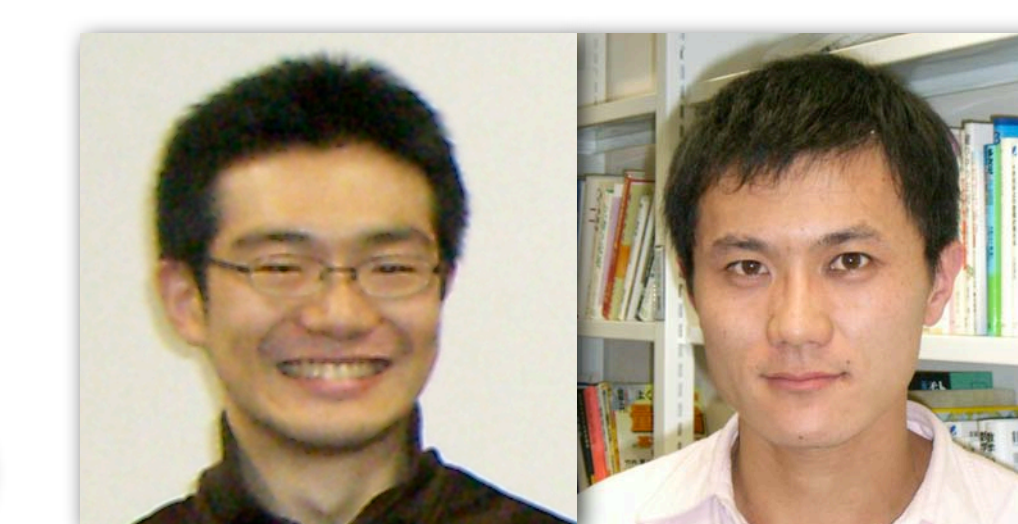# PRONUNCIATION CLINIC

## Nobuaki Minematsu, Max Takazawa, Xuebin Ma @ The University of Tokyo
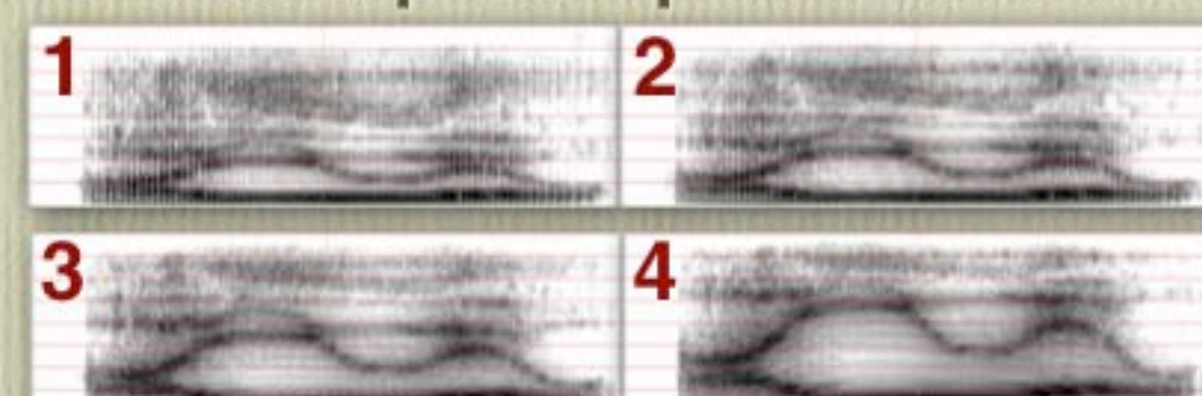
### !! A big problem !!

