

# Binaural Techniques with Cross-Talk Canceled Loudspeakers

Yangang Cao

March 6, 2019

Binaural recordings are meant to be played back in such a way that the sound which originates from the left ear is played back only to the left ear, and correspondingly with the right ear. If such a recording is played back with stereophonic setup of loudspeakers, the sound from the left loudspeaker also travels to the right ear, and vice versa, called cross-talk, which ruins the spatial audio quality.

In order to be able to listen to binaural recordings over two loudspeakers, some methods have been proposed. In these methods, the loudspeakers are driven in such a way that in practice the cross-talk is canceled as much as possible.

A system can be formed to deliver binaurally recorded signals to the listener's ears using two closely spaced loudspeakers with cross-talk cancellation. The binaural signals are represented as a  $2 \times 1$  vector in  $\mathbf{x}(n)$ , and the produced ear canal signals also as  $2 \times 1$  vector  $\mathbf{d}(n)$ . The system can be formulated in the  $z$ -domain

$$\mathbf{d}(z) = \mathbf{C}(z)\mathbf{H}(z)\mathbf{x}(z),$$

where  $\mathbf{C}(z) = \begin{bmatrix} C_{11}(z) & C_{12}(z) \\ C_{21}(z) & C_{22}(z) \end{bmatrix}$  contains the electro-acoustical responses of the loudspeakers measured in the ear canals, and  $\mathbf{H}(z) = \begin{bmatrix} H_{11}(z) & H_{12}(z) \\ H_{21}(z) & H_{22}(z) \end{bmatrix}$  contains the responses for performing inverse filtering to minimize the cross-talk. Ideally,  $\mathbf{x}(z) = \mathbf{d}(z)$ , which can be obtained if  $\mathbf{H}(z) = \mathbf{C}(z)^{-1}$ . Unfortunately, the direct inversion is not feasible due to unidealities of the loudspeakers and the listening conditions. A regularized method to find an optimal  $\mathbf{H}_{opt}(z)$  has been proposed

$$\mathbf{H}_{opt}(z) = [\mathbf{C}^T(z^{-1})\mathbf{C}(z) + \beta\mathbf{I}]^{-1}\mathbf{C}^T(z^{-1})z^{-m},$$

where  $\beta$  is a positive scalar regularization factor, and  $z^{-m}$  models the time delay due to the sound reproduction system. If  $\beta$  is selected very low, there will be sharp peaks in the resulting time-domain inverse filters, which may exceed the dynamic range of the loudspeakers. If  $\beta$  is selected to be higher, the inverse filter will have longer duration in time, which is less demanding on the loudspeakers, but unfortunately the inversion is also less accurate.

A **MATLAB** example is provided in the following to compute inverse filters for a cross-talk canceling system:

- The responses in  $\mathbf{C}$  are moved into the frequency domain with discrete Fourier transform (DFT) with the desired length of time window.
- The filter responses are computed by  $\mathbf{H}_{opt}(k) = [\mathbf{C}^H(k)\mathbf{C}(k) + \beta\mathbf{I}]^{-1} \mathbf{C}^H(k)$ , where  $k$  presents the frequency bin indexes and  $\mathbf{H}$  Hermitian transposition.
- The inverse DFT is taken of  $\mathbf{H}$ , resulting in the inverse filters for cross-talk cancellation.
- A circular shift of half of the applied time-window length is implemented on the inverse filters.

```

1  % Simplified cross-talk canceler
2
3  theta = 10; % spacing of stereo loudspeakers in azimuth
4  Fs = 44100; % sample rate
5  b = 10^-5; % regularization factor
6  % loudspeaker HRIRs for both ears (ear_num,loudspeaker_num)
7  % If more realistic HRIRs are available, pls use them
8  HRIRs(1,1,:) = simpleHRIR(theta/2,Fs);
9  HRIRs(1,2,:) = simpleHRIR(-theta/2,Fs);
10 HRIRs(2,1,:) = HRIRs(1,2,:);
11 HRIRs(2,2,:) = HRIRs(1,1,:);
12 Nh = length(HRIRs(1,1,:));
13 %transfer to frequency domain
14 for i = 1:2
15     for j = 1:2
16         C_f(i,j,:) = fft(HRIRs(i,j,:),Nh)
17     end
18 end
19 % Regularized inversion of matrix C
20 H_f = zeros(2,2,Nh);
21 for k = 1:Nh
22     H_f(:, :, k) = inv((C_f(:, :, k)'*C_f(:, :, k)+eye(2)*b))*C_f(:, :, k)';
23 end
24 % Moving back to time domain
25 for k = 1:2
26     for m = 1:2
27         H_n(k,m,:) = real(ifft(H_f(k,m,:)));
28         H_n(k,m,:) = fftshift(H_n(k,m,:));
29     end
30 end
31 % Generate binaural signals. Any binaural recording should also ...
    be ok
32 binauralsignal = simplehrtfconv(70);
33 %binauralsignal = wavread('road_binaural.wav');
34 % Convolve the loudspeaker signals
35 loudpsig = [conv(reshape(H_n(1,1,:),Nh,1), binauralsignal(:,1)) ...
    + ...
36 conv(reshape(H_n(1,2,:),Nh,1), binauralsignal(:,2)) ...
37 conv(reshape(H_n(2,1,:),Nh,1), binauralsignal(:,1)) + ...
38 conv(reshape(H_n(2,2,:),Nh,1), binauralsignal(:,2))];
39 soundsc(loudpsig,Fs) % play sound for loudspeakers

```

In practice, this method works best with loudspeakers close to each other, as a larger loud- speaker base angle would lead to coloration at lower frequencies. The listening area in which the effect is audible is very small, as if the listener

departs from the mid line between the loudspeakers by about 1–2 cm, the effect is lost.

A nice feature of this technique is that the sound is typically externalized. This may be due to the fact that head movements of the listener produce somewhat relevant cues, and since the sound is reproduced using far-field loudspeakers generating plausible monaural spectral cues. However, although the sound is externalized, a surrounding spatial effect is hard to obtain with this technique. With a stereo dipole in the front, the reproduced sound scene is typically perceived only at the front.

The technique also is affected by the reflections and reverberation of the listening room. It works best only in spaces without prominent reflections. To get the best results, the HRTFs of the listener should be known, however already very plausible results can be obtained with generic responses.