

Structure-based Pronunciation Assessment

Nobuaki Minematsu, Masayuki Suzuki @ The University of Tokyo

Assessment of Japanese Learners' Pronunciation of English Vowels

!! A big problem !!

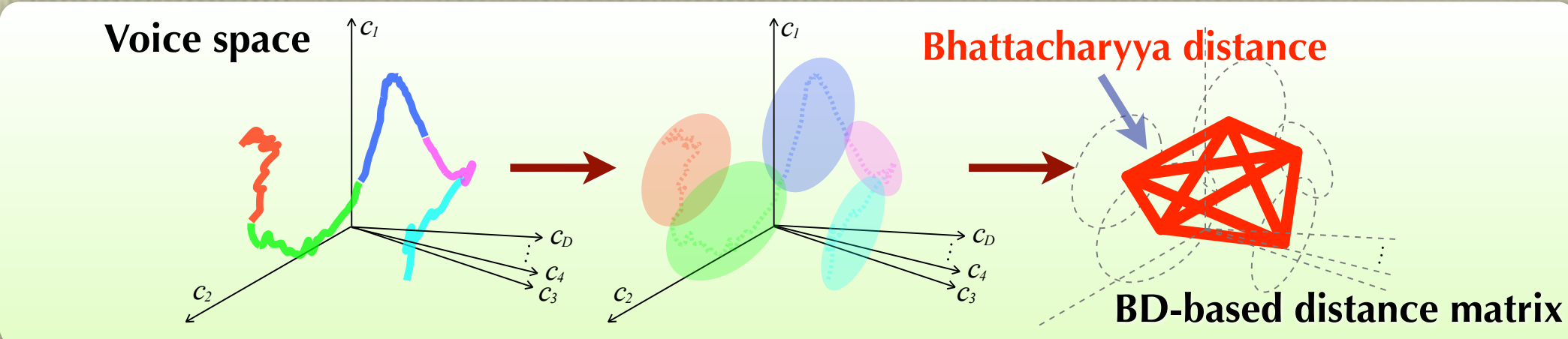
A very important and requisite function for CALL systems

- The system has to be able to ignore speaker individualities.
 - Age and gender (the size and length of the vocal tube)
 - But no current system can ignore speaker individualities well enough.
- Requirement of "no acoustic mismatch" bet. HMMs and learners
 - Collection of children's speech or speaker adaptation of adult HMMs
 - Q: Learning to pronounce is learning to impersonate?



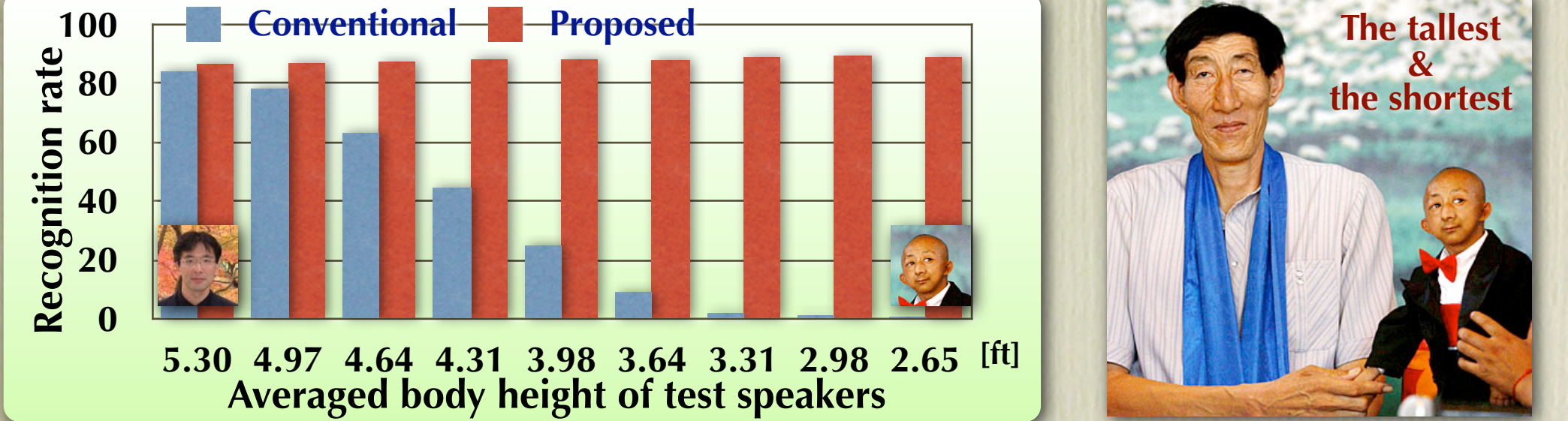
!! Our solution for that problem !!

