# Toward Spontaneous Speech Recognition and Understanding

*Sadaoki Furui*

**Tokyo Institute of Technology**
**Department of Computer Science**
**2-12-1, O-okayama, Meguro-ku, Tokyo, 152-8552 Japan**
**Tel/Fax: +81-3-5734-3480**
**furui@cs.titech.ac.jp**
**http://www.furui.cs.titech.ac.jp/**

---

# Outline

- Fundamentals of automatic speech recognition

- Acoustic modeling

- Language modeling

- Database (corpus) and task evaluation

- Transcription and dialogue systems

- Spontaneous speech recognition

- Speech understanding

- Speech summarization

# Outline

- Fundamentals of automatic speech recognition

- Acoustic modeling

- Language modeling

- Database (corpus) and task evaluation

- Transcription and dialogue systems

- Spontaneous speech recognition

- Speech understanding

- Speech summarization

# Speech recognition technology

## Categorization of speech recognition tasks

|  | Dialogue | Monologue |
|---|---|---|
| Human to human | (Category I)<br><br>Switchboard,<br>Call Home (Hub 5),<br>meeting task | (Category II)<br><br>Broadcasts news (Hub 4),<br>lecture, presentation,<br>voice mail |
| Human to machine | (Category III)<br><br>ATIS, Communicator,<br>information retrieval,<br>reservation | (Category IV)<br><br>Dictation |

## Major speech recognition applications

- **Conversational systems** for accessing information services
  - Robust conversation using wireless handheld/hands-free devices in the real mobile computing environment
  - Multimodal speech recognition technology
- Systems for **transcribing, understanding and summarizing** ubiquitous speech documents such as broadcast news, meetings, lectures, presentations and voicemails

## Mechanism of state-of-the-art speech recognizers

**Speech input**

**Acoustic analysis**

$x_1 \cdots x_T$

**Global search:**
**Maximize**

$P(x_1 \ldots x_T \mid w_1 \ldots w_k) \; P(w_1 \ldots w_k)$

**over** $w_1 \ldots w_k$

$P(x_1 \ldots x_T \mid w_1 \ldots w_k)$

$P(w_1 \ldots w_k)$

**Phoneme inventory**

**Pronunciation lexicon**

**Language model**

**Recognized word sequence**

---

## State-of-the-art algorithms in speech recognition

**LPC or mel cepstrum, time derivatives, auditory models**

**Speech input**

**Context-dependent, tied mixture sub-word HMMs, learning from speech data**

**Acoustic analysis**

**SBR, MLLR**

**Cepstrum subtraction**

**Phoneme inventory**

**Pronunciation lexicon**

**Global search**

**Frame synchronous, beam search, stack search, fast match, A* search**

**Language model**

**Recognized word sequence**

**Bigram, trigram, FSN, CFG**

# Outline

- Fundamentals of automatic speech recognition

- Acoustic modeling

- Language modeling

- Database (corpus) and task evaluation

- Transcription and dialogue systems

- Spontaneous speech recognition

- Speech understanding

- Speech summarization

**Digital sound spectrogram**

**Feature vector (short-time spectrum) extraction from speech**

**Spectral structure of speech**

**Logarithmic spectrum**

**Cepstrum**

**Spectral fine structure**

log

Fast periodical function of $f$

$\longrightarrow f$

**+**

**Spectral envelope**

log

Slow periodical function of $f$

$\longrightarrow f$

**IDFT**

**0**

$\longrightarrow t$

**+**

**Concentration at different positions**

0.

$\longrightarrow t$

**Relationship between logarithmic spectrum and cepstrum**

Log amplitude

10

8

6

4

2

0

**Spectral envelope by LPC**

**Spectral envelope by LPC cepstrum**

**Short-time spectrum**

**Spectral envelope by FFT cepstrum**

0          1          2          3

**Frequency [kHz]**

**Comparison of spectral envelopes by LPC, LPC cepstrum, and FFT cepstrum methods**

**Parameter (vector) trajectory**

**Instantaneous vector
(Cepstrum)**

**Transitional (velocity) vector
(Delta-cepstrum)**

**Cepstrum and delta-cepstrum coefficients**

**Speech**

**FFT**

**FFT based
spectrum**

**Mel scale
triangular filters**

**Log**

**DCT**

$\Delta$

$\Delta^2$

**Acoustic
vector**

**MFCC-based front-end processor**

$b_1(x)$  $b_2(x)$  $b_3(x)$

$x$  $x$  $x$

**Output probabilities**

0.2  0.4  0.7

1  0.5  2  0.6  3  0.3

0.3

**Phoneme models**

**Feature vectors**

time

**Phoneme $k$-1**  **Phoneme $k$**  **Phoneme $k$+1**

## Structure of phoneme HMMs

**Words**  grey  whales

**Phonemes**  ~ g  r  eʸ  w  eʸ  l  z  ~

**Allophones**  ~[g] r  g[r]eʸ r[eʸ]  w eʸ[w]eʸ  w[eʸ]l eʸ[l]z  l[z]~

**Allophone models**

**Spectrogram**

Frequency (kHz)

**Speech signal**

Amplitude

0.0  0.1  0.2  0.3  0.4  0.5  0.6

**Times (seconds)**

## Units of speech (after J.Makhoul & R. Schwartz)

# Outline

- Fundamentals of automatic speech recognition
- Acoustic modeling
- Language modeling
- Database (corpus) and task evaluation
- Transcription and dialogue systems
- Spontaneous speech recognition
- Speech understanding
- Speech summarization

---

# Language model is crucial !

- Rudolph the red nose reindeer.
- Rudolph the Red knows rain, dear.
- Rudolph the Red Nose reigned here.

