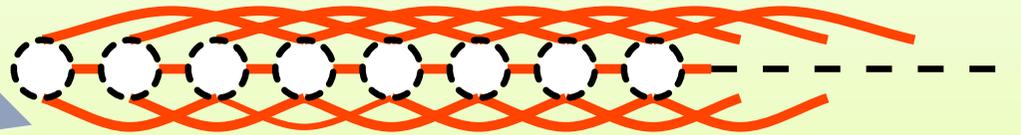
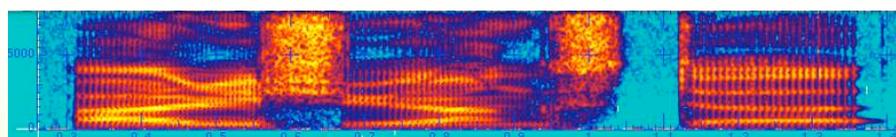
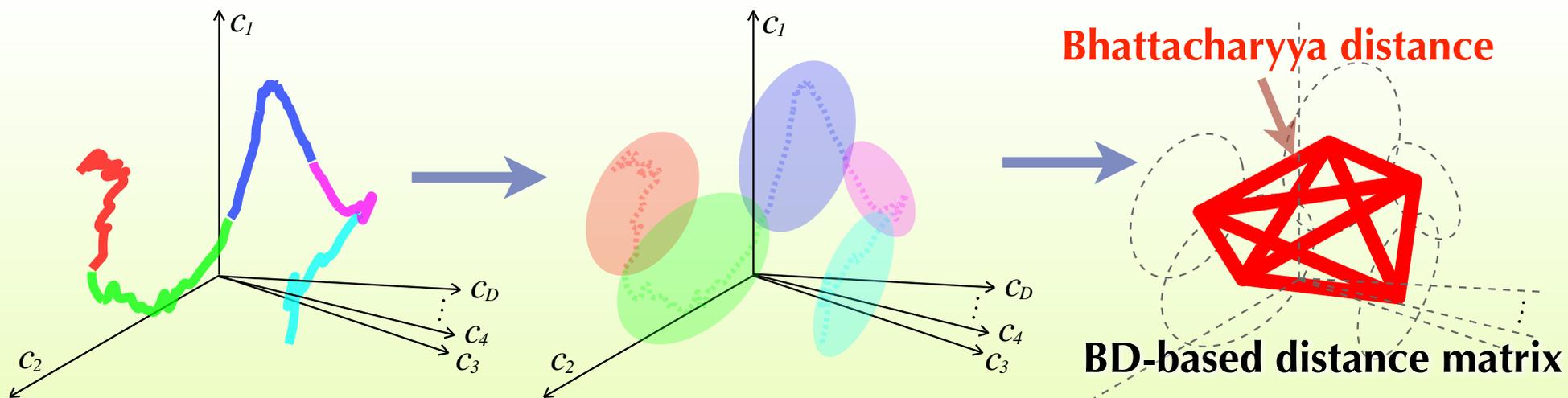


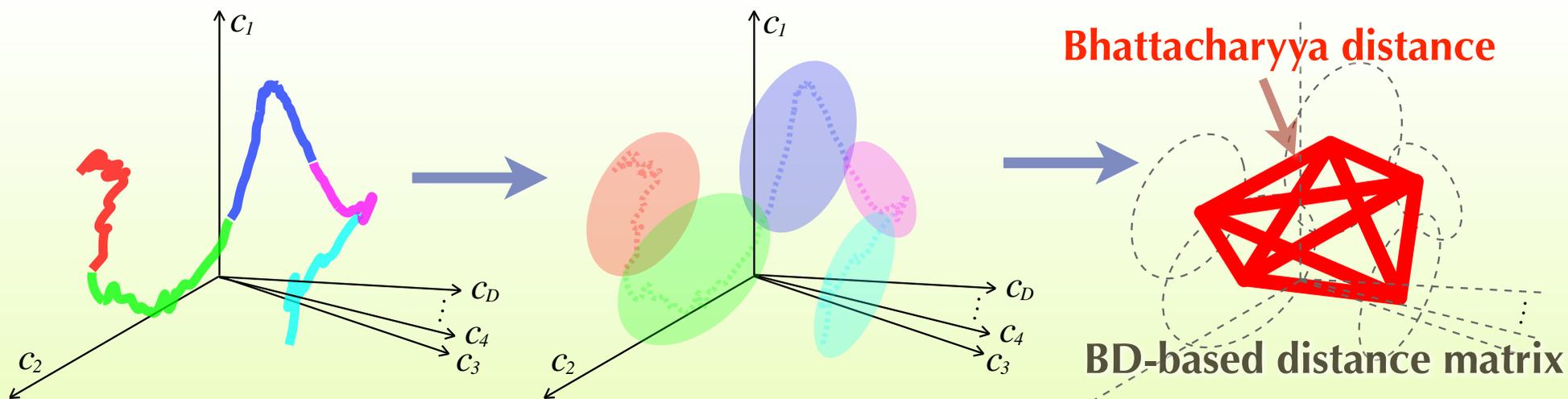
音声の構造的表象の工学的・実験的検証

f-div. (BD)に基づく一発声の構造化



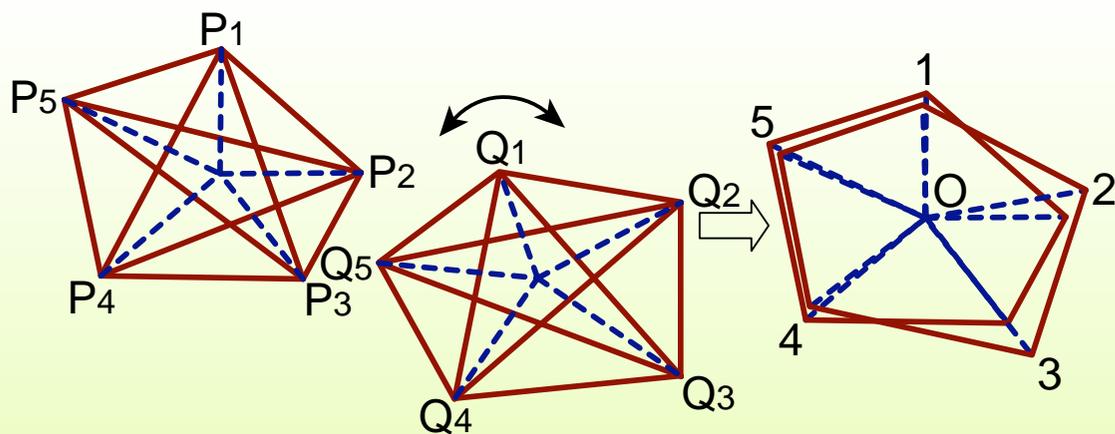
音声の構造的表象の工学的・実験的検証

f-div. (BD)に基づく一発声の構造化



2発声 (= 2距離行列) 間の音響照合

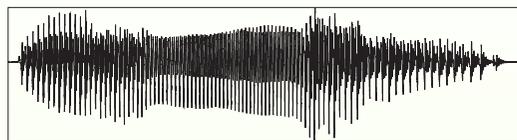
2距離行列間のユークリッド距離



- 回転：声道長差異
- シフト：マイク差異
- 話者適応・環境適応後のスコアが適応処理無しで算出
- 話者性を削除した音声表象

音声の構造的表象の工学的・実験的検証

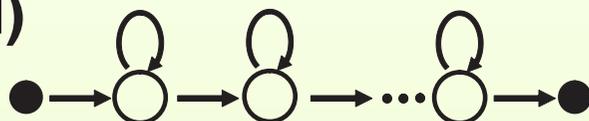
Speech signal



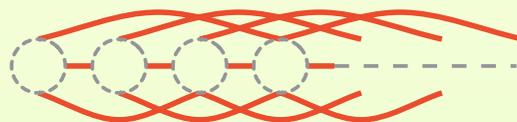
Cepstrum vector sequence



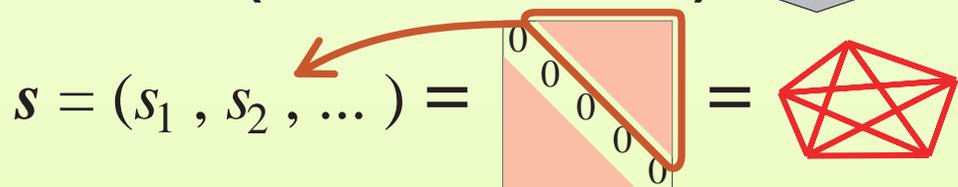
Cepstrum distribution sequence (HMM)



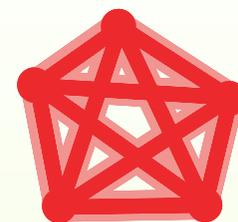
Distances of distributions



Structure (distance matrix)



Statistical structure model

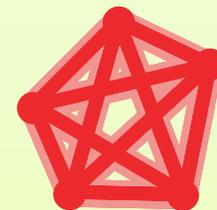


Word 1



Word 2

⋮

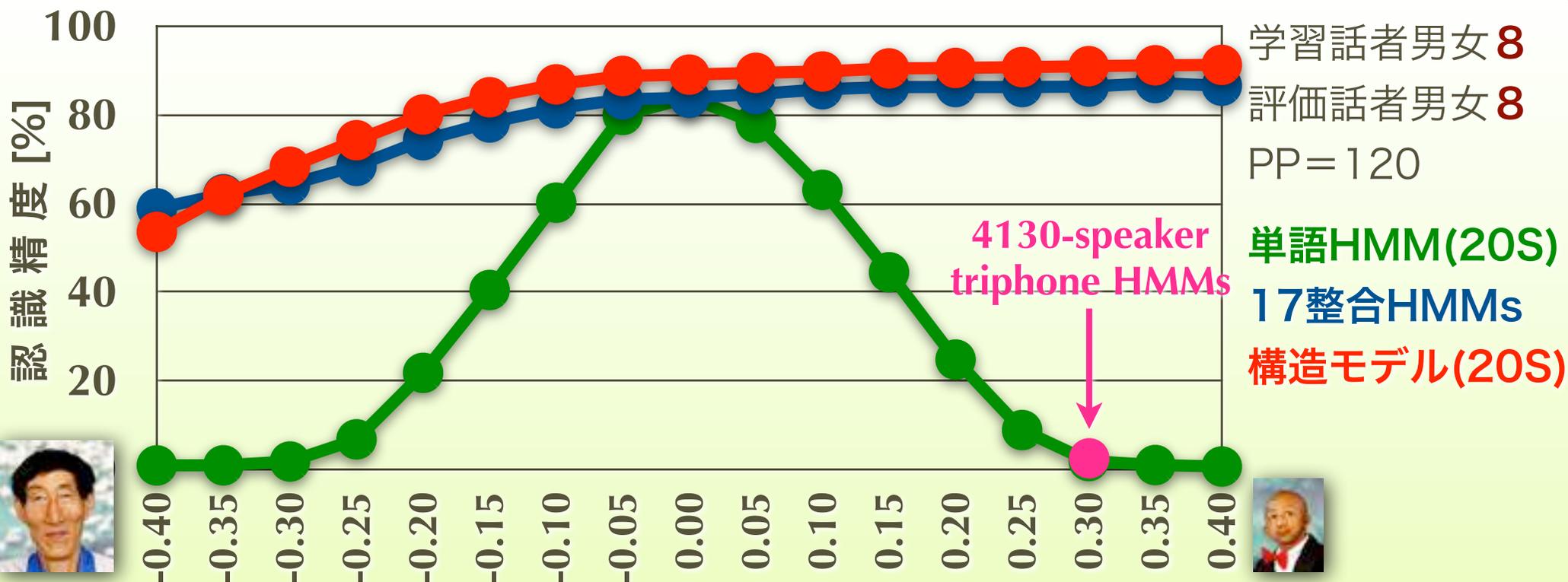


Word N

音声の構造的表象の工学的・実験的検証

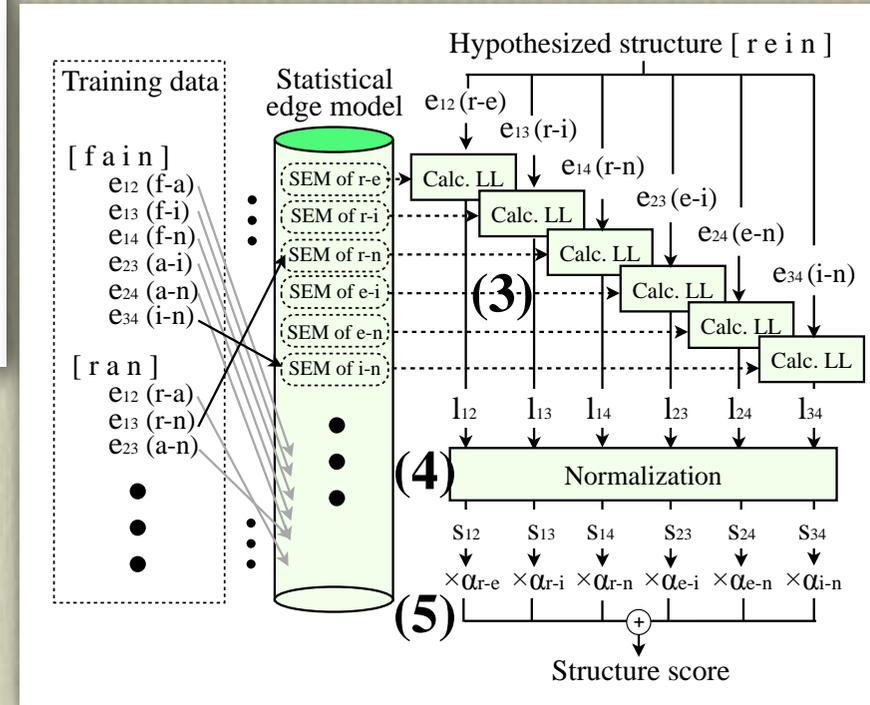
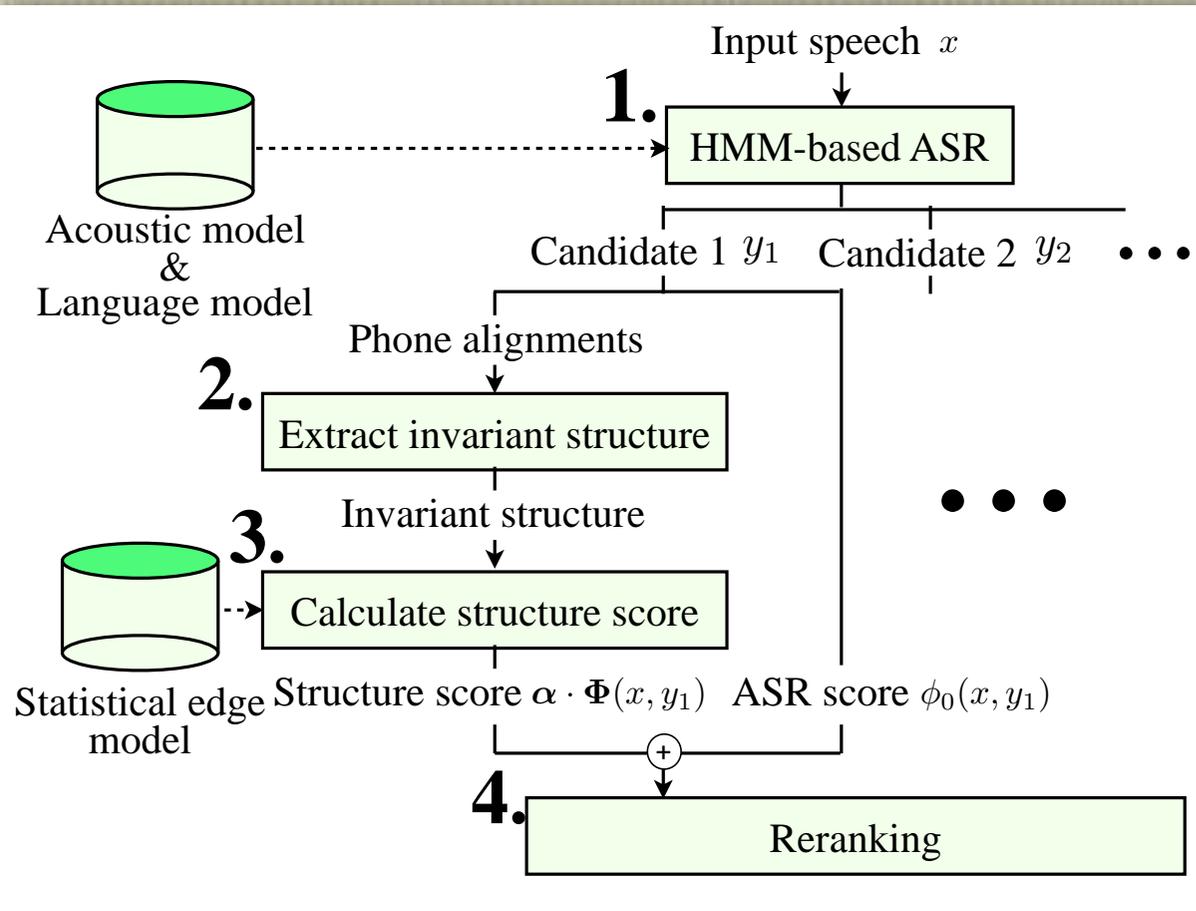
孤立単語音声の認識実験

