

Progress in audible (normal) and inaudible speech recognition based on NAM microphones

Panikos Heracleous

Speech and Acoustics Processing Lab.

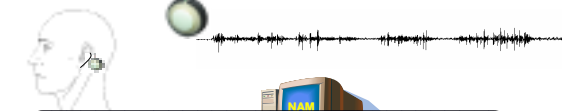
2004/12/22

NAIST

Introduction

NAM microphone is a special device (sensor) which is attached behind the talker's ear and perceives very quietly uttered speech

NAM microphone [Nakajima et al., NAIST, 2003-2004]



NAM plays an important role in Ubiquitous Networking

It extends the speech media applications to cover:

- Applications with privacy (e.g. speech recognition, telephony)
- Applications under several (e.g. noisy) environments
- Applications for specific people groups (e.g. speech-impaired)

2004/12/22

NAM microphones

1. Stethoscopic NAM microphone [2003]

- Based on stethoscope used in medical science
- Able to capture audible and inaudible speech
- Only low frequency band
- Low speech intelligibility



Male

2. Silicon NAM microphone [2004]

- Silicon is used to wrap a sensitive microphone
- Able to capture audible and inaudible speech
- Wide frequency band
- Higher speech intelligibility



Male

Female

3. First pre-commercial prototype NAM [2004]

- Based on the Silicon NAM microphone
- Still under development and investigation



2004/12/22

NAIST

Some advantages of NAM

■ Able to capture audible and inaudible speech

- Privacy in communication (human machine, human-human)
- Direct attachment to the skin – Expected noise robustness

Music is heard

■ No additional techniques are required

- Only a microphone pre-amplifier is required
- Directly connected to PC for recording, or speech recognition

■ Size, weight



Silicon NAM microphone



A throat microphone

2004/12/22

NAIST

Some similar approaches

Bone-microphone



Fig. 1 The Ear and Bone-Conductive Integrated Microphone (BICM)

Microsoft

PARAT



Norwegian University of Science and Technology

Throat microphones



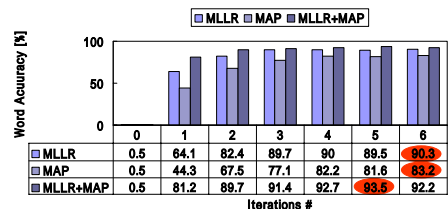
- Special applications
- Military, Police, mobile telephone
- Usually normal speech recognition
- Integration with other sensors
- Big size and uncomfortable usage
- Results show high noise robustness

2004/12/22

NAIST

Stethoscopic NAM

Speaker-dependent experiments (Clean)



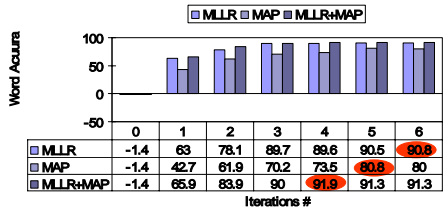
700 training utterances, 24 test utterances
Julius recognizer, 20k vocabulary Japanese dictation task
High performance! Comparable to normal-speech recognition

2004/12/22

NAIST

Stethoscopic NAM

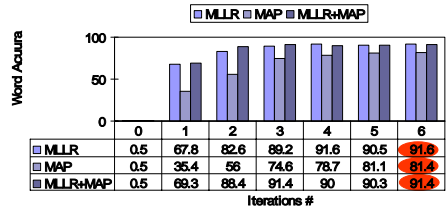
Speaker-dependent experiments (Noise) - I



700 training utterances, 24 test utterances
Room noise is superimposed to clean data
Noise level: About 50dBa - HMMs noisy

Stethoscopic NAM

Speaker-dependent experiments (Noise) - II



700 training utterances, 24 test utterances
Room noise is superimposed to clean data
Noise level: About 50dBa - HMMs clean

Stethoscopic NAM

Speaker-independent experiments

Speaker	Training method			
	MLLR	MAP	MLLR+MAP	EM
Male	73.5	62.4	77.8	75.9
Female	70.0	58.4	71.1	65.4

Reasonable results

Recording conditions

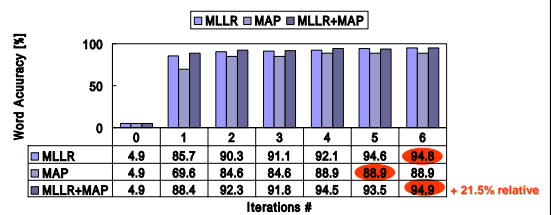
- 16-bit SONY ZA5ES SUPER BIT MAPPING DAT tape digital recorder
- Anechoic room 30dBa, Sampling frequency: 44100kHz -> 16kHz

Recording data

- 25 speakers (14 male, 11 female), 100 newspaper, 50 phoneme balanced utterances
- 3189 training utterances, 48 training utterances

Silicon NAM

Speaker-dependent experiments (Clean)

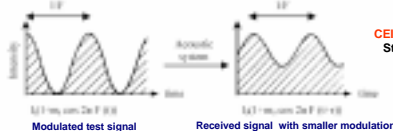


+ 21.5% relative

350 training utterances, 24 test utterances
Julius recognizer, 20k vocabulary Japanese dictation task

Intelligibility objective measurement

The intelligibility of speech refers to the accuracy with which a normal listener can understand a spoken word or phrase



CE-IEC standard n. 60269-216
Steeneken and Houtgast

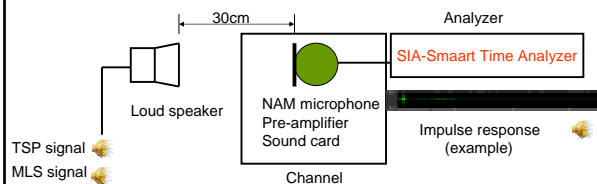
Speech Transmission Index (STI): Based on the reduction of m_i modulation index

Modulation function: Shows the reduction of the modulation index

$$m(F) = \frac{m_o}{m_i} \quad m(F) = \frac{\int_0^{\infty} h^2_f(\tau) \cdot \exp(-j \cdot 2\pi \cdot F \cdot \tau) d\tau}{\int_0^{\infty} h^2_f(\tau) d\tau}$$

Base on impulse response (Schroeder, 1981)

NAM microphone measurement



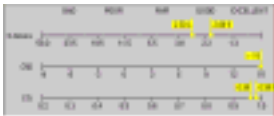
Speech Intelligibility objective measures

- STI Speech transmission index - Optimal 1.0 / Cannot hear 0.0
- RSTI Rapid Transmission Index - Data from narrow band
- Clarity Early to late energy ratio in impulse response
- % AICons Articulation Loss of consonants

NAM microphones measurement

Investigating effects of design (e.g. wrapping), quality of NAM and recording setup correctness

1



Reference microphone (Simple Desktop PC mic)

TSP signal

Impulse response

2



Pre-commercial Silicon NAM microphone

3



Silicon NAM microphone

4



Stethoscope microphone

Silicon microphones do not show significant degradations

Conclusions

- NAM microphones introduction
 - Able to capture audible and inaudible speech
 - Can be applied in speech recognition and communication
 - Easy to use – Small, light
- Experiments on audible speech recognition
 - Stethoscope NAM
 - Clean: 93.5% word accuracy
 - Noisy: 93.3% word accuracy (50dBA)
 - Silicon NAM
 - Clean: 94.86% word accuracy
- NAM quality and intelligibility investigation
 - NAM microphones show similar intelligibility to other sensors
- Data collection for speaker-independent experiments



Merry Christmas !