

CHANNEL SELECTION BASED ON MULTICHANNEL CROSS-CORRELATION COEFFICIENTS FOR DISTANT SPEECH RECOGNITION

Kenichi Kumatani¹, John McDonough², Jill Fain Lehman^{1,2}, and Bhiksha Raj²

¹Disney Research, Pittsburgh
Pittsburgh, PA 15213, USA

²Carnegie Mellon University
5000 Forbes Avenue, Pittsburgh, PA 15213, USA

ABSTRACT

In theory, beamforming performance can be improved by using as many microphones as possible, but in practice it has been shown that using all possible channels does not always improve speech recognition performance [1, 2, 3, 4, 5]. In this work, we present a new channel selection method in order to increase the computational efficiency of beamforming for distant speech recognition (DSR) without sacrificing performance.

To achieve better performance, we treat a channel that is uncorrelated with the others as unreliable and choose a subset of microphones whose signals are most highly correlated with each other. We use the *multichannel cross-correlation coefficient* (MCCC) [6] as a measure for selecting the reliable channels. The selected channels are then used for beamforming.

We evaluate our channel selection technique with DSR experiments on real children’s speech data captured using a linear array with 64 microphones. A single distant microphone provided a word error rate (WER) of 15.4%, which was reduced to 8.5% by super-directive beamforming with all the sensors. The experimental results suggest that almost the same recognition performance can be obtained with half the number of sensors in the case of super-directive beamforming. Maximum kurtosis beamforming [7] with 48 sensors out of a total of 64 achieved a WER of 5.7%, which is very comparable to the 5.2% WER obtained with a close-talking microphone.

Index Terms— channel selection, microphone arrays, beamforming, speech recognition

1. INTRODUCTION

There has been a great and growing interest in distant speech recognition (DSR) [8] within the research community, as this technology offers the possibility of relieving users from the necessity of donning close talking microphones (CTMs) before interacting with automatic speech recognition (ASR) systems. Moreover, DSR may be especially useful for young children who may find CTMs too cumbersome and intrusive to use in interactive attractions.

First of all, the authors would like to thank Prof. Jessica Hodgins for giving us the opportunity to study this work. The authors would also like to thank Cedrick Rochet for his support in developing the Mark IV microphone array. Also due thanks are Wei Chu, Spencer Diaz, Jerry Feng, Ishita Kapur, and Moshe Mahler for their assistance in collecting the audio-visual material used for the experiments described in this work.

The presence of noise and reverberation effects in real environments severely degrades the performance of DSR systems. Depending on the distance between each microphone and the noise source, some channels will have lower signal-to-noise ratios (SNR) than others, especially when a large microphone array is used. The reverberation effects also differ among the sensors. Therefore, the performance of speech enhancement might not always be improved by using as many microphones as possible in a real environment. Moreover, it is generally assumed in microphone array processing that all the microphones have the same gain and phase characteristics. This assumption may not hold due to variations in system response introduced by the microphone and analog-to-digital converter (ADC) [9, 10].

Various methods have been proposed for selecting a suitable channel or using a cluster of microphones. These methods can be categorized into the following approaches:

- selecting a channel with a high SNR [2];
- choosing a channel to which a speech recognizer assigns the maximum likelihood [11];
- measuring how much the system’s outputs are changed by a noise adaptation technique based on the comparison of word hypotheses of uncompensated and compensated features, and choosing the one with the smallest change [1];
- calculating the class separability measure of feature vectors and selecting the channel which maximizes the separation measure [3]; and
- clustering microphones based on the distance between two microphones and choosing the cluster of microphones according to the proximity measure to a speaker that considers the distance between the reference microphone and speaker as well as the size of the cluster [4, 5].

The SNR-based method is simple and can be calculated efficiently, but requires voice activity detection which often fails in noisy environments. Moreover, the SNR measure does not consider any information about ASR.