- **A very important and requisite function for CALL systems**
  - The system has to be able to ignore speaker individualities.
    - Age and gender (the size and length of the vocal tube)
    - A desirable system must not be able to see differences in age and gender.
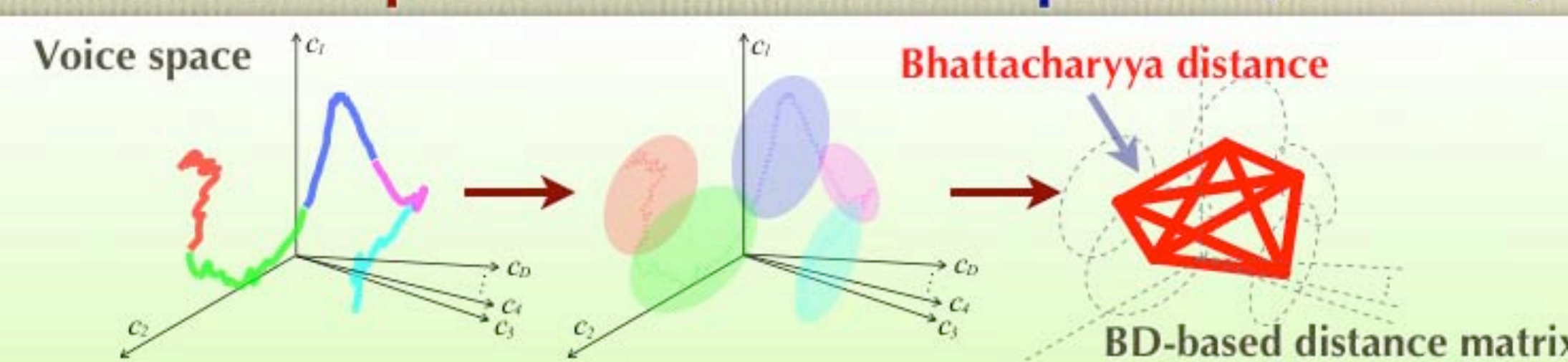  - Some examples of speaker differences
  - Source + filter model
    - Separation between source and filter
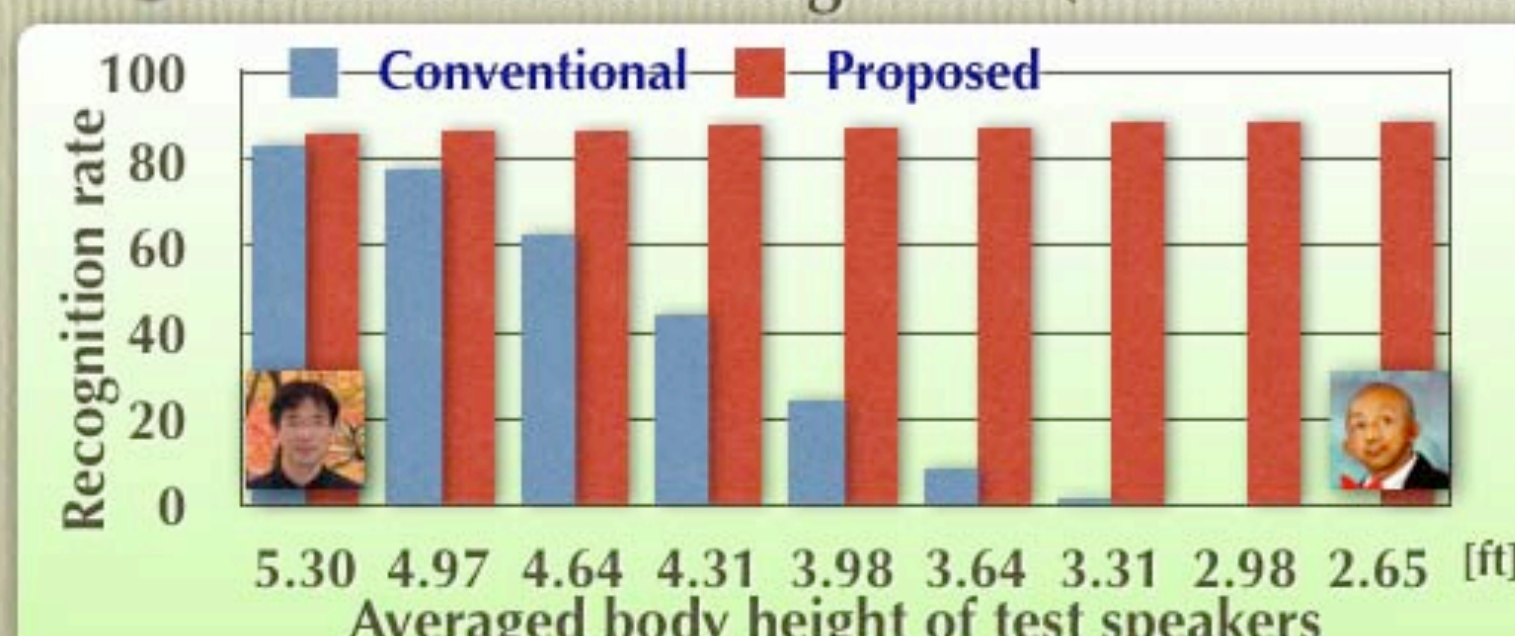    - Separation between ling. and extra-ling.

