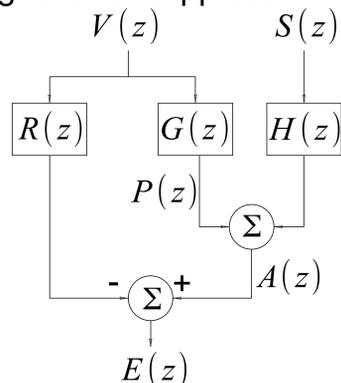


## Abstract

In earlier work, we proposed a voice prompt suppression (VPS) algorithm based on a Kalman filter, in which the temporal update or correction step is performed in information space. The advantage of this approach is that the information matrix can be diagonally loaded in order to control the magnitude of the subband filter coefficients, which provides for better robustness. In this work, we propose a square root implementation of the information filter VPS algorithm, and a technique for diagonally loading the Cholesky factor of the error covariance matrix used in this implementation. We also investigate the effectiveness of cascading VPS after maximum kurtosis beamforming. In a set of distant speech recognition experiments we demonstrate that VPS can reduce word error rate from 19.9% to 16.1% for an adult speaker, and from 44.4% to 40.0% for a child.

## Voice Prompt Suppression

- $V(z)$  denotes the transform of the known voice prompt;
- $S(z)$  denotes the transform of the unknown desired speech;
- $R(z)$  denotes the FIR filter simulating the *room impulse response* (RIR);
- $G(z)$  is the transform of the RIR for the voice prompt  $V(z)$ ;
- $H(z)$  is the transform of the actual, unknown RIR for the speech  $S(z)$ ;
- $A(z)$  is the combined signal at single channel of the microphone array;
- $E(z)$  is the residual signal after suppression of the voice prompt.

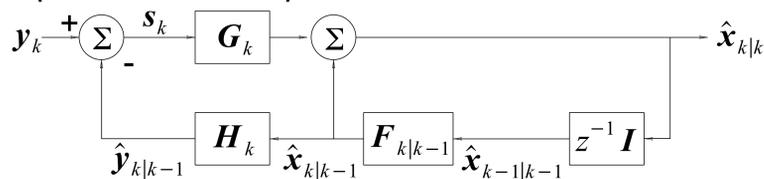


## Kalman Filter

The Kalman filter is governed by a *state* and an *observation equation*

$$\begin{aligned} \mathbf{x}_k &= \mathbf{F}_{k|k-1} \mathbf{x}_{k-1} + \mathbf{u}_k, \\ \mathbf{y}_k &= \mathbf{H}_k \mathbf{x}_k + \mathbf{v}_k, \end{aligned}$$

The state update involves a *prediction* and a *correction*:



The all important *Kalman gain* is calculated through the recursion

$$\begin{aligned} \mathbf{S}_k &= \mathbf{H}_k \mathbf{K}_{k|k-1} \mathbf{H}_k^H + \mathbf{V}_k \\ \mathbf{G}_k &= \mathbf{K}_{k|k-1} \mathbf{H}_k^H \mathbf{S}_k^{-1} \\ \mathbf{K}_{k|k-1} &= \mathbf{F}_{k|k-1} \mathbf{K}_{k-1} \mathbf{F}_{k|k-1}^H + \mathbf{U}_{k-1} \\ \mathbf{K}_k &= (\mathbf{I} - \mathbf{G}_k \mathbf{H}_k) \mathbf{K}_{k|k-1} \end{aligned}$$

## Information Filter

The *Fisher information matrix* and *vector* are defined as

$$\begin{aligned} \mathbf{Z}_k &\equiv \mathbf{K}_k^{-1}, \\ \hat{\mathbf{d}}_{k|k-1} &\equiv \mathbf{Z}_{k|k-1} \hat{\mathbf{x}}_{k|k-1}. \end{aligned}$$

The *temporal update* or *prediction* in information space is given by

$$\begin{aligned} \mathbf{A}_k &= \mathbf{F}_k^{-H} \mathbf{Z}_k \mathbf{F}_k \\ \mathbf{Z}_{k|k-1} &= \left[ \mathbf{I} - \mathbf{A}_k (\mathbf{A}_k + \mathbf{U}_{k-1}^{-1})^{-1} \right] \mathbf{A}_k \\ \hat{\mathbf{d}}_{k|k-1} &= \left[ \mathbf{I} - \mathbf{A}_k (\mathbf{A}_k + \mathbf{U}_{k-1}^{-1})^{-1} \right] \mathbf{F}_k^{-H} \hat{\mathbf{d}}_{k-1|k-1} \end{aligned}$$