- 二つの問題とその解決
 - 強すぎる不変性→マルチストリーム構造化による都合のよい不変性へ
 - 高すぎる次元数→線形判別分析 (LDA) による次元数削減
- 孤立単語認識実験による提案手法の評価
 - 日本語五母音を並び替えて作成される120単語の孤立単語認識



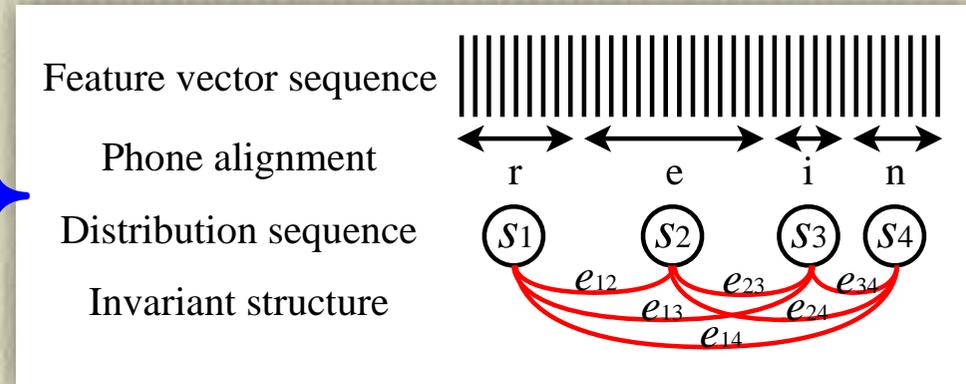
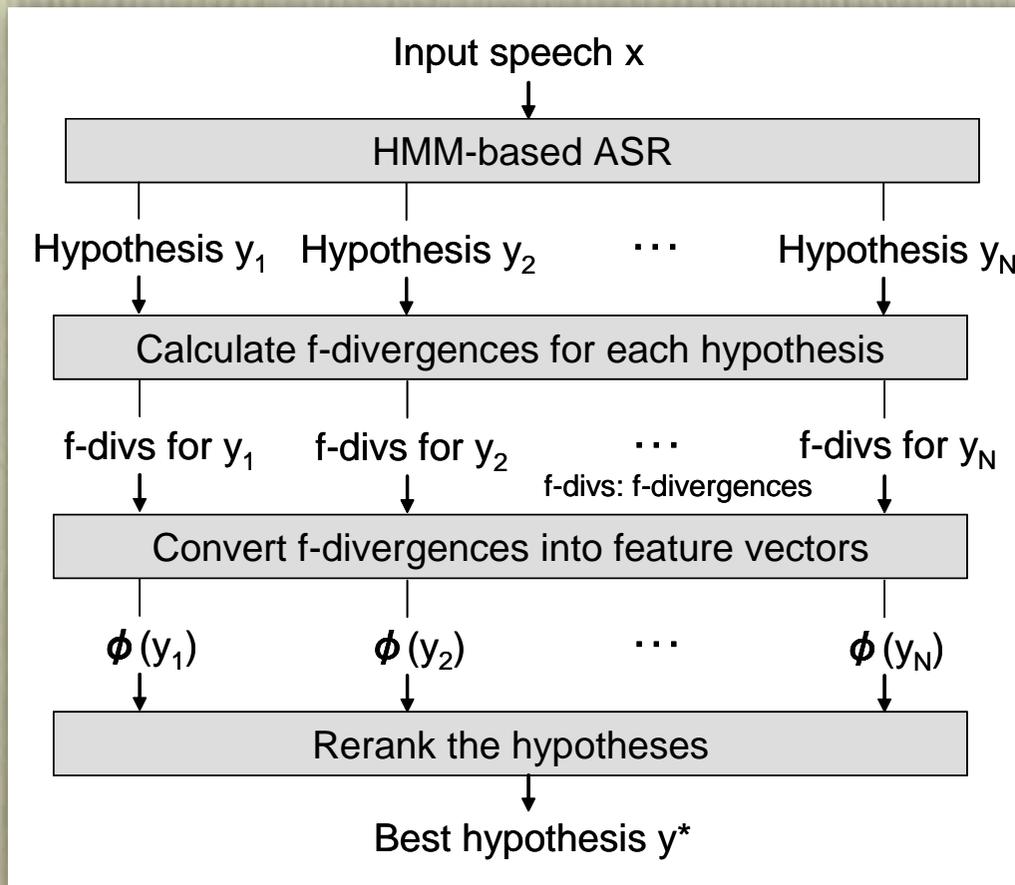
大語彙連続音声認識への応用

構造表象を複数仮説のリランキング処理に応用





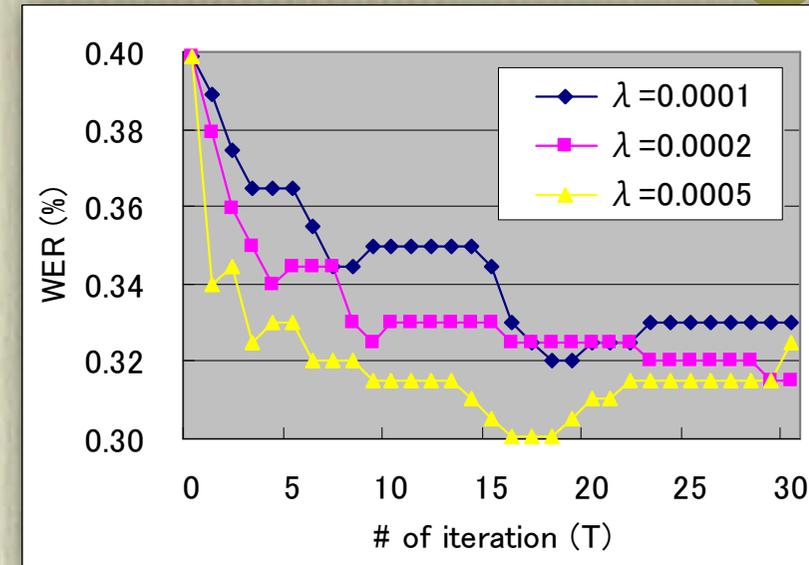
- Application to more realistic ASR tasks [Suzuki+'15]
 - Digits recognition and LVCSR (dictation)
 - Use of structural features in discriminative reranking**
 - Str. scores and ASR scores are combined with average perceptron.





Continuous digits recognition

- Language = Japanese
- Baseline = GMM-HMM ASR
- Reranking = averaged perceptron
- Error reduction rate = 30%



Large vocabulary continuous speech recognition

- Language = Japanese
- Baseline = DNN-HMM ASR
- Reranking = averaged perceptron
- Error reduction rate = 5%

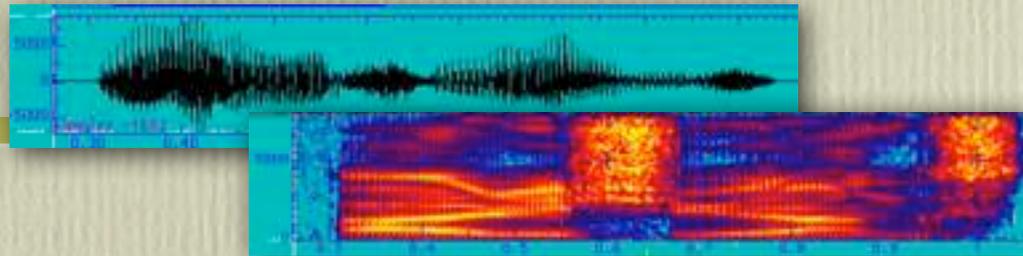
Many errors are due to a large number of homonyms in Japanese.

Table 6: CERs of the LVCSR experiment.

Baseline	Proposed	Relative improvement
2.67%	2.53%	5.24%

音響音声学

(Topics in Acoustic Phonetics)



峯松 信明

工学系研究科電気系工学専攻

本発表の流れ

刺激の物理的多様性とその認知的不変性

- 見え／色み／音高の多様性と自然・進化が編み出した解決方法

音声の物理的多様性とその認知的不変性

- 音色の多様性と工学者が編み出した解決方法

音声の構造的表象とそれに関する様々な考察

- 常識を覆すことで、違和感の解消を試みしてみる。

音声の構造的表象と数学的表現と技術的実装

- 体格・性別に不変な音声波形・スペクトルの表現とは？

音声の構造的表象を用いた音声アプリケーション

- 音声認識, 音声合成, 発音分析, etc

音声の構造的表象の言語学的妥当性

- 何故, こうしてこなかったのか? 観測技術の功罪?

そもそもの出発点はこちら

学習者の英語発音を評価する技術の構築

- 発音の習熟度を定量的にスコア化する。
- 個々の発声における、不適切な発音部位を検出する。

発音評価技術構築における根本的な問題

- 教師音声データベースから個々の音素の音響モデルを構築する。
- 学習者は小学生かもしれない。
- 音韻の違い = 話者の違い = 管形状の違い = 音色の違い

学習者にあった音響モデルの構築

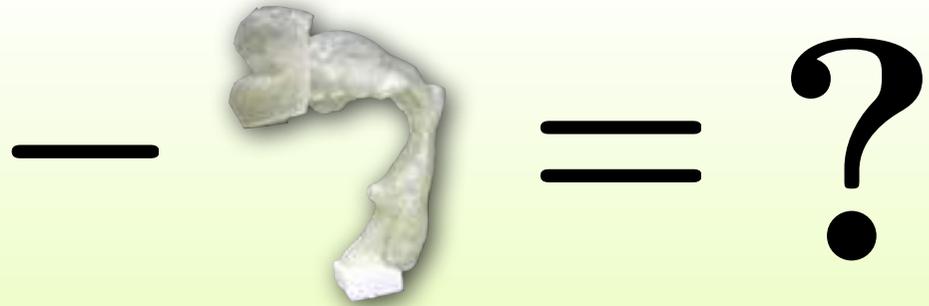
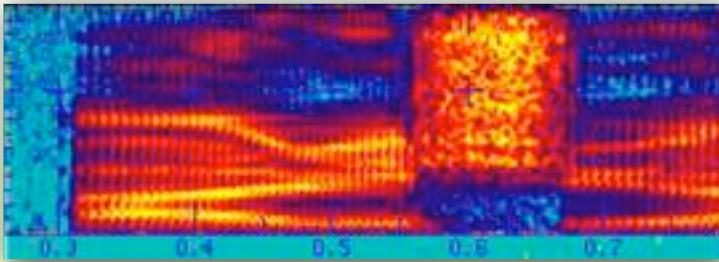
- 話者適応技術
 - 大人用モデル / 子供用モデル
- 発音評定？ 声帯模写評定？

峯松の声から峯松を消したい

- そんなことができるのか？



Mismatch
problem



- 個々の発声における，不適切な発音部位を検出する。

● 発音評価技術構築における根本的な問題

- 教師音声データベースから個々の音素の音響モデルを構築する。
- 学習者は小学生かもしれない。
- 音韻の違い = 話者の違い = 管形状の違い = 音色の違い

● 学習者にあった音響モデルの構築

- 話者適応技術
 - 大人用モデル / 子供用モデル
- 発音評定？ 声帯模写評定？

● 峯松の声から峯松を消したい

- そんなことができるのか？



**Mismatch
problem**

で、こういうのができました。

峯松の普通の英語発音はどちらに近いのか？

speaker	USA/F12	✗	Minematsu	○	Minematsu
gender	female	✗	male	○	male
age	??	✗	36	○	36
mic	Sennheiser	✗	cheap mic	○	cheap mic
room	recording room	✗	living room	○	living room
AD	SONY DAT	✗	PowerBook	○	PowerBook
proficiency	perfect	△	good	✗	Japanized

Those answers will be straight forward if you think them through carefully first.

で、こういうのができました。

母語話者発音と日本語的峯松発音の音響モデルを利用

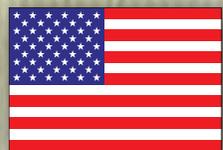
各モデルから峯松英語発声が観測される確率 $P(o|M)$



USA/F12



Minematsu
(Japanized)



USA/M08



Minematsu
(Japanized)



で、こういうのができました。

母語話者発音と日本語的峯松発音の音響モデルを利用

峯松発音が 所望の文章列の読み上げによる発音 $P(M|o)$

$$\begin{aligned} P(M|o) &= P(p_1, \dots, p_N|o) \\ &= \frac{P(o|p_1, \dots, p_N)P(p_1, \dots, p_N)}{\sum_{p_i} P(o|p_1, \dots, p_N)P(p_1, \dots, p_N)} \\ &\approx \frac{P(o|p_1, \dots, p_N)}{\sum_{p_i} P(o|p_1, \dots, p_N)} \\ &\approx \frac{P(o|p_1, \dots, p_N)}{\max_{p_i} P(o|p_1, \dots, p_N)} \\ &= \frac{P(o|M)}{\max_M P(o|M)} \\ &= \text{GOP (Goodness Of Pronunciation)} \end{aligned}$$



USA



USA



matsu
nized)



matsu
nized)

で、こういうのができました。

母語話者発音と日本語的峯松発音の音響モデルを利用

● 峯松発声が、所望の音素列の読み上げである確率 $P(M|o)$



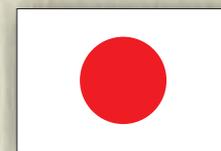
USA/F12



Minematsu
(Japanized)



