

App development with audio applications from m-file to app

Daniel JANKOVIC

September 2, 2015

1 Part 1

Voice Activity Detector (VAD) is a technique used in signal processing to detect the presence of human voice in a signal. It can be an energy detector that indicates speech when the energy of the filtered signal exceeds a predefined threshold. Considered an important technology in speech based communication, today there are various types of applications that use it. Therefore a wide variety of VAD algorithms have been developed to provide the needed features.

There are different kind of stand-alone commercial baby monitors on the market today. From the most basic, that use one-way radio communication, to advance two-way communication monitors that use signal processing to transmit audio when a predefined threshold has been reached. It is also possible to find baby video monitors that broadcast both audio and video when the sensors notice movement. Since most of the monitor applications rely on radio signals to communicate between the units there is a probability that the signal will weaken or possibly not even reach the receiver because it needs to pass through multiple walls of varying thickness. The signal could also be effected by other applications. As the stand-alone monitor focuses on reliability (among other important sale strategies), little is known about the security features. It is possible to assume that the communication is unencrypted, at least in some products, and therefore introduces a potential risk for intrusion of peoples privacy.

To resolve the issues brought up above, an application such as the *Baby Activity Detector* (BAD) can be made more portable, versatile and secure with the help of todays smart-phone technology and VAD. There are many VAD algorithms to choose from and they all have their strengths and weaknesses. Complex algorithms such as Linear Predictive coding (LPC), mel-frequency cepstrum (MFC) are very powerful but quite difficult to grasp and also to implement, they can be considered out of scope for this course. The following VAD algorithms are easy to implement and can be, when combined, quite robust for the task of a basic BAD. The simple short-time energy algorithm calculates the energy levels for each frame to detect voice, unvoiced or silenced regions. Voiced regions will have higher energy levels, however, the algorithm does not

take unwanted noise into account which means that we can have false indication of voice detection. In order to remove the noise from the signal, spectral subtraction can be preformed. In the case of BAD, the threshold needs to be adjusted so that unforeseeable sound is not interpret as the infants cry. Zero-crossing rate (ZCR), is the rate at which a signal changes from plus to minus and back. The higher the rate the higher the frequency which indicates possible voice activity. According to [1] the cry sound that an infant makes has a fundamental frequency of 250-600 Hz (pitch). To be able to use ZCR together with the information above, it is necessary to extract the pitch from the signal in order to match the frequency interval.

The main task of BAD is to detect infant activity, an alternate algorithm is proposed in [1]. It describes an cry detection algorithm that is build up by three main stages. i) *VAD*, a statistical model-based detector [5] is used for detecting sections with sufficient audio acitivity. It also helps to reduce the power consumption. ii) *Classification*, uses k-nearest neighbours (k-NN) algorithm [6] to label each frame as either 'cry' (1), close enough, or 'no cry' (0). iii) Post-processing, validation stage in order to reduce false-negative errors. The idea of having devoted algorithm to detect infant cry is a winning concept for a BAD application, according to the authors it even had promising results in low SNR. Despite simplicity of the algorithm many of the features required to implement were mentioned earlier to be out of scope for this project.

An algorithm that might be of interest [3] suggests a new approach to speech enhancement, without the help of VAD technology. The signal is divided into multiple sub-bands and an noise floor level estimate is calculated simultaneously as the short-time average. The goal is to boost the sub-bands with high Signal-to-Noise Ration (SNR) instead of to suppress the lower. This algorithm has great potential to reduce the noise levels when analyzing incoming signals to the BAD application.

A quick search on the net gives significant amount of hits for smart-phone based BAD applications. The techniques vary, from bluetooth to Wi-Fi and 3-/4G solutions. The award winning application *Baby Monitor 3G* [7] is a feature rich cross platform application that solves most the of issues brought up in this report. It supports both Wi-Fi and 3G/LTE networks, ability to transfer high quality live video, adjustable microphone sensitivity, talk-back functionality and guarantees both reliability and privacy. It can be assumed that the Android based BAD application *Dormi* [8] offers similar features as Baby Monitor 3G, even if Baby Monitor 3G's feature list provides barely any deeper information. It is noteworthy that *Domri* have both *Smart noise detection* and *Adaptive audio enhancement* as sales pitch, which from a engineering point of view is very attractive.

Following is a purposed algorithm for an BAD application

```

while(true){
    % Record audio can place it into the register %
    Get frame from the register
    Divide the signal into different sub-bands with FFT
    Calculate the total short-time energy average
    Calculate noise level for each sub-bands
    if energy average above threshold
        Calculate the gain for each sub-band
        Boost the sub-band with high SNR
        % Extract fundamental frequency (pitch)
        Count the ZCR under 1 seconds
        if the ZCR is within 250-600 Hz
            Possible infant activity detected!
            (Occurs only first time, and needs to be reseted)
    % Send frames %
        Start broadcasting the sound to the receiver
    end
    else
        Reset ZCR
        Dismiss and get next frame from register
    end
}

```

2 Part 2

Three various baby and noise sounds were provided by the institute. The sample frequency for all of the recorded files are 8 kHz. In Figure 1 and Figure 3 the frequency spectrum is plotted for the baby respective noise sounds. Because the energy levels are unproportionally low for two out of the three audio recordings, an enhanced plot is displayed in Figure 2. From the plots it can be seen that the baby frequency interval is between ~ 300 -2500 Hz, including cry and talk, and the noise frequency interval is between ~ 0 -200 Hz and ~ 1300 -2100 Hz. One way to work with only the desired frequencies is to create a filter with a bandpass between 300-1300 Hz.