Holistic and speaker-invariant sound pattern (structure)



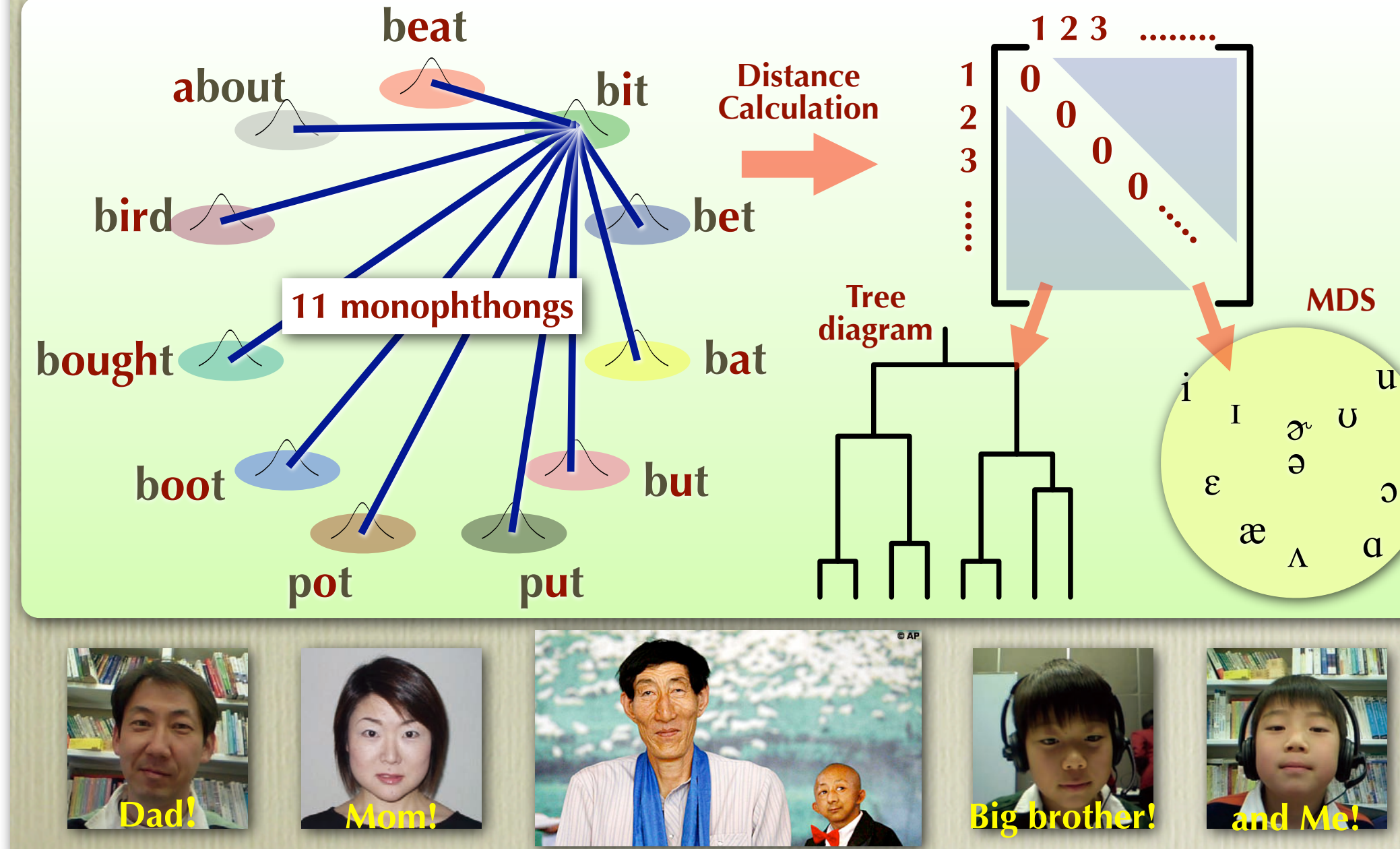
Use of structures for automatic speech recognition

- Isolated word recognition (word = 5 vowel sequence, e.g. /aeoui/)