USA/M08



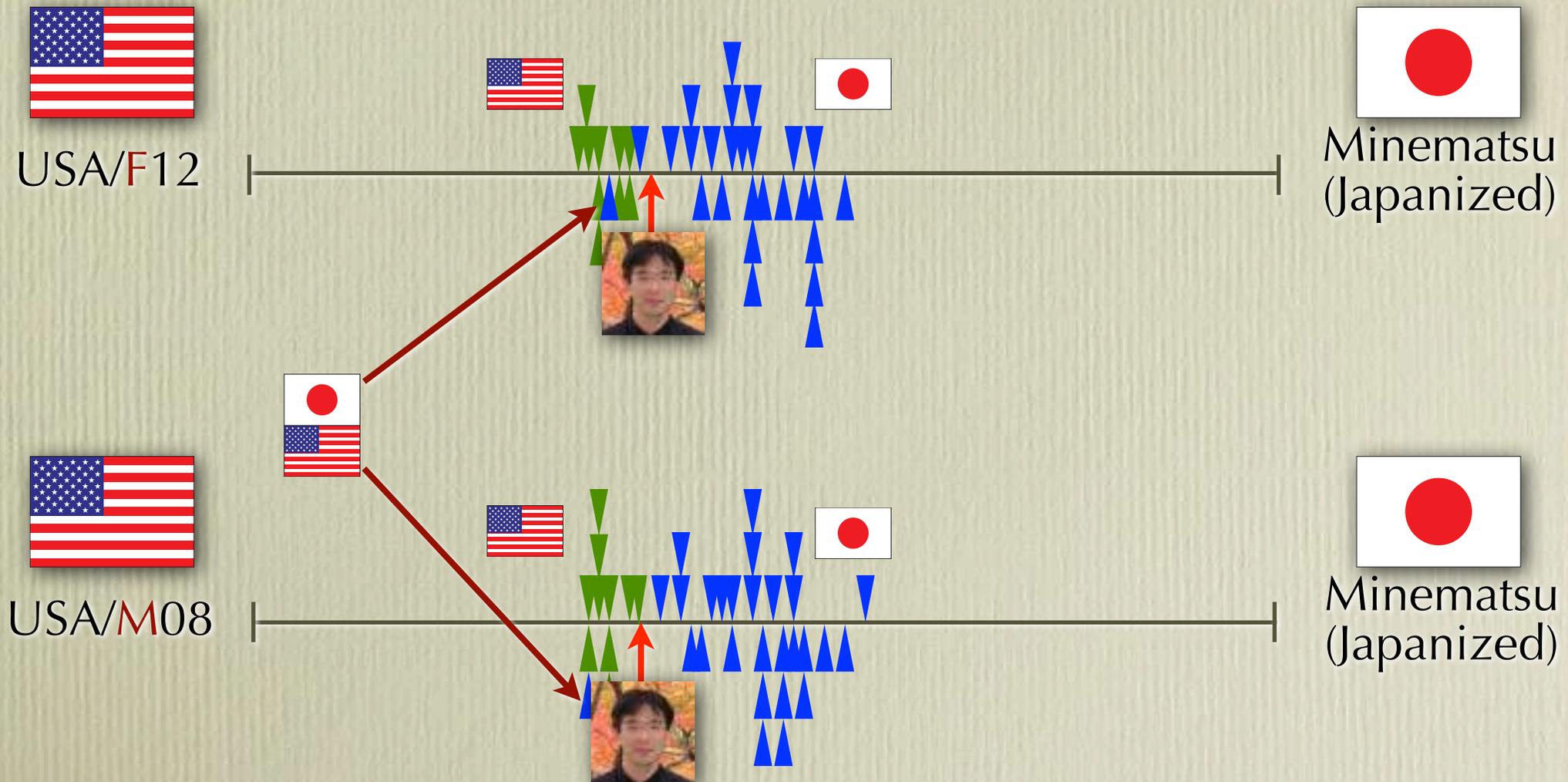
Minematsu
(Japanized)



で、こういうのができました。

母語話者発音と日本語的峯松発音の構造モデルを利用

母語発音構造, 日本語的峯松発音構造と, 峯松発音構造を比較



speaker	USA/F12	Minematsu	Minematsu
gender	female	male	male
age	??	36	36
mic	Sennheiser	cheap mic	cheap mic
room	recording room	living room	living room
AD	SONY DAT	PowerBook	PowerBook
proficiency	perfect	good	Japanized