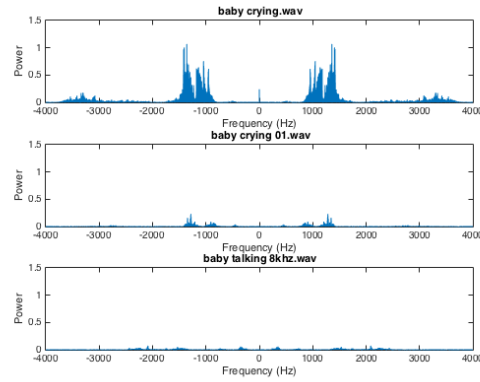


Figure 1: Frequency spectrum for baby sound

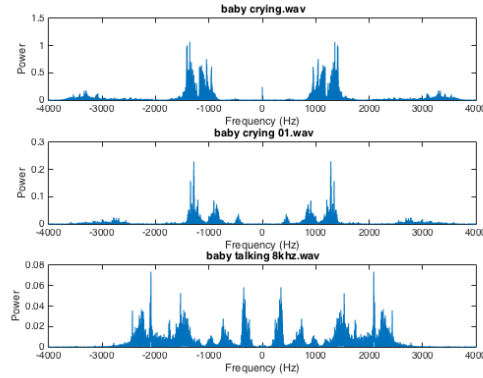


Figure 2: Frequency spectrum for baby sound, scaled y-axis

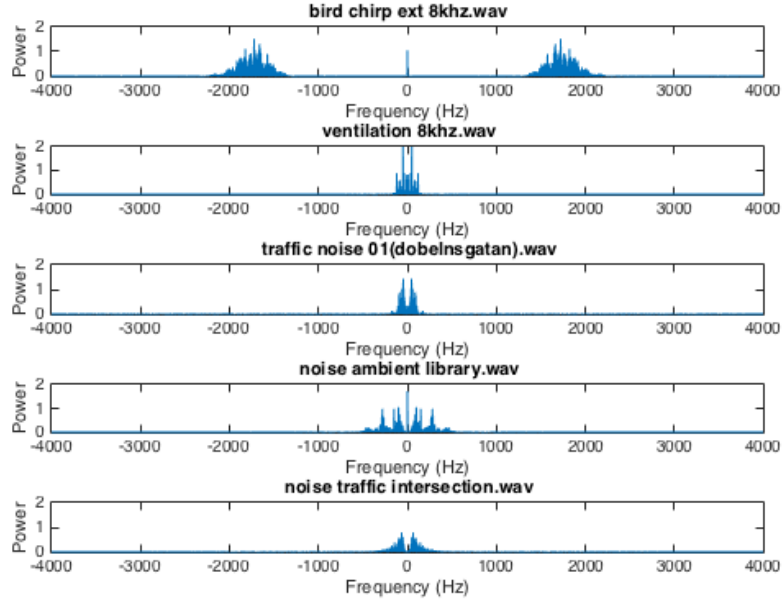


Figure 3: Frequency spectrum for noise files

2.1 The algorithms

A quick and simple algorithm for a BAD application is preferably an algorithm that measures the energy from the input sound. Another benefit of using this technique is that the implementation in Android OS may be easier for a novice application developer. Mathematically described, the power equation, which is close relationship to energy, is given by the following formula:

$$P(n) = \frac{1}{N} \sum_{k=0}^{N-1} x^2(n - k)$$

Since the BAD application will be performing these power calculations of the input sound in real-time, it is undoubtedly impossible to implement the equation above. A solution to the problem is to use the *recursive averaging* algorithm. It is small and hardware friendly algorithm that calculates the power of a given input without using too much memory. Given the equation below,

$$P(n) = \alpha P(n - 1) + (1 - \alpha)x^2(n)$$

instead of performing the calculation for one sample at the time, a predefined number of samples in blocks, referred to as frames, are squared and summed.

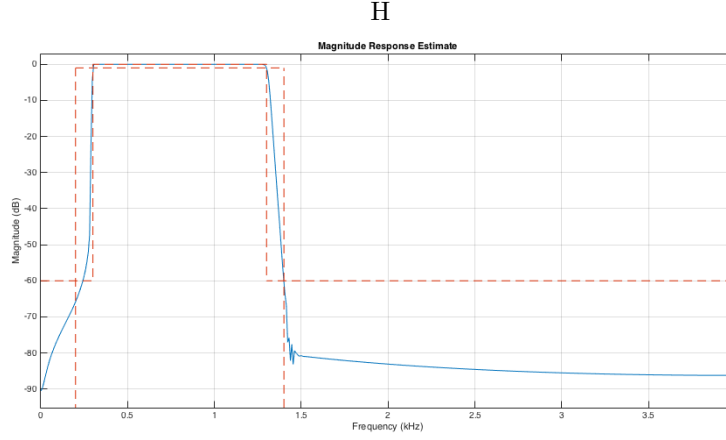


Figure 4: Butterworth 300-1300 Hz bandpass filter

Each frame represents $x^2(n)$. Every result, $P(n)$, is saved to be used as $P(n-1)$. The α is a constant between 0 and 1 and is related to the following formula:

$$\alpha = \frac{1}{T_s F_s}$$

Despite that the α can be derived from the equation, to get the better results further tweaking and testing is required. The α value used in this research, 0.5, suggests that the new $P(n)$ is equally weighted between old and new calculations. The MATLAB code for the simple algorithm can be found in Appendix B.

