of a stationary signal is estimated by formula 5 and calculation of the Shannon entropy of the ( $m^{\text{th}}$ ) segment is given by formula 6.

$$\phi_x(k) = \frac{1}{N - |k|} \sum_{n=0}^{N-|k|-1} x(n) \cdot x(n+k) \quad (5)$$

$$\epsilon_n(Xan) = \sum_{i=0}^{N-1} Xan^2(i) \cdot \log[Xan^2(i)] \quad (6)$$

The discrete samples of the speech signal are normalized between -1 and 1, The  $m^{\text{th}}$  segment is assumed silent or unvoiced if  $\epsilon_n < 0.1$ . We then compute the energy of the approximation and the detail coefficients using:

$$E_n(A_j) = \sum_{i=0}^{M-1} A_j^2(i) \quad (7)$$

$A_j$  is the  $j^{\text{th}}$  approximation coefficient of the  $Xan$  segment (same for  $D_j$ )

## IV. SYSTEM IMPLEMENTATION

In order to estimate the hardware performance of the system, the design has been prototyped targeting DE2 development board kit of Altera containing the FPGACyclone II EP2C70F896 [18]. The proposed system is presented in figure 2 showing the codec interface. This latter has been designed using Qsysto act as the input/output interface to the system [19].

# International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 8, August 2014

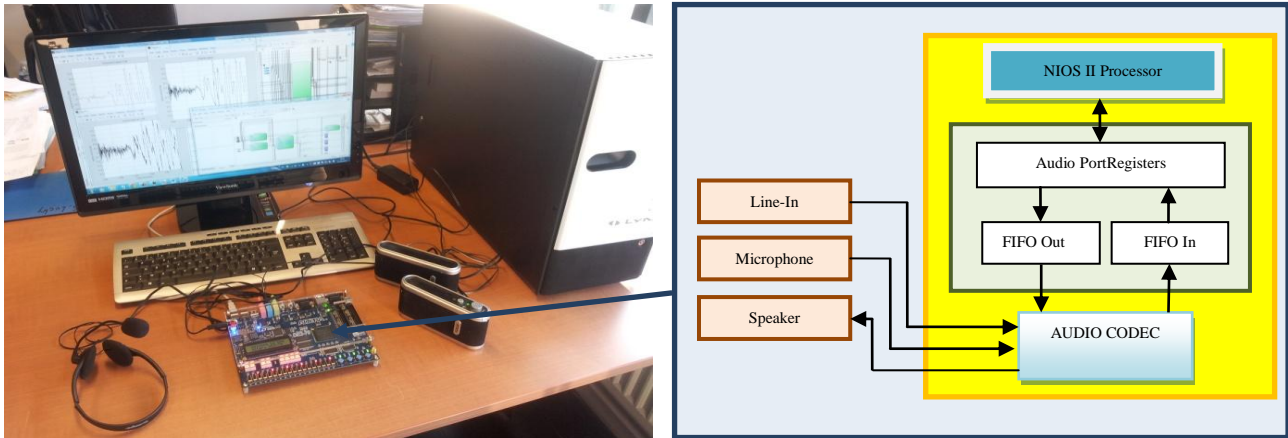


Figure 2. System prototype targeting the DE2 Board and Codec interface

The platform is composed of the I/O block and the processing block containing the CONV-OLA components. The samples of the speech signal at the input are directly taken from the computer via the line-in connector or from the microphone and sampled at 16 KHz using the analog to digital converter (ADC). The output samples are collected through the Fifo-Out port; they are converted to analog signal and send to the speaker. The OLA module makes possible the continuous behavior of the treated signal by overlapping the contiguous segments. As the input samples change at each instant, zero padding to the input data is applied. In the calculation process, only the overlapped data is temporarily stored yielding to a gain of memory space. The output data is obtained by adding the neighbored segments. Using QMF, figure 3, the signal at the input of each level is split by a low pass (h) and a high pass (g) filters given by formulas 9 and 10 respectively.