proficiency

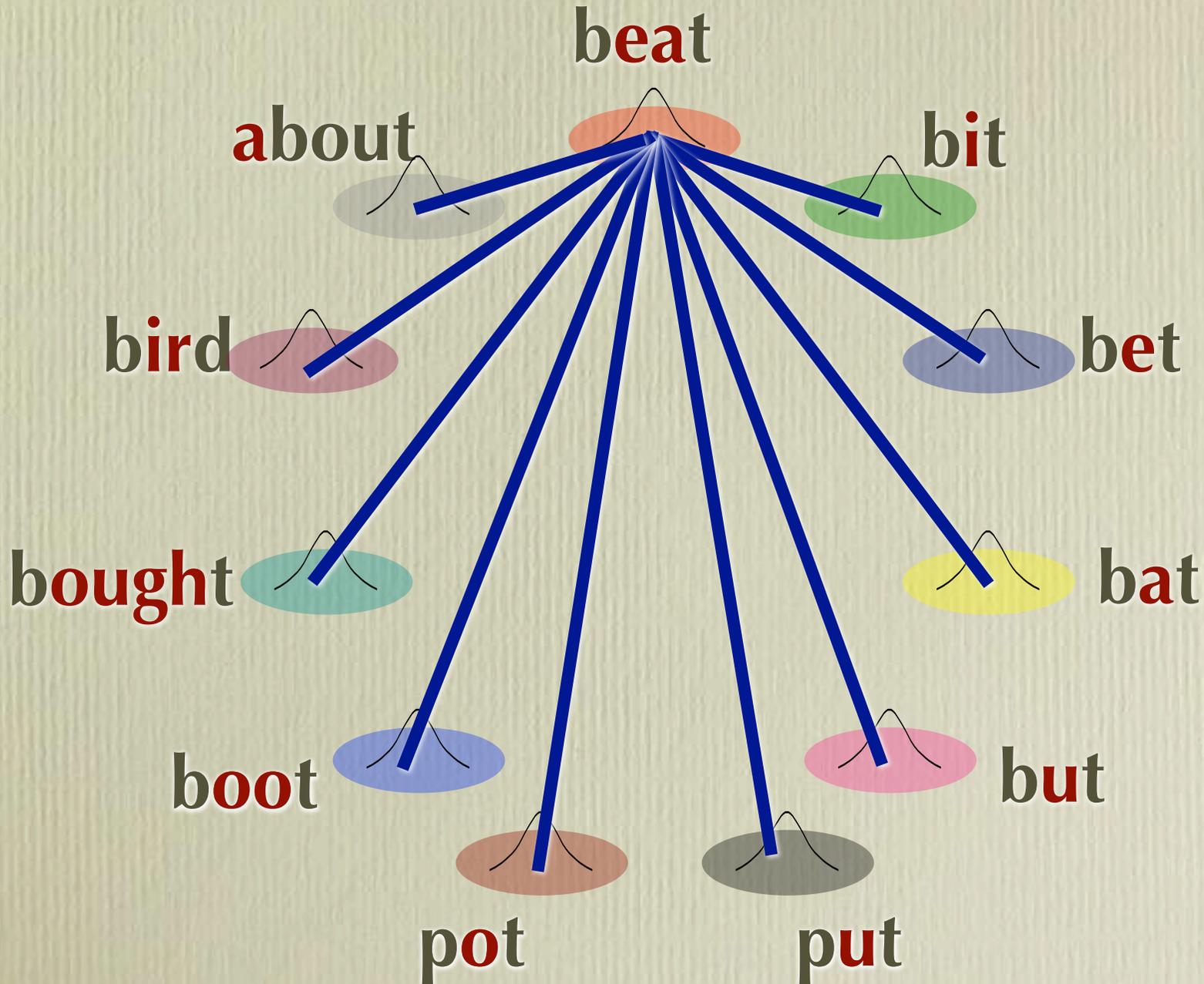
perfect

good

Japanized

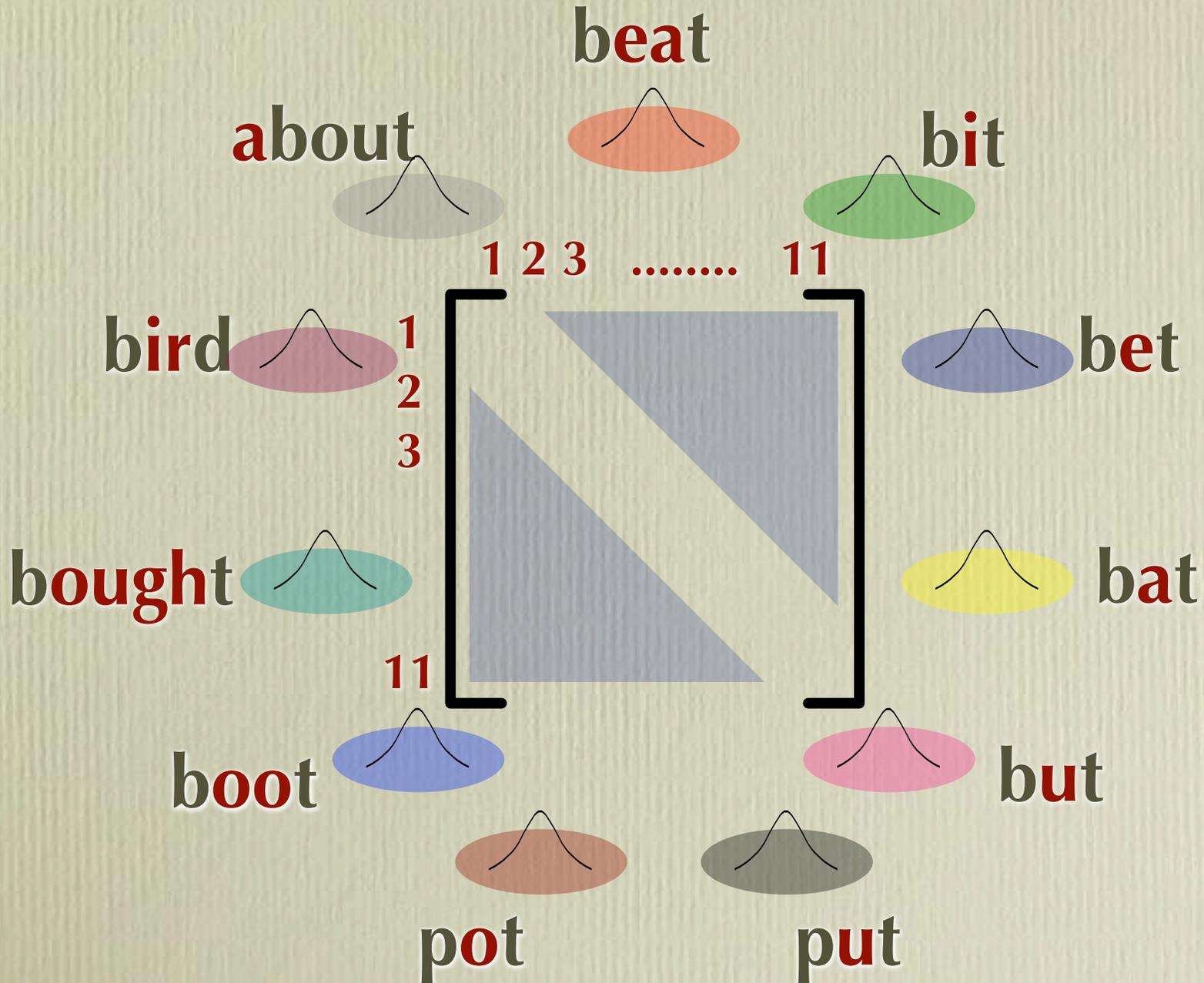
英語母音の発音トレーニング

個々の母音ではなく，母音群の体系を評価対象



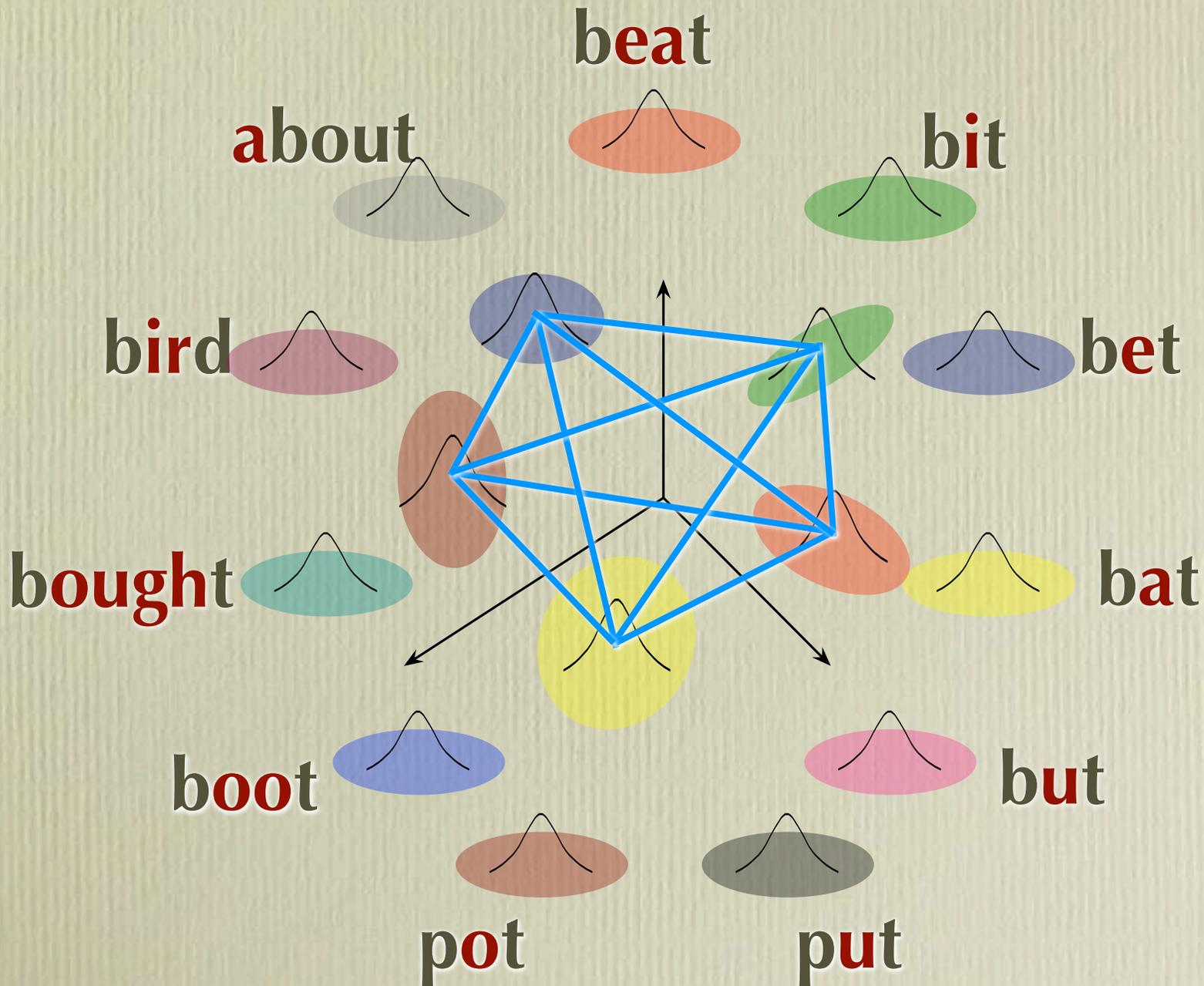
英語母音の発音トレーニング

個々の母音ではなく，母音群の体系を評価対象



英語母音の発音トレーニング

個々の母音ではなく，母音群の体系を評価対象



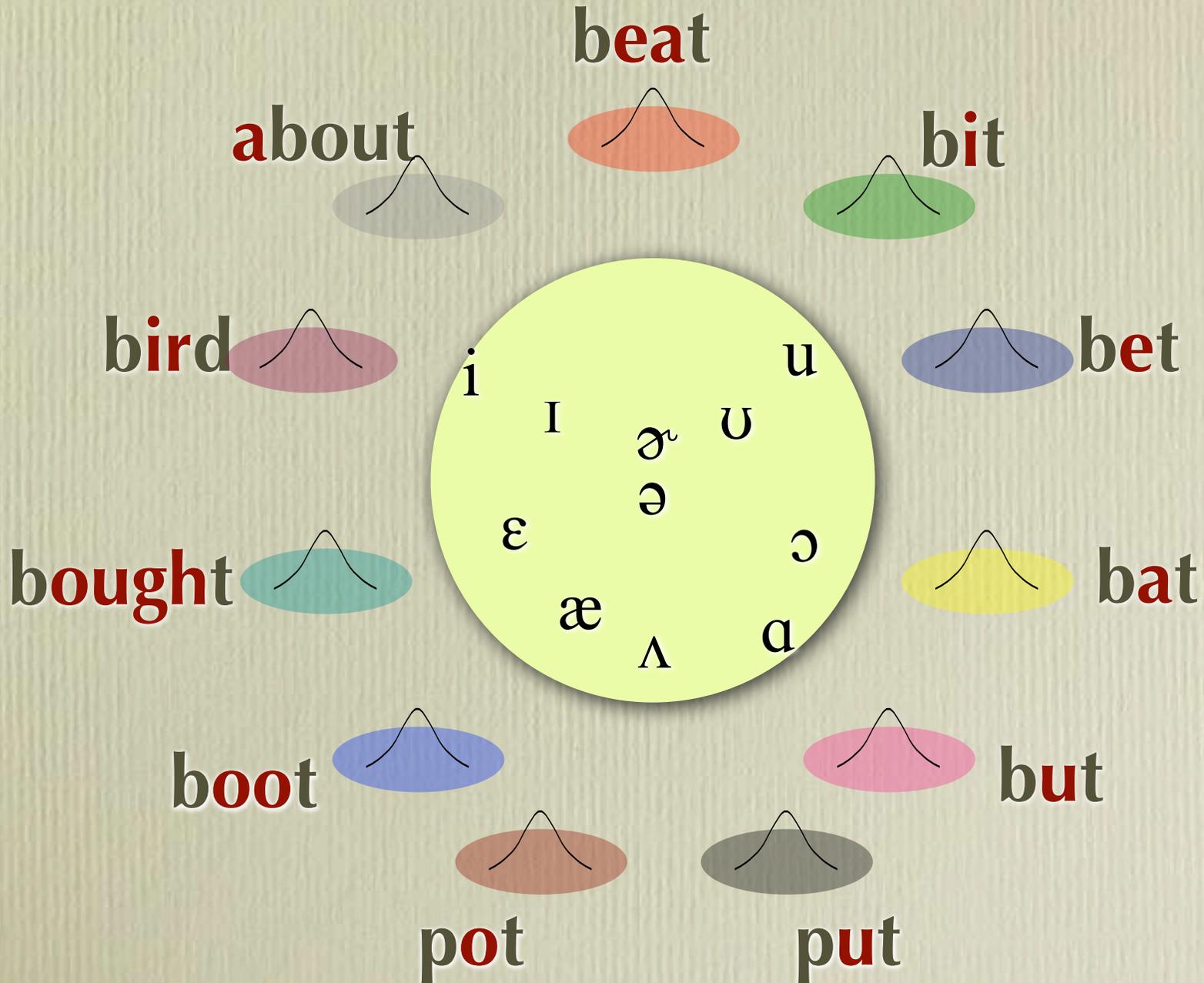
英語母音の発音トレーニング

個々の母音ではなく，母音群の体系を評価対象



英語母音の発音トレーニング

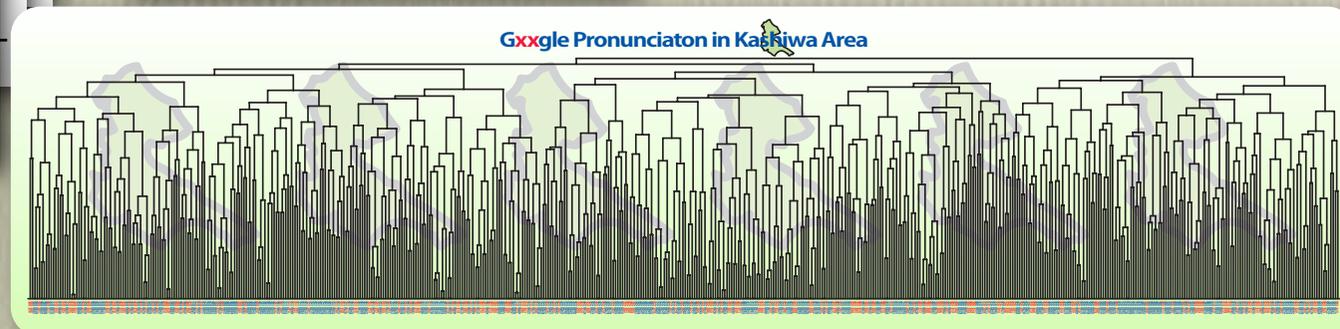
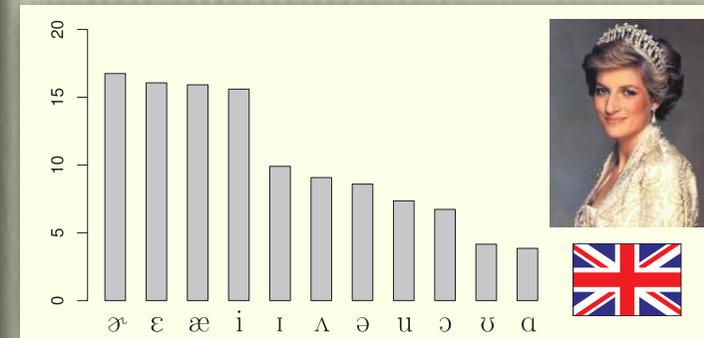
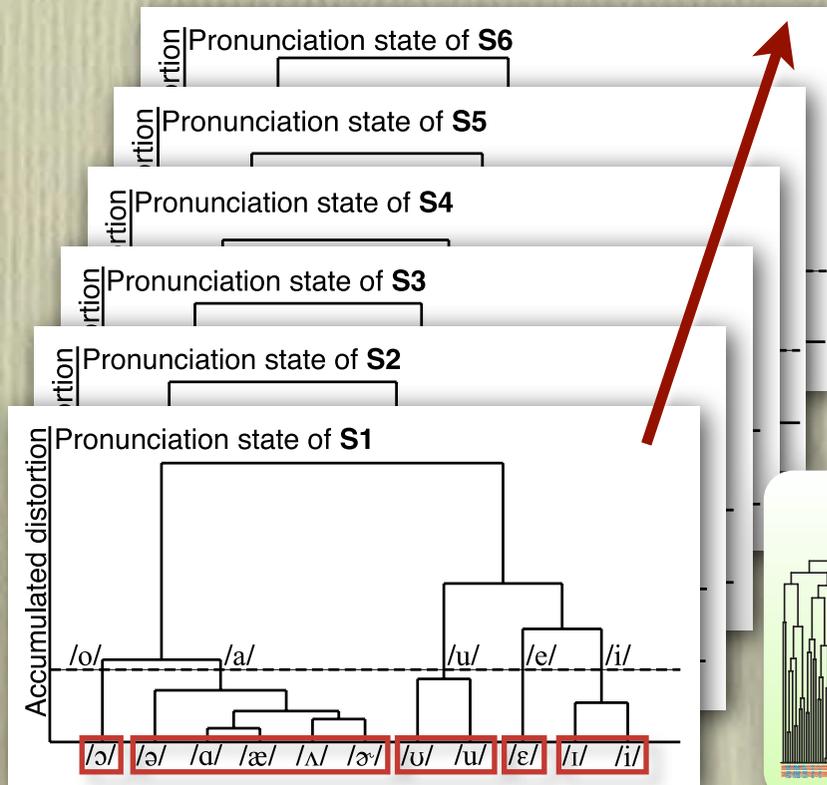
個々の母音ではなく，母音群の体系を評価対象



母音体系に見る発音習熟度

発音の構造表象に基づいて構築した技術・機能

- 学習者の母音体系がどのような状態にあるのかを記録する
- お手本の母音構造と比較して、矯正対象の母音を教示する
- 複数の教師からお手本としたい教師（相手）を選ぶ
- 年齢，性別などには影響されない学習者発音の分類



学習者の発音状態の記録とその遷移

分析対象の音声資料

● 英語劇経験者による英語&日本語単語音声

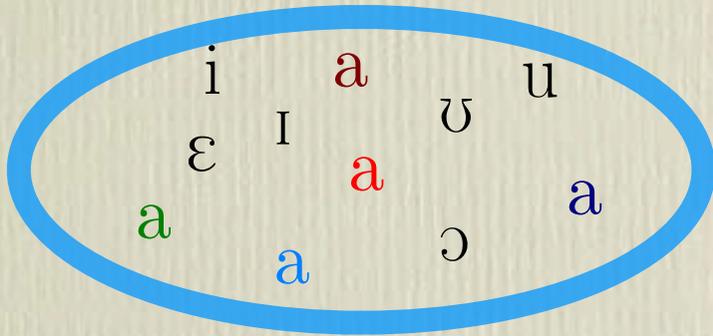
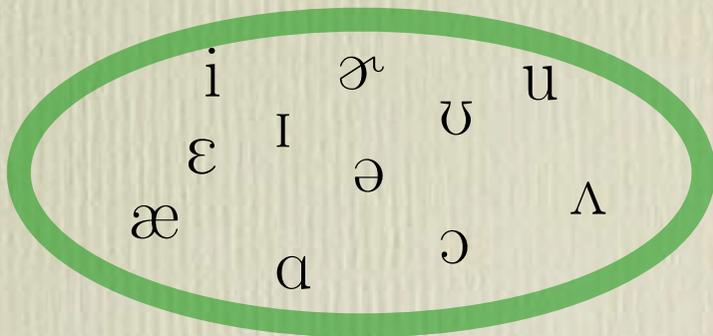
● beat, bit, bet, bat, but, pot, bought, boot, put, bird, about x 1

● ばと, びと, ぶと, べと, ぼと x 5

日本人学習者の母音発音の模擬的生成

● 幾つかの米語母音と日本語母音の置き換え

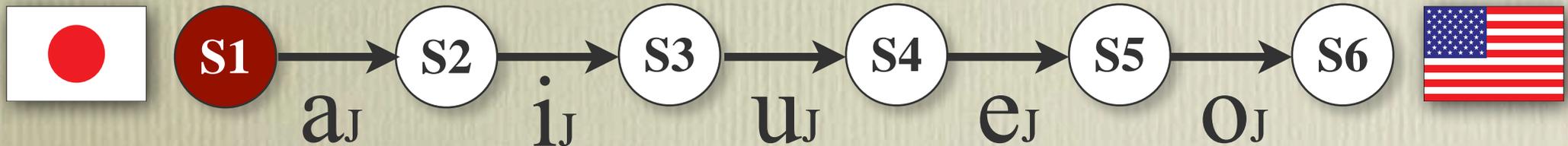
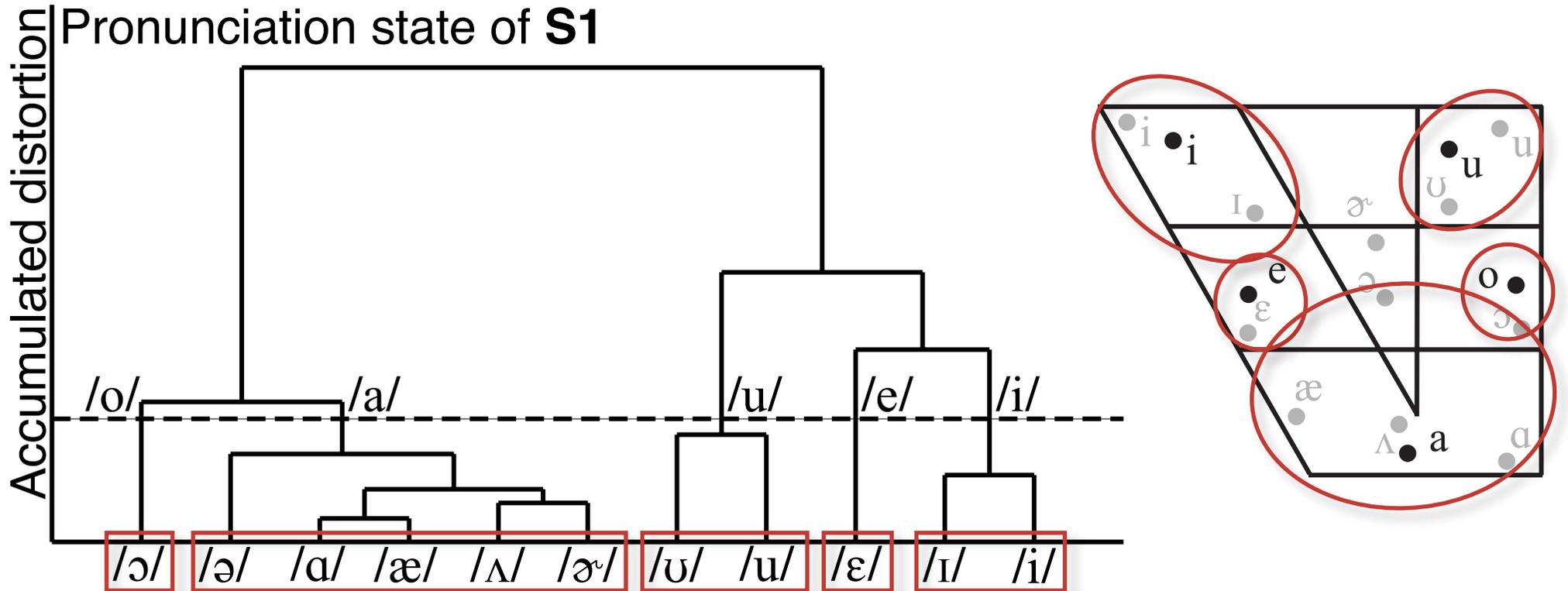
a	ɑ, æ, ʌ, ə, ɚ
i	ɪ, i
u	ʊ, u
e	ɛ
o	ɔ



- 「あ」 的な母音の置換
- 「い」 的な母音の置換
- :