Mismatch problem

### !! Our solution of that problem !!

- **Holistic and speaker-invariant sound pattern (structure)**
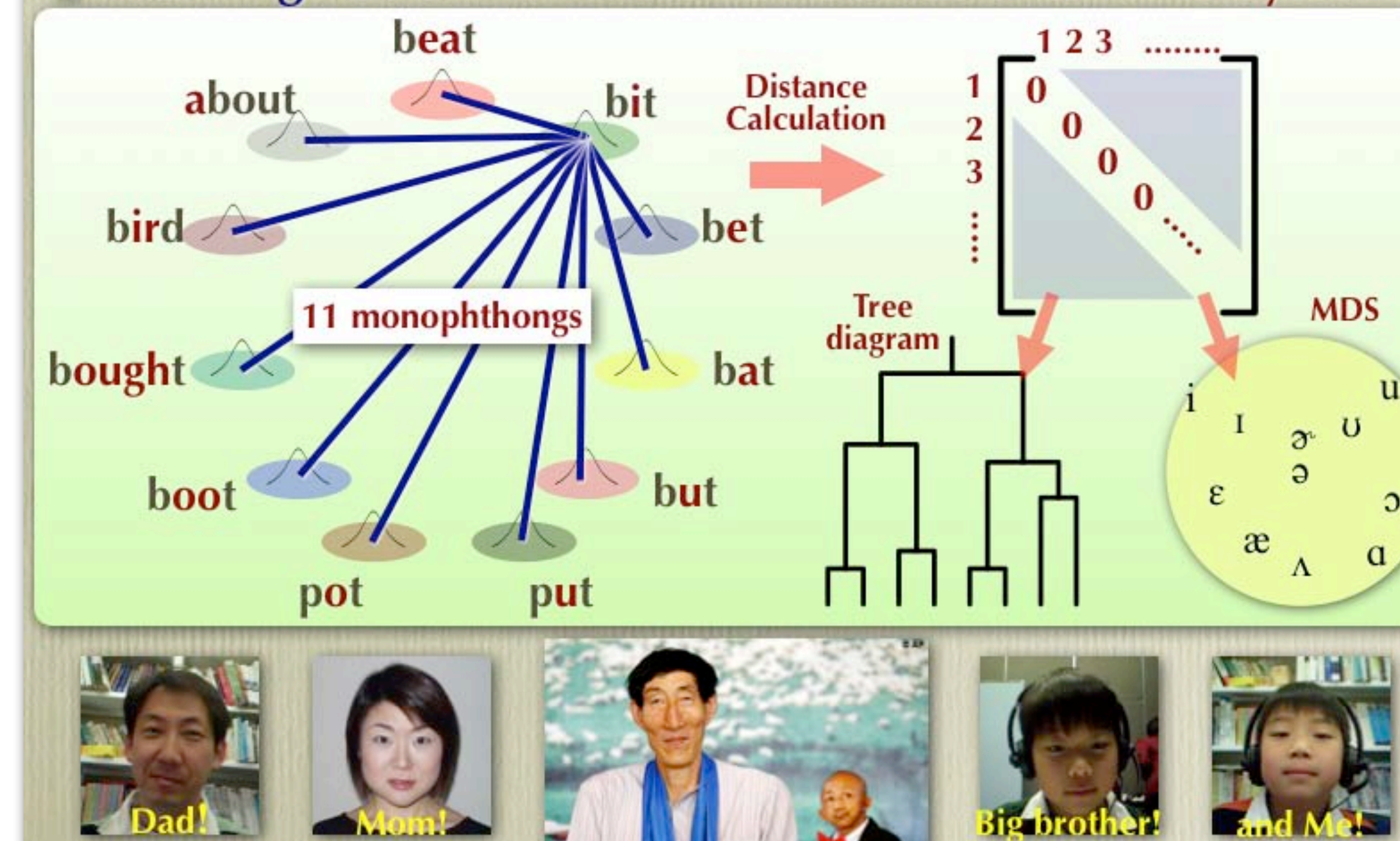
  Voice space

  Bhattacharyya distance

  BD-based distance matrix

- **Use of structures for automatic speech recognition**
  - Isolated word recognition (word = 5 vowel sequence, e.g. /aeoui/)

  Recognition rate — Conventional / Proposed
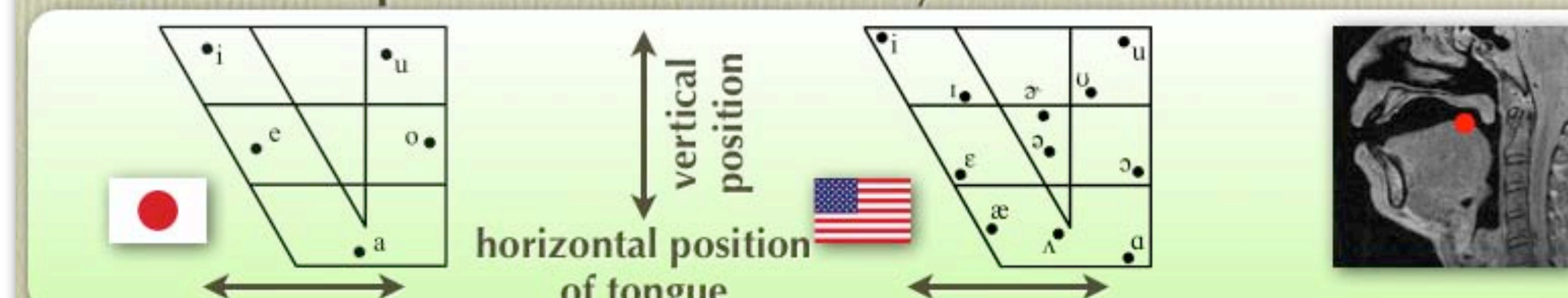
  5.30 4.97 4.64 4.31 3.98 3.64 3.31 2.98 2.65 [ft]
  Averaged body height of test speakers

  The tallest & the shortest

### A vowel training system for everybody!!

- **Learning not of individual vowels but of a vowel system**

  beat, bit, about, bird, bet, bought, bat, boot, but, boot, pot, put

  11 monophthongs

  Distance Calculation

  Tree diagram

  MDS

  Dad! Mom! Big brother! and Me!