In terms of ASR, it might be straightforward to use outputs from the speech recognizer for channel selection. As Wölfel noted in [3], however, the disadvantage of this approach is that at least one decoding process is required for each channel in order to avoid mismatch between different channels. Such additional calculation leads to a drastic increase in computational complexity.

In contrast to the SNR measure, the class separability criterion can take into account speech features for ASR and requires less

computation than the decoder-based methods. Wölfel demonstrated in [3] that the channel selection method based on the class separability criterion provided better recognition performance than the SNR-based approach. However, Wölfel selected a single channel and thus did not consider using beamforming, which can drastically decrease word error rate (WER). Moreover, the computation required by his method is still significant in the case of multi-channel processing. In contrast, we propose here a technique which selects a subset of all channels for microphone array processing.

Himawan *et al.* [5], addressed the situation where microphones are placed on an ad hoc basis. Accordingly, clustering of microphones must be done without any knowledge of microphone positions. In contrast, we consider the situation where the microphones are regularly spaced and whose positions are known a priori. This assumption simplifies the problem significantly.

In essence, we consider the *multichannel cross-correlation coefficient* (MCCC) [6] as a measure for selecting the reliable channels. The MCCC represents correlation among more than two channels and the cross-correlation coefficient can be viewed as the special case where the MCCC is calculated with two channels. Although Benesty *et al.* [6] originally proposed the MCCC for the speaker localization problem, we use the maximum MCCC criterion for channel selection.

The basic idea behind the algorithm is that signals of unreliable channels are uncorrelated with most others. For the sake of computational efficiency, we first compensate for the delays of the signals based on the *phase transform* (PHAT) [8, §10.1]. After the multi-channel signal is aligned, we compute the MCCC and then choose a set of channels with the maximum MCCC. Finally, beamforming and post-filtering are performed on the selected channels. We demonstrate the effectiveness of our channel selection technique through a series of DSR experiments on real data captured with real microphones. In these experiments we used both traditional *super-directive beamforming* [8, §13.3.4] and state-of-the-art *maximum kurtosis* beamforming; the latter adapts the subband filter coefficients on each channel so as to maximize the kurtosis of the beamformer's output subject to a distortionless constraint in the look direction [7].

We also investigated other microphone array design methods [12, 13] in order to reduce the number of microphones for beamforming. Logarithmically spaced and non-redundant linear array design methods were evaluated in terms of recognition performance.

The balance of this paper is organized as follows. Section 2 describes the formulation of the problem for microphone array processing and defines the notation used in this work. Section 4 reviews the MCCC. Section 5 presents our channel selection method based on maximizing MCCC. Recognition experiments are described in Section 6. Our conclusions about this work and future plans are summarized in Section 7.

2. PROBLEM FORMULATION

Consider the anechoic situation shown in Figure 1 where a single source signal is captured with a microphone array.

In the time domain, a vector of the M -channel signal captured with M microphones at discrete time n can be denoted as

$$\mathbf{x}_M[n] = [x_1[n] \quad x_2[n] \quad \cdots \quad x_M[n]]^T. \quad (1)$$

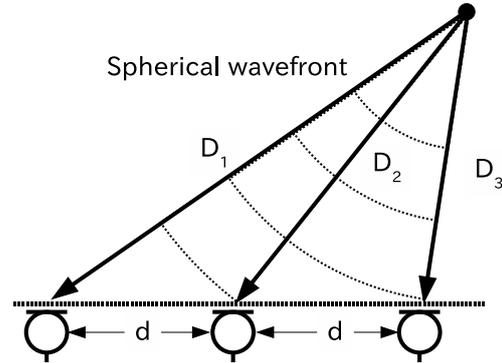


Fig. 1. Illustration of the single source signal model under the near-field assumption.