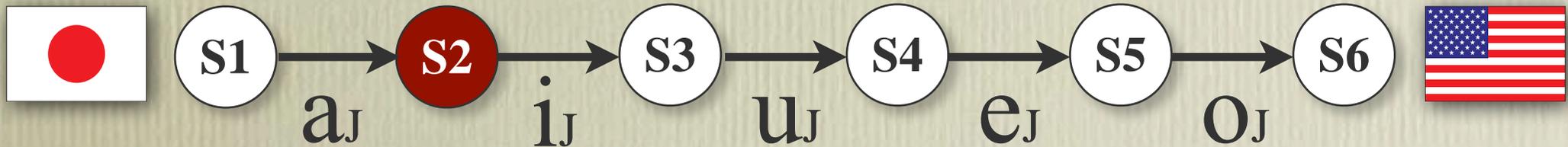
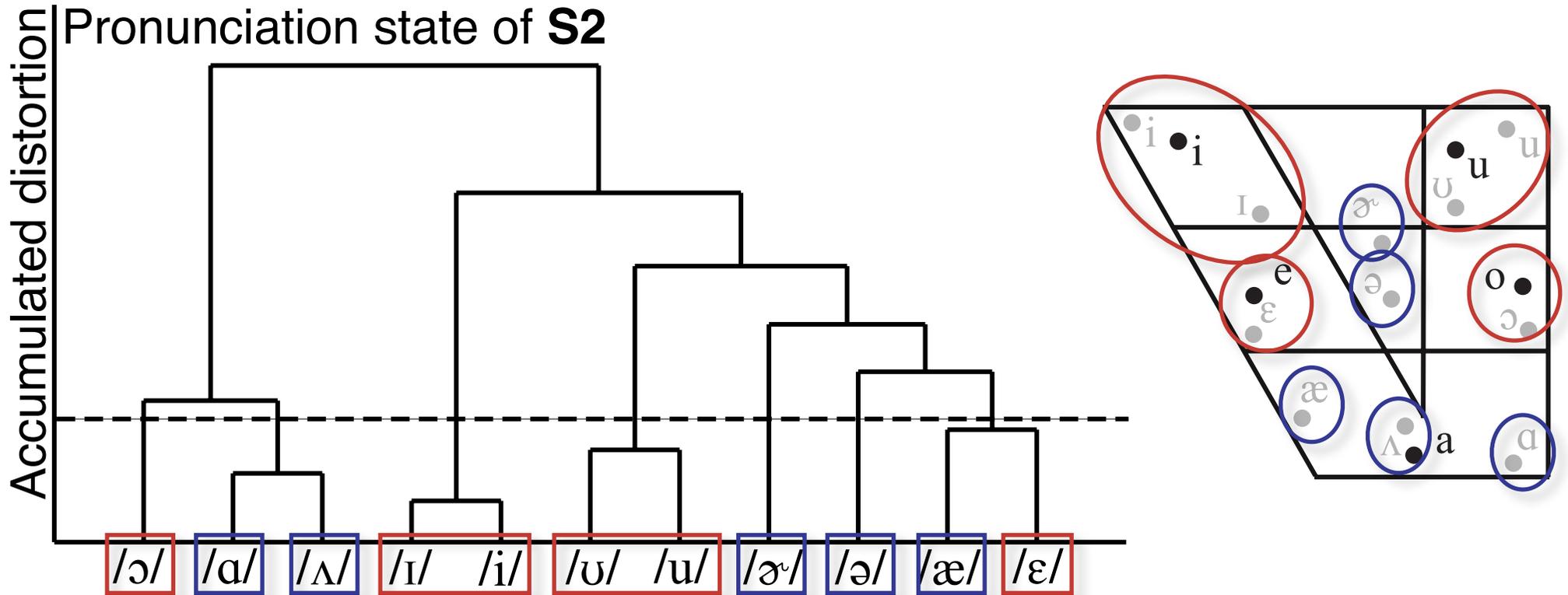
学習者の発音状態の記録とその遷移

S1からS6への遷移の様子の記事



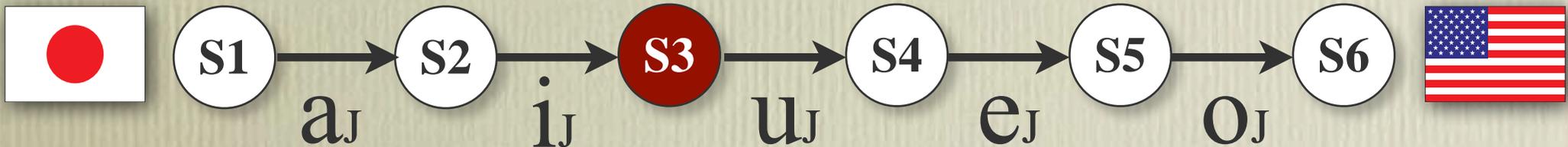
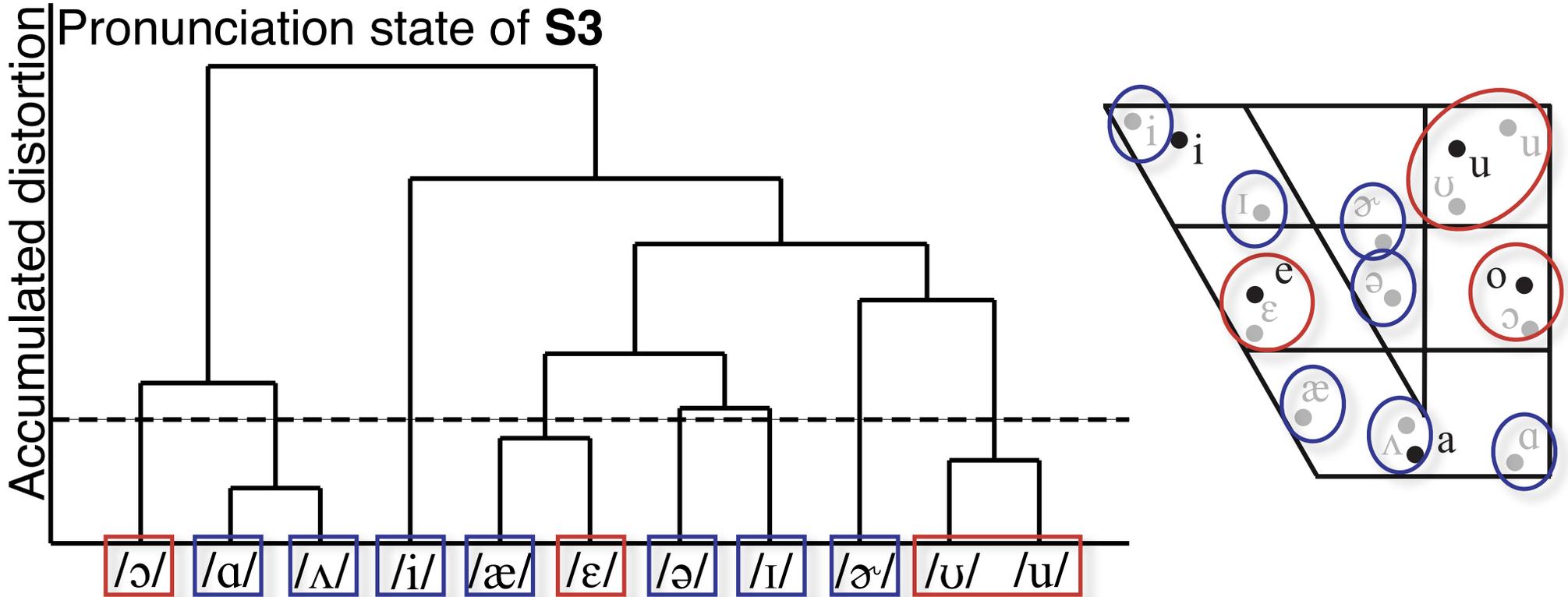
学習者の発音状態の記録とその遷移

S1からS6への遷移の様子の記事



学習者の発音状態の記録とその遷移

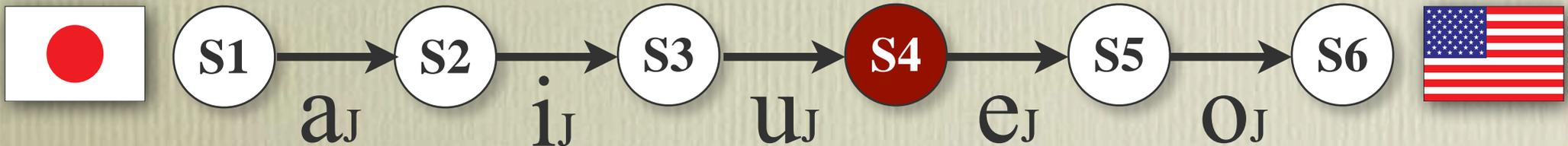
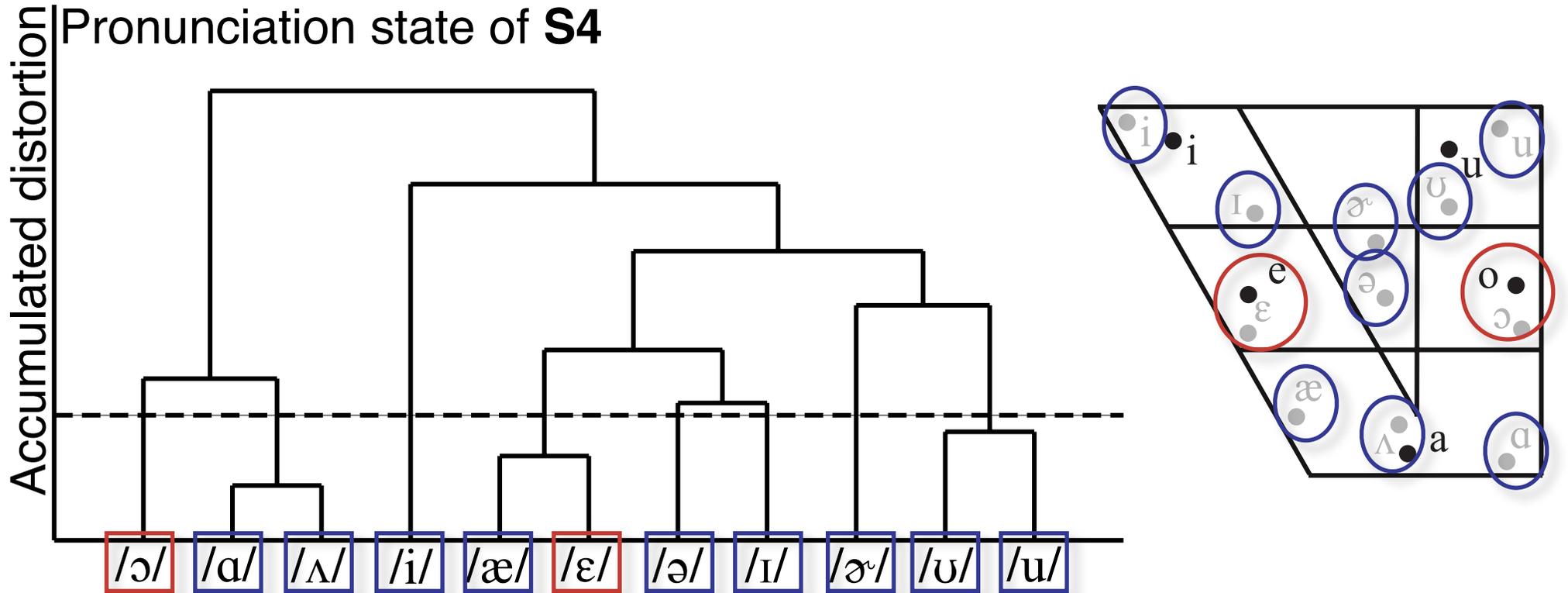
S1からS6への遷移の様子の記事



学習者の発音状態の記録とその遷移

S1からS6への遷移の様子の記事

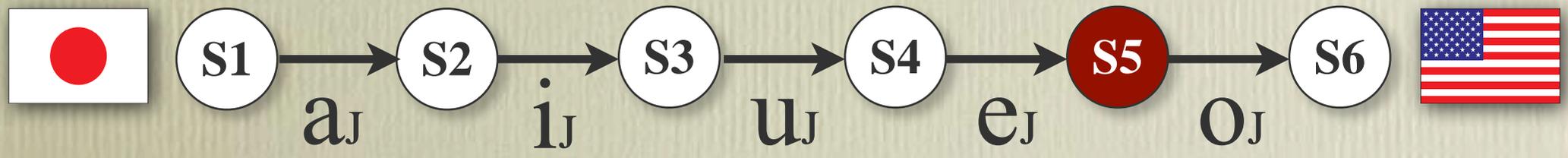
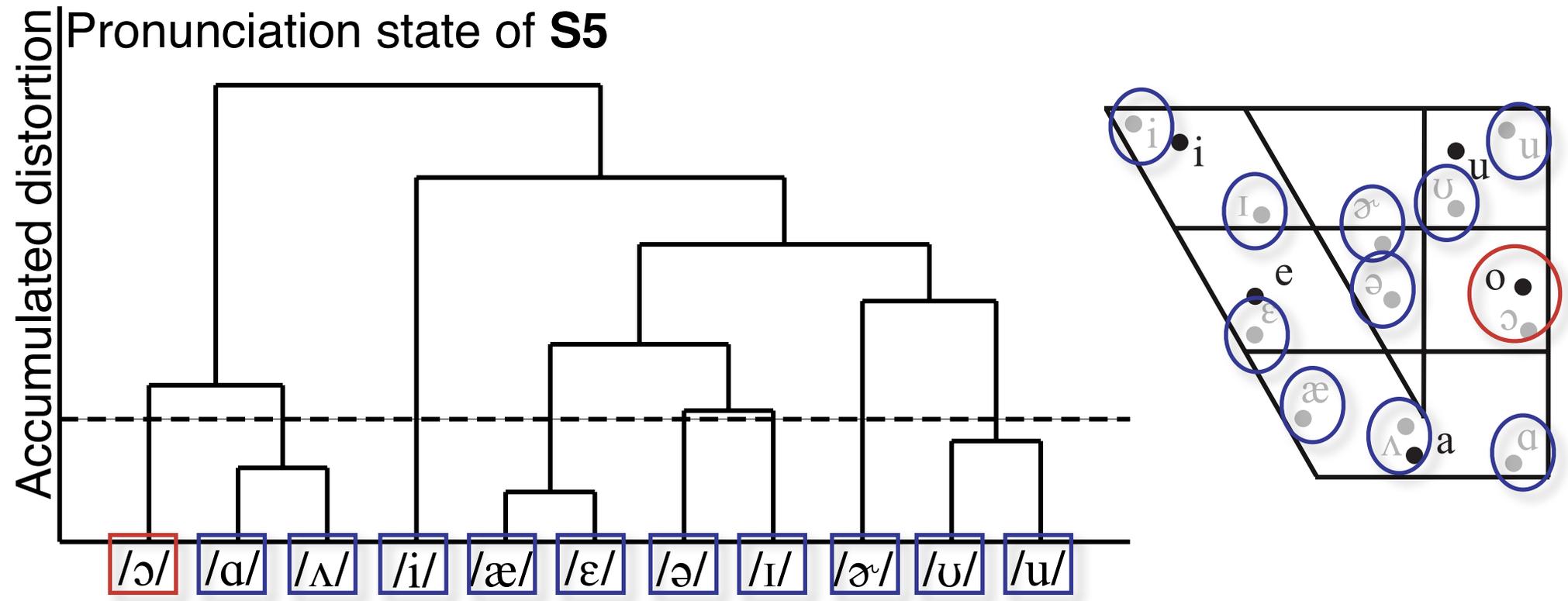
Pronunciation state of S4