The advanced algorithm is based upon the same, *recursive averaging*, algorithm but before calculating the power of the input sound, the signal is first filtered through a Butterworth bandpass filter to remove unwanted frequencies. Butterworth filter, because it gives the least amount of ripple in the bandpass. By having the noise frequencies suppressed, more precision is acquired for finding baby activity. The MATLAB code for the filter can be found in Appendix A.

2.2 Evaluation and performance

To evaluate and test the performance of the simple and advance algorithm new sound files were created, the three provided baby sound files created the based. The goal was to make the alarm go off with the sound that the baby caused, in clean and noisy configurations. For each baby sound file, three new sound files were created to mislead the algorithms. The MATLAB code for this can be found in Appendix C. Giving three sets of 4 different test configurations, each set contained the following files:

1. clean baby sound without any noise

2. simulate early mornings, bird and ventilation noise added to the base sound
3. simulate daily environment, all noise files added to the base sound
4. simulate *extremely* noisy environment, all noise files amplified and added to the base sound

As can be seen from Figure 5, the green horizontal line indicates the threshold value. The red circle indicates where the alarm has been set off, when the algorithm has notified the user of baby activity. If the algorithm failed to set off the alarm, the red circle is placed in the origo. In this test environment the alarm was set off after five frames successfully breached the threshold value.

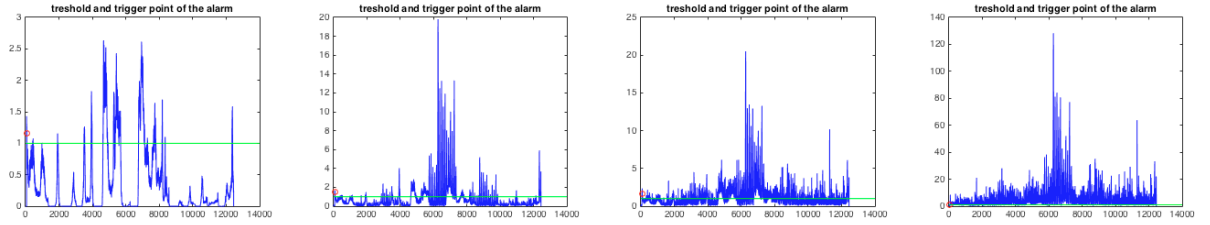


Figure 5: Baby crying.wav, simple algorithm

The test results from the *Baby crying.wav* set, Figure 6, show that the alarm was set off in approximate same time in every configuration and the reason for this could be that there is a high energy baby sound in the beginning of all of the files. The *clean* and the *bird and ventilation* configurations can be considered as plausibly successful results since the noise levels are not interfering as much as the two, to the right, configurations.

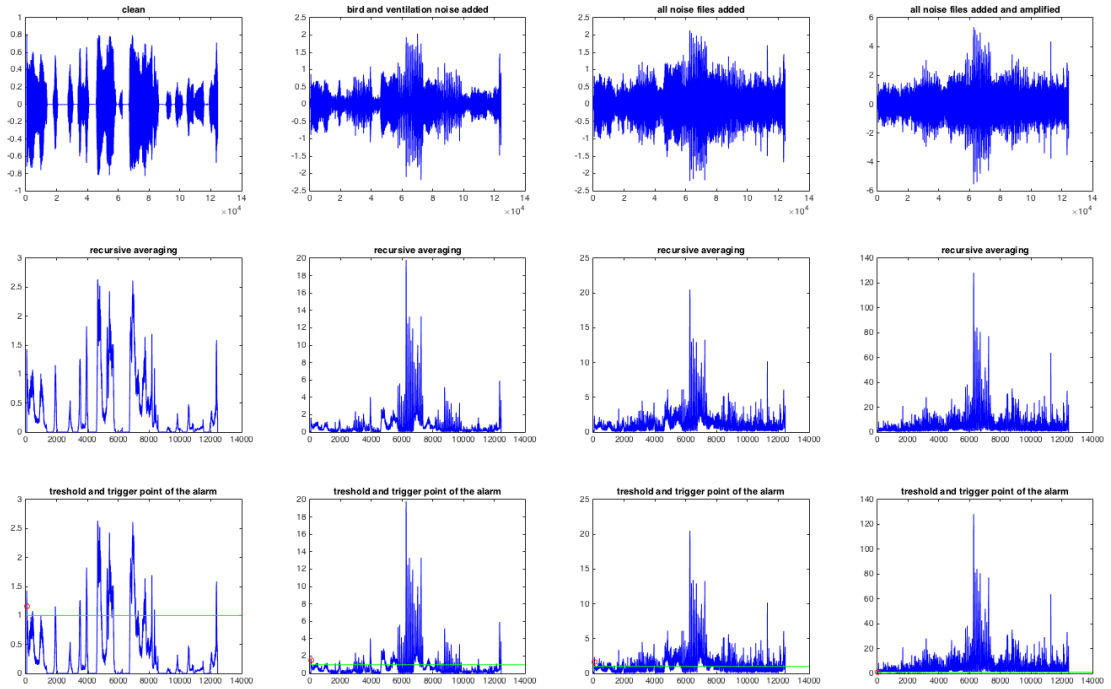


Figure 6: Baby crying.wav, simple algorithm

Both *Baby crying1.wav* and *Baby talking.wav* gave unsatisfactory results without any amplification, as can be seen in Figure 7 and Figure 8. The success of setting of the alarm in the other configurations were due to the noise.. Not remotely close to trigger the alarm in clean configuration, illustrated by the left most plots, the test were performed a second time with amplified baby sound.