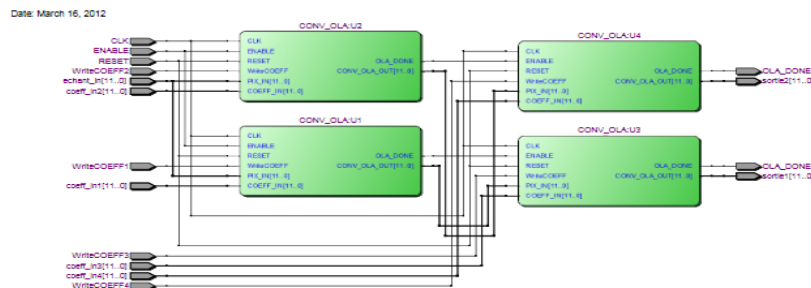


Figure 3. Quadratic Mirror Filter at RTL

$$Y_L(n) = \sum_{k=0}^{L-1} h(k) \cdot X(2n - k) \quad (9)$$

$$Y_H(n) = \sum_{k=0}^{L-1} g(k) \cdot X(2n - k) \quad (10)$$

The advantages in using QMF are the delayed but perfect reconstruction of the signal and the no aliasing. We apply time alignment between successive windows with respect to signal similarities in order to remove the phase discontinuities. The whole synchronization for the system to process the data acquisition, calculation and transfer is insured by a state machine,

## V. SYSTEM PERFORMANCE

We decreased the number of logic registers by using embedded DSP48A1 as MACslices. We can also observe from Table I that the number of resources we obtained by the Transpose Form is much less than that of the Direct Form.



# International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 8, August 2014

TABLE I. : EDA NETLIST FINAL REPORT

Usage	Logic elements	With DSP48		Without DSP48
	Slices	TF	DF	DF
	LUTs	1668 (2%)	6157 (9%)	10944 (10%)
	Flip Flops	1868 (3%)	3420 (5%)	10944 (18%)
	Memory	240 (< 1%)	80640 (7%)	-
	Multipliers	32 (21%)	96 (64%)	-
	IOBs	-	-	178/240 (74%)
Timing	Clockperiod		Frequency	
	10.774 ns		92.816 Mhz	

We can observe from figure 4 that the I/O power is high (38.05 %) because of the throughput and the memory dissipates less power (4.41 %). This is explained by the fact that the architecture is fully parameterized and pipelined.

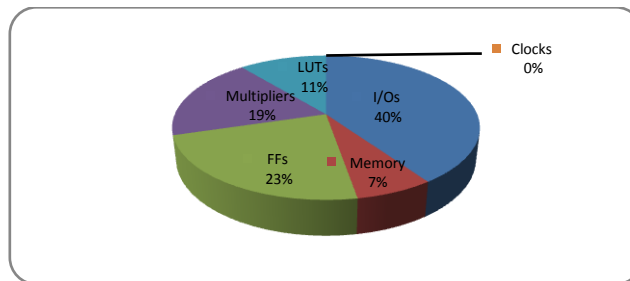


Figure 4. Thermal Power Block

## VI. EXPERIMENTAL RESULTS AND EVALUATION

**EXPERIMENTAL SETUP:**  
**Listeners:** The group is composed of 12 male and female listeners. 10 are normal hearing (NHP) subjects and 2 are hearing-impaired (HIP) subjects of 27 and 45 years old. The age is between 24 and 57 with mean age of 34 year.  
**Speakers:** Utterances are made by native male and female English speakers.  
**Sentence materials:** In total 11 English sentences taken from Arctic speech corpus data base have been.  
**Environment:** The tests were held during an English session. The classroom is located in an isolated place with a background noise (45 Db).  
**Procedures:** The phrases are presented through loudspeakers from computer. We recall entire sentences by both male and female speakers. The corresponding loudness was first adjusted to allow the participants a stable quality of perception and no change of the volume was applied. Each listener was given an Opinion Score Table (OST) to put his own records of the listening quality. The subjects were asked to evaluate sounds on five-point scale (1–5).

Using DWT-OLA, the signal is first denoised as in figure 5(a), (b). Hard threshold technique has been applied because of the high frequency components of the speech which are corrupted by the noise. The pitches are then detected and manipulated in order to make shift of some frequencies before speech reconstruction.

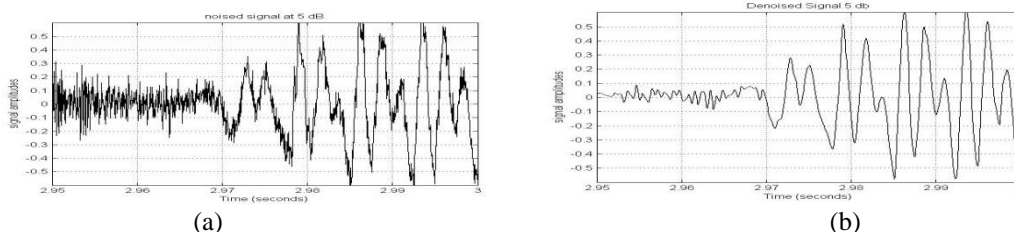


Figure 5. (a): Noised speech signal, (b): Denoised speech signal.