In this case, the observation vector of source signal $s[n]$ can be expressed as

$$\mathbf{x}_M[n] = \begin{bmatrix} a_1 s[n - T_p - \tau_{1r}] \\ \vdots \\ a_m s[n - T_p - \tau_{mr}] \\ \vdots \\ a_M s[n - T_p - \tau_{Mr}] \end{bmatrix} + \begin{bmatrix} v_1[n] \\ \vdots \\ v_m[n] \\ \vdots \\ v_M[n] \end{bmatrix} \quad (2)$$

where a_m denotes the attenuation factor from the source to microphone m , T_p denotes the propagation time to the reference microphone r , τ_{mr} denotes the *time delay of arrival* (TDOA) between two microphones m and r , and $v_m[n]$ is an additive noise signal.

We denote the signal model of (2) in the subband or frequency domain as

$$\mathbf{X}_M(e^{j\omega n}) = \begin{bmatrix} a_1 S_1 e^{j\omega(n - T_p - \tau_{1r})} \\ \vdots \\ a_M S_M e^{j\omega(n - T_p - \tau_{Mr})} \end{bmatrix} + \begin{bmatrix} V_1(e^{j\omega n}) \\ \vdots \\ V_M(e^{j\omega n}) \end{bmatrix}. \quad (3)$$

In our channel selection algorithm, the TDOA τ_{mr} is first estimated in order to align the signals and calculate the correlation measure among the multiple microphones more accurately. This is not a straightforward task in real acoustic environments, as each microphone captures multiple attenuated and delayed replicas of the source signal due to reflections from, for example, tables and walls.

3. TIME DELAY ESTIMATION

In this work, we use the phase transform (PHAT) for time delay estimation. It is a variant of generalized cross-correlation (GCC) and, is perhaps, the most widely used method due to its computational efficiency and robustness in the presence of noise and reverberation [14, 8]. The PHAT between two microphones m and n can be expressed as

$$\rho_{mn}(\tau) = \int_{-\pi}^{\pi} \frac{X_m(e^{j\omega\tau}) X_n^*(e^{j\omega\tau})}{|X_m(e^{j\omega\tau}) X_n^*(e^{j\omega\tau})|} e^{j\omega\tau} d\omega, \quad (4)$$

where $X_m(e^{j\omega\tau})$ denotes the spectrum of the signal captured with by the m -th sensor. We use a Hamming window for analysis in order to calculate these short-time. The normalization term in the denominator of (4) is intended to weight all frequencies equally; it has been shown that such a weighting conduces to more robust time delay estimation [14]. The TDOA between the m th and n th channels is then estimated from

$$\hat{\tau}_{mn} = \max_{\tau} \rho_{mn}(\tau). \quad (5)$$

Thereafter, an interpolation is performed to overcome the granularity in the estimate corresponding to the sampling interval.

4. MULTICHANNEL CROSS-CORRELATION COEFFICIENT

Once the time delays of the M signals are estimated based on the PHAT, time-aligned signal can be obtained according to

$$\mathbf{x}_{d,M}[n] = [x_1[n + \hat{\tau}_{1r}], x_2[n + \hat{\tau}_{2r}], \dots, x_M[n + \hat{\tau}_{Mr}]]^T. \quad (6)$$

In order to calculate the MCCC, we first need a spatial correlation (covariance) matrix of the observations. The spatial correlation matrix can be expressed as

$$\mathbf{R}_M = E \left\{ \mathbf{x}_{d,M}[n] \mathbf{x}_{d,M}^T[n] \right\}. \quad (7)$$

Then, given the TDOA estimates, the MCCC can be computed as

$$\varrho_M^2 = 1 - \frac{\det[\mathbf{R}_M]}{\prod_{i=1}^M \sigma_i^2}, \quad (8)$$

where $\det[\cdot]$ denotes the determinant and σ_i^2 is the i th diagonal component of the spatial correlation matrix \mathbf{R}_M . It can be readily confirmed that the MCCC is equivalent to the cross-correlation coefficient normalized by the energy in the case of $M = 2$ [6].