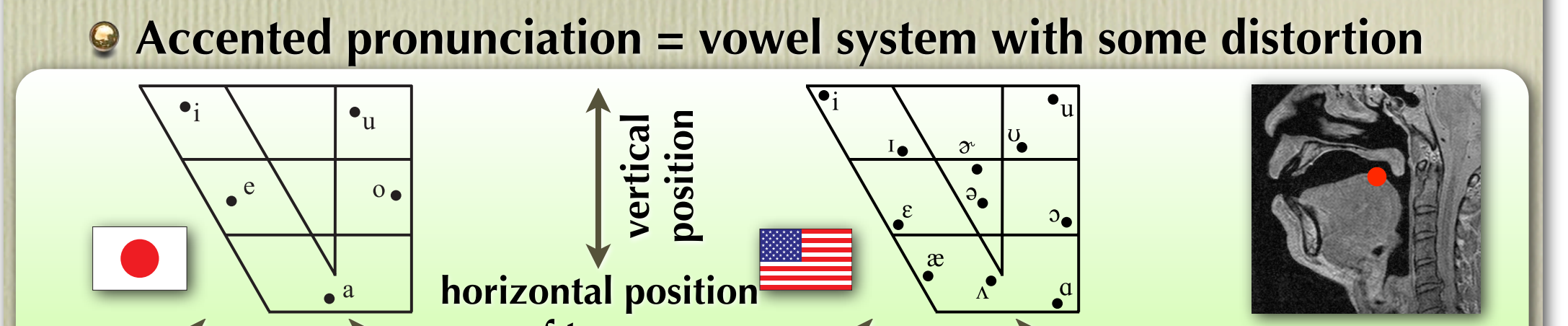
A vowel training system for everybody!!

