

Robust Speech Recognition and its ROBOT implementation

Yoshikazu Miyanaga

Hokkaido University Laboratory of Information Communication Networks Graduate School of Information Science and Technology Sapporo 060-0814, Hokkaido Japan



Part 1

NOISE ROBUST TECHNOLOGY



Conditions for Speech Recognition



Robot Implementation

- Autonomous Speech Recognition
- Speech Synthesis
- Quick Response
- Control to Consumer Electronics and Machines





Hokkaido University Speech Communication System (HU-SCS)









Spectrum including noise can be modeled as,

$$X(n,\omega) = S(n,\omega)H(\omega) + A(\omega)$$

Clean	Multiplicative	Additive nois
spectrum	noise	

e

$$\log \int E(n,\omega) = S(n,\omega)H(\omega)$$

 $\log X(n,\omega) = \log(E(n,\omega) + A(\omega))$

All right reserved. Copyright ©2015- Yoshikazu Miyanaga

 \vee



 Noise corruptions make differences on gains and DC components.





Running Spectrum

Running spectrum is obtained by accumulating short-time spectrum





Modulation Spectrum

RSF focuses on modulation spectrum



Modulation spectrum: spectrum versus time trajectory of frequency.





Mod-F of Clean and Noisy Speech

 Speech components are dominant around 4 Hz in modulation spectrum.



Lower modulation frequency components can be assumed as noise because of little changes in noise components.

All right reserved. Copyright ©2015- Yoshikazu Miyanaga



RSF (Running Spectrum Filtering)

Speech components are dominant around 4 Hz in modulation spectrum.



All right reserved. Copyright ©2015- Yoshikazu Miyanaga



Green shows the 4th MFCC before DRA and Blue shows after DRA



Running Spectrum Filtering



The 4th MFCC by only DRA in the left hand side, and by both DRA and RSF in the right hand side

Evaluation on Likelihoods



RSF/DRA reduce noise factors.

MFCC Likelihood of MFCC into HMMs

HMM



The maximum likelihood

is selected and its label is recognized as the result.



Phase

Candidates of Recognition Results

- (1) Good Morning
- (2) See you
- (3) How are you ?

Automatic Speech Rejection



Recognition Result Good Morning

Confidential Phase

Evaluation on Likelihoods





HMM

Training HMM



The maximum likelihood is selected and its label is recognized as the result.



The result is correct, isn't it ?

All right reserved. Copyright ©2015- Yoshikazu Miyanaga



The result of the top score is trusted.

The result of the top score is **NOT** trusted.



Rejection Method







All right reserved. Copyright ©2015- Yoshikazu Miyanaga

First SCS HW •LSI IP Mobile Fine Advantage Intelligent Consumer Electronics etc (1) Mobile Appli HW with Small Low Power Low Power (2) PC free

- Super Low-Power Consumption Design
- Real-Time SCS
 - →180nsec/word (10MHz clk) Recognition Time
- Small Scale Design with Special Designed LSI
- Noise Reduction by Array Microphone



Current HU-SCS v4

PC Interface with HU-SCS Board





HU-SCS Board

55mm × 44 mm

Comparison on Performance

Environment	Noise Level	Correctness	
Environment		Current	Previous
Meeting Room	50dB	96.4%	90.0%
Elevator	50dB	95.0%	84.4%
Stairs	45dB	85.1%	50.5%
Car A (Idling, No-Moving)	50dB	99.4%	95.6%
Car B (High Speed, Open Window)	75dB	93.3%	85.0%
Car C(High Speed, Audio ON(FM))	75dB	88.9%	65.6%
Total		93.0%	78.5%
%Cruiser Board (Outside, high speed)	80dB	82.7%	-

Comparisons between HU-SCS v4 and v3





Part 2

ADVANCED NOISE ROBUST TECHNOLOGY

Sequential Word Stream

- Only keywords are recognized in a continuous speech under noisy environments.
 - Any words except keywords should not be recognized.
- The speech is given as a sequential word stream like voice commands.





Issues in This System

- Intelligent Keyword Rejection
 Many of non-Keywords are recorded.
- High Precise Speech Detection
 - Keywords and non-Keywords are sequentially recorded into a system like continuous speech waveform.

SP-ASR Recognition Processing

Set a segment after dividing the short-time frame



All right reserved. Copyright ©2015- Yoshikazu Miyanaga



2nd Rejection Methods



Note that the vector P_i might be similar to the vector P_j

All right reserved. Copyright ©2015- Yoshikazu Miyanaga



Rejection Method



k: Center segment number

1.
$$p_{k,1} \in P_r \ (r = k - 2, \dots, k + 2)$$

2. $L(p_{r,i}) > L(p_{r,j})$

{
$$i \mid p_{k,1} = p_{r,i}$$
}, { $j \mid p_{k,2} = p_{r,j}$ }
 $r = k - 2, \dots, k + 2$

 $p_{m,i} = HMM_p(m)$: Word /p/ at the *m*-th segment. The index *i* shows the *i*-th largest likelihood.

 $L(p_{m,i})$: Likelihood Value It is the i-th largest ML in the a-th segment. m: segment index, i: ranking index

Experiment - Conditions

- Short-time frame
 - Width: 23.2ms(256point) Shift: 11.6ms(128point)
- Segment
 - Length(L): 1.0s(11025point)Shift: 8[frame]
- Speech data:
 - Connect two words.
 - Add a non-speech section of 1.4 sec to both ends and 0.5 sec between the words.
 - 100 data by random two words combination.
- White Noise : SNR=∞,20,10dB



Results of recognition

Recognizer/SNR	00	20d B	10d B
Isolated Word	99.4	97.6	84.6
Recognition (1 word)	%	%	%
SP-ASR with 2 nd Rejection (2 sequential words)	98.5	96.5	88.5
	%	%	%

The SP-ASR recognizes keywords with the same level accuracy as the past word recognition.





Small, Fast and Low Power

Hokkaido University Speech Communication System Integrated Architecture of Speech Detection, Robust Speech Analysis, Speech Recognition, Speech Rejection

Analysis, Speech Recognition, Speech Rejection Higher Speed Processing than DSP and Software Superior in Energy Saving than DSP Solutions Improving Noise Robustness by RSF/DRA Technique