Chen, Benesty and Huang originally used the MCCC for estimating the direction of arrival (DOA) based on the far-field assumption [15, 16]. In their work, the MCCC was viewed as a function of the time delays. In contrast to their work, we estimate the TDOA based on the PHAT which leads to a drastic computational reduction in the case of the near field assumption and calculate the MCCC with fixed time delays for channel selection.

In the context of source localization, Chen *et al.* [15, 16], showed that

$$0 \leq \frac{\det[\mathbf{R}_M]}{\prod_{i=1}^M \sigma_i^2} \leq 1, \quad (9)$$

and noted that the MCCC has the following properties:

- $0 \leq \varrho_M^2 \leq 1$;
- $\varrho_M^2 = 1$ if two or more signals are perfectly correlated;
- $\varrho_M^2 = 0$ if all the signals are completely uncorrelated with one another; and
- if one of the signals is completely uncorrelated with the $M - 1$ other signals, the MCCC of all the signals will be equal to that of those $M - 1$ remaining signals.

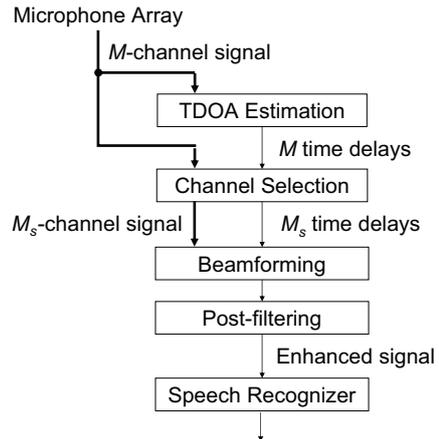


Fig. 2. A flow chart of our distant speech recognition system.

5. CHANNEL SELECTION

Here we describe our channel selection method. Let us assume that we select M_s channels with the maximum MCCC out of M microphones. We ideally want to find a set of channels \mathcal{C}_{M_s} which provides the largest MCCC among all the possible combinations as follows:

$$\hat{\mathcal{C}}_{M_s} = \operatorname{argmax}_{\mathcal{C}_{M_s}} \varrho_{M_s}^2. \quad (10)$$

An exhaustive search requires computing the MCCC $M C_{M_s}$ times. If we have a large number of microphones, this computation is intractable.

We avoid this problem by iteratively reducing the number of the search candidates from M to M_s . More specifically, we ignore the channel that provides the smallest MCCC and keep the remaining channels for the next step. This process is repeated until we obtain the desired number of channels, M_s . By doing so, the computation for the MCCC is reduced from $M C_{M_s}$ to $\sum_{i=0}^{M-M_s} M - i$.

Our channel selection algorithm is summarized as follows:

1. Estimate the time delays of the M -channel signal with (5) and align the signals.
2. Push all the M channels onto a search stack.
3. Denoting the number of the candidates in the search stack as M_c , find a set of the $M_c - 1$ channels with the largest MCCC.
4. Remove the channel which provides the smallest MCCC in Step 2 from the stack.
5. Go to Step 3 if $M_c > M_s$.

Clearly at least two channel must be retained so that the correlation can be evaluated.

6. EXPERIMENTS

Figure 2 shows a block diagram of the *distant speech recognition* (DSR) system used to generate the experimental results reported here. Our DSR system involves the time delay estimation step described in Section 3, the channel selection method depicted in Section 5, beamforming, post-filtering and automatic speech recognition (ASR) components which we will now describe.

In our experiments, beamforming is performed on the channels selected by the algorithm proposed above. We consider both the widely used super-directive beamforming [8, §13.3.4] and one of the state-of-art techniques, maximum kurtosis beamforming [7]. As the experimental results presented in Section 6.1 show, the computation required for beamforming can be significantly decreased by reducing the number of channels without degrading recognition performance. Following beamforming, Zelinski post-filtering [17], a variant of Wiener filtering, is carried out in order to remove the uncorrelated noise among the sensors.