Learning not of individual vowels but of a vowel system



Structural representation of the vowel system

Vowel system and vowel chart

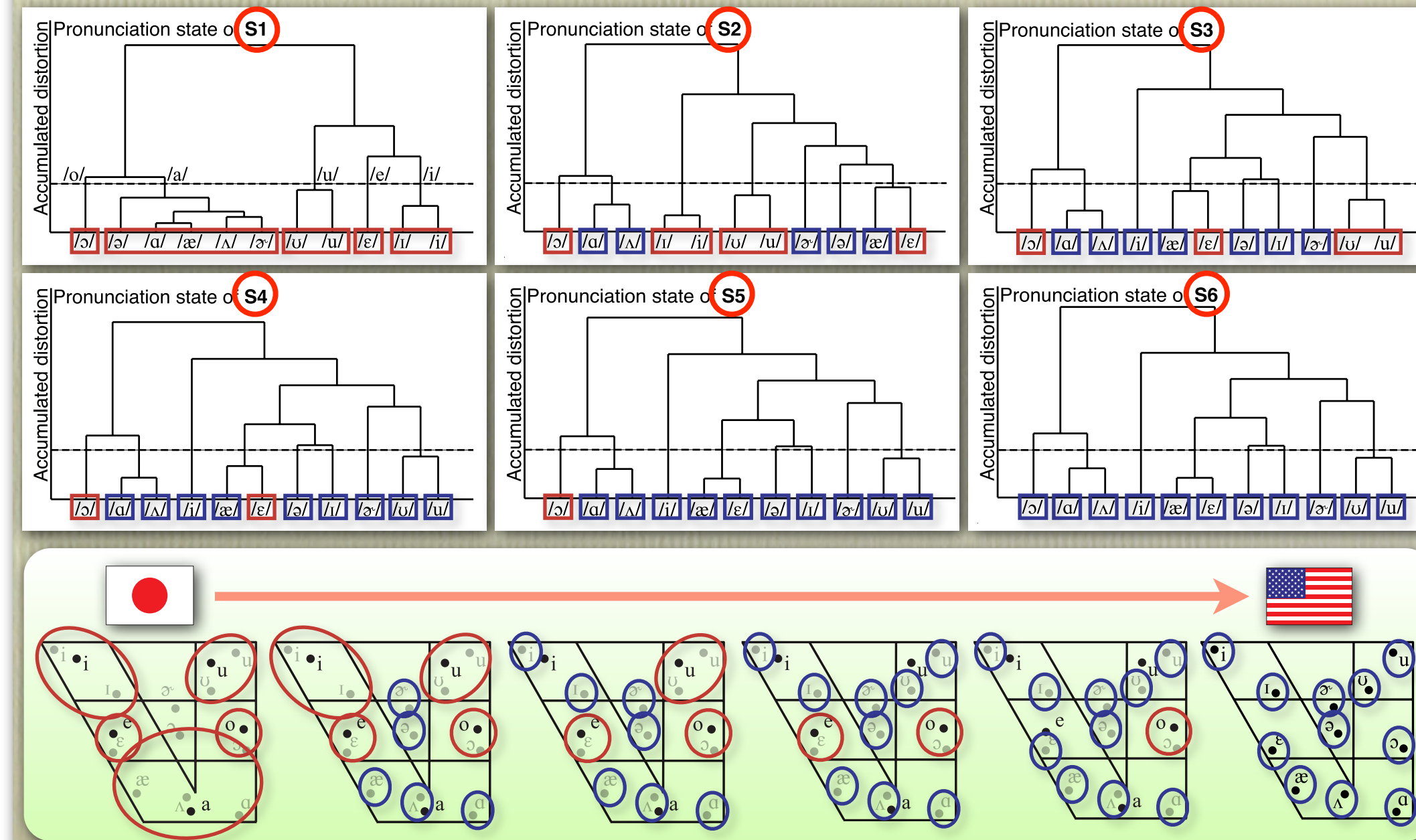


What's possible in the proposed demo system

- 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

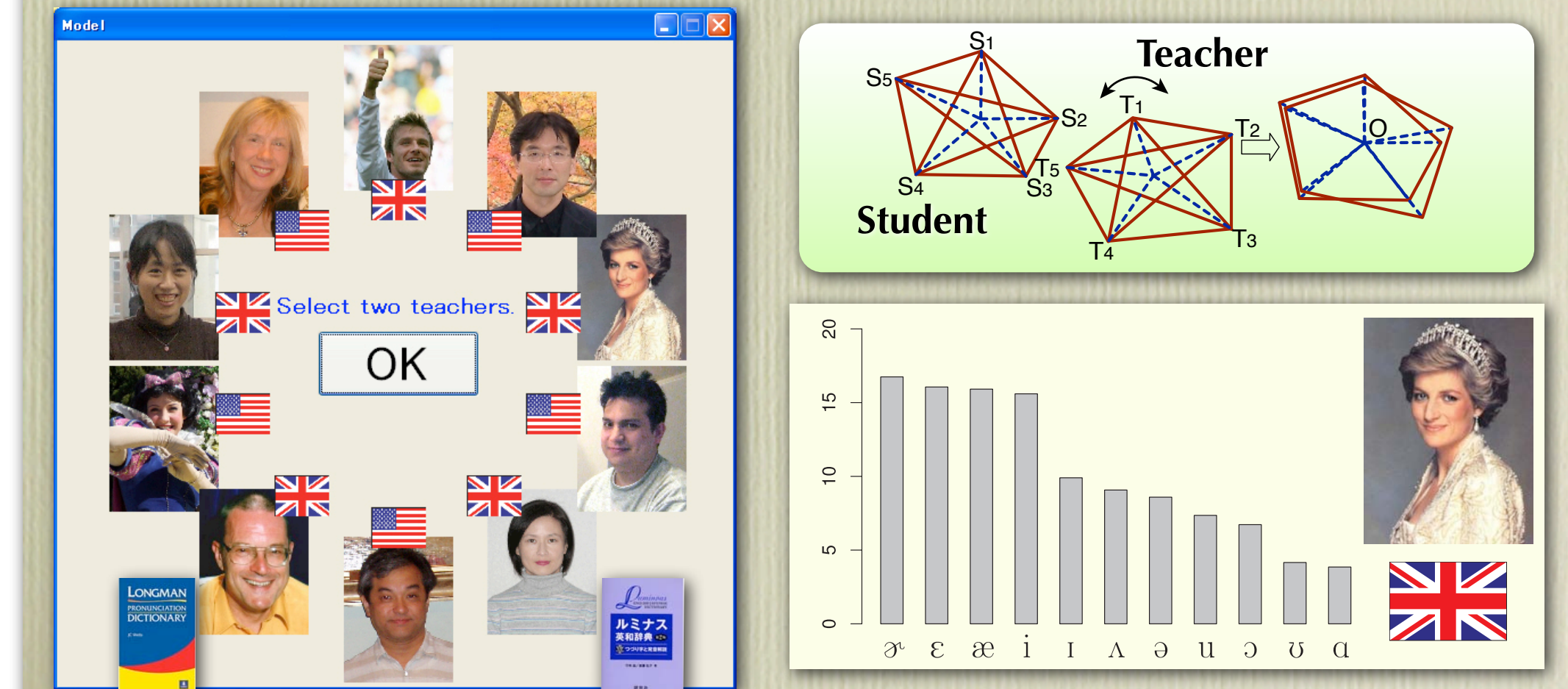
Completely Japanized pronunciation to AE pronunciation



Which vowels to correct at first in your case?