### Structural representation of the vowel system

- **Vowel system and vowel chart**
  - Accented pronunciation = vowel system with some distortion

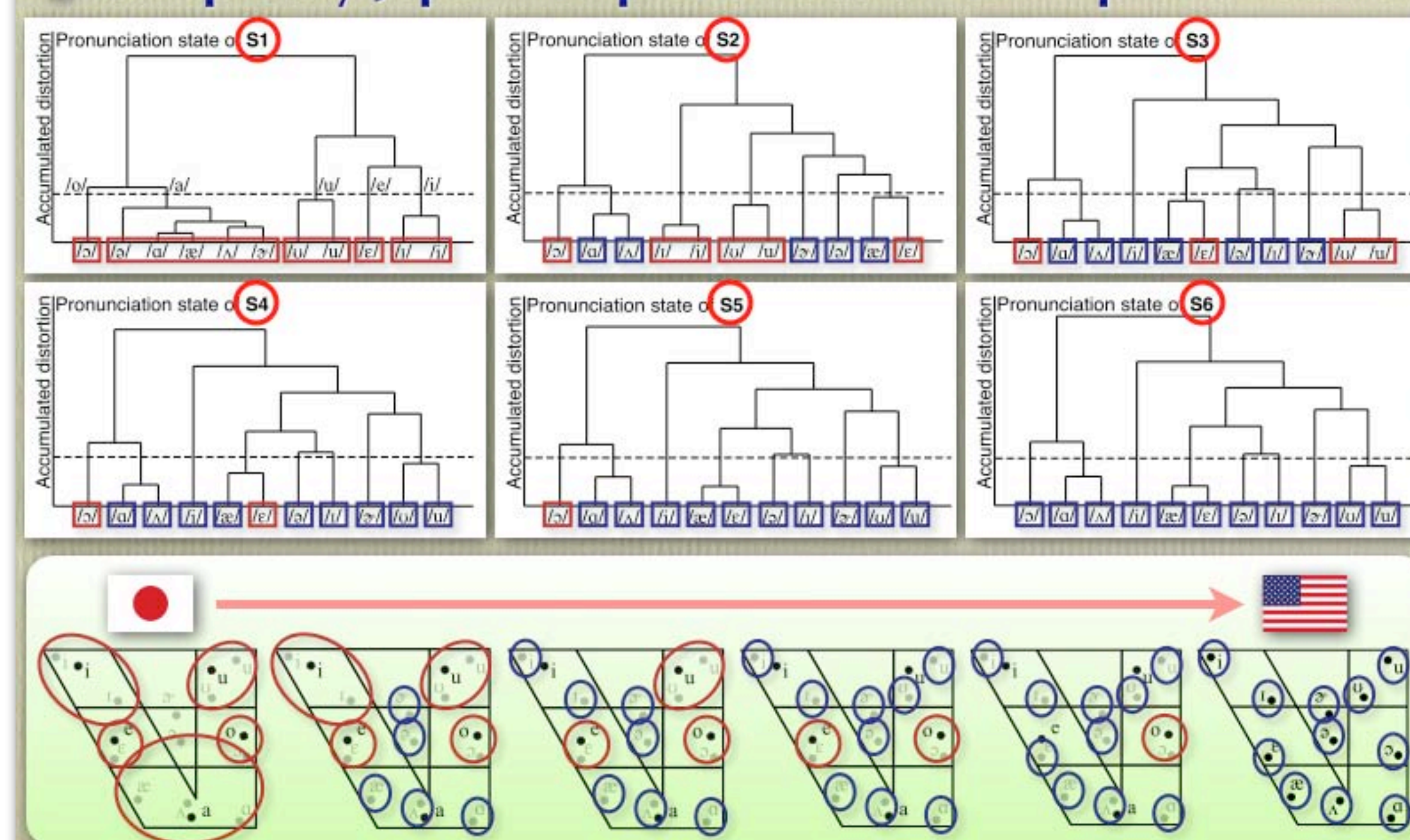  vertical position

  horizontal position of tongue

- **What's possible in the proposed demo system**
  - The demo system can
    1. record or log a history of vowel pronunciation training of each learner.
    2. provide for learners a window of "favorite teacher selection".
    3. show which vowel to correct first to become like the selected teacher.