Our basic DSR system was trained on three corpora of children’s speech:

1. the CMU Kids’ Corpus, which contains 9.1 hours of speech from 76 speakers;
2. the Center for Speech and Language Understanding (CSLU) Kids’ Corpus, which contains 4.9 hours of speech from 174 speakers.
3. A set of *Copycat* data collected at the Carnegie Mellon Children’s School in June, 2010.

The feature extraction used for the ASR experiments reported here was based on cepstral features estimated with a warped *minimum variance distortionless response* (MVDR) spectral envelope of model order 30 [8, §5.3]. Front-end analysis involved extracting 20 cepstral coefficients per frame of speech, and then performing *cepstral mean normalization* (CMN). The final features were obtained by concatenating 15 consecutive frames of cepstral coefficients together, then performing *linear discriminant analysis* (LDA), to obtain a feature of length 42. The LDA transformation was followed by a second CMN step, then a global semi-tied covariance transform estimated with a maximum likelihood criterion [18].

HMM training was conducted initializing a context independent model with three states per phone with the global mean and variance of the training data. Thereafter, five iterations of Viterbi training [8, §8.1.5] were conducted. This was followed by an additional five iterations whereby optional silences and optional breath phones were allowed between words. The next step was to treat all triphones in the training set as distinct and train three-state single-Gaussian models for each. Then state clustering was conducted as in [19]. In the final stage of conventional training, the context-dependent state-clustered model was initialized with a single Gaussian per codebook from the context-independent model; three iterations of Viterbi training followed by splitting the Gaussian with the model training steps. These steps were repeated until no more Gaussians had sufficient training counts to allow for splitting. The conventional model had 1,200 states and a total of 25,702 Gaussian components. Conventional training was followed by *speaker-adapted training* (SAT) as described in [8, §8.1.3].

In our experiments, the ASR system consisted of three passes:

1. Recognize with the unadapted conventionally trained model;
2. Estimate *vocal tract length normalization* (VTLN) [20], *maximum likelihood linear regression* (MLLR) [21] and *constrained maximum likelihood linear regression* (CM-LLR) [22] parameters, then recognize once more with the adapted conventionally trained model;
3. Estimate VTLN, MLLR and CMLLR parameters for the SAT model, then recognize with same.

For all but the first unadapted pass, unsupervised speaker adaptation was performed based on word lattices from the previous pass.

Algorithm	Pass (%WER)		
	1	2	3
Single distant microphone	38.1	19.8	15.4
SD beamforming with CS	24.1	11.3	8.6
MK beamforming with CS	21.8	8.3	5.7
SD beamforming without CS	32.5	11.4	8.5
MK beamforming without CS	33.6	11.3	7.1
Lapel microphone	19.3	5.9	5.2

Table 1. Word error rates (WERs) for each decoding pass.

6.1. Recognition results

Test data for experiments were collected at the Carnegie Mellon University Children’s School over weeks. The database consists of 4 sessions which were recorded on different dates. The speech material in this corpus was captured with a 64-channel Mark IV microphone array; the elements of the Mark IV were arranged linearly with a 2 cm intersensor spacing. In order to provide a reference for the DSR experiments, the subjects of the study were also equipped with Shure lavelier microphones with a wireless connection to an RME Hammerfall Octamic II preamp and ADC. The Octamic II was connected via an ADAT optical cable to a RME Hammerfall HDSPe AIO sound card. A PNC coaxial connection between the Mark IV and the Octamic II ensured that *all* audio capture was sample synchronous. This was required to enable voice prompt suppression experiments. All the audio data were captured at 41.1 kHz with a 24-bit per sample resolution.

The test set consists of 354 utterances (1,297 words) spoken by nine children. The children were native-English speakers (aged four to six). They were asked to play *Copycat*, a listen-and-repeat paradigm in which an adult experimenter speaks a phrase and the child tries to copy both pronunciation and intonation. As is typical for children in this age group, pronunciation was quite variable and the words themselves sometimes indistinct.