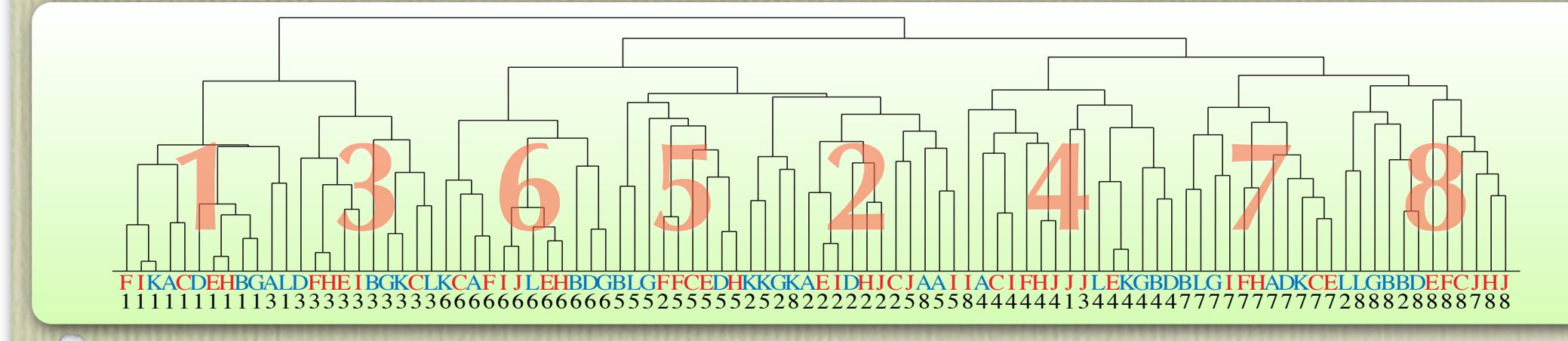
Who is your model speaker?

- 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.