    4. classify all the registered learners only wrt pronunciation proficiency by ignoring gender, age, etc. very effectively.
    5. give a very motivating user-interface for pronunciation training.
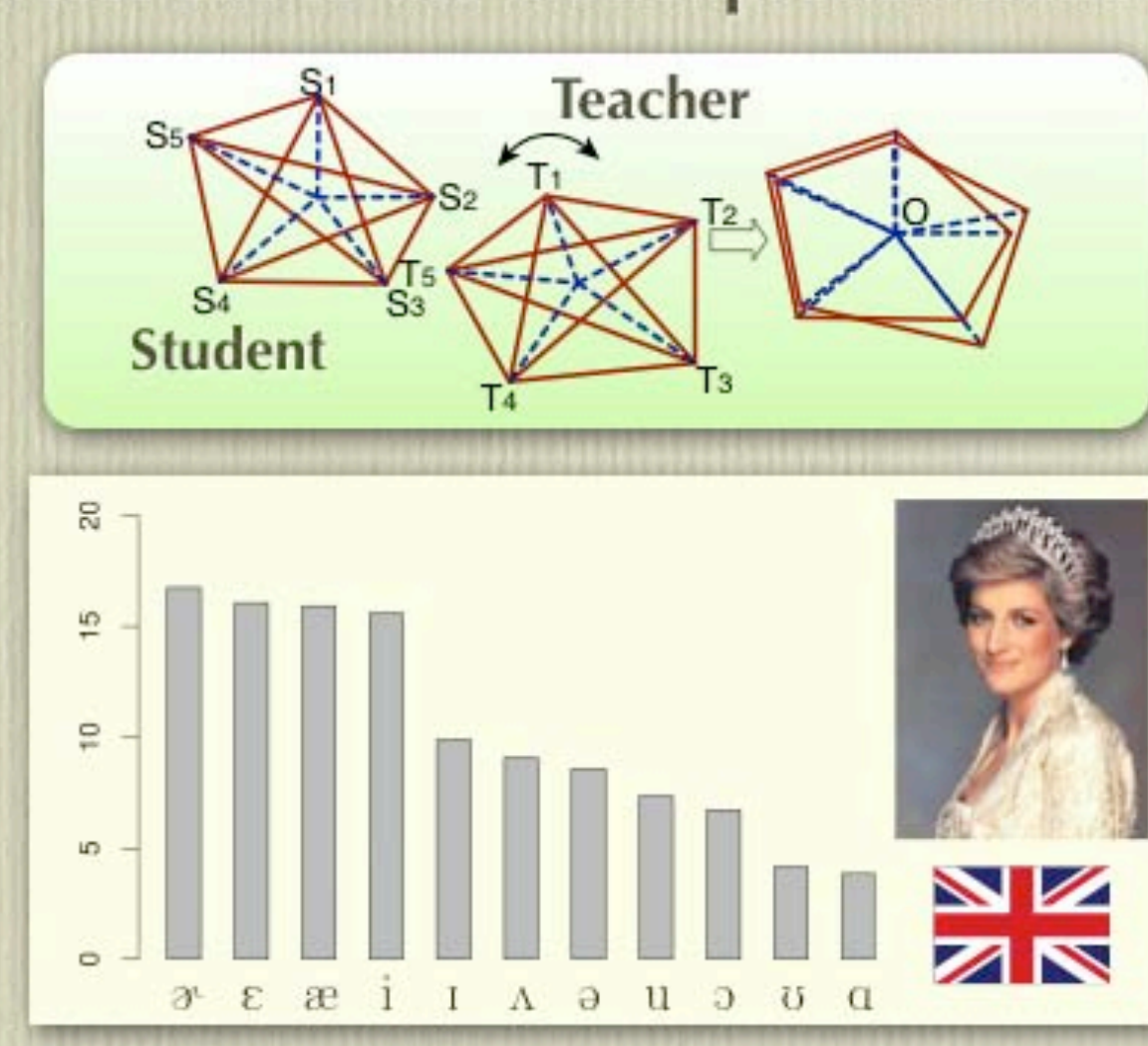
### Developmental changes in vowel training ①

- **Completely Japanized pronunciation to AE pronunciation**

  Pronunciation state of S1, S2, S3, S4, S5, S6

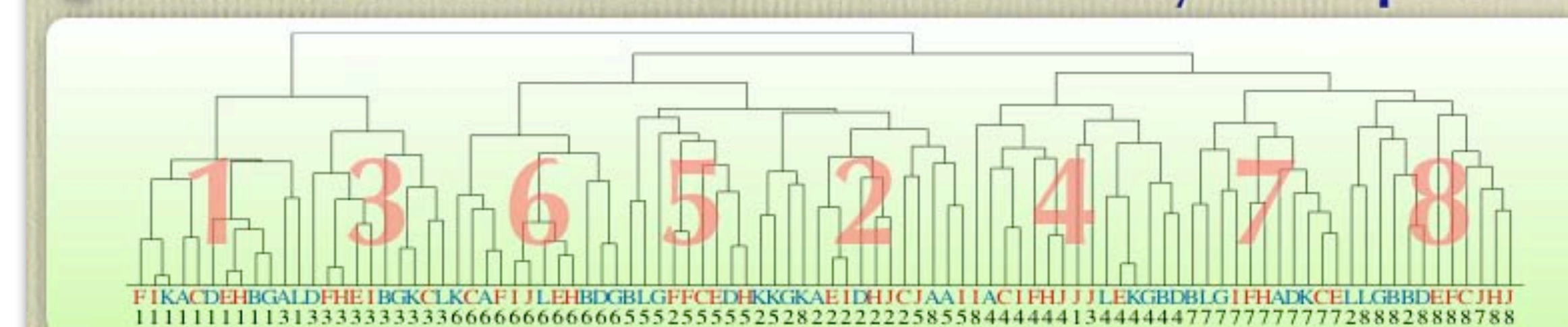### Which vowels to correct at first in your case? ②③

