# Prologue: Introduction to Intelligent Speech Technology
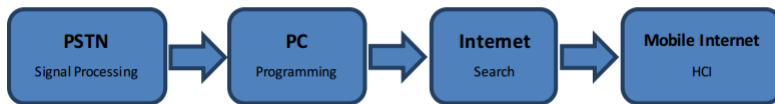
Yanmin Qian

SpeechLab
Department of Computer Science & Engineering
Shanghai Jiao Tong University

Autumn 2017

# Syllabus

- **Prologue: Introduction to Intelligent Speech Technology**
- **Part 1: Basic Concepts and Theories**
  - Probability and stochastic process
  - Pattern recognition and machine learning
- **Part 2: Fundamental of Speech Recognition**
  - Speech signal processing
  - Acoustic modelling (Hidden Markov Models)
  - Language modelling ($n$-grams)
  - Decoding algorithm
  - Large vocabulary continuous speech recognition (LVCSR)
- **Part 3: Advanced Topics of Speech Recognition**
  - Deep neural network for speech recognition
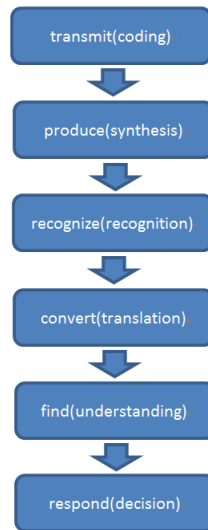  - Discriminative training and adaptation

# Intelligent Speech Technology



**Speech interaction** is one of the most important forms of Human Computer Interaction (HCI).

# Speech and Language Processing

**Speech and language processing** aims at modelling and manipulating information from speech and text to

- **transmit (coding)** speech signal efficiently

- **produce (synthesis)** human-like natural text and/or speech

- **recognize (recognition)** underlying text and/or other information from speech signals

- **convert (translation)** text from language to another language

- **find (understanding)** semantic and syntactic content from recognized text and other info.

- **respond (decision)** to incoming semantic content to form conversation
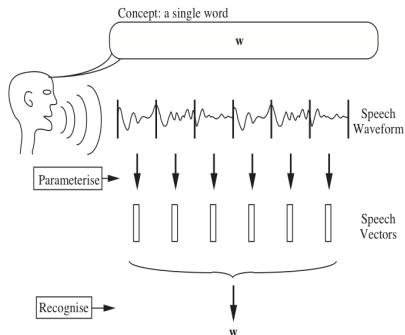
Name a few applications?

transmit(coding)

produce(synthesis)

recognize(recognition)

convert(translation)

find(understanding)

respond(decision)

# Speech Recognition

# Statistical Speech Recognition

## Diagram of Isolated Word Recognition

$$\hat{\mathbf{W}} = \arg\max_{\mathbf{W}} p(\mathbf{A}|\mathbf{O})p(\mathbf{O}|\mathbf{L})P(\mathbf{L}|\mathbf{W})P(\mathbf{W})$$

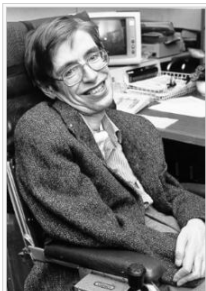# Speech Synthesis



中文  英语  德语  检测语言

语音识别                                              ×
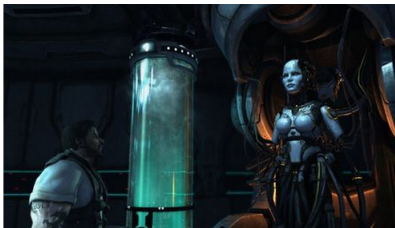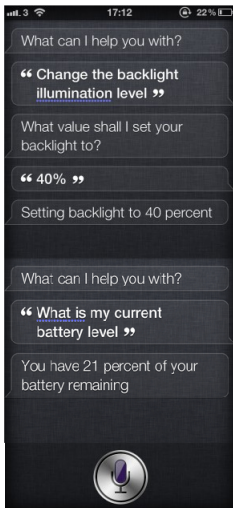
🎤 拼 ▾                                    🔊 💬 Ā