学習者の発音状態の記録とその遷移

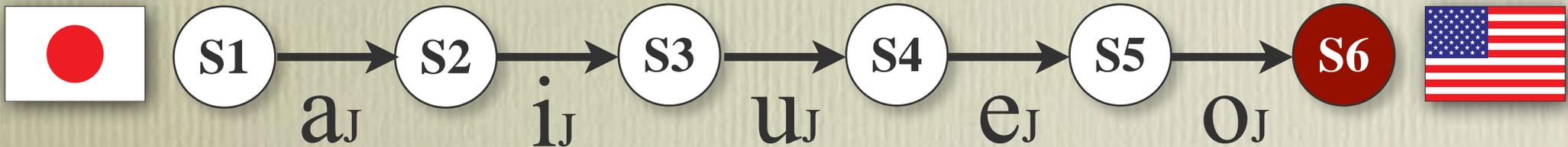
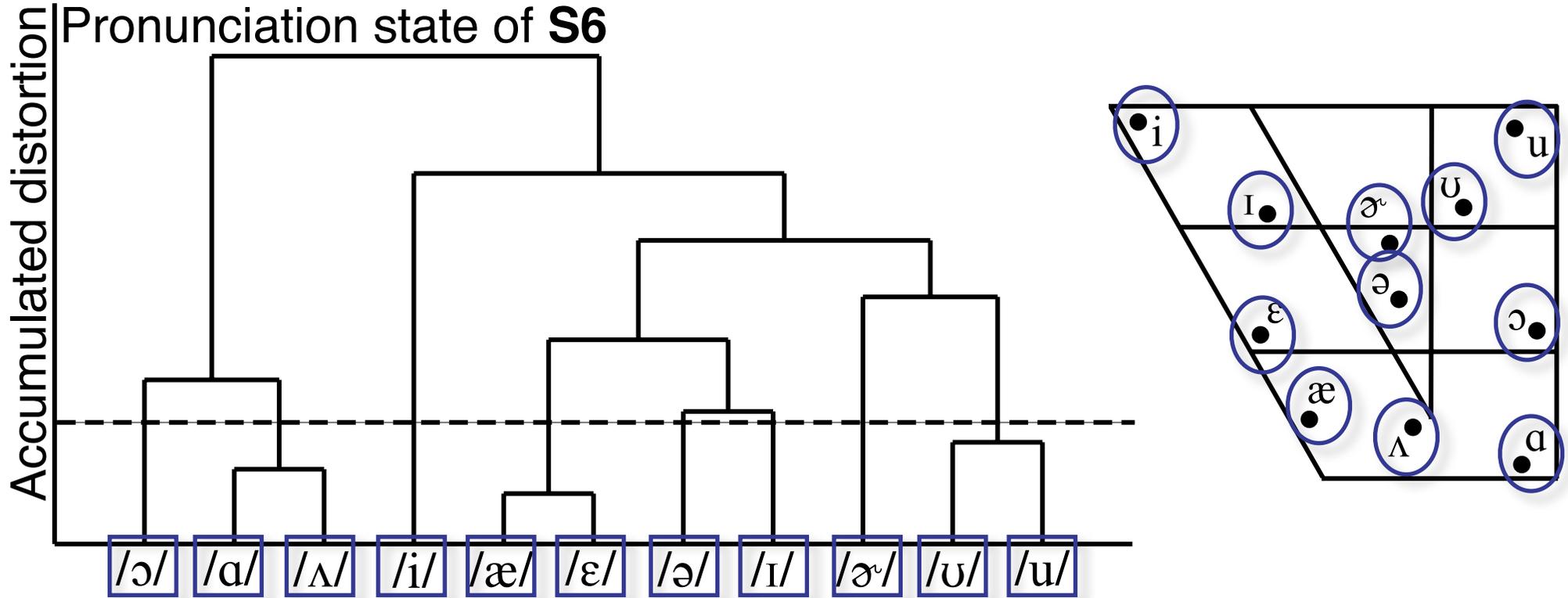
S1からS6への遷移の様子の記事

Pronunciation state of S5



学習者の発音状態の記録とその遷移

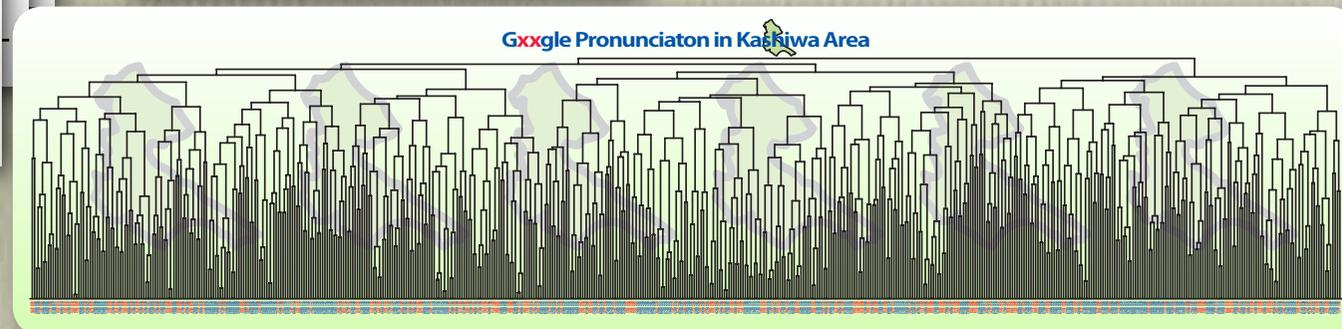
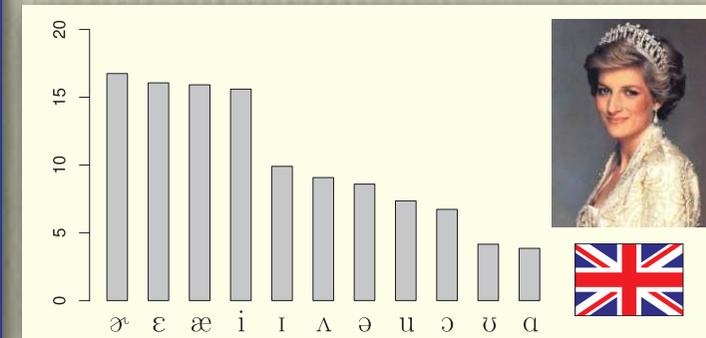
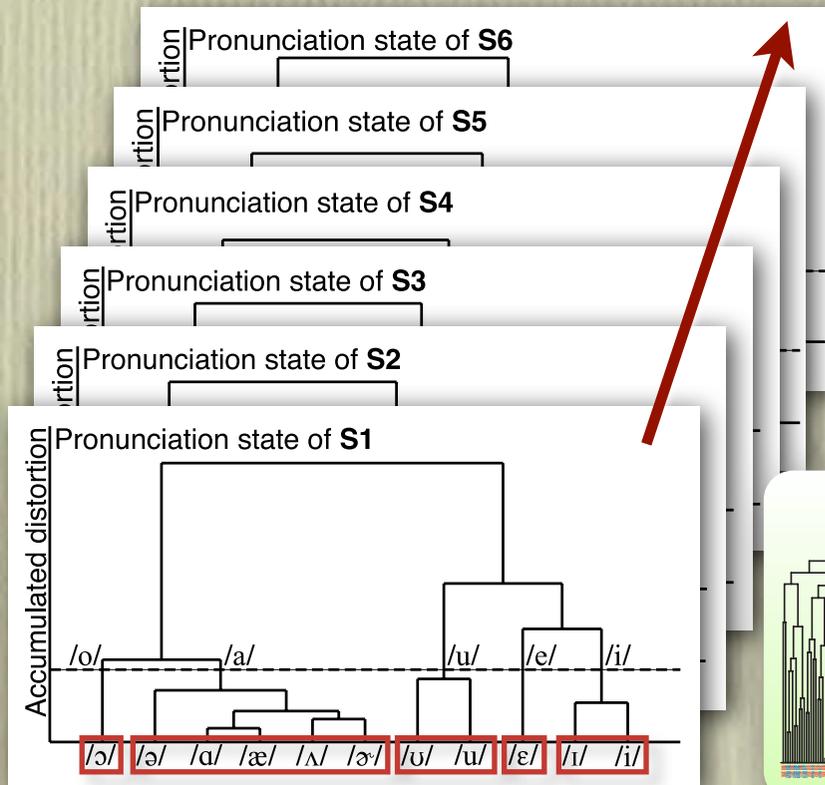
S1からS6への遷移の様子の記事



母音体系に見る発音習熟度

発音の構造表象に基づいて構築した技術・機能

- 学習者の母音体系がどのような状態にあるのかを記録する
- お手本の母音構造と比較して、矯正対象の母音を教示する
- 複数の教師からお手本としたい教師（相手）を選ぶ
- 年齢，性別などには影響されない学習者発音の分類



優先的に矯正すべき母音の教示

全体的構造歪みと局所的構造歪み

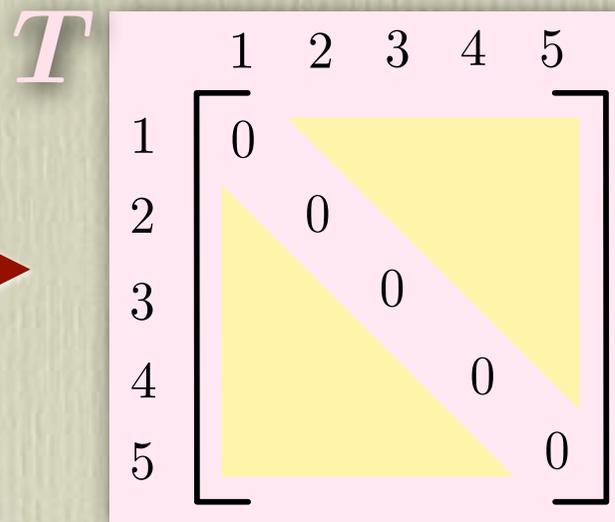
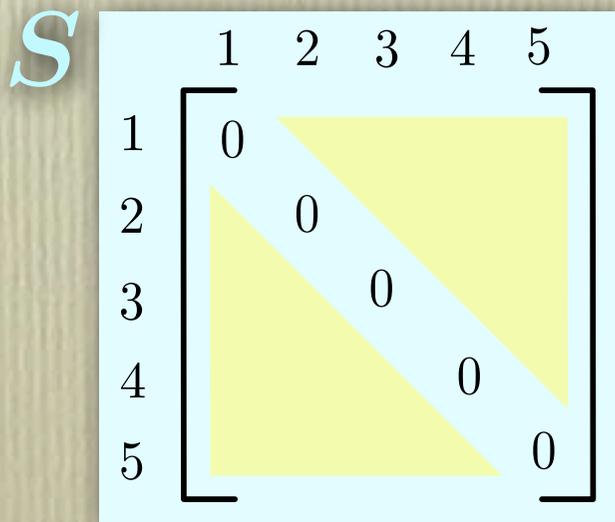
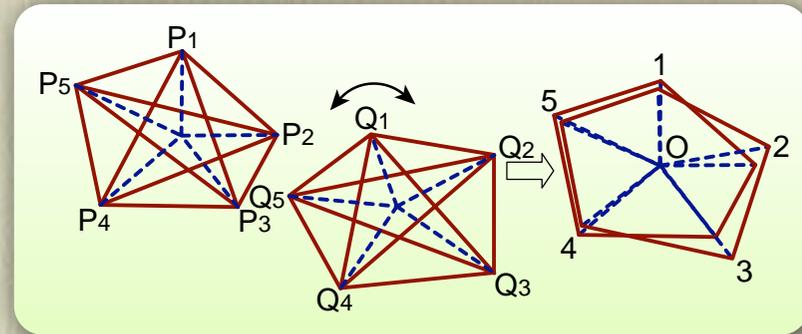
- 全体的構造歪み → 発音習熟度の推定

$$TD(S, T) = \sqrt{\frac{1}{M} \sum_{i < j} (S_{ij} - T_{ij})^2}$$

- 局所的構造歪み → 矯正すべき発音部位の推定

$$LD(S, T, i) = \sum_{j=1}^M |S_{ij} - T_{ij}|$$

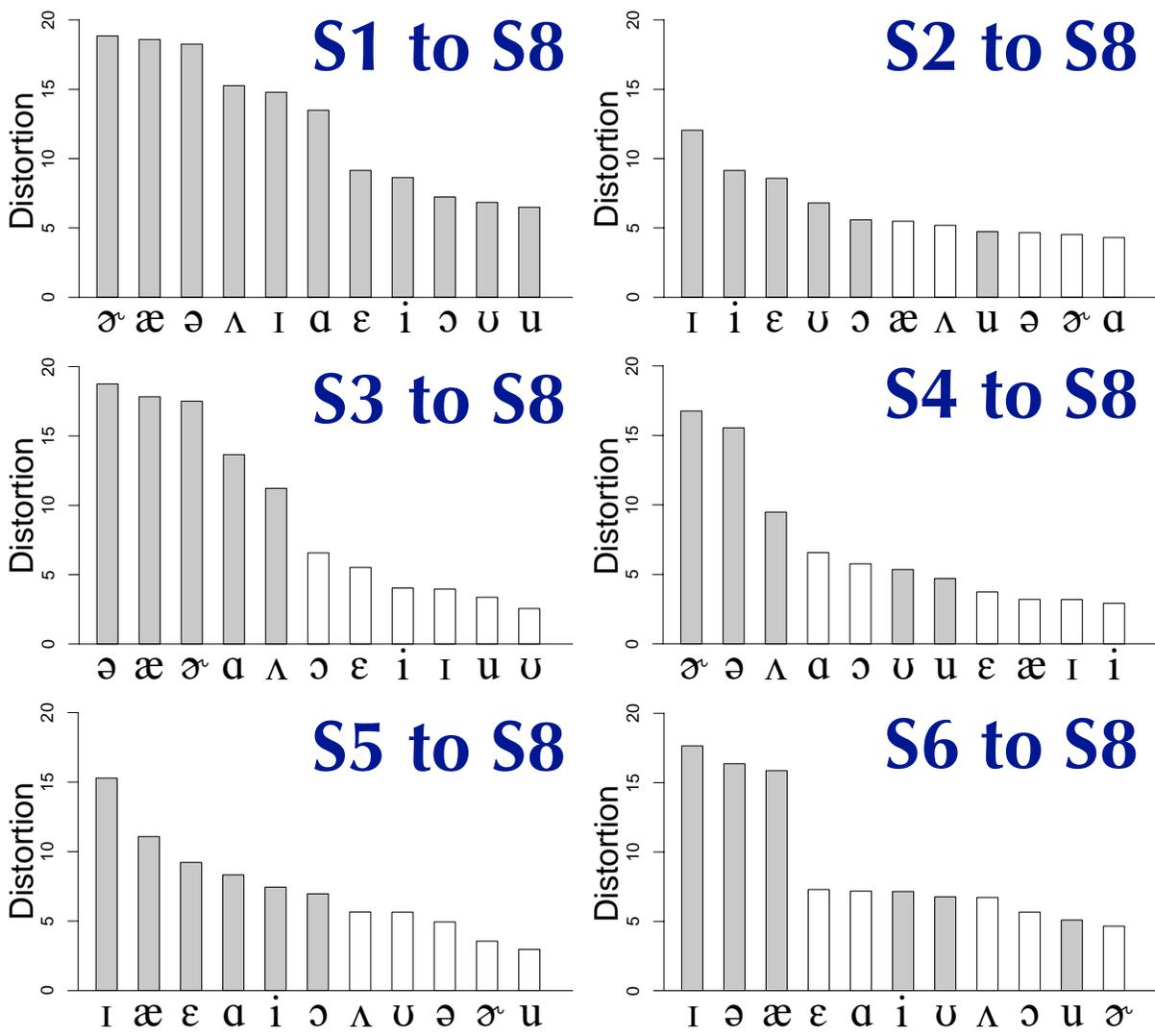
- LD = 大 ↔ 当該母音の矯正度 = 大



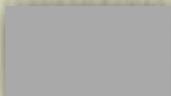
優先的に矯正すべき母音の教示

S8 : 米語, S1-S7 : 日本語混入米語

教師と学習者の母音距離行列のみから, 優先的矯正対象を推定



	ɑ	æ	ʌ	ə	ɝ	ɪ	ɪ	ʊ	u	ε	ɔ
S1	J	J	J	J	J	J	J	J	J	J	J
S2	E	E	E	E	E	J	J	J	J	J	J
S3	J	J	J	J	J	E	E	E	E	E	E
S4	E	E	J	J	J	E	E	J	J	E	E
S5	J	J	E	E	E	J	J	E	E	J	J
S6	E	J	E	J	E	J	J	J	J	E	E
S7	J	E	J	E	J	E	E	E	E	J	J
S8	E	E	E	E	E	E	E	E	E	E	E

 : 日本語母音で置換
 : 置換無し

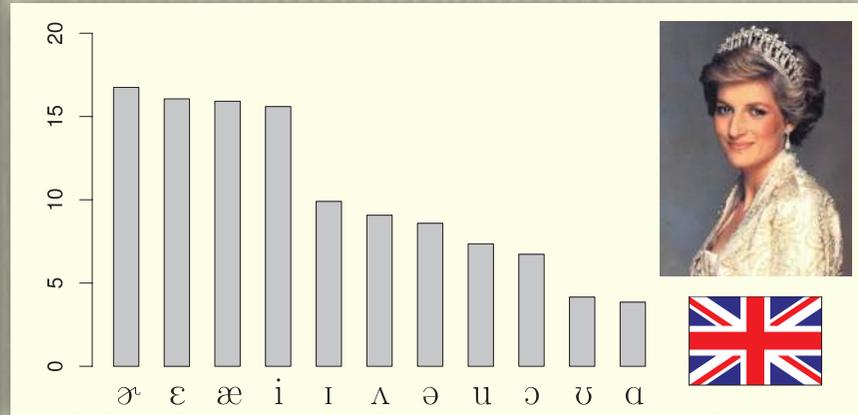


Figure 15: The estimated order of vowel correction

お手本となる相手を選ぶインタフェース



Model

Select two teachers.