- **Who is your model speaker?**
  - A famous phonetician, a movie star (character) or a sport player??
- **Which vowels to correct at first to become like him/her?**
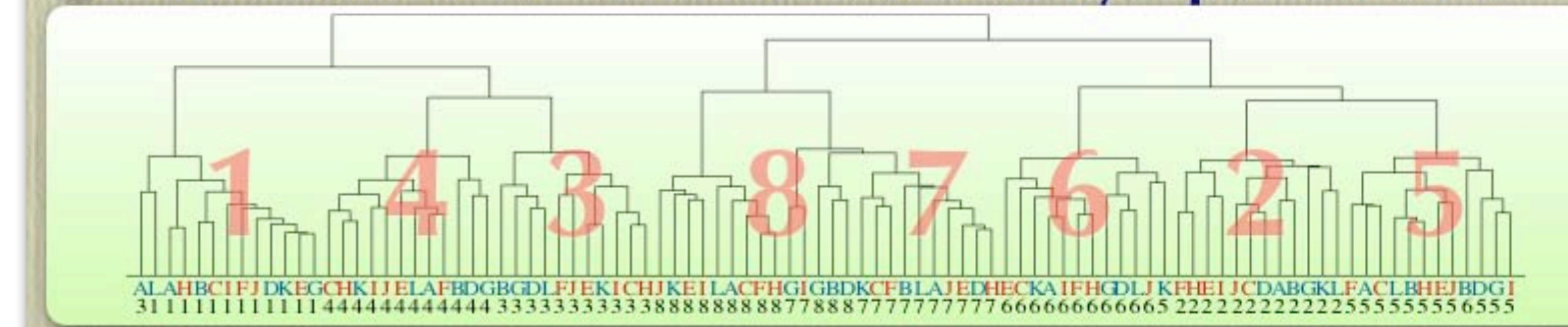  - The system can show the shortest cut to the model pronunciation.

  Select two teachers — OK

  Teacher

  Student

### Classification of learners ④

- **Automatic classification of 96 learners by a computer**

  1 3 6 5 2 4 7 8

- **Manual classification of 96 learners by a phonetician**

  1 4 3 8 7 6 2 5

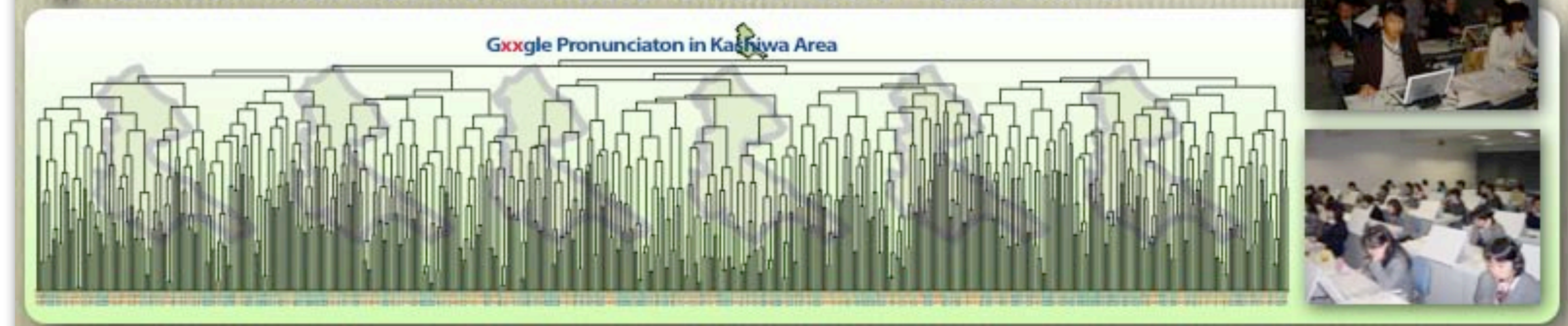  - 8 speakers x 12 pronunciations = 96 simulated learners
  - 1 to 8 = pronunciations, A to L = speakers
  - By substituting J vowels for some E ones, 12 v-systems are defined.

### Classification of learners ④⑤

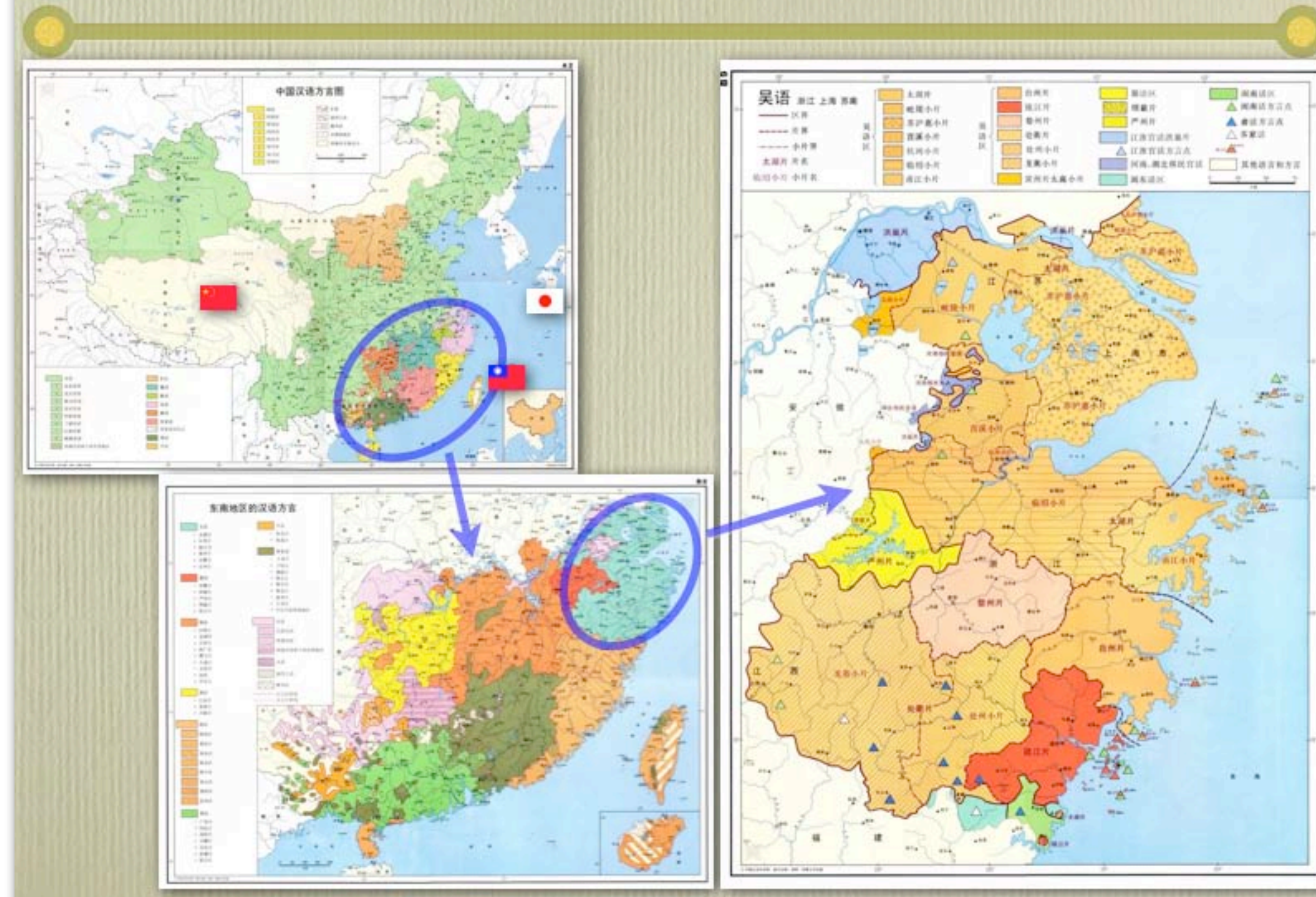- **Changes of students in a class before and after training**