The output signal becomes synthetic for the normal-hearing but is more comprehensible by the impaired-hearing persons. Figures 6 show the generated speech used in the conducted experiments.

# International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 8, August 2014

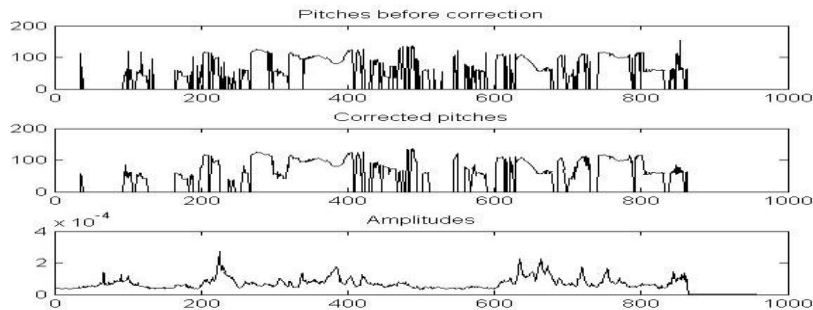


Figure 6. Spectrogram of the synthetic speech signal

Figures 7 show the MOS results obtained from the first conducted experiment. In the graph, we can observe that for the same conditions the Impaired Hearing Persons (IHP) have deficiency in understanding than Normal Hearing Persons (NHP). Also, it is obvious that difficulties appear when female speaking (pitches : 200-300) than with men speaking (pitches 300-500).

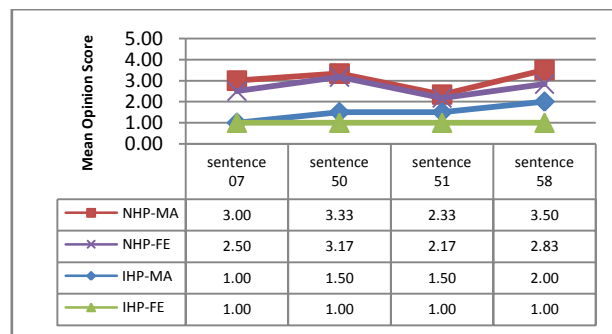


Figure 7. MOS comparison of female versus male listening.

When the speech signal is treated by means of denoising and modification, the gain obtained in the intelligibility for the impaired participants which reaches the 70%. We can see from the graph of figure 8 that the HPL can reach the normal hearings in normal conditions.

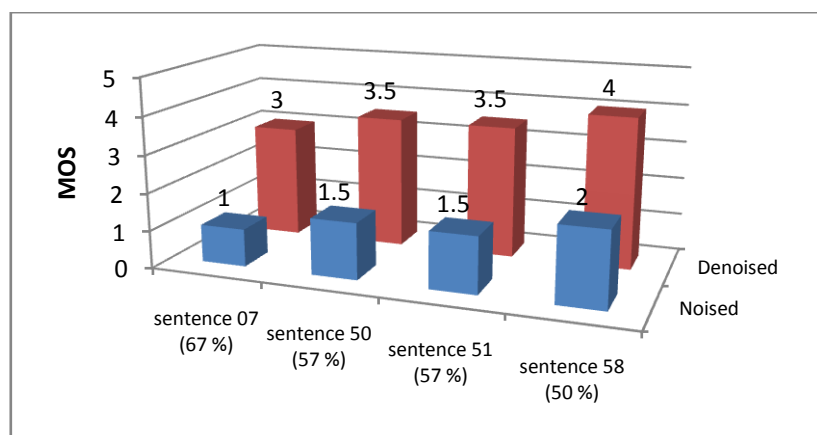


Figure 8. Noised versus denoised speech for impaired persons

**Acknowledgement:** Authors would like to address their great thanks to Professor Pedro OSSES who accepted the conduction of the experiments in his English class session and to the volunteers who participated to these experiments.



# International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 8, August 2014

## VII. CONCLUSION

In this paper, we implemented a platform on FPGA for Hearing-Aid and showed the possible ways to get efficient design using DSP techniques. Using the DWT-OLA, the speech signal is segmented without any distortion. The efficacy of the algorithm was evaluated using subjects with and without hearing deficiency. Listening tests showed that the proposed algorithm increases the quality and intelligibility of the denoised speech. The comparative experiments for the capacity to perceive speech between the normal-hearing and the hearing-impaired people had shown that under identical conditions, the hearing-impaired people have generally very low scores of speech recognition and requires raised Sound Pressure Levels (SPL) to reach the performances of a good hearing. Our aim is to provide an efficient system and the proposed architecture gives satisfactory results based on the evaluation by the Mean Opinion Scores (MOS). Since this embedded device should be portable, the work has also focused on some optimizations namely the reduction of FPGA resources and power consumption. The reconfigurability of the FPGA made possible the use of the DWT algorithm with different parameters so as to meet the specifications for different hearing pathologies. We are currently pursuing our research to design hybrid architecture for noise reduction and echo cancellation.