- This new display can recognize speech.
- This nudist play can wreck a nice beach.

| 1. I | 5. ONE | 9. BOOKS | 13. OLD |
|------|--------|----------|---------|
| 2. WANT | 6. A | 10. COAT | |
| 3. NEED | 7. AN | 11. COATS | |
| 4. THREE | 8. BOOK | 12. NEW | |

**An example of FSN (Finite State Network) grammar**

# Syntactic language models

- **Rely on a formal grammar of a language.**

- **Syntactic sentence structures are defined by rules that can represent global constraints on word sequences.**

- **Mainly based on context-free grammars.**

- **Very difficult to extend to spontaneous speech.**

# Problems of context-free grammars

- **Over generation problem: not only generates correct sentences but also many incorrect sentences.**

- **Ambiguity problem: the number of syntactic ambiguities in one sentence becomes increasingly unmanageable with the number of phrases.**

- **Suitability for spontaneous speech is arguable.**

→**Stochastic context-free grammars**

---

# Statistical language modeling

Probability of the word sequence $w_1^k = w_1 w_2 ... w_k$ :

$$P(w_1^k) = \prod_{i=1}^{k} P(w_i | w_1 w_2 \ldots w_{i-1}) = \prod_{i=1}^{k} P(w_i | w_1^{i-1})$$

$$P(w_i | w_1^{i-1}) = N(w_1^i) / N(w_1^{i-1})$$

where $N(w_1^i)$ is the number of occurrences of the string $w_1^i$ in the given training corpus.

Approximation by Markov processes:
    **Bigram** model       $P(w_i | w_1^{i-1}) = P(w_i | w_{i-1})$
    **Trigram** model     $P(w_i | w_1^{i-1}) = P(w_i | w_{i-2} w_{i-1})$

Smoothing of trigram by the deleted interpolation method:
    $P(w_i | w_{i-2} w_{i-1}) = \lambda_1 P(w_i | w_{i-2} w_{i-1}) + \lambda_2 P(w_i | w_{i-1}) + \lambda_3 P(w_i)$

$P(w_t = \text{YES} \mid w_{t-1} = \text{sil}) = 0.2$

$P(s_t \mid s_{t-1})$

$S^{(1)}$   $S^{(2)}$   $S^{(3)}$   $S^{(4)}$   $S^{(5)}$   $S^{(6)}$

**Phoneme 'YE'**    **Phoneme 'S'**

**$w$ = YES**

$P(w_t = \text{sil} \mid w_{t-1} = \text{YES}) = 1$

**0.6**

$S^{(0)}$ **Silence**

$P(w_t = \text{sil} \mid w_{t-1} = \text{NO}) = 1$

**Start**

$S^{(7)}$   $S^{(8)}$   $S^{(9)}$   $S^{(10)}$   $S^{(11)}$   $S^{(12)}$

**Phoneme 'N'**    **Phoneme 'O'**

**$w$ = NO**

$P(Y \mid s_t = s^{(12)})$

$Y$

$P(w_t = \text{NO} \mid w_{t-1} = \text{sil}) = 0.2$

**Complete Hidden Markov Model of a simple grammar**

---

# Variations of *N*-gram language models

- **Various smoothing techniques**

- **Word class language models**

- ***N*-pos (parts of speech) models**

- **Combination with a context-free grammar**

- **Extend *N*-grams to the processing of long-range dependencies**

- **Cache-based adaptive/dynamic language models**

# Good-Turing estimate

**For any *n*-gram that occurs *r* times, we should pretend that it occurs *r\** times as follows:**

$$r^* = (r+1)\,\frac{n_{r+1}}{n_r}$$

**where $n_r$ is the number of *n*-grams that occur exactly *r* times in the training data.**

---

# Katz smoothing algorithm

Katz smoothing extends the intuitions of the Good-Turing estimate by adding the combination of higher-order models with lower-order models.

$$P_{Kats}(w_i|w_{i\text{-}1}) = \begin{cases} C(w_{i-1}w_i)/C(w_{i-1}) & \text{if } r > k \\ d_r C(w_{i-1}w_i)/C(w_{i-1}) & \text{if } k \geq r > 0 \\ \alpha(w_{i-1})/P(w_i) & \text{if } r = 0 \end{cases}$$

where $d_r = \dfrac{\dfrac{r^*}{r} - \dfrac{(k+1)\,n_{k+1}}{n_1}}{1 - \dfrac{(k+1)\,n_{k+1}}{n_1}}$ and $\alpha(w_{i-1}) = \dfrac{1 - \sum\limits_{w_i\,:\,r>0} P_{Kats}(w_i|w_{i\text{-}1})}{1 - \sum\limits_{w_i\,:\,r>0} P(w_i)}$