  1 week later...

  Family members & univ. students

  Father, Younger Son, Older Son, Mother

  American English, British English

  English teachers

- **Classification of more than 600 learners**
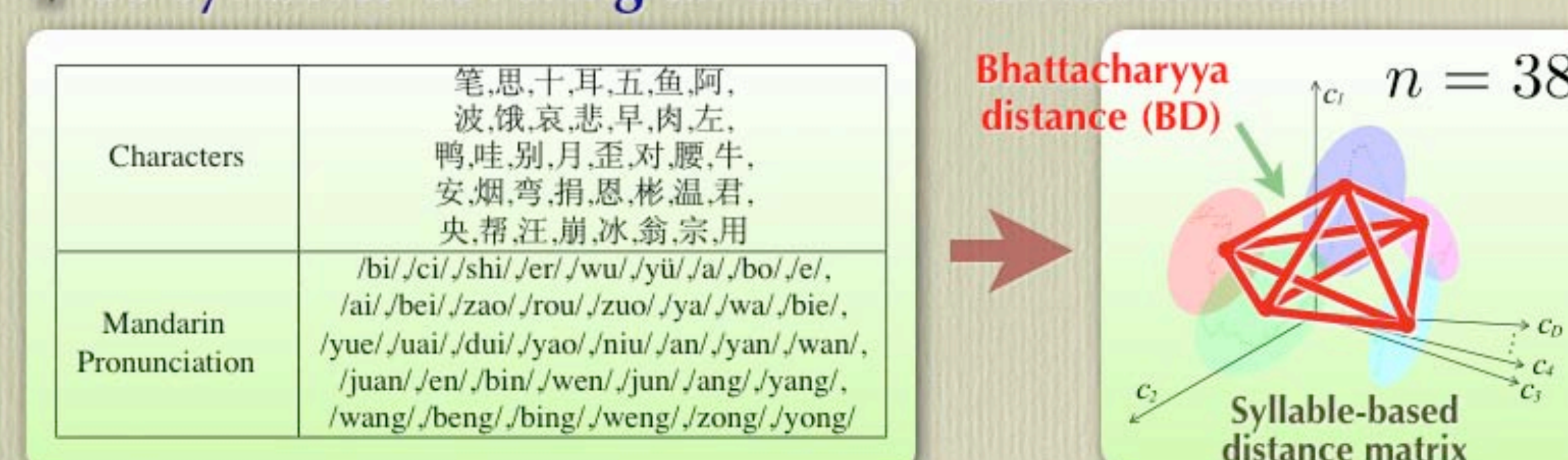
  Google Pronunciation in Kagawa Area

# DIALECT-BASED SPEAKER CLASSIFICATION

### Chinese dialects & sub-dialects

### Pronunciation structure of the 38 syllables

- **38 syllables covering all the 38 Mandarin finals**

  | | |
  |---|---|
  | Characters | 笔思十耳五鱼阿, 波饿哀悲早肉左, 鸭娃别月盂对腿牛, 安烟弯捐恩彬温君, 央帮正崩冰翁宗用 |
  | Mandarin Pronunciation | /bi/ /ci/ /shi/ /er/ /wu/ /yü/ /a/ /bo/ /e/, /ai/ /bei/ /zao/ /rou/ /zuo/ /ya/ /wa/ /bie/, /yue/ /uai/ /dui/ /tui/ /niu/ /an/ /yan/ /wan/, /juan/ /en/ /bin/ /wen/ /jun/ /ang/ /yang/, /wang/ /beng/ /bing/ /weng/ /zong/ /yong/ |