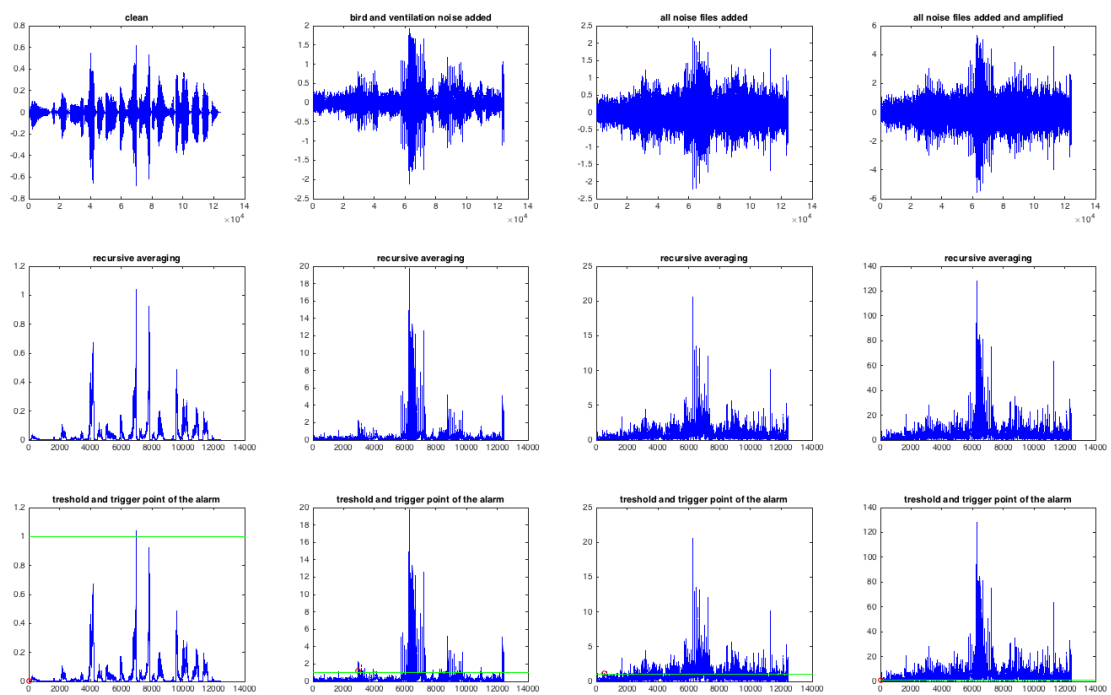


Figure 7: Baby crying1.wav, simple algorithm

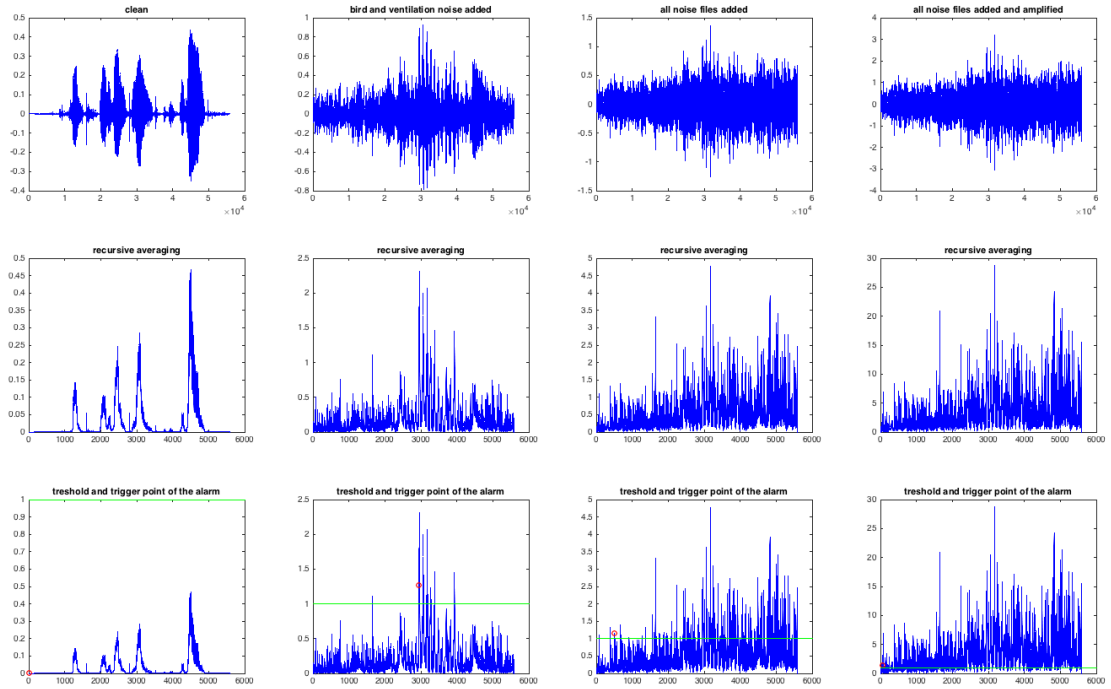


Figure 8: Baby talking.wav, simple algorithm

By having the baby sounds amplified before the noise is added, making sure that the alarm sets off in clean configuration, gives a possibility for the baby sounds to trigger the alarm in other configurations as well, see Figure 9 and Figure 10. Unfortunately, the simple algorithm performed poorly for the *Baby crying1.wav* set, Figure 9, but in the *Baby talking.wav* set, Figure 10, the first two configurations, *clean* and *bird and ventilation*, performed surprisingly well.