日语  中文(简体)  英语

Speech Recognition



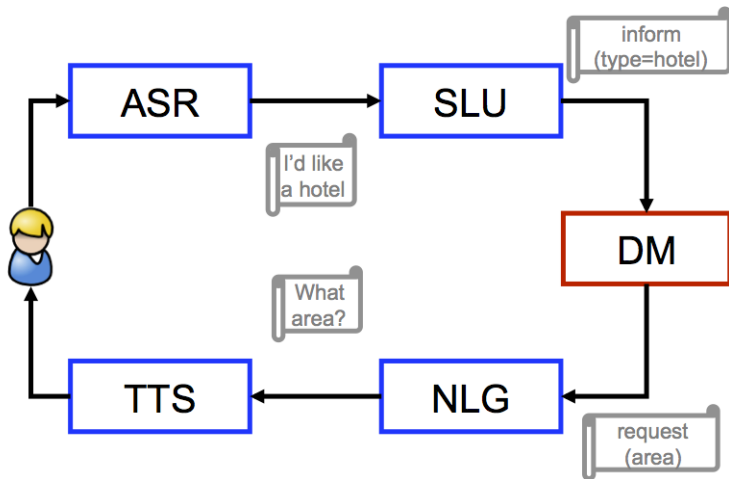Stephen Hawking is one of the most famous people using speech synthesis to communicate

Raynor talking to Adjutant in Starcraft 2

In-car spoken dialogue system

# Spoken Dialogue System

## Task-oriented spoken dialogue system

# Recognition of Non-text Information (Speaker, Language, Emotion, Humming etc.)

- Prologue: Introduction to Intelligent Speech Technology
- **Part 1: Basic Concepts and Theories**
    - Probability and stochastic process
    - Pattern recognition and machine learning
- **Part 2: Fundamental of Speech Recognition**
    - Speech signal processing
    - Acoustic modelling (Hidden Markov Models)
    - Language modelling ($n$-grams)
    - Decoding algorithm
    - Large vocabulary continuous speech recognition (LVCSR)
- **Part 3: Advanced Topics of Speech Recognition**
    - Neural network in speech recognition
    - Discriminative training and adaptation

# Key Actions and Course Works

- **Evaluation**: No Examination
  - Attendance or request for leave (10%)
  - Projects + Talks (90%)

- **Part 1-Project 1: Basic concepts and theories (40%)**
  - Memerize and Instantiate
  - Implement EM algorithm for GMM – Individual
- **Part 2-Project 2: Automatic Speaker Verification Spoofing Detection (30%)**
  - Derive and Practise
  - Design the complete system – Challenge (Group)
- **Part 3-Talk: Advances on Speech Processing (20%)**
  - Read and Summarize
  - Paper/Derivation/Tool presentation – Group

- **EM algorithm for GMM (40%)**
  - A simple binary classification task
    - 2-dimensional feature are provided
    - Three sets are provided: training, dev and test
    - GMM training and evaluation
  - Detailed report + Classification results + Source code

# Project 2 - Automatic Speaker Verification Spoofing Detection

- **Automatic Speaker Verification Spoofing Detection Challenge (30%)**
    - 1st ASVspoofing Challenge in SJTU
    - A detection and recognition task
        - Training set and test set are provided
    - The complete system design
        - Feature extraction
        - Model training
        - Recognition or classification
    - The data and rules are the same as the ASVspoof 2017 Challenge
        - http://www.asvspoof.org/
    - Detailed technical report with results (formal paper style)
        - The formal ICASSP template will be provided
    - Rank & Reward

# Talk - Advances on Speech Processing

- **Advanced Talks on Speech Processing (20%)**
    - New Progress in Speech & Language processing
    - First Read and then Present
    - Any topics are welcomed
        - Speech Enhancement
        - Speech Recognition
        - Speaker / Language Identification
        - Speech Emotion Recognition
        - Language Modeling
        - ......

# References

- **Probability and pattern recognition basics**

  Pattern Recognition and Machine Learning.
  Christopher M. Bishop, Springer.

- **Speech recognition theory and tools**

  HTK Book, Steve Young, et al. Cambridge University

- **Speech and language technology**

  Spoken Language Processing.
  A Guide to Theory, Algorithm and System Development
  Xuedong Huang, Alex Acero, Hsiao-Wuen Hon

- **Deep Learning for Speech Processing**

  Automatic Speech Recognition-A Deep Learning Approach.
  Dong Yu & Li Deng. Springer

# Information

- Yanmin Qian
  - Computer Science Department-SEIEE 3#515
  - Email: yanminqian@sjtu.edu.cn
  - Cell-Phone: +86-18516056597
- Teaching Assistant: Heinrich Dinkel
  - Computer Science Department-SEIEE 3#225
  - Email: richman@sjtu.edu.cn
  - Cell-Phone: +86-13262830583
- Speech-Lab in SJTU (Speech-Lab in SJTU)
  - Homepage: http://speechlab.sjtu.edu.cn/

# FTP

- Materials for download and upload
  - FTP: ftp://202.120.38.125
  - Port: 8821
  - Username: speech2017m
  - Password: speech2017m