  Bhattacharyya distance (BD)

  $n = 38$

  Syllable-based distance matrix

- **Two definitions of structure-to-structure distance**

  Contrast-based (proposed)

  $$D_1(A, B) = \sqrt{\frac{1}{38} \sum_{i<j} (A_{ij} - B_{ij})^2}$$

  Substance-based (conventional)

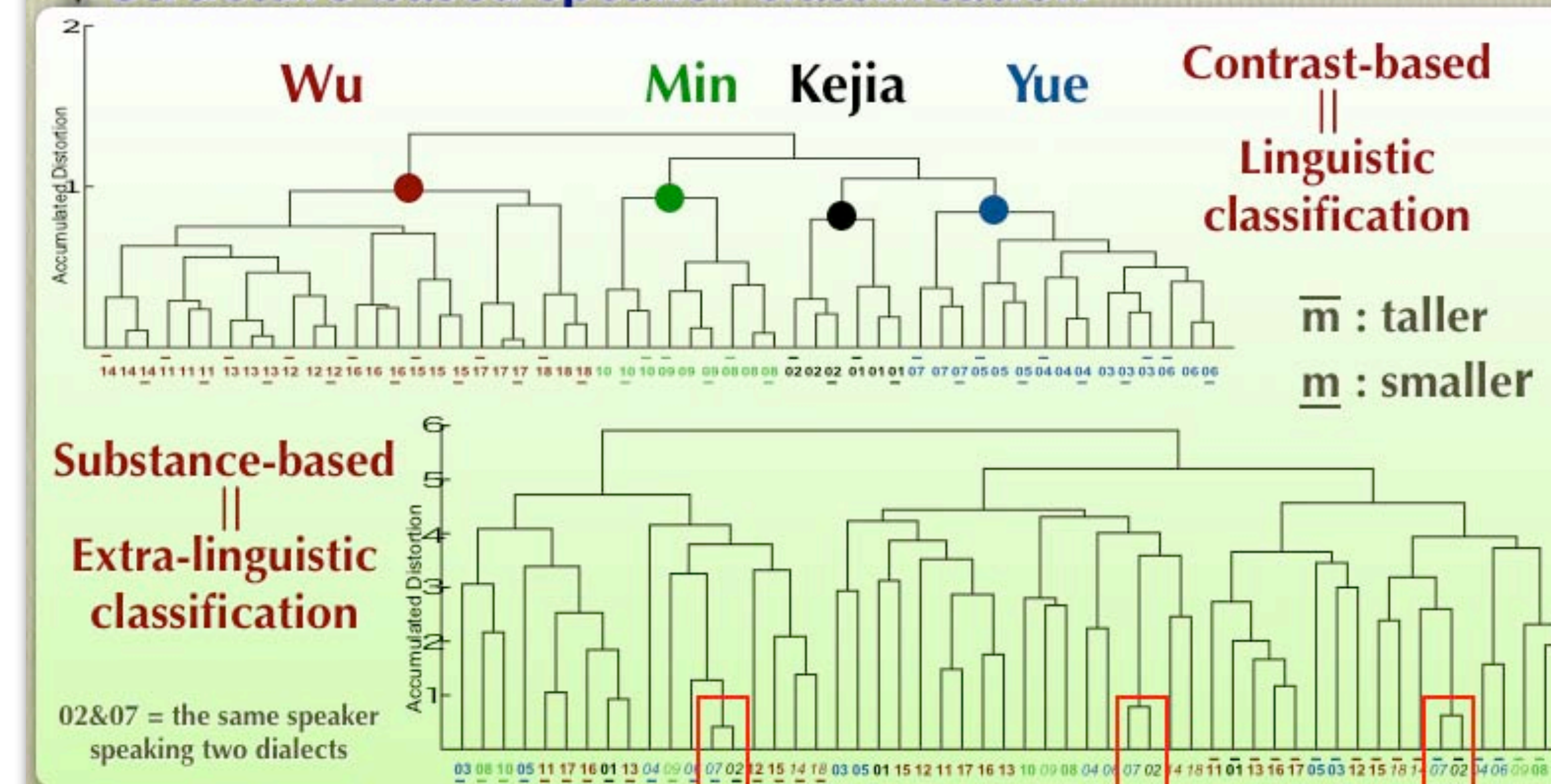  $$D_2(A, B) = \sqrt{\frac{1}{38} \sum_{i} BD(S_i^A, S_i^B)}$$

  $A, B$ = speakers, $A_{ij}$ = distance matrix for $A$
  $S$ = syllable, $i, j$ = syllable index

### Classification of Chinese speakers

- **Spectrum warping done to increase speaker variability**
  - Original speakers(m=18) → taller version(18) + smaller version(18)
- **Structure-based speaker classification**

  Wu   Min   Kejia   Yue

  Contrast-based = Linguistic classification

  $\overline{m}$ : taller
  $\underline{m}$ : smaller

  Substance-based = Extra-linguistic classification

  02&07 = the same speaker speaking two dialects

### The ultimate goal of these studies

  G00gle Pronunciation!?