

# CCIT Report 810 EE 1767

## Example-based Cross-Modal Denoising: Supplementary Document

Dana Segev and Yoav Y. Schechner  
 Dept. Electrical Engineering  
 Technion - Israel Inst. Technology  
 Haifa 32000, ISRAEL  
 sdanone@tx.technion.ac.il  
 yoav@ee.technion.ac.il

Michael Elad  
 Dept. Computer Science  
 Technion - Israel Inst. Technology  
 Haifa 32000, ISRAEL  
 elad@cs.technion.ac.il

### Abstract

In Ref. [1] Example-based Cross-Modal Denoising we refer to temporal cropping of audio-visual data and fine alignment during synthesis of a denoised audio. In this additional document we give more details.

### 1. Temporal Density

The input segments are extracted from the input video. A simple way is to extract segments using a fixed period of  $p_F$  frames, as illustrated in Fig. 1. Then, the initial frame in segment  $m$  is

$$f_m^0 = 1 + (m - 1)p_F . \quad (1)$$

Segment  $m$  comprises frames  $[f_m^0, \dots, (f_m^0 + N_F - 1)]$  of the input video, where  $N_F$  is the number of frames per segment. Note that  $p_F \leq N_F$ , so that each video frame is extracted in at least one segment. If  $p_F < N_F$ , there is temporal overlap between consecutive video segments. Similarly, the segment encapsulates  $N_S$  input audio samples, indexed

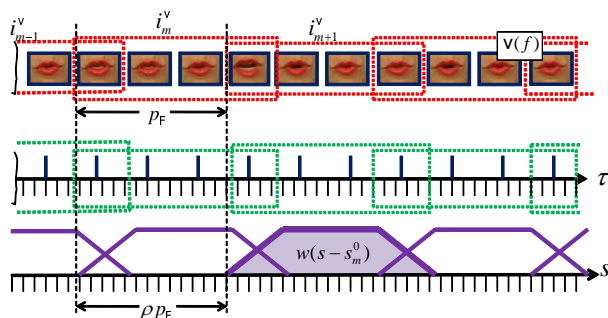


Figure 1. Periodic extraction of video and audio segments. Consecutive segments partially overlap in time. [Bottom]: Mosaicing output audio segments is done by fading in and out each segment, using a weighting function  $w(s - s_m^0)$ .

by  $[s_m^0, \dots, (s_m^0 + N_S - 1)]$ . Here,

$$s_m^0 = 1 + (m - 1)\rho p_F , \quad (2)$$

where

$$\rho = \Delta\tau^V / \Delta\tau^A \quad (3)$$

is the ratio between the temporal sampling periods  $\Delta\tau^A$  and  $\Delta\tau^V$  of the continuous-time audio and inputs. We used  $\rho p_F < N_S$  in order to have temporal overlap between audio segments.

There is a difference between *input* and *example* cropping of data. To fully exploit the short training, it is sometimes useful to extract as many examples as possible from the training video. The highest yield of segments to the example set is achieved if a new example segment  $k$  starts after each frame (as if  $p_F = 1$ ):

$$f_k^0 = k, \quad s_k^0 = 1 + (k - 1)\rho . \quad (4)$$

The hidden Markov model, described in Section. 4 of [1] generally does *not* use  $p_F = 1$ . Rather the value of  $p_F$  for the model is the one used in the *input* segments, as we describe now.

For an *input speech* sequence to be denoised,  $\rho p_F$  is in the order of  $N_S$ : we used  $p_F = 7$ , corresponding to 0.28sec. There are two benefits for this diluted extraction. First, this time span is comparable the duration of a syllable. Therefore, we can use directly the hidden Markov model, which we describe in Section. 4 of [1]. Had the segments been much shorter, capturing the relation between consecutive segments would have diminished. The second benefit is reduction of processing time.

This said, a periodic extraction as in Eqs. (1,2,4) is generally *not necessary*, in neither input nor training. It should often be more efficient to extract segments adaptively, focusing on audio onsets. Alternatively, the example set  $\mathbf{E}$  can be automatically *pruned*, to eliminate useless example segments. This is what we do for music sequences:  $p_F = 1$  at the input, while examples correspond only to onsets and their lingering sound.