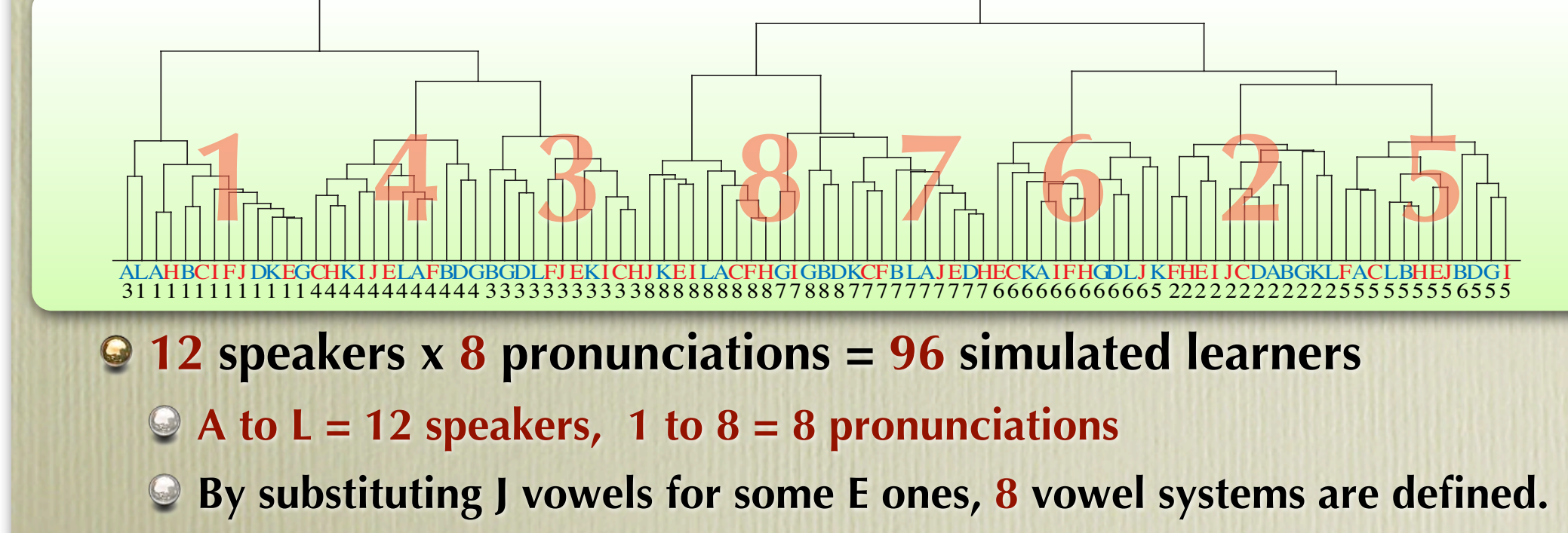


Classification of learners

Automatic classification of 96 learners by a computer

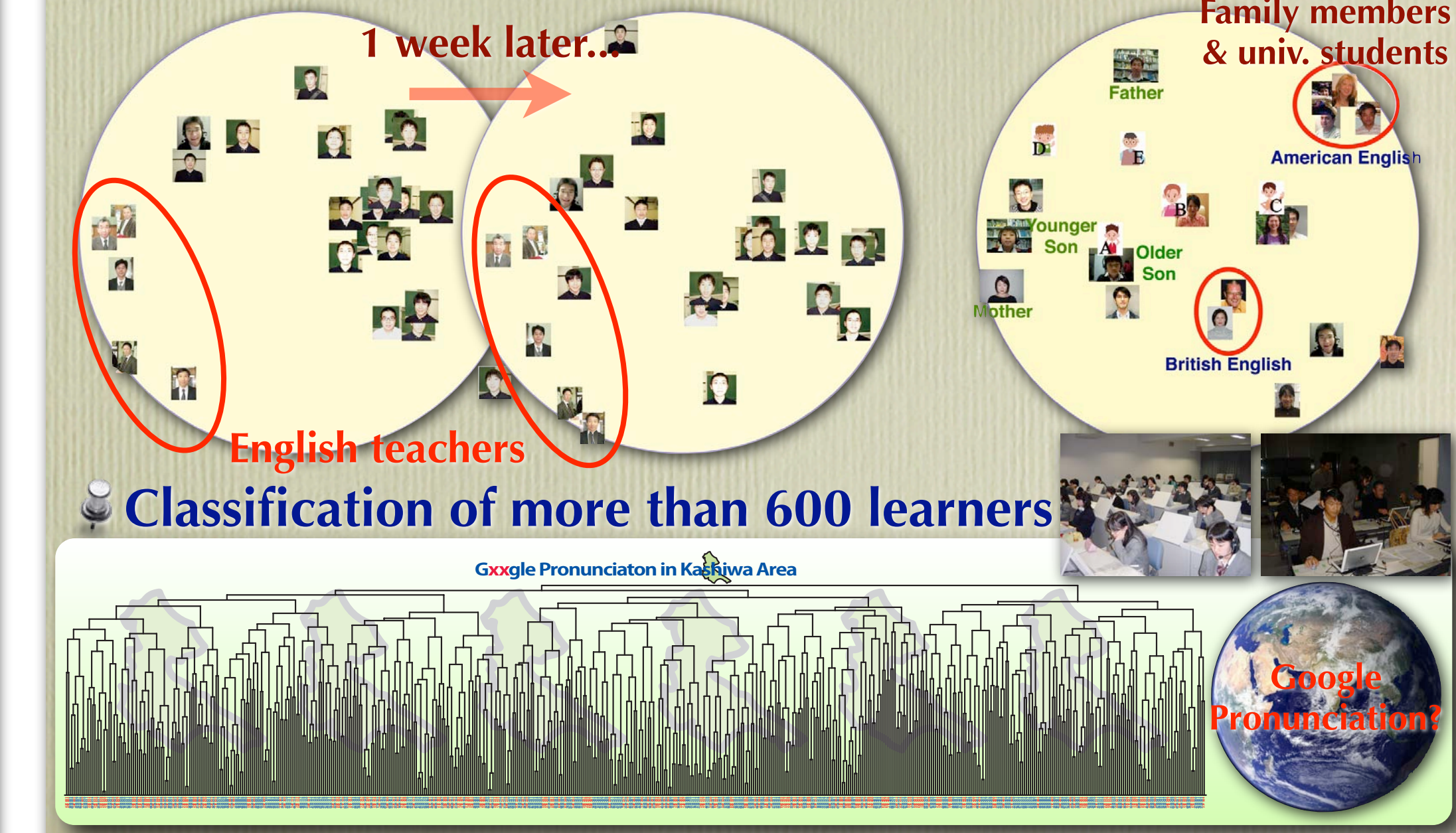


Manual classification of 96 learners by a phonetician