# Outline

- Fundamentals of automatic speech recognition

- Acoustic modeling

- Language modeling

- Database (corpus) and task evaluation

- Transcription and dialogue systems

- Spontaneous speech recognition

- Speech understanding

- Speech summarization

# Spontaneous speech corpora

- **Spontaneous speech variations: extraneous words, out-of-vocabulary words, ungrammatical sentences, disfluency, partial words, repairs, hesitations, repetitions, style shifting, ….**

- **"There's no data like more data" – Large structured collection of speech is essential.**

- **How to collect *natural* data?**

- **Labeling and annotation of spontaneous speech is difficult; how do we annotate the variations, how do the phonetic transcribers reach a consensus when there is ambiguity, and how do we represent a semantic notion?**

# Spontaneous speech corpora (cont.)

- **How to ensure the corpus quality?**

- **Research in automating or creating tools to assist the verification procedure is by itself an interesting subject.**

- **Task dependency: It is desirable to design a task-independent data set and an adaptation method for new domains → Benefit of a reduced application development cost.**

### Main database characteristics (Becchetti & Ricotti)

| Name | Quantities | | | | No. units |
|---|---|---|---|---|---|
| | No.CD | No.hours | Giga-bytes | No. speakers | |
| TI Digits | 3 | ~14 | 2 | 326 | >2,500 numbers |
| TIMIT | 1 | 5.3 | 0.65 | 630 | 6,300 sentences |
| NTIMIT | 2 | 5.3 | 0.65 | 630 | 6,300 sentences |
| RM1 | 4 | 11.3 | 1.65 | 144 | 15,024 sentences |
| RM2 | 2 | 7.7 | 1.13 | 4 | 10,608 sentences |
| ATIS0 | 6 | 20.2 | 2.38 | 36 | 10,722 utterances |
| Switchboard (Credit Card) | 1 | 3.8 | 0.23 | 69 | 35 dialogues |
| TI-46 | 1 | 5 | 0.58 | 16 | 19,136 isol. words |
| Road Rally | 1 | ~10 | ~0.6 | 136 | Dialogues/sentences |
| Switchboard (Complete) | 30 | 250 | 15 | 550 | 2,500 dialogues |
| ATC | 9 | 65 | 5.0 | 100 | 30,000 dialogues |
| Map Task | 8 | 34 | 5.1 | <256 | 128 dialogues |
| MARSEC | 1 | 5.5 | 0.62 | - | 53 monologues |
| ATIS2 | 6 | ~37 | ~5 | 351 | 12,000 utterances |
| WSJ-CSR1 | 18 | 80 | 9.2 | >124 | 38,000 utterances |

## Further database characteristics (Becchetti & Ricotti)

| Name | Transcription | | Speech style | Recording environment | SR kHz | Sponsor |
|---|---|---|---|---|---|---|
| | Based on: | TA | | | | |
| TI Digits | Word | No | Reading | QR | 20 | TI |
| TIMIT | Phone | Yes | Reading | QR | 16 | DARPA |
| NTIMIT | Phone | Yes | Reading | Tel | 8 | NYNEX |
| RM1 | Sentence | No | Reading | QR | 20 | DARPA |
| RM2 | Sentence | No | Reading | QR | 20 | DARPA |
| ATIS0 | Sentence | No | Reading spon. | Ofc | 16 | DARPA |
| Switchboard (Credit Card) | Word | Yes | Conv. spon. | Tel | 8 | DARPA |
| TI-46 | Word | No | Reading | QR | 16 | TI |
| Road Rally | Word | Yes | Reading spon. | Tel | 8 | DoD |
| Switchboard (Complete) | Word | Yes | Conv. spon. | Tel | 8 | DARPA |
| ATC | Sentence | Yes | Spon. | RF | 8 | DARPA |
| Map Task | Sentence | Yes | Conv. spon. | Ofc | 20 | HCRC |
| MARSEC | Phone | - | Spon. | Various | 16 | ESRC |
| ATIS2 | Sentence | No | Spon. | Ofc | 16 | DARPA |
| WSJ-CSR1 | Sentence | Yes | Reading | Ofc | 16 | DARPA |

---

# Entropy: Amount of information
# (Task difficulty)

- **Yes, No**            **1 bit (log 2)**

- **Digits**

  * **0 ~ 9 : 0.1**          **3.32 bits (log 0.1)**

  * $\begin{cases} \mathbf{0 : 0.91} \\ \mathbf{1 \sim 9 : 0.01} \end{cases}$   **0.722 bits**