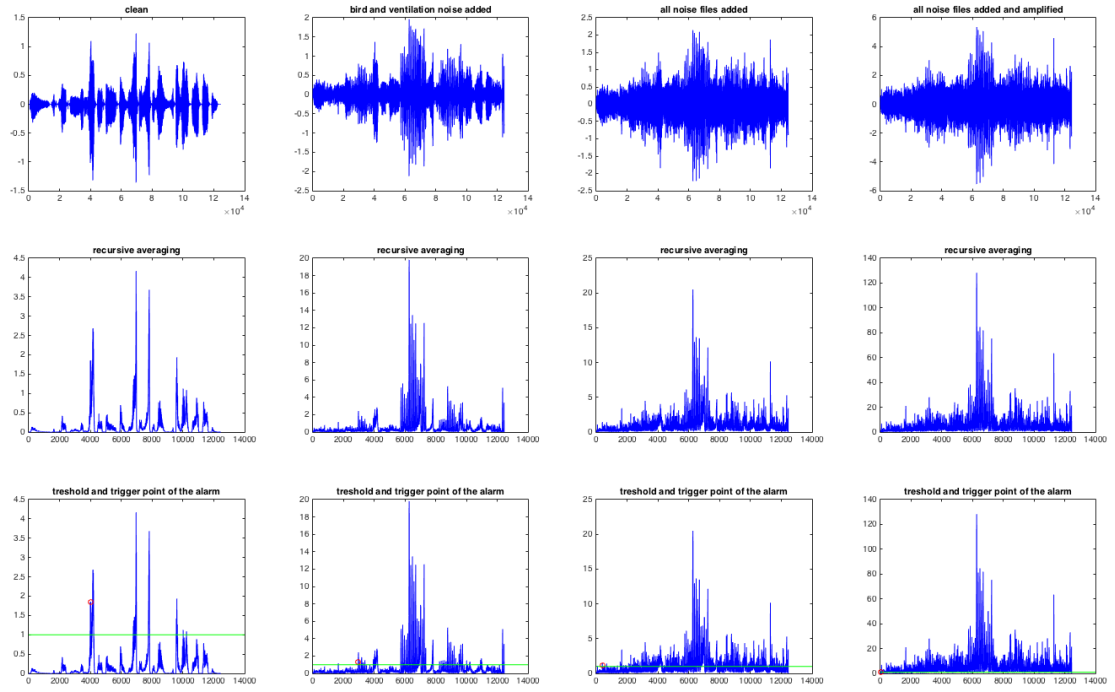


Figure 9: Amplified Baby crying1.wav, simple algorithm

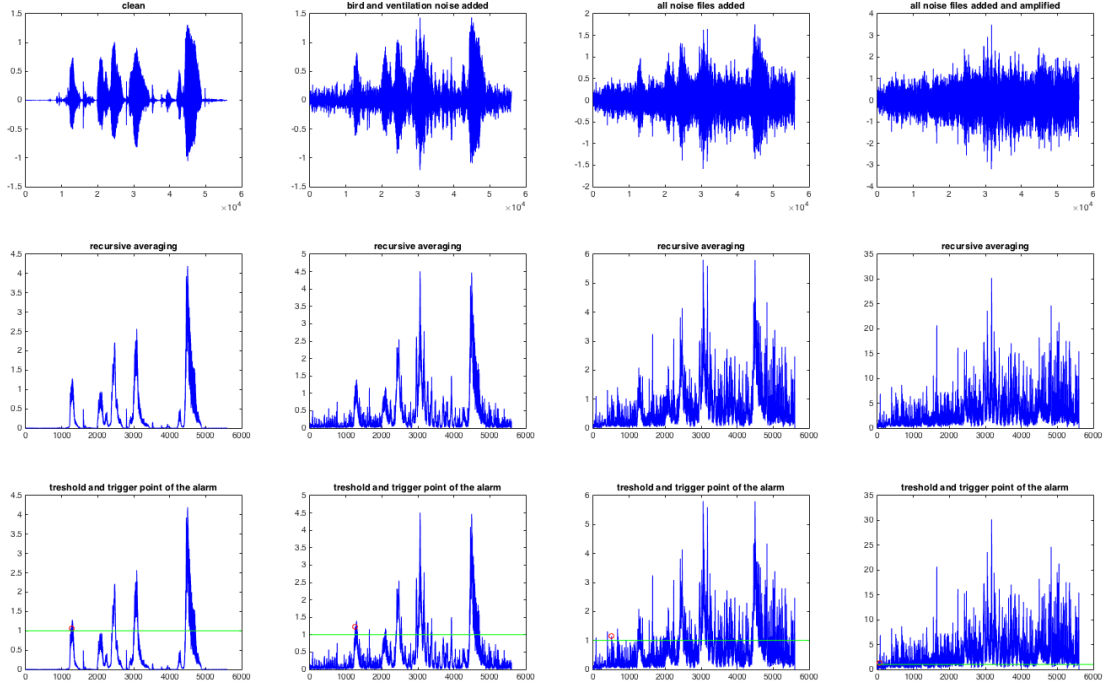


Figure 10: Amplified Baby talking.wav, simple algorithm

The advanced algorithm performed, as expected, exceptionally well due to unwanted frequencies being suppressed, not interfering with the power calculations. Figure 11, Figure 12 and Figure 13 show great results, except for the early trigger in the last configuration of the *Baby talking.wav* set. The alarm was set off earlier because the signal was infected with noise. Despite the early trigger of the alarm, it was executed during baby activity and not by the added noise. The alarm could have been triggered during the same time in other configurations of that set as well.