The *observational update* or *correction* can be expressed as

$$\begin{aligned} \mathbf{Z}_k &= \mathbf{Z}_{k|k-1} + \mathbf{H}_k^H \mathbf{V}_k^{-1} \mathbf{H}_k, \\ \hat{\mathbf{d}}_{k|k} &= \hat{\mathbf{d}}_{k|k-1} + \mathbf{H}_k^H \mathbf{V}_k^{-1} \mathbf{y}_k. \end{aligned}$$

## Square Root Formulation

The *square root* implementation of the information filter is based on definitions

$$\begin{aligned} \mathbf{Z}_k &= \mathbf{Z}_k^{1/2} \mathbf{Z}_k^{H/2} \\ \mathbf{z}_k &\equiv \mathbf{Z}_k^{H/2} \mathbf{x}_k. \end{aligned}$$

The temporal update involves a unitary *transformation*

$$\begin{bmatrix} \mathbf{U}_{k-1}^{-1/2} & -\mathbf{F}_{k|k-1}^{-H} \mathbf{Z}_{k-1}^{1/2} \\ \mathbf{0} & \mathbf{F}_{k|k-1}^{-H} \mathbf{Z}_{k-1}^{1/2} \\ \mathbf{0} & \hat{\mathbf{z}}_{k-1|k-1}^H \end{bmatrix} \Theta_{pred} = \begin{bmatrix} \mathbf{B}_{11} & \mathbf{0} \\ \mathbf{B}_{21} & \mathbf{Z}_{k|k-1}^{1/2} \\ \mathbf{b}_{31}^H & \hat{\mathbf{z}}_{k|k-1}^H \end{bmatrix}$$

Proving the equivalence of the direct and square root implementations involves applying the matrix factorization lemma then showing that

$$\begin{aligned} \mathbf{B}_{11} &= (\mathbf{A}_k + \mathbf{U}_{k-1}^{-1})^{1/2}, \\ \mathbf{B}_{21} &= -\mathbf{A}_k (\mathbf{A}_k + \mathbf{U}_{k-1}^{-1})^{1/2}, \\ \mathbf{b}_{31} &= -(\mathbf{A}_k + \mathbf{U}_{k-1}^{-1})^{-1/2} \mathbf{F}_{k|k-1}^{-1} \hat{\mathbf{d}}_{k-1|k-1}. \end{aligned}$$

The observational update involves a second unitary transformation

$$\begin{bmatrix} \mathbf{Z}_{k|k-1}^{1/2} & \mathbf{H}_k^H \mathbf{V}_k^{-1/2} \\ \hat{\mathbf{z}}_{k|k-1}^H & \mathbf{y}_k^H \mathbf{V}_k^{-1/2} \end{bmatrix} \Theta_{corr} = \begin{bmatrix} \mathbf{Z}_k^{1/2} & \mathbf{0} \\ \hat{\mathbf{z}}_{k|k}^H & \beta \end{bmatrix}$$

Finally, diagonal loading can be applied to control the magnitude of the weight vector

$$\begin{bmatrix} \mathbf{Z}_k^{1/2} & \sigma_D \mathbf{e}_n \end{bmatrix} \Theta_{diag} = \begin{bmatrix} (\mathbf{Z}'_k)^{1/2} & \mathbf{0} \end{bmatrix}$$

## Experiments and Results

The data collection scenario used for the DSR experiments described here was a simple listen-and-repeat task known as *Copycat*, in which children were shown an illustration of an object and asked to repeat the referring phrase spoken by the experimenter (e.g., "I want the dragon's tail," or "Give her the crown"). To obtain a large number of segments of high overlap between a voice prompt and speech of the subjects, the former was artificially mixed with the latter after capture with far-field microphones. All far-field data capture was conducted with a 64 channel linear microphone array with an intersensor spacing of 2 cm.

Type	$L_{sub}$	Pass					
		1		2		3	
		Exp.	Child	Exp.	Child	Exp.	Child
None		37.7	60.0	19.9	44.4	18.4	41.9
Direct	4	30.5	57.3	15.9	41.2	16.9	41.7
	8	31.2	57.6	17.1	42.3	16.7	43.0
	16	31.1	58.1	17.5	40.0	17.1	43.2
S/R	4	30.3	55.9	16.3	40.9	15.2	44.0
	8	31.0	56.7	16.1	40.0	16.7	42.3
	16	31.4	56.8	17.3	40.9	17.3	43.4

Table 1: Word error rates (WERs) for several subband filter lengths using both the direct form and square root (S/R) implementation of VPS after MK beamforming; also shown are comparable results with no VPS.

## Conclusions

We have proposed a voice prompt suppression algorithm based on a square root implementation of the information filter. This formulation enables diagonal loading to be applied information matrix to control the magnitude weight vector. Much like in beamforming, diagonal loading provides for superior robustness. Further work is need to directly compare the proposed algorithm to convention techniques based on normalized LMS algorithms.