         **(0.91 log 0.91 + 0.09 log 0.01)**

## Test-set perplexity

**Vocabulary: *A*, *B*    Test sentence: *ABB***

$$A\ 1/8 \qquad A\ 3/4 \qquad A\ 3/4$$
$$B\ 7/8 \qquad B\ 1/4 \qquad B\ 1/4$$

Test sentence entropy: $\log 8 + \log 4 + \log 4 = 7\,\mathrm{bit}$
  (*ABB*)

Entropy per word: 7 / 3 = 2.33 bits

Test-set perplexity (branching factor): $2^{2.33} = 5.01$

# Outline

- Fundamentals of automatic speech recognition
- Acoustic modeling
- Language modeling
- Database (corpus) and task evaluation
- Transcription and dialogue systems
- Spontaneous speech recognition
- Speech understanding
- Speech summarization

**A unigram grammar network where the unigram probability is attached as the transition probability from starting state S to the first state of each word HMM.**

**A bigram grammar network where the bigram probability $P(w_j|w_i)$ is attached as the transition probability from word $w_i$ to $w_j$.**

$$P(w_1|w_1, w_1)$$

$w_1$

$$P(w_1|w_2, w_1)$$

$$P(w_2|w_1, w_1)$$

$w_2$

$$P(w_2|w_2, w_1)$$
$$P(w_1|w_1, w_2)$$

$w_1$

$$P(w_2|w_1, w_2)$$

$$P(w_1|w_2, w_2)$$

$w_2$

$$P(w_1|w_2, w_2)$$

**A trigram grammar network where the trigram probability**
**$P(w_k|w_i, w_j)$ is attached to transition from grammar state $w_i$, $w_j$**
**to the next word $w_k$. Illustrated here is a two-word vocabulary,**
**so there are four grammar states in the trigram network.**

---

| **Text DB** | | **Speech** | | **Speech DB** |

**Language model training**

**Acoustic analysis**

**Acoustic model training**

**Bigram**

**Trigram**

**Beam-search decoder (First path)**

**Phone models (Tied-state Gaussian mixture HMM)**

**N-best hypotheses with acoustic score**

**Rescoring (Second path)**

**Recognition results**

**Two-pass search structure used in the Japanese broadcast-news transcription system**

20

# Outline

- Fundamentals of automatic speech recognition

- Acoustic modeling

- Language modeling

- Database (corpus) and task evaluation

- Transcription and dialogue systems

- Spontaneous speech recognition

- Speech understanding

- Speech summarization

**The main factors influencing speech communication**
The possible distortions in the communication and the redundancy of the speech signal interact, yielding to a message intelligible or unintelligible to a listener

# Difficulties in (spontaneous) speech recognition

- **Lack of systematic understanding in variability**

  **Structural or functional variability**

  **Parametric variability**

- **Lack of complete structural representations of (spontaneous) speech**

- **Lack of data for understanding non-structural variability**

---

**Overview of the Science and Technology Agency Priority Program "Spontaneous Speech: Corpus and Processing Technology"**



Large-scale spontaneous speech corpus → World knowledge / Linguistic information / Para-linguistic information / Discourse information

Spontaneous speech → Speech recognition → Tran-scription → Understanding / Information extraction / Summarization → Summarized text / Keywords / Synthesized voice

**For training
a morphological
analysis and POS
tagging program**

**C S J**

**Spontaneous
monologue**

*Core*

**Manually tagged
with segmental
and prosodic
information**

**500k words**

**Digitized
speech,
transcription,
POS and
speaker
information**

**7M words**

**Overall design of
the Corpus of Spontaneous Japanese (CSJ)**

**Test-set perplexity and OOV rate for the two language models**

| Acous. model | RdA | SpnA | RdA | SpnA | |
|---|---|---|---|---|---|
| Ling. model | SpnA | | | | |
| Speaker adapt. | WebL | WebL | SpnL | SpnL | SpnL |
| | w/o | w/o | w/o | w/o | with |

**Word accuracy for each combination of models**

---

0204-13

## Mean and standard deviation
## for each attribute of presentation speech

| | Acc | AL | SR | PP | OR | FR | RR |
|---|---|---|---|---|---|---|---|
| Mean | 68.6 | -53.1 | 15.0 | 224 | 2.09 | 8.59 | 1.56 |
| Standard deviation | 7.5 | 2.2 | 1.2 | 61 | 1.18 | 3.67 | 0.72 |

Acc: word accuracy (%),    AL: averaged acoustic frame likelihood,
SR: speaking rate (number of phonemes/sec),   PP: word perplexity,
OR: out of vocabulary rate,    FR: filled pause rate (%),
RR: repair rate (%)

- ━━ **Correlation** ━━ **Spurious correlation**

Acc: word accuracy, OR: out of vocabulary rate,
RR: repair rate, FR: filled pause rate,
SR: speaking rate, AL: averaged acoustic frame likelihood,
PP: word perplexity