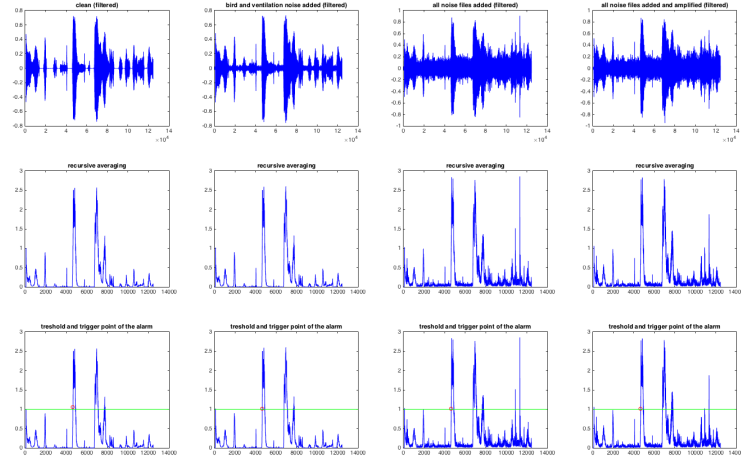


Figure 11: Baby crying.wav, advanced algorithm

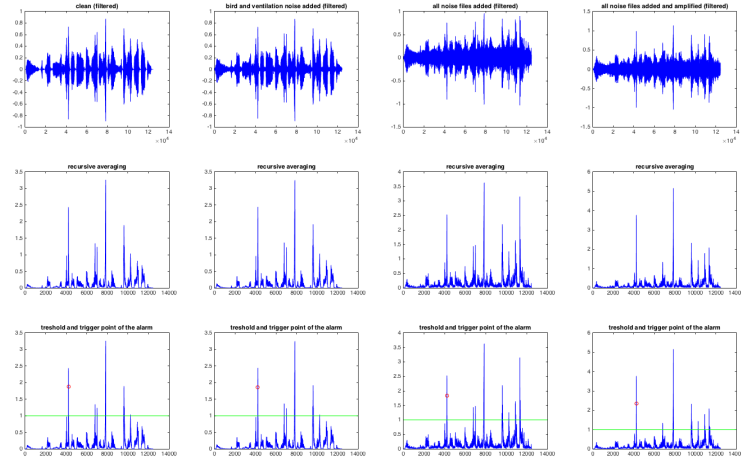


Figure 12: Amplified Baby crying1.wav, advanced algorithm

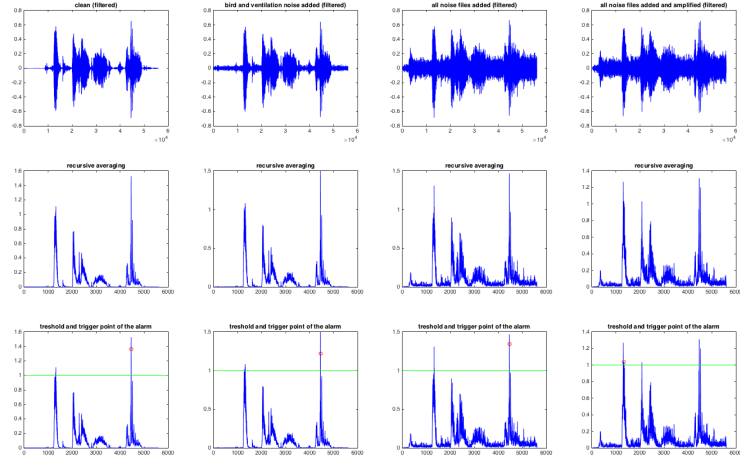


Figure 13: Amplified Baby talking.wav, advanced algorithm

2.3 Conclusion

Giving the results from the two different algorithms, it is fair to state the simple algorithm was unreliable since the noise levels could trigger the alarm. If the design goal is accuracy, reliability and good performance for the BAD application then implementing the advanced algorithm would be the best option. However, both algorithms rely on the *recursive averaging* algorithm. Implementing it must be done either way and should this be done within short amount of time, an attempt should be preformed to implement the Butterworth bandpass filter. If the filter issue can not be overcome with built-in Java classes, third-party libraries, or guidelines from the institute then the advanced algorithm will not be implemented. This should be considered as bonus task for Part 3 of the course.

3 Part 3 (4)

The last section of this report, *Part 3*, actually represents *Part 4* in the "step-by-step" instructions provided from the institute.