OK

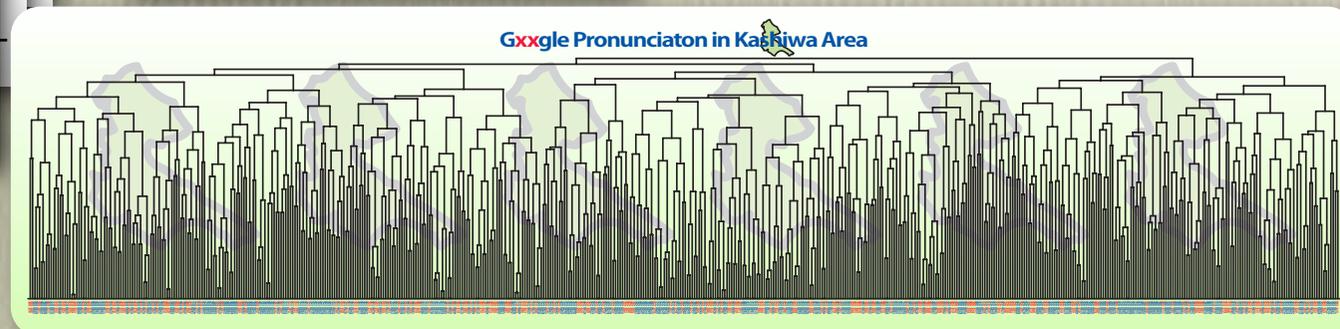
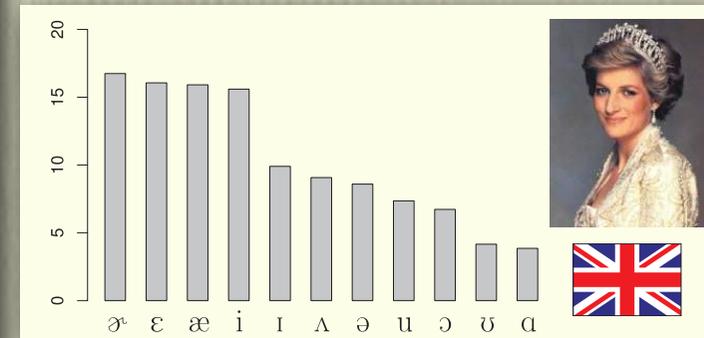
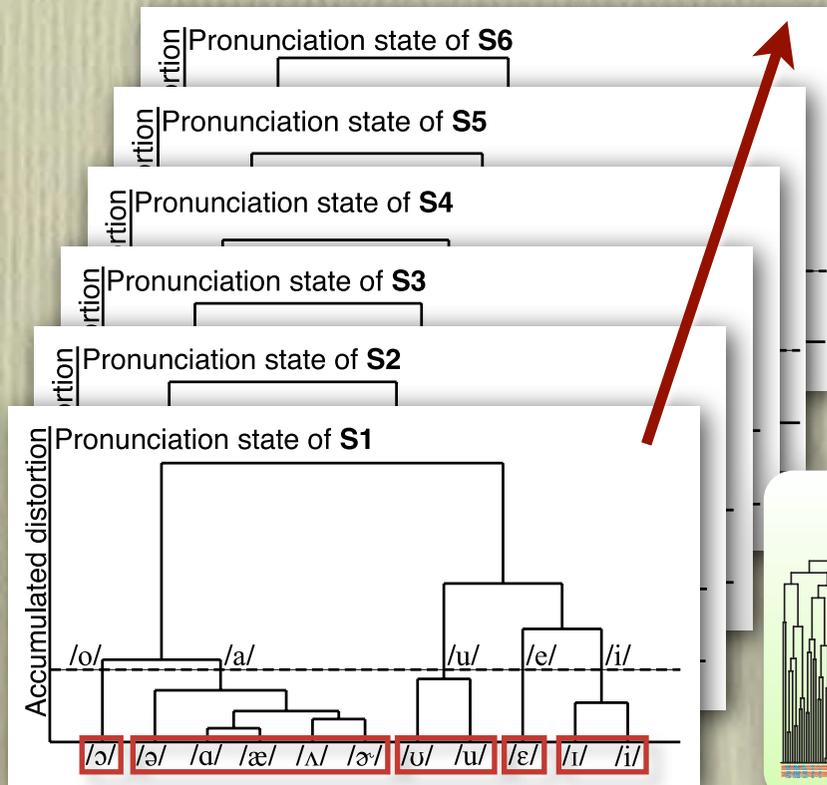
LONGMAN
PRONUNCIATION
DICTIONARY
C. Wells



母音体系に見る発音習熟度

発音の構造表象に基づいて構築した技術・機能

- 学習者の母音体系がどのような状態にあるのかを記録する
- お手本の母音構造と比較して、矯正対象の母音を教示する
- 複数の教師からお手本としたい教師（相手）を選ぶ
- 年齢，性別などには影響されない学習者発音の分類

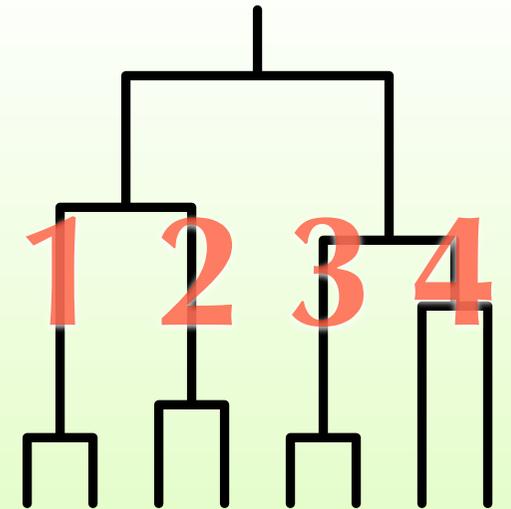
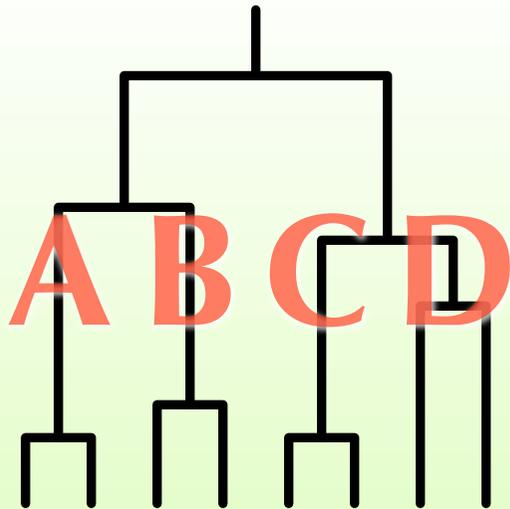
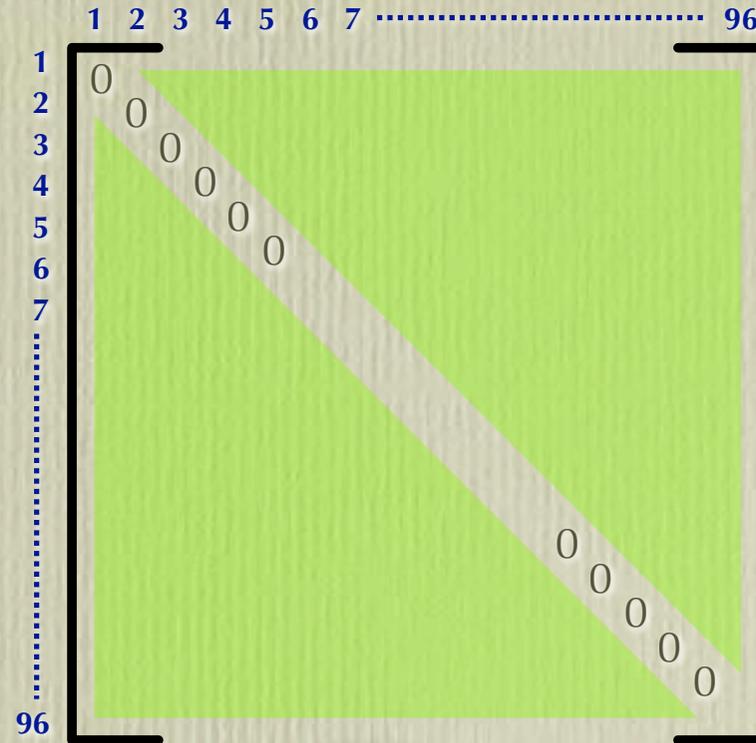


2通りの学習者分類

96 x 96 の距離行列に基づく学習者分類

● 話者：A~L

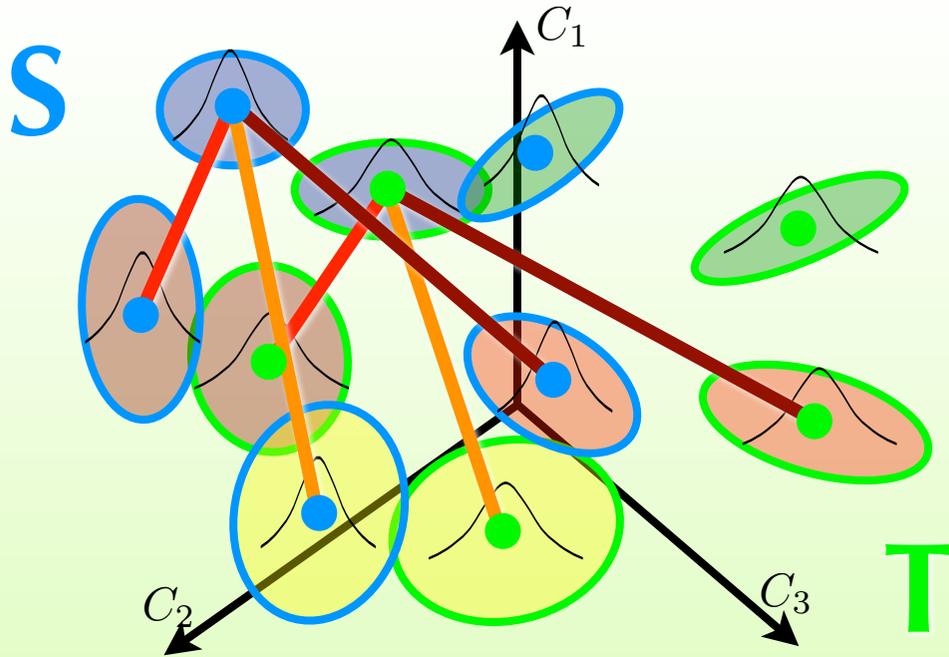
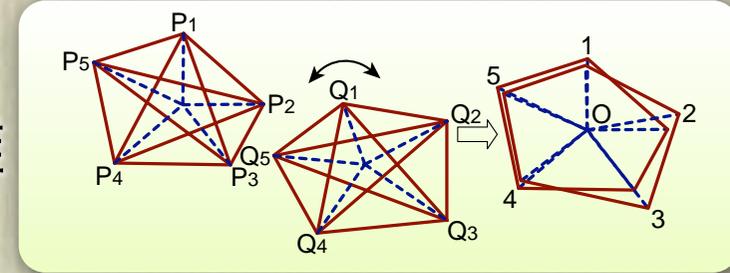
● 発音：1~8



2通りの学習者分類

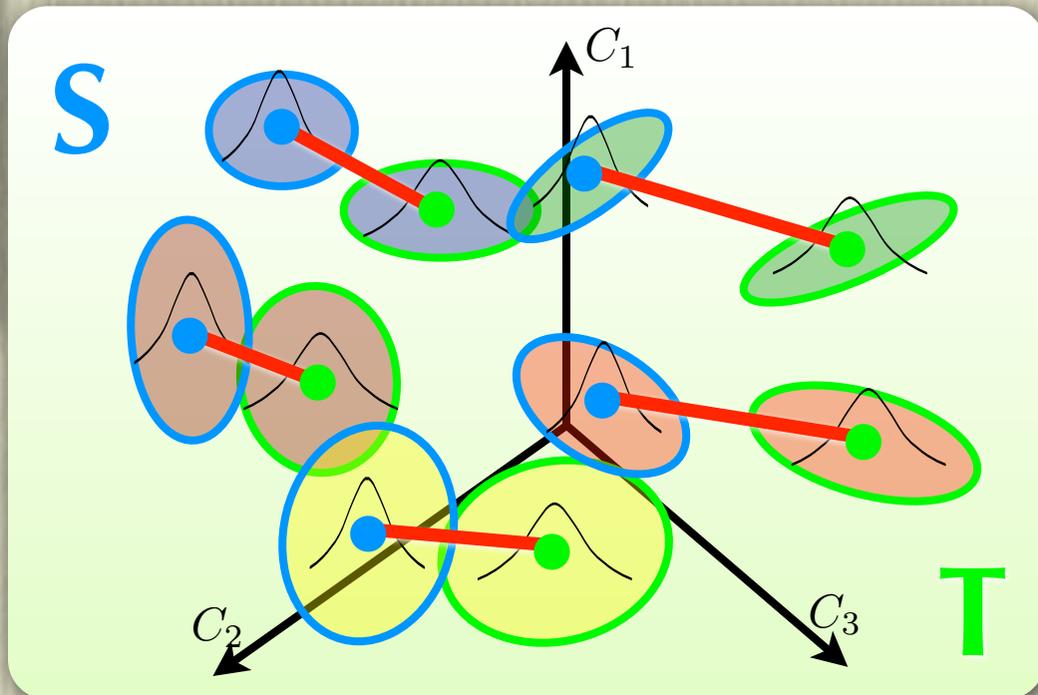
2通りの構造間距離

- コントラストの比較に基づく構造間距離
- 音の実体の比較に基づく構造間距離



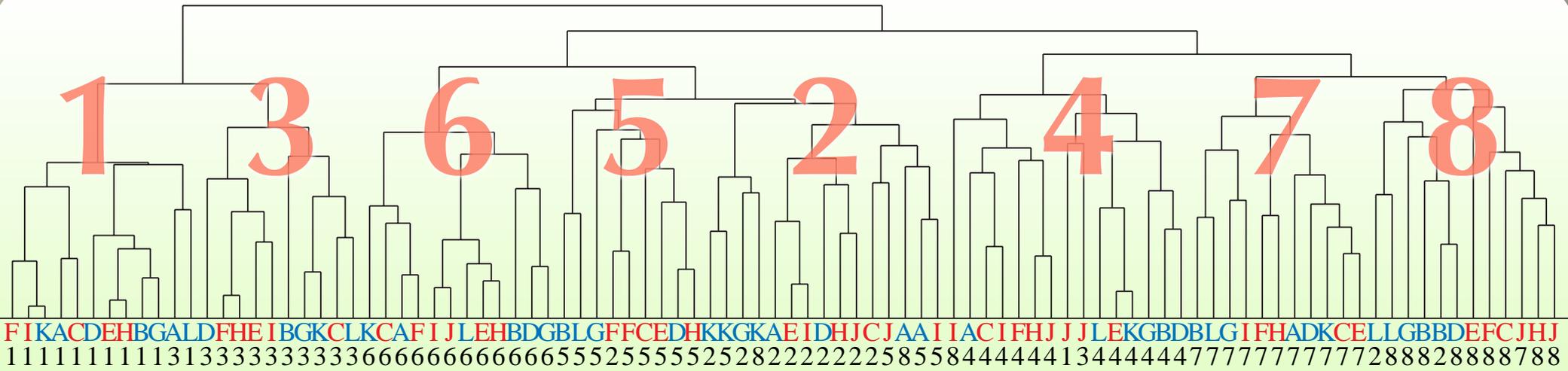
$$\sqrt{\frac{1}{M} \sum_{i < j} (S_{ij} - T_{ij})^2}$$