The search graph for the recognition experiments was created by initially constructing a finite-state automaton by stringing *Copycat* utterances in parallel between a start and end state. This acceptor was convolved together with a finite-state transducer representing the phonetic transcriptions of the 147 words in the *Copycat* vocabulary. Thereafter this transducer was convolved with the *HC* transducer representing the context-dependency decision tree estimated during state-clustering [8, §7.3.4].

The channel selection algorithm is performed with 460 milliseconds of speech data from the beginning of each session. After that, we perform beamforming on the same channel set consistently. In this data set, we do not need to select the channels in the online manner since a speaker does not move significantly in each session.

Table 1 shows word error rates (WERs) of every decoding pass obtained with one of 64 microphones, super-directive (SD) beamforming and maximum kurtosis (MK) beamforming with channel selection (CS) and without it. In the experiments with channel selection, the numbers of channels for SD and MK beamforming are 32 and 48 respectively because those settings provided the best results. As a reference, the WERs of the lapel microphone are also depicted in Table 1.

Table 1 demonstrates that the improvement from the adaptation techniques is dramatic. The reduction in the WER from the first pass to the third is approximately four-fold in the case of MK beamforming. It is also clear that the performance of far-field speech recog-

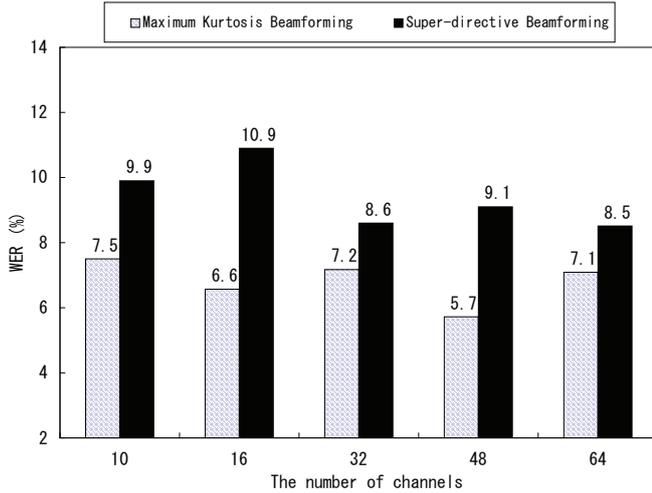


Fig. 3. WERs of the third pass as a function of the number of channels.

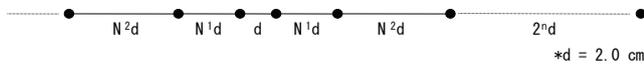


Fig. 4. Logarithmically spaced linear microphone array.

tion can be improved by beamforming techniques, and the MK beamforming algorithm achieves the best performance in the experiments. The MK beamforming technique provides almost the same recognition performance as the lapel microphone.

We also investigated the WERs as a function of the number of channels used for beamforming. Figure 3 shows the WERs of the third pass for the number of channels when SD and MK beamforming algorithms were applied. The MK beamformer provides better recognition performance than SD beamforming when the same number of channels is used. Using all the microphones does not provide the best recognition performance because several channels are distorted by reverberation and noise. The results in Figure 3 suggest that we could improve recognition performance by automatically finding the optimum number of channels although the effect would be relatively small. In practice, the number of channels could be empirically decided based on the computer resources available for the application.

Another interesting result comes from a comparison of our channel selection algorithm and the microphone array design methods [12, 13]. In our case, due to the fixed geometry of the Mark IV, our adoption is to select among the channels with the 2 cm inter-sensor spacing. In other words, the microphone array design method can be viewed as a channel selection method.