**Summary of correlation between various attributes**

---

## Linear regression models of the word accuracy (%) with the six presentation attributes

**Speaker-independent recognition**
$$Acc = 0.12AL - 0.88SR - 0.020PP - 2.2OR + 0.32FR - 3.0RR + 95$$

**Speaker-adaptive recognition**
$$Acc = 0.024AL - 1.3SR - 0.014PP - 2.1OR + 0.32FR - 3.2RR + 99$$

Acc: word accuracy,  SR: speaking rate,
PP: word perplexity,  OR: out of vocabulary rate,
FR: filled pause rate,  RR: repair rate

# Outline

- Fundamentals of automatic speech recognition
- Acoustic modeling
- Language modeling
- Database (corpus) and task evaluation
- Transcription and dialogue systems
- Spontaneous speech recognition
- Speech understanding
- Speech summarization

---

*Speech generation*

(*Text generation*)

(*Speech production*)

Intension

Phonemes, prosody

Articulatory motions

| Message formulation | → | Language code | → | Neuro-muscular controls | → | Vocal tract system |

**Discrete signal**

**Continuous signal**

| 50 bps | 200 bps | 2,000 bps | 30,000 ~ 50,000 bps | **Acoustic waveform** |

*Information rate*

Semantics

Phonemes, words, syntax

Feature extraction, coding

Spectrum analysis

| Message understanding | ← | Language translation | ← | Neural transduction | ← | Basilar membrane motion |

**Discrete signal**

**Continuous signal**

(*Linguistic decoding*)

(*Acoustic processing*)

*Speech recognition*

**Human speech generation and recognition process**

26

## A communication - theoretic view of speech generation & recognition



$P(M)$ — Message source — $M$ → $P(W/M)$ — Linguistic channel — $W$ → $P(X/W)$ — Acoustic channel — $X$ → Speech recognizer

Language
Vocabulary
Grammar
Semantics
Context
Habits

Speaker
Reverberation
Noise
Transmission
   characteristics
Microphone

---

## Message-driven speech recognition

Maximization of *a posteriori* probability,

$$\max_{M} P(M \mid X) = \max_{M} \sum_{W} P(M \mid W) P(W \mid X) \tag{1}$$

Using Bayes' rule, it can be expressed as

$$\max_{M} P(M \mid X) = \max_{M} \sum_{W} P(X \mid W) P(W \mid M) P(M) / P(X) \tag{2}$$

For simplicity, we can approximate the equation as

$$\max_{M} P(M \mid X) \approx \max_{M,W} P(X \mid W) P(W \mid M) P(M) / P(X) \tag{3}$$

$P(W|M)$: hidden Markov models, $P(M)$: uniform probability for all $M$.

We assume that $P(W|M)$ can be expressed as follows.

$$P(W \mid M) \approx P(W)^{1-\lambda} P(W \mid M)^{\lambda} \tag{4}$$

where $\lambda$, $0 \le \lambda \le 1$, is a weighting factor.

# Word co-occurrence score

$P(W|M)$ is represented by word co-occurrences of nouns;

$$CoScore(w_i, w_j) = \log \frac{p(w_i, w_j)}{(p(w_i) p(w_j))^{1/2}}$$

$p(w_i, w_j)$: probability of observing words $w_i$ and $w_j$
       in the same news article

$p(w_i), p(w_j)$: probabilities of observing word $w_i$ and $w_j$
       in all the articles

A square root term was employed to compensate
the probabilities of the words with very low frequency.

---

**Speech**

↓

**Large vocabulary or keyword recognizer** ← **Language model**

↓

**Interface between recognition and natural language (e.g. N-best, word-graphs)**

↓

**Parser dealing with speech fragments** ← **Language model**

↓

**Additional semantic pragmatic or dialogue constraints**

↓

**Database query**

**Generic block diagram for spontaneous speech understanding**

**Data collection**

**Training**

**Decoding**

Speaker ↔ Speaker or Machine

Annotated data

Annotated learning

Syntax
**(short and long term structural relationships)**

Lexicon

Semantics

Discourse

Pragmatics

Integrated understanding framework

**Speech input**

**Meaning**

**Generic semi-automatic language acquisition for speech understanding**

## An architecture of a detection-based speech understanding system

**Speech input**

**Detector 1**

**Detector 2**

**Detector 3**

**Detector N-1**

**Detector N**

**Integrated search & confidence evaluation**

**Understanding & response**

**Partial language model**

**Solved by discriminative training, consistent with Neymann-Pearson Lemma**

**Partial language modeling and detection-based search still need to be solved**

# Outline

- Fundamentals of automatic speech recognition

- Acoustic modeling

- Language modeling

- Database (corpus) and task evaluation

- Transcription and dialogue systems

- Spontaneous speech recognition

- Speech understanding

- Speech summarization

# Sayings

- **The shortest complete description is the best understanding –** *Ockham*
- **If I had more time I could write a shorter letter –** *B. Pascal*
- **Make everything as simple as possible – *A. Einstein***

# From speech recognition to summarization

**LVCSR** (**L**arge **V**ocabulary **C**ontinuous **S**peech **R**ecognition) **systems** can transcribe **read speech** with 90% word accuracy or higher.