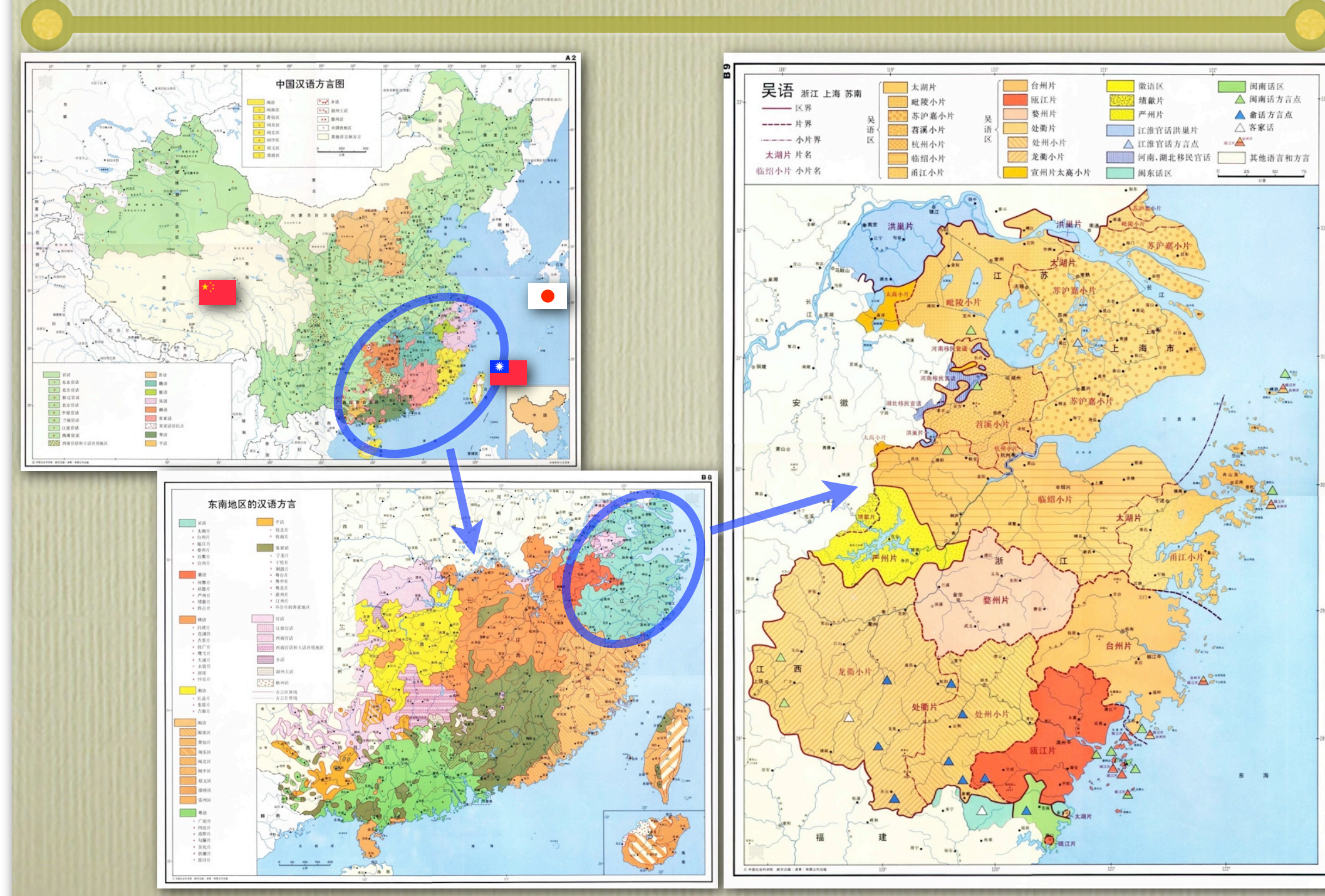
Classification of all the learners on earth!?

Changes of students in a class before and after training



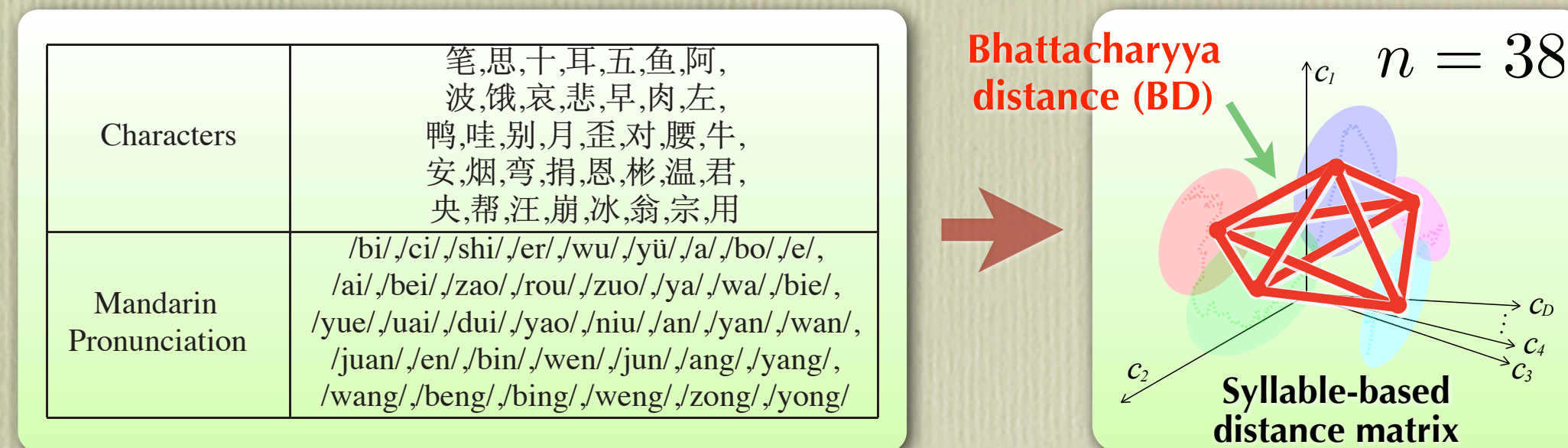
Dialect- and Sub-dialect-based Speaker Classification of Chinese