As mentioned in *Part 2*, the *recursive averaging* algorithm was crucial for both the simple and advanced algorithm. The implementation required comprehensive study of the framework that was provided since the author of this report had no previous experience of Android application development. The first modification was performed in *startFragment.java*, *run()* method.

```
public void run()
{
    ...
    try
    {
        int N = AudioRecord.getMinBufferSize(8000,
            AudioFormat.CHANNEL_IN_MONO,
            AudioFormat.ENCODING_PCM_16BIT); /*N:640*/
        recorder = new AudioRecord(AudioSource.MIC, 8000,
            AudioFormat.CHANNEL_IN_MONO,
            AudioFormat.ENCODING_PCM_16BIT, N*10);

        recorder.getState();
        recorder.startRecording();

        /*
         * Loops until something outside of this
         * thread stops it. Reads the data from the
         * recorder and calculates an average to use
         * in the baby detector.
         */
        while(!stopped)
        {
            //was 160 now 10
            short[] buffer = new short[10];
            N = recorder.read(buffer,0,buffer.length);
            double sum = 0;
            for (short val : buffer) {
                sum = sum + Math.abs(val);
            }

            /**
             * To scale it down, dive by buffer length
             */
            double average = sum / buffer.length;
            recursiveSum =
                (ALPHA * oldRecursiveSum) + ((1-ALPHA)*average);
            oldRecursiveSum = recursiveSum;

            final BabyState babyState =
                detector.updateState(recursiveSum);

            ...
        }
    }
}
```

}

4 Appendix

4.1 A *butterFilter.m*

```
% Values for the bandpass
A_stop1 = 50;
F_stop1 = 200;
F_pass1 = 300;
F_pass2 = 1300;
F_stop2 = 1400;
A_stop2 = 50;
A_pass = 1;
Fs = 8000;
filtering = 1;

% Design and create Bandpass filter
BandPassSpecObj = fdesign.bandpass(F_stop1, F_pass1, ...
    F_pass2, F_stop2, A_stop1, A_pass, A_stop2, Fs);
BandPassFilt = design(BandPassSpecObj, 'butter');

% Graphical representaion of the bandpass filter
% Uncomment for usage, requires DSP toolbox for MATLAB
fvtool(BandPassFilt)
```

4.2 B *recursiveAverg.m*

```
function [ P ] = recursiveAverg( buffer, alpha )
% SIMPLE.RECURSIVEAVERG function calculates reverse averaging
% given a squared buffer. The buffer is squared and the frames
% are summed up and placed in a cell in vec-array. Each cell
% represent a frame and is later used in the recursive averaging
% algorithm. The alpha value is 0.5.
%
% Input: buffSquared, alpha
% Output: P-array with recursive averaging power calculations

% The buffer-matrix is squared and summed to fit in
% vec-array
n=size(buffer,2);
buffSquared=buffer.^2;
for i=1:n
    vec(i,1)=sum(buffSquared(1:end,i));
end

% The vec-array is used to perform recursive averaging, which
% is stored in P-array
rows = size(vec,1);
P=zeros(rows,1);
P(1,1)=0;
for i=2:size(vec,1)-1
    P(i,1)=alpha*P(i-1,1) + (1-alpha)*vec(i,1);
end
end
```