**Current target**
**LVCSR systems** for **spontaneous speech** recognition
to generate closed captions, abstracts, etc.

○ **Spontaneous speech features**
**filled pauses, disfluency,**
**repetition, deletion, repair,**
**etc.**

○ **Outputs from LVCSR systems**
**include recognition errors**

**Automatic speech summarization**

**Important information extraction**

---

# Summarization levels

**Summarization** ⟹ **Information extraction**
**Understanding speech**

**Indicative** summarization

**Informative** summarization

**Topics**

**Sentence(s)**

**Abstract**

**Summarized utterance(s)**

**Raw utterance(s)**

**Target**
**Closed captioning**
**Lecture summarization**

# Automatic speech summarization system

**Language database** → **Language model**

**Speech database**

**Speech** → **Acoustic model**

**Language model**, **Acoustic model** → **LVCSR module**

**Spontaneous speech**
**News, lecture, meeting etc.**

**Speaker adaptation**

**Knowledge database** → **Context model**

**Summary corpus** → **Summarization model**

**Context model**, **Summarization model** → **Summarization module (Understanding)**

**Captioning**

**Taking minutes**

**Making abstracts**

**Indexing**

---

# Approach to speech summarization utterance by utterance

**Each transcribed utterance**

1  2  3  4  5  6  7  8  9  10

**A set of words is extracted (sentence compaction)**

**Specified ratio**
**e.g. Extracting 7 words from 10 words: 70%**

1  2  3      6  7  8  9

**Summarized (compressed) sentence**

# Target of summarized speech

**Maintaining original meanings of speech
as much as possible**

| Information extraction | $\Rightarrow$ | **Significant (topic) words** |
| Linguistic correctness | $\Rightarrow$ | **Linguistic likelihood** |
| Semantic correctness | $\Rightarrow$ | **Dependency structure between words** |
| Word error exclusion | $\Rightarrow$ | **Acoustic and linguistic reliability** |

---

# Summarization score

**Summarized sentence with $M$ words $V = v_1, v_2, \ldots, v_M$**

**Summarization score**

$$S(V^M) = \sum_{m=1}^{M} \left\{ L(v_m | \cdots v_{m-1}) \right.$$

$\Rightarrow$ **Linguistic score**
Trigram

$$+ \lambda_I \, I(v_m)$$

$\Rightarrow$ **Significance (topic) score**
Amount of information

$$+ \lambda_C \, C(v_m)$$

$\Rightarrow$ **Confidence score**
Acoustic & linguistic reliability

$$\left. + \lambda_T \, T_r(v_m) \right\}$$

$\Rightarrow$ **Word concatenation score**
Word dependency probability

33

# Linguistic score

**Linguistic likelihood of word strings**
**(bigram/trigram)**
**in a summarized sentence**

$$\log P\,(v_m|\,v_{m\text{-}2}\,v_{m\text{-}1})$$

**Linguistic score is trained using a summarization corpus.**

# Word significance score

**Amount of information**

$$f_i \log \frac{F_A}{F_i}$$

$f_i$ : Number of occurrences of $u_i$ in the transcribed speech
$u_i$ : Topic word in the transcribed speech
$F_i$ : Number of occurrences of $u_i$ in all the training articles
$F_A$ : Summation of all $F_i$ over all the training articles

$$(F_A = \sum_i F_i\,)$$

•**Significance scores of words other than topic words and reappearing topic words are fixed.**

# Confidence score

**Acoustic and linguistic reliability of a word hypothesis**

**Posterior probability**

$$C(w_{k,l}) = \log \frac{\alpha_k \; P_{ac}(w_{k,l}) \; P_{lg}(w_{k,l}) \; \beta_l}{P_G}$$

$C(w_{k,l})$ : Log posterior probability of $w_{k,l}$

$k,l$ : Node index in a graph

$w_{k,l}$ : Word hypothesis between node $k$ and node $l$

$\alpha$ : Forward probability from the beginning node $S$ to node $k$

$\beta$ : Backward probability from node $l$ to the end node $T$

$P_{ac}$ : Acoustic likelihood of $w_{k,l}$

$P_{lg}$ : Linguistic likelihood of $w_{k,l}$

$P_G$ : Forward probability from the beginning node $S$ to the end node $T$

# Word concatenation score

**A penalty for word concatenation with no dependency in the original sentence**

**Inter-phrase**

**Intra-phrase**

**Intra-phrase**

| the | beautiful | cherry | blossoms |

| in | Japan |

*Phrase* 1

*Phrase* 2

**"the beautiful Japan"** ⟹ **Grammatically correct but incorrect as a summary**

35

# Dependency structure

**Dependency Grammar**

Dependency

modifier → head

Right-headed dependency

Left-headed dependency

The cherry blossoms bloom in spring

**Phrase structure grammar for dependency**

DCFG **(Dependency Context Free Grammar)**

$\alpha \rightarrow \beta\alpha$ **(Right-headed dependency)**
$\alpha \rightarrow \alpha\beta$ **(Left-headed dependency)**
$\alpha \rightarrow w$