First, we compare our channel method with the logarithmically spaced linear array shown in Figure 4. In the logarithmically spaced linear array, the sensor is symmetrically placed in the center of the linear array on a logarithmic scale. Figure 5 shows the WERs obtained by selecting the channels in order to form a logarithmically spaced array. In Figure 5, SD beamforming is performed for the sake of efficiency. Due to the physical restrictions imposed by the Mark IV, it was not possible to change the channel spacing. Hence, the microphones were chosen so as to conform to a logarithmic design as closely as possible. It is clear from Figure 5 that our channel selection method provides lower WERs than the logarithmically

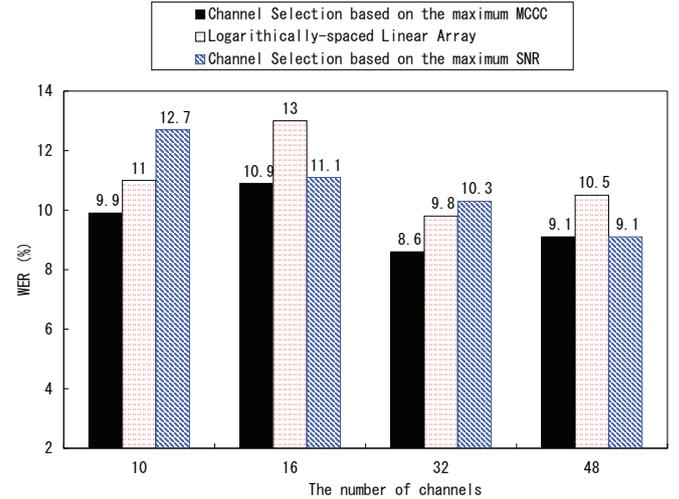


Fig. 5. WERs of the third pass as a function of the number of channels in the case of super-directive beamforming.

Method	WER of the third pass
Logarithmically spaced array	11.0
Non-redundant array	9.6
Channel selection method	9.9

Table 2. WERs for each array design method from the super-directive beamformer with 10 sensors.

spaced linear array. This improvement occurs because our channel selection method can adaptively choose the channels based on signal characteristics as opposed to the static logarithmic design.

Figure 5 also shows the WERs obtained by a channel selection algorithm based on the maximum SNR criterion as contrast condition. The maximum SNR-based algorithm used here first measures the SNR of each channel from the noise and speech segments aligned by the speech recognizer and then selects the channels with the best SNRs. Figure 5 illustrates that the maximum SNR-based algorithm performs worse than the method based on the maximum MCCC criterion. The increases in the WERs occur mainly because it is not feasible to precisely measure the SNR in noisy acoustic environments due to the absence of perfect speech activity detection. The results might also suggest that the SNR is not related to the WER.

Finally, we tabulated the WERs of our channel selection method and two array design methods in Table 2 in the case of SD beamforming with 10 sensors. Again, because of the uniform spacing of the MarkIV, we cannot compare our channel selection method with the non-redundant linear array design [13, §3.9] in the case of more than 10 sensors. We can, however, observe from Table 2 that the non-redundant array and our channel selection method provide almost the same recognition performance in the experiment with 10 microphones. This result is promising because these techniques could be combined if we had the freedom to choose the actual geometry of the array. For instance, we could select the channel of the non-redundant microphone array based on the maximum MCCC criterion.

7. CONCLUSIONS

In this work, we have proposed a new channel selection algorithm for distant speech recognition (DSR) based on acoustic beamforming. We have demonstrated through a series of DSR experiments that our algorithm can reduce the number of channels for beamforming effectively. Our channel selection method can also improve recognition performance.

In future, we plan to combine our channel selection and array design methods as well as other conventional channel selection methods. We also plan to extend the algorithm proposed here to the situation where multiple sources are active. We also plan to investigate the eigenvalues of the spatial covariance matrix and develop an automatic method to determine the optimum number of channels.