Chinese dialects & sub-dialects



Pronunciation structure of the 38 syllables

38 syllables covering all the 38 Mandarin finals



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

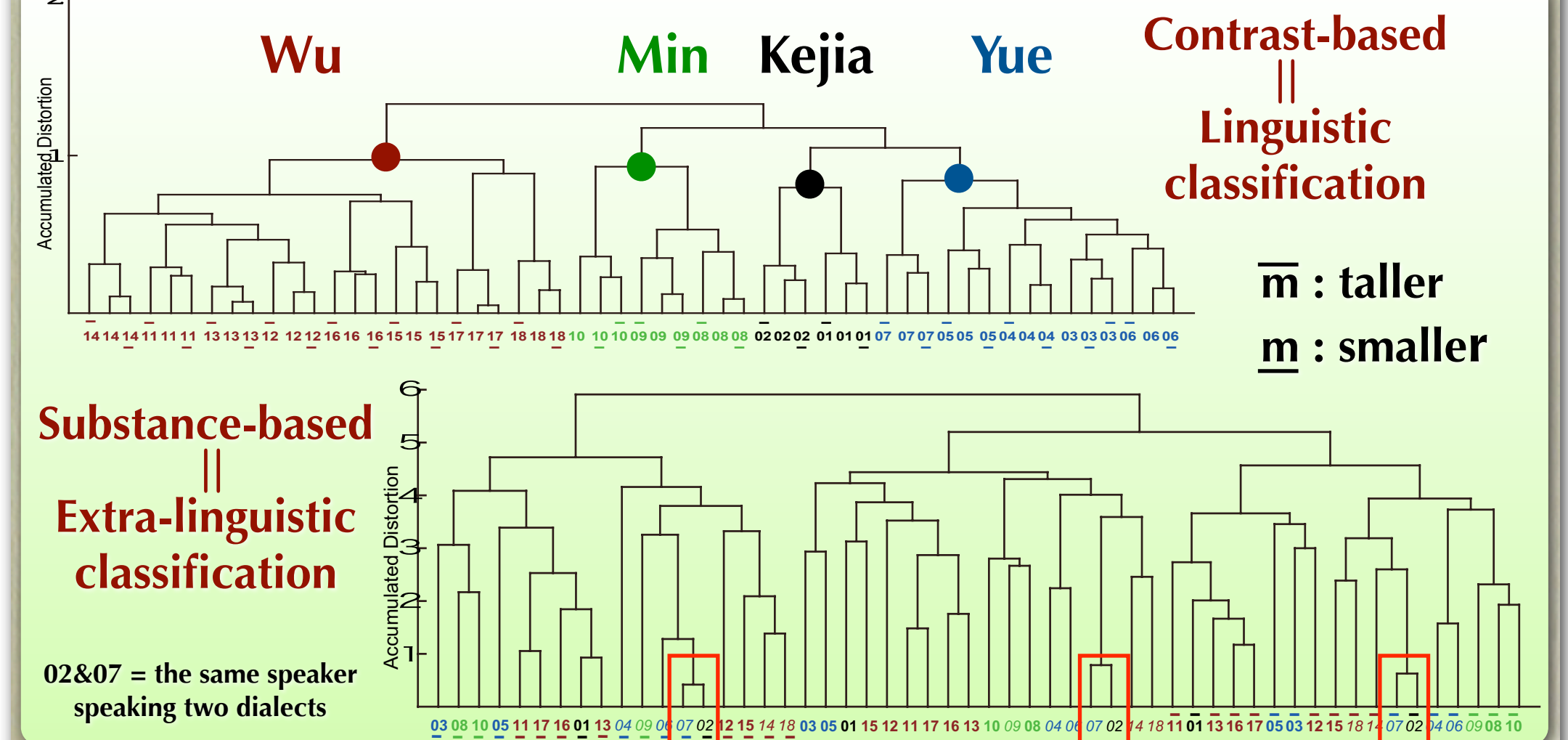
A, B = speakers, A_{ij} = distance matrix for A
 S = syllable, i, j = syllable index

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



Which syllables are more different?

Comparison between Min and Mandarin