$\alpha, \beta$: **Non-terminal symbols,**
$w$: **Terminal symbols**

VP
 ├ NP
 │  └ NP
 │     ├ DET ADJ NP
 │        The cherry blossoms
 └ VP
    ├ VP PP
       bloom  ├ PP NP
                in spring

# Word concatenation score based on SDCFG

## Word dependency probability

**If the dependency structure between words is deterministic,**

⇩

**0 or 1**

**If the dependency structure between words is ambiguous,**

⇩

**SDCFG**
(Stochastic DCFG)

**The dependency probability between $w_m$ and $w_l$, $d(w_m, w_l, i, k, j)$ is calculated using Inside-Outside probability based on SDCFG.**

$T(w_m, w_n)$

$= \log \sum_{i=1}^{m} \sum_{k=m}^{n-1} \sum_{j=n}^{L} \sum_{l=n}^{j} d(w_m, w_l, i, k, j)$

S
α
β α
β (Inside probability) α (Outside probability)

$w_1 \ldots\ldots w_{i-1} w_i \ldots w_m \ldots w_k w_{k+1} \ldots w_n \ldots w_l \ w_j \ w_{j+1} \ldots\ldots w_L$

**S**: Initial symbol, $\alpha, \beta$: Non-terminal symbol, $w$: Word

36

# Dynamic programming for summarizing each utterance

**Selecting a set of words maximizing the summarization score**

**Ex.:** &lt;s&gt; $w_2$ $w_4$ $w_5$ $w_7$ $w_{10}$ &lt;/s&gt;

Transcription result (vertical axis): &lt;/s&gt;, $w_{10}$, $w_9$, $w_8$, $w_7$, $w_6$, $w_5$, $w_4$, $w_3$, $w_2$, $w_1$, &lt;s&gt;

Summarized sentence (horizontal axis): &lt;s&gt; $v_1$ $v_2$ $v_3$ $v_4$ $v_5$ &lt;/s&gt;

$$g(m, l, n) = \max_{k < l}\{g(m-1, k, l)$$
$$+ \log P(w_n | w_k w_l)$$
$$+ \lambda_I I(w_n) + \lambda_C C(w_n)$$
$$+ \lambda_T T_r(w_l, w_n)\}$$

$n$, $l$, $k$ with $w_n$, $w_l$, $w_k$; columns $m-2$, $m-1$, $m$

---

# Summarization of multiple utterances

**The method of summarizing each utterance is extended to summarize a set of multiple utterances by adding a rule giving a restriction at utterance boundaries.**

**Original utterances having many informative words are preserved and utterances having few informative words are deleted or shortened.**

⟹

**Important sentence extraction**

$+$

**Summarization of each utterance**

37

# Dynamic programming
# for summarizing multiple utterances

**\* Initial and terminal symbols cannot be skipped.**
**\* Word concatenation score is not applied to the utterance boundaries.**



0%    **Summarization ratio**    100%

Recognition result

Utterance 3

Utterance 2

Utterance 1

$<s>\ v_{1,1}\ v_{1,2}\ ...</s><s>\ v_{2,1}\ v_{2,2}\ ...</s><s>\ v_{3,1}\ v_{3,2}\ ...</s>$

**Summarization hypothesis**

---

# Evaluation experiments



**40% or 70%**
**summarization ratio**

Random word selection → Random word selection of REC

LVCSR system → Automatic transcription (REC) → Summarization Module (Training corpora)

→ Automatic summarization of REC

→ Automatic summarization of TRS

Broadcast news utterances

Manual transcription (TRS) → Sentence compaction by human subjects → Manual summarization of TRS

Random word selection → Random word selection of TRS

## Word network of manual summarization results for evaluation

**Manual summarization results are merged into a network.**

• **The network approximately expresses all possible correct summarization including subjective variations.**



For 5 subjects

Summarization accuracy is defined as the word accuracy based on
the word string, extracted from the word network, that is most similar
to the automatic summarization result.

Summarization accuracy = {Len-(Sub+Ins+Del) }/Len*100 [%]
   Len: number of words in the most similar word string in the network
   Sub: number of substitution errors
   Ins: number of insertion errors
   Del: number of deletion errors

---

## Examples of automatic summarization for manually transcribed CNN news - 1

• **Transcription:**
*It's sulfur, and as Ed Garsten reports in today's edition of tech trends, the petroleum industry is proposing a cleanup.*

• **Automatic summarization (30-40% summarization ratio) :**
*sulfur Ed Garsten reports tech petroleum __ proposing cleanup.*

• **The most similar word string in the manual summarization network:**
*Ed Garsten reports tech trends industry proposing cleanup.*

• **Automatic summarization (50-70% summarization ratio):**
*sulfur Ed Garsten reports in today's edition tech trends petroleum industry is proposing cleanup.*

• **The most similar word string in the manual summarization network :**
*Sulfur, Garsten reports in today's tech trends the industry is proposing cleanup.*