## 2. Tuning to Audio Resolution

As written in Section 6 of [1], the digital audio track example  $\mathbf{e}_{\hat{k}_m}^A$  is a clean version of the noisy input  $\mathbf{i}_m^A$ . This is one example out of the example set  $\{\mathbf{e}_k\}_{k=1}^{N_E}$ . As described above (Sec. 1), this set is extracted in increments that cannot be smaller than a single-frame, i.e, temporal lag of  $\Delta\tau^V$ . This is coarse, relative to the audio dynamics. The audio we seek from  $\mathbf{e}_{\hat{k}_m}$  might not be  $\mathbf{e}_{\hat{k}_m}^A$ , but a slightly shifted track, with a temporal lag smaller than  $\Delta\tau^V$ . Such a slight temporal shift in the training audio is not resolvable by the set of examples  $\mathbf{E}$ .

To handle this issue, we perform a finer association of inputs to examples. It proceeds after the cross-modal process described in [1]. However, in the temporally-fine association to examples, only the audio modality is used. Recall from [1] that the data in *audio-example*  $k$  is

$$\mathbf{e}_k^A = [a_e(s_k^0), a_e(s_k^0 + 1), \dots, a_e(s_k^0 + N_S - 1)]. \quad (5)$$

Therefore,  $\mathbf{e}_{\hat{k}_m}^A$  is created by extracting samples  $[s_{\hat{k}_m}^0, \dots, (s_{\hat{k}_m}^0 + N_S - 1)]$  from a long soundtrack. Similarly, a temporally shifted sequence  $\mathbf{e}_{\hat{k}_m, \delta}^A$  can be extracted:

$$\mathbf{e}_{\hat{k}_m, \delta}^A = [a_e(s_{\hat{k}_m}^0 + \delta), \dots, a_e(s_{\hat{k}_m}^0 + \delta + N_S - 1)]. \quad (6)$$

Here  $\delta \in [-\rho N_F, \dots, \rho N_F]$ , i.e, the shifts are much subtler than in  $\mathbf{E}$ . Each of these finely shifted audio examples yields a corresponding feature vector  $\tilde{\mathbf{e}}_{\hat{k}_m, \delta}^A = \mathcal{P}(\mathbf{e}_{\hat{k}_m, \delta}^A)$ . The shift that optimizes the match of the audio example to the input audio is

$$\hat{\delta}_m = \arg \min_{\delta \in [-\rho N_F, \dots, \rho N_F]} d_A(\tilde{\mathbf{i}}_m^A, \tilde{\mathbf{e}}_{\hat{k}_m, \delta}^A). \quad (7)$$

Substituting  $\hat{\delta}_m$  derived in Eq. (7) for  $\delta$  in Eq. (6) results in the estimated clear audio segment  $m$ ,  $\mathbf{e}_{\hat{k}_m, \hat{\delta}_m}^A$ . This is the denoised version of  $\mathbf{i}_m^A$ .

The vector  $\mathbf{e}_{\hat{k}_m, \hat{\delta}_m}^A$  is simply mapped to segment  $m$  of the output audio. Define  $h_m = (s_{\hat{k}_m}^0 + \hat{\delta}_m) - s_m^0$ . Then, create an output soundtrack  $\mathbf{o}_m$ , whose samples are

$$o_m(s) = \begin{cases} a_e(h_m + s), & s_m^0 \leq s \leq s_m^0 + N_S - 1 \\ 0, & \text{otherwise.} \end{cases} \quad (8)$$

This soundtrack is silent, except for  $[s_m^0, \dots, (s_m^0 + N_S - 1)]$ , which includes the denoised content  $\mathbf{e}_{\hat{k}_m, \hat{\delta}_m}^A$ . Finally, as we describe in Section 6 of [1],  $\mathbf{o}_m$  is used to render the final audio output, by *mosaicing*.

## References

- [1] D. Segev, Y.Y. Schechner, M. Elad: ‘‘Example-based Cross-Modal Denoising’’ Proc. IEEE CVPR 2012.