8. REFERENCES

- [1] Yasunari Obuchi, "Multiple-microphone robust speech recognition using decoder-based channel selection," in *Proc. ISCA Tutorial and Research Workshop on Statistical and Perceptual Audio Processing (SAPA)*, Jeju, Korea, 2004.
- [2] Matthias Wölfel, Christian Fügen, Shajith Ikbali, and John W. McDonough, "Multi-source far-distance microphone selection and combination for automatic transcription of lectures," in *Proc. Interspeech*, Pittsburgh, Pennsylvania, 2006.
- [3] Matthias Wölfel, "Channel selection by class separability measures for automatic transcriptions on distant microphones," in *Proc. Interspeech*, Antwerp, Belgium, 2007.
- [4] Ivan Himawan, Iain McCowan, and Sridha Sridharan, "Clustering of ad-hoc microphone arrays for robust blind beamforming," in *ICASSP*, Dallas, Texas, 2010.
- [5] Ivan Himawan, Iain McCowan, and Sridha Sridharan, "Clustered blind beamforming from ad-hoc microphone arrays," *IEEE Trans. Speech Audio Processing*, vol. 18, pp. –, 2010.
- [6] Jacob Benesty, Jingdong Chen, and Yiteng Huang, *Microphone Array Signal Processing*, Springer, 2008.
- [7] Kenichi Kumatani, John McDonough, Barbara Rauch, Philip N. Garner, Weifeng Li, and John Dines, "Maximum kurtosis beamforming with the generalized sidelobe canceller," in *Proc. Interspeech*, Brisbane, Australia, September 2008.
- [8] Matthias Wölfel and John McDonough, *Distant Speech Recognition*, Wiley, New York, 2009.
- [9] Ivan Jelev Tashev, *Sound Capture and Processing: Practical Approaches*, Wiley, 2009.
- [10] Carsten Sydow, "Broadband beamforming for a microphone array," *Journal of the Acoustical Society of America*, vol. 96, pp. 845–849, 1994.
- [11] Y. Shimizu, S. Kajita, K. Takeda, and F. Itakura, "Speech recognition based on space diversity using distributed multi-microphone," in *ICASSP*, Istanbul, Turkey, 2000.
- [12] Saeed Gazor and Yves Grenier, "Criteria for positioning of sensors for a microphone array," *IEEE Trans. Speech Audio Processing*, vol. 3, pp. 294–303, 1995.
- [13] H. L. Van Trees, *Optimum Array Processing*, Wiley-Interscience, New York, 2002.
- [14] M. Brandstein and D. Ward, Eds., *Microphone Arrays*, Springer Verlag, Heidelberg, Germany, 2001.
- [15] Jingdong Chen, Jacob Benesty, and Yiteng Huang, "Robust time delay estimation exploiting redundancy among multiple microphones," *IEEE Trans. Speech Audio Processing*, vol. 11, pp. 549–557, 2003.
- [16] Jacob Benesty, Jingdong Chen, and Yiteng Huang, "Time delay estimation via linear interpolation and cross-correlation," *IEEE Trans. Speech Audio Processing*, vol. 12, pp. 509–519, 2004.
- [17] Claude Marro, Yannick Mahieux, and K. Uwe Simmer, "Analysis of noise reduction and dereverberation techniques based on microphone arrays with postfiltering," *IEEE Transactions on Speech and Audio Processing*, vol. 6, pp. 240–259, 1998.
- [18] M. J. F. Gales, "Semi-tied covariance matrices for hidden Markov models," *IEEE Transactions Speech and Audio Processing*, vol. 7, pp. 272–281, 1999.
- [19] S. J. Young, J. J. Odell, and P. C. Woodland, "Tree-based state tying for high accuracy acoustic modelling," in *Proc. of HLT*, Plainsboro, NJ, USA, 1994, pp. 307–312.
- [20] Ellen Eide and Herbert Gish, "A parametric approach to vocal tract length normalization," in *Proc. of ICASSP*, 1996, vol. I, pp. 346–8.
- [21] C.J. Leggetter and P.C. Woodland, "Maximum likelihood linear regression for speaker adaptation of continuous density hidden markov models," *Jour. on CSL*, vol. 9, no. 2, pp. 171–185, 1995.
- [22] M. J. F. Gales, "Maximum likelihood linear transformations for hmm-based speech recognition," *Jour. on CSL*, vol. 12, pp. 75–98, 1998.