## REFERENCES

- [1] R. Plomp “Auditory handicap of hearing impairment and the limited benefit of hearing aids”, Journ. of Acoust. Amer. Soc., 63, 533-49, 1978.
- [2] F. Marino, D. Guevorkian and J.T. Astola, “Highly efficient high speed/low power for the 1-D discret wavelet transform”, IEEE Transactions on Analog and Digital Signal Processing Circuits and Systems, Vol. 47, pp. 1492-1502, 2000
- [3] X. Hu, L. DeBrunner and V. DeBrunner, “An efficient design for FIR filters with variable precision”, Proceeding of the IEEE International Symposium on Circuits and Systems Vol. 4, pp. 365-368, May 2002
- [4] S. Edward, and S. Rogers, “FPGA Architecture: Survey and challenges”, Journ. of Found. & Trends in Elect. Desi. Autom. Vol. 2, N° 2, 2007.
- [5] R. Hourani, W Alexander, T. Raithatha, “Automated design space exploration for DSP applications” Journ. of Sign. Proc. Sys. Springer, 2009
- [6] S. Chan, W. Liu and K. Ho “Multiplier less perfect reconstruction modulated filter banks with sum of powers of two coefficients” IEEE Signal Processing Letters, Vol. 8, N° 6, pp. 163-166, June 2001.
- [7] S. Powell and P. Chan “Reduced complexity programmable FIR Filters” IEEE Int. Symposium on Circuits and Systems pp 561-564 May. 1992
- [8] L. Bendaouia et al. “Fast DWT based FPGA implementation for medical application”, IEEE Intern.Conf. on Phealth, Lyon, France, June 2010.
- [9] L. Bendaouia, SM. Karabernou, L. Kessal, H. Salhi and F. Ykhlef, “DWT based FPGA implementation of a reconfigurable platform for a bio-inspired medical hearing aid” International Conference on Systems, Modeling and Design, Istanbul, Turkey Feb. 3<sup>rd</sup>-5<sup>th</sup> 2012
- [10] J.B. Allen et al. “Modelling the noise damaged cochlea”, The mechanics and biophysics of hearing, Springer, pp. 321-332, 1991.
- [11] Y.M. Cheng and D.O. Shaughnessy, “Automatic and reliable estimation of glottal closure instant and period ” IEEE Transaction on Acoustics, Speech and Signal Processing, pp. 1805-1815, 1989.
- [12] S. Roucoux and A. Wilgus “High quality time scale modification of speech ”, IEEE Int. Conf. on Acous. Speech & Sig., pp. 493-496, 1985
- [13] X. Hung et al. “Spoken language processing, a guide to theory, algorithm and system development ”, Prentice Hall Inc 1<sup>st</sup> ed, 2001.
- [14] Y. Laprie and V. Colotte “Automatic Pitch Marking for speech transformations via TDPSOLA”, Proceeding of the European Signal Processing Conference, pp. 1133-1136, 2011.
- [15] J. Flanagan and M. Saslow “Speech analysis, synthesis and perception” Springer, 2<sup>nd</sup> edition, New York 1972.
- [16] J.O. Smith and J.S. Abel “Bark and ERB bilinear transforms”, IEEE Trans. On speech and audio Processing Vol. 7, N° 6, Nov. 1999.
- [17] O. Rioul et al. “Fast algorithms for discrete and continuous wavelet transform”, IEEE Trans. Info. Theory Vol. 38, pp. 569-753, Oct. 1999.
- [18] Altera “DE2 Development and Education Board: User Manual”, Copyright 2006, Altera Corporation V1.4.
- [19] Tim Erjavec, “Introducing the Xilinx targetd design platform”, [www.eetimes.com](http://www.eetimes.com) [Retrieved , February 2<sup>nd</sup>, 2009].

## BIOGRAPHY



**Lotfi Bendaouia** is currently working for his Ph.D. degree at the ETIS laboratory. The main subject concerns the implementation of bio-inspired medical systems for impaired persons. He has completed his Magister in Cybernetics and worked as researcher in the Center of Advanced Technologies (CDTA-Algeria) and as Assistant Professor in ENSEA-France.