0205-26

## Examples of automatic summarization for manually transcribed CNN news - 2

- **Transcription:**
  *We are dealing with something of such a massive uh size and potential impact, um that a lot of people wisely are saying hands off.*

- **Automatic summarization (20-40% summarization ratio) :**
  *We're dealing something __ impact lot of people saying hands __.*

- **The most similar word string in the manual summarization network :**
  *We're dealing something such impact lot of people saying hands off.*

- **Automatic summarization (50-70% summarization ratio) :**
  *We're dealing with something  of a size and impact, a lot of people wisely are saying hands __.*

- **The most similar word string in the manual summarization network :**
  *We're dealing with something of such size and impact, a lot of people wisely are saying hands off.*

---

0205-27

## Examples of automatic summarization for recognized CNN news (80% recognition accuracy)

- **Recognition result:**
  *Vice president Al Gore says the government has a plan to avoid the inevitable prospect of increased airplane crashes and fatality is*

- **Automatic summarization (40% summarization ratio) :**
  *Gore the government has a plan to avoid the increased airplane crashes*

- **The most similar word string in the word network:**
  *<INS> the government has a plan to avoid the increased airplane crashes*

- **Automatic summarization (70% summarization ratio) :**
  *Vice president Al Gore says the government has a plan to avoid <DEL> increased airplane crashes*

- **The most similar word string in the word network:**
  *Vice president Al Gore says the government has a plan to avoid the increased airplane crashes*

**English news speech summarization**
**(Each utterance summarization)**

Summarization accuracy [%]

REC  TRS  REC  TRS
40%  70%

RDM: Random word selection
SIG: Word significance score
2gram: Linguistic score
CM: Confidence score

SDCFG: Word concatenation score
TRS_SUB: Manual summarization
of manual transcription

**English news speech summarization**
**(Multiple utterance summarization)**

Summarization accuracy [%]

REC  TRS  REC  TRS
40%  70%

RDM: Random word selection
SIG: Word significance score
2gram: Linguistic score
CM: Confidence score

SDCFG: Word concatenation score
TRS_SUB: Manual summarization
of manual transcription

# Recognition error reduction

**Each utterance summarization**

40%    70%

Number of word errors in summarization results

200
160
120
80
40
0

**Exclusion of out of context words**

**Contribution by CM**

**Multiple utterance summarization**

40%    70%

350
300
250
200
150
100
50
0

| ☐ Word error | ☐ SIG |
| ■ RDM | ☐ SIG+2gram |

| ☐ SIG+2gram+CM |
| ☐ SIG+2gram+SDCFG |
| ☐ SIG+2gram+SDCFG+CM |

---

**Speech recognition /understanding**
• Speaker-independent
• Spontaneous speech

**Speech coding**
• Wide/narrow-band
• Very-low-bit-rate

**Speech synthesis**
• Synthesis by rule
• Text-to-speech

**Robustness**
• Noise/distortion

**Individuality**
• Speaker recognition
• Speaker adaptation/normalization
• Voice conversion

**Human-machine interface**
• Ergonomics
• Subjective/objective evaluation

Database

**Feature extraction (dynamics)**

*Speech analysis*

Database

Psychology    Speech perception    Nerve system    Speech production    Acoustic phonetics    Memory/learning    Artificial Intelligence    Signal processing

Physiology    Articulation    Acoustics

**Speech information processing "tree"**, consisting of present and future speech information processing technologies supported by scientific and technological areas serving as the foundations of speech research.

## Summary

- *Speech recognition technology* **has made significant progress with many potential applications.**
- **How to model and recognize** *spontaneous speech* **is one of the most important issues.**
- **Construction of a large-scale** *spontaneous speech corpus* **is crucial.**
- **Paradigm shift from recognition to** *understanding* **is needed.**
- *Speech summarization* **is attractive as information extraction and speech understanding.**

## References

- C. Becchetti and L. P. Ricotti: *Speech Recognition*, John Wiley & Sons, Ltd., New York, 2000
- S. Furui: *Digital Speech Processing, Synthesis, and Recognition, Second Edition*, Signal Processing and Communications Series, Marcel Dekker, New York, 2000
- D. Gibbon, I. Mertins and R. K. Moore (Eds.): *Handbook of Multimodal and Spoken Dialogue Systems*, Kluwer Academic Publishers, Boston, 2000
- B.-H. Juang and S. Furui: "Automatic recognition and understanding of spoken language processing – A first step toward natural human-machine communication," Proc. IEEE, 88, 8, pp. 1142-1165, 2000
- C. Hori, S. Furui, R. Malkin, H. Yu and A. Waibel: "Automatic summarization of English broadcast news speech," Proc. Human Language Technology 2002, San Diego, pp. 228-233, 2002
- X.-D. Huang, A. Acero and H.-W. Hon: *Spoken Language Processing*, Prentice Hall PTR, New Jersey, 2001
- S. Young and G. Bloothooft (Eds.): *Corpus-based Methods in language and Speech Processing*, Kluwer Academic Publishers, Boston, 1997