4.3 C *makeSound.m*

```
% Clear and close all
close all
clear all

% Read all audio files and extract fs
% Baby sounds
baby_crying_1 = ('baby-signals/baby-crying.wav');
[x_baby_crying_1, fs_baby_crying_1] = audioread(baby_crying_1);
baby_crying_2 = ('baby-signals/baby-crying-01.wav');
[x_baby_crying_2, fs_baby_crying_2] = audioread(baby_crying_2);
baby_talking = ('baby-signals/baby-talking-8khz.wav');
[x_baby_talking, fs_baby_talking] = audioread(baby_talking);
% Noise sounds
noise_bird = ('noise-signals/bird-chirp-ext-8khz.wav');
[x_noise_bird, fs_noise_bird] = audioread(noise_bird);
noise_traffic = ('noise-signals/traffic-noise-01(dobelnskatan).wav');
[x_noise_traffic, fs_noise_traffic] = audioread(noise_traffic);
noise_ventilation = ('noise-signals/ventilation-8khz.wav');
[x_noise_ventilation, fs_noise_ventilation] = audioread(noise_ventilation);
noise_amb_lib = ('noise-signals/noise-ambient-library-2.wav');
[x_noise_amb_lib, fs_noise_amb_lib] = audioread(noise_amb_lib);
noise_traff_inter = ('noise-signals/noise-traffic-intersection.wav');
[x_noise_traff_inter, fs_noise_traff_inter] = audioread(noise_traff_inter);

% Add noise to baby recordings
L_baby_crying_1=length(x_baby_crying_1);
L_baby_crying_2=length(x_baby_crying_2);
L_baby_talking=length(x_baby_talking);

L_noise_bird=length(x_noise_bird);
L_noise_traffic=length(x_noise_traffic);
L_noise_ventilation=length(x_noise_ventilation);
L_noise_amb_lib=length(x_noise_amb_lib);
L_noise_traff_inter=length(x_noise_traff_inter);

xL_baby_crying1=min(L_baby_crying_1,min(L_noise_bird,min(L_noise_traffic,...
    min(L_noise_ventilation,min(L_noise_amb_lib,L_noise_traff_inter)))));
xL_baby_crying2=min(L_baby_crying_2,min(L_noise_bird,min(L_noise_traffic,...
    min(L_noise_ventilation,min(L_noise_amb_lib,L_noise_traff_inter)))));
xL_baby_talking=min(L_baby_talking,min(L_noise_bird,min(L_noise_traffic,...
    min(L_noise_ventilation,min(L_noise_amb_lib,L_noise_traff_inter)))));

% Version:
% 0 - clean, without any noise
% 1 - slightly amp bird & vent noise added      Nstrength=2
% 2 - slightly amp (all) noise files added      Nstrength=2
% 3 - highly amp (all) noise files added        Nstrength=5

% Baby crying 1
x_BC10=x_baby_crying_1(1:xL_baby_crying1);
Nstrength = 2;
x_BC11=x_baby_crying_1(1:xL_baby_crying1)+...
    Nstrength*x_noise_bird(1:xL_baby_crying1)+...
    Nstrength*x_noise_ventilation(1:xL_baby_crying1);
```

```

x_BC12=x_baby_crying_1(1:xL_baby_crying1)+...
    Nstrength*x_noise_bird(1:xL_baby_crying1)+...
    Nstrength*x_noise_traffic(1:xL_baby_crying1)+...
    Nstrength*x_noise_ventilation(1:xL_baby_crying1)+...
    Nstrength*x_noise_amb_lib(1:xL_baby_crying1)+...
    Nstrength*x_noise_traff_inter(1:xL_baby_crying1);
Nstrength = 5;
x_BC13=x_baby_crying_1(1:xL_baby_crying1)+...
    Nstrength*x_noise_bird(1:xL_baby_crying1)+...
    Nstrength*x_noise_traffic(1:xL_baby_crying1)+...
    Nstrength*x_noise_ventilation(1:xL_baby_crying1)+...
    Nstrength*x_noise_amb_lib(1:xL_baby_crying1)+...
    Nstrength*x_noise_traff_inter(1:xL_baby_crying1);

% Baby Crying 2
N=2;
x_BC20=N*x_baby_crying_2(1:xL_baby_crying2);
Nstrength = 2;
x_BC21=N*x_baby_crying_2(1:xL_baby_crying2)+...
    Nstrength*x_noise_bird(1:xL_baby_crying2)+...
    Nstrength*x_noise_ventilation(1:xL_baby_crying2);
x_BC22=N*x_baby_crying_2(1:xL_baby_crying2)+...
    Nstrength*x_noise_bird(1:xL_baby_crying2)+...
    Nstrength*x_noise_traffic(1:xL_baby_crying2)+...
    Nstrength*x_noise_ventilation(1:xL_baby_crying2)+...
    Nstrength*x_noise_amb_lib(1:xL_baby_crying2)+...
    Nstrength*x_noise_traff_inter(1:xL_baby_crying2);
Nstrength = 5;
x_BC23=N*x_baby_crying_2(1:xL_baby_crying2)+...
    Nstrength*x_noise_bird(1:xL_baby_crying2)+...
    Nstrength*x_noise_traffic(1:xL_baby_crying2)+...
    Nstrength*x_noise_ventilation(1:xL_baby_crying2)+...
    Nstrength*x_noise_amb_lib(1:xL_baby_crying2)+...
    Nstrength*x_noise_traff_inter(1:xL_baby_crying2);

% Baby talking
N=3;
x_BT0=N*x_baby_talking(1:xL_baby_talking);
Nstrength = 2;
x_BT1=N*x_baby_talking(1:xL_baby_talking)+...
    Nstrength*x_noise_bird(1:xL_baby_talking)+...
    Nstrength*x_noise_ventilation(1:xL_baby_talking);
x_BT2=N*x_baby_talking(1:xL_baby_talking)+...
    Nstrength*x_noise_bird(1:xL_baby_talking)+...
    Nstrength*x_noise_traffic(1:xL_baby_talking)+...
    Nstrength*x_noise_ventilation(1:xL_baby_talking)+...
    Nstrength*x_noise_amb_lib(1:xL_baby_talking)+...
    Nstrength*x_noise_traff_inter(1:xL_baby_talking);
Nstrength = 5;
x_BT3=N*x_baby_talking(1:xL_baby_talking)+...
    Nstrength*x_noise_bird(1:xL_baby_talking)+...
    Nstrength*x_noise_traffic(1:xL_baby_talking)+...
    Nstrength*x_noise_ventilation(1:xL_baby_talking)+...
    Nstrength*x_noise_amb_lib(1:xL_baby_talking)+...
    Nstrength*x_noise_traff_inter(1:xL_baby_talking);

```

References

- [1] R. Cohen, Y. Lavner, *Infant Cry Analysis and Detection*, 2012
- [2] E. Verteletskaya, K. Sakhnov, *Voice Activity Detection for Speech Enhancement Applications*, 2010
- [3] N. Westerlund, M. Dahl, *Speech Enhancement using an Adaptive Gain Equalizer*, 2003
- [4] R. Narayanam, *An Efficient Peak Valley Detection based VAD Algorithm for Robust Detection of Speech Auditory Brainstem Responses*, 2013
- [5] J. Sohn, N.S. Kim, W. Sung, *A Statistical Model-Based Voice Activity Detection*, 1999
- [6] O. Sutton, *A Statistical Model-Based Voice Activity Detection*, 2012
- [7] TappyTaps, <https://www.babymonitor3g.com>,
- [8] Sleekbit, <http://dormi.sleekbit.com/index.html>,
- [9] D. Jankovic, M. Johansson, M. Lichota, *Adaptive Gain Control in Digital Signal Processors*, 2015