$$\sqrt{\frac{1}{M} \sum_i BD(v_i^S, v_i^T)}$$

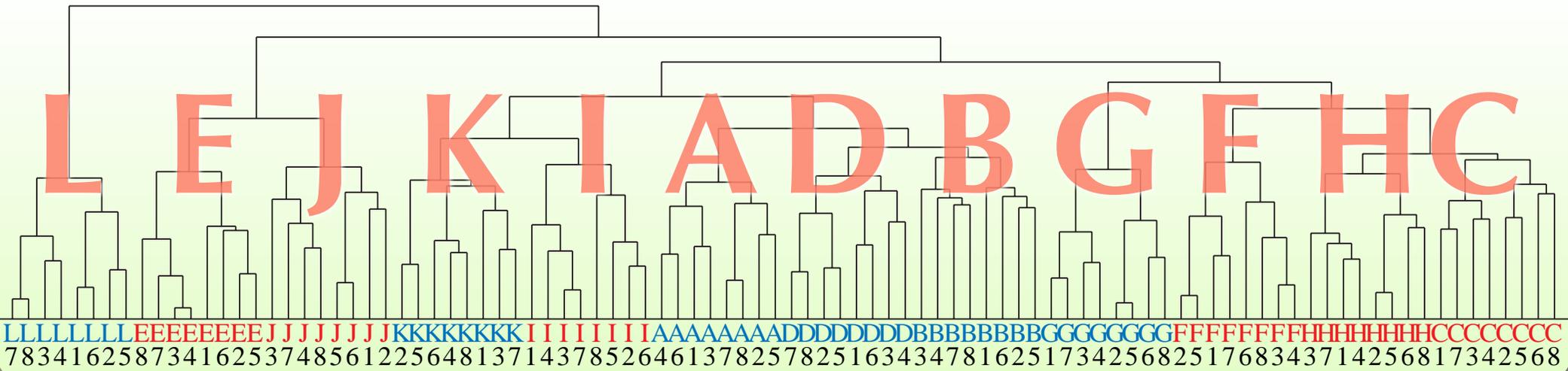


2通りの学習者分類

コントラストの比較に基づく構造間距離

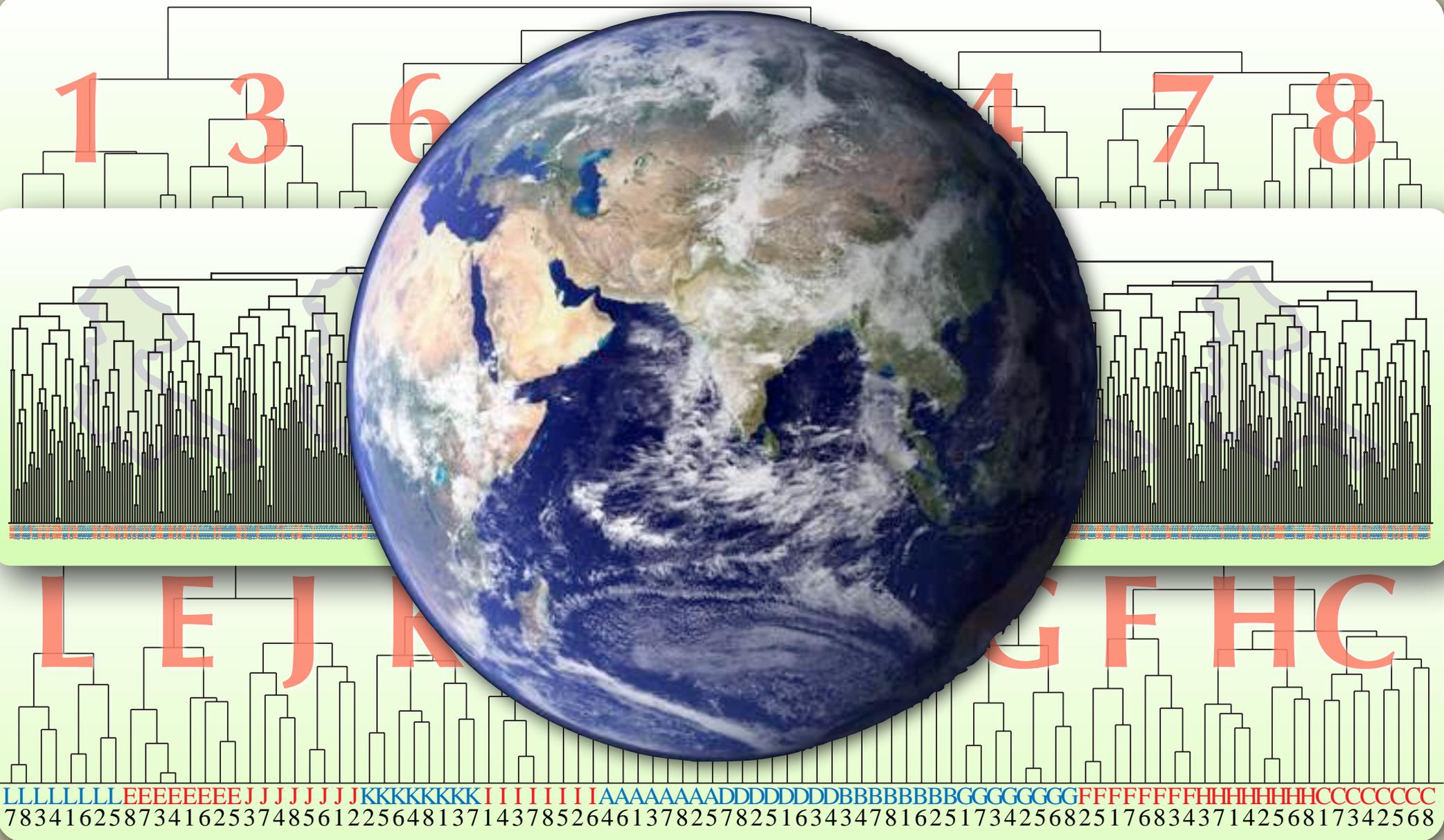


音の実体の比較に基づく構造間距離



2通りの学習者分類

コントラストの比較に基づく構造間距離



The current state of English

It is the only language used for global communication.

- About **1.5 billion** users on earth

It has the largest diversity in its form.

- Internationalization of a thing inevitably alters its form.

- English is not exceptional.

 - Syntax, pragmatics, lexical choice, spelling, **pronunciation**, etc

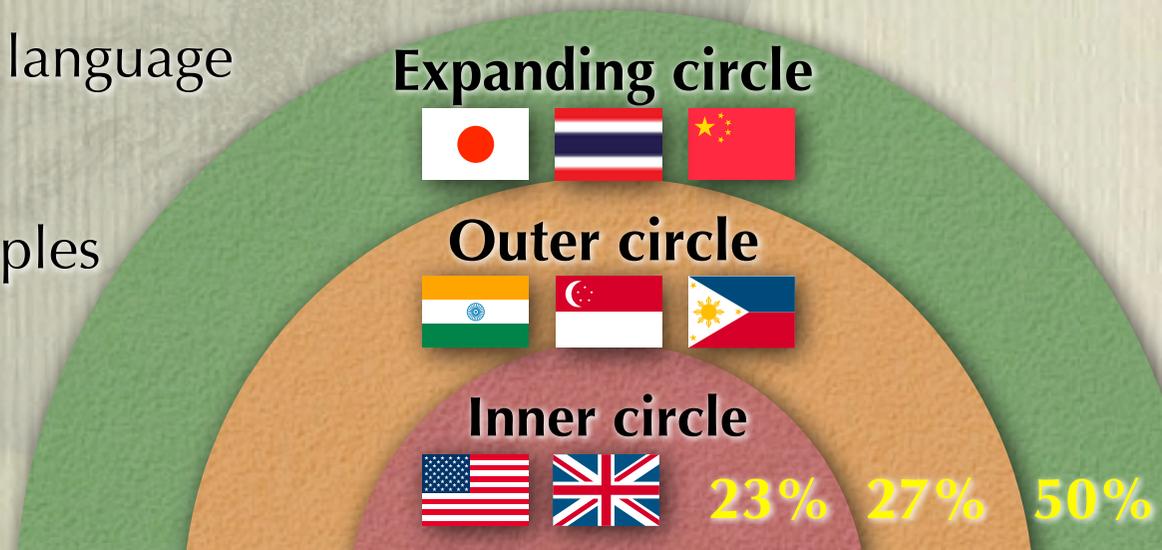
World Englishes (WE)

- Three circles model [Kachru1992]

 - E as native / official / foreign language

- No standard pronunciation

 - AE and BE are just two examples of **accented** Englishes.



Diversity of pronunciation in WE

What is the minimal unit and how many units?



[Kachru 1992]

Diversity of pronunciation in WE

What is the minimal unit and how many units?



[Kachru 1992]



1. British - Southern English - East London - Cockney

9. British - Scottish (unsure of specific type)

3. British - Southern English - Formal RP (Received Pronunciation)

1) native language, 2) official language



Requires a speaker-basis pronunciation distance matrix

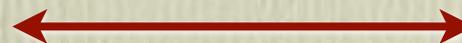
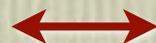


$$\begin{bmatrix} d_{11} & d_{12} & \dots & d_{1N} \\ d_{21} & d_{22} & \dots & d_{2N} \\ d_{31} & & & \\ \vdots & & & \\ d_{N1} & d_{N2} & \dots & d_{NN} \end{bmatrix}$$



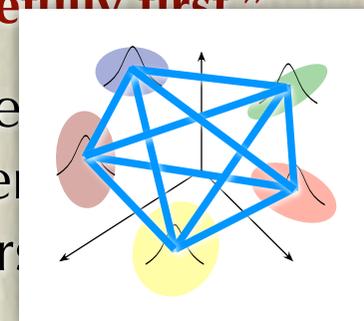
What is technically challenging?

To which is Minematsu's natural pronunciation closer?



“Those answers will be straightforward if you think them through carefully first”

- Pronunciation distance = phonetic distance between speakers
 \neq **acoustic** distance between speakers
 \neq **spectral** distance between speakers



Speech Accent Archive (SAA) [Weinberger'13]

- A common paragraph read by about 1.8K international speakers
- The paragraph is designed to achieve high phonemic coverage of AE.
- Speech samples and **their narrow IPA transcripts** are provided.

Please call Stella. Ask her to bring these things with her from the store: Six spoons of fresh snow peas, five thick slabs of blue cheese, and maybe a snack for her brother Bob. We also need a small plastic snake and a big toy frog for the kids. She can scoop these things into three red bags, and we will go meet her Wednesday at the train station.



Intius kol stela ask ha tu bñij
[plis kol stela ask ha tu bñij] tor
di: ðigks wð hs fðrm ða stor lts
siks spñns of fæf ə sno pils u
fæf ðik' ðik' sglaps ov blu su
fjz ?æn mæbi e snæk fo hau
buaðe bap wi ?olsð nid ?e
smol plæstik snæk ?æn e bik
toi frog fðrm ðe kits fj kæn
skyp ðoʒ ðigs ðnu jri: æt'
baks ?æn wi wil go mit hs
wenzdeis ?að ðe trejn stei[ʒ]





Speech Accent Archive (SAA) [Weinberger'13]

A common paragraph read by about 1.8K international speakers

[p^hlis kəl steɪlɑ ask^ɹ ɜ tə bɪŋ ðiz θɪŋz wɪf hɜ
 fɪlŋ ðə stɔɪ sɪks spunz əv fɪf snəu pi:s faɪf
 θɪk slæʃs əv blu tʃi:z en məɪbi ɜ snæk^ɹ foɹ
 hɜ bɪɑɹə ʔə brʌðə bɑp wɪ ɔl^ʷsə nid ə smɔl^ʷ
 plæstɪk sneɪk en ə bɪk tʷɪ fɪɔg fɛ ðə kɪdʒ ʃɪ
 kɛn skəp ðiz θɪŋs ɪntu fɪɪ ɹed bægz en wɪ
 wɪl goʊ mɪd ʒ wɛnzdeɪ et ɔðə tɹeɪn steɪʃən]

Please
 the s
 chee
 smal
 scoo
 Wed



smol plæstɪk snæk ʔɛn ə bɪk
 tɔɪ fɪɔg fɪɔm ðə kɪts ʃɪ kɛn
 skɪp ðəz θɪŋs ɪntu ɹi: ɹɛɹ^ɹ
 bɑks ʔɛn wɪ wɪl go mɪt hɜ
 wɛnzdeɪs ʔəð ɔðə tɹeɪn steɪʃən]





Pron.

of WE



Speech Ac

- A commo
- The para
- Speech sa

Please call Stella. Ask
the store: Six spoons o
cheese, and maybe a s
small plastic snake an
scoop these things into
Wednesday at the train



1. i	2. ĩ	3. i:	4. j	5. ĩ	6. ĩ
7. y	8. ɹ	9. ɪ	10. ɪ:	11. ɹ	12. ĩ
13. e	14. ě	15. ě	16. ε	17. ě	18. ě
19. æ	20. æ	21. æ:	22. æ̃	23. a	24. ā
25. i	26. j	27. ĩ	28. u	29. u	30. ø
31. ɜ	32. ə	33. e	34. ē	35. ũ	36. o
37. ɵ	38. ə	39. ə	40. ə	41. ə	42. ə
43. ɯ	44. ü	45. ũ	46. u	47. ü	48. u:
49. ü	50. ũ	51. ũ:	52. u	53. ɣ	54. o
55. ɵ	56. ɵ	57. ʌ	58. ʌ̃	59. ɔ	60. ɔ:
61. ə	62. ə	63. a	64. a:	65. ä	66. ä
67. p	68. p ^h	69. p̄	70. b	71. b̄	72. b
73. φ	74. β	75. β̄	76. β̄	77. f	78. v
79. ɣ	80. v	81. m	82. m̄	83. m̄	84. n
85. ŋ	86. ŋ	87. ŋ	88. ŋ	89. ŋ	90. n
91. t	92. t ^h	93. t̄	94. t̄	95. t'	96. t̄
97. d	98. d̄	99. d̄	100. d̄	101. s	102. ʒ
103. s ^j	104. z	105. z	106. ɹ	107. ɹ	108. ɹ
109. r	110. r	111. ɸ	112. l	113. l	114. l ^v
115. θ	116. ð	117. e	118. z	119. z	120. ʃ
121. ɜ	122. ç	123. j	124. j	125. k	126. k ^h
127. k̄	128. k'	129. k ^h	130. k	131. g	132. g
133. ġ	134. ġ	135. x	136. ɣ	137. ɣ	138. ɰ
139. ʔ	140. h	141. fi	142. w	143. ɥ	144. pφ
145. tθ	146. dð	147. ts	148. dz	149. tɕ	150. dz
151. tʃ	152. dʒ	153. kx			

nal speakers
verage of AE.
provided.



Pron. clustering only based on SAA

N speakers



$$\begin{matrix} & \begin{matrix} 1 & 2 & \dots & N \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ \vdots \\ N \end{matrix} & \left[\begin{array}{cccc} d_{11} & d_{12} & \dots & d_{1N} \\ d_{21} & d_{22} & \dots & d_{2N} \\ d_{31} & & & \\ \vdots & \vdots & & \\ d_{N1} & d_{N2} & \dots & d_{NN} \end{array} \right] \end{matrix}$$

[plis kol stala ask ha tu biftj
di: 0fjks w0 h3 f30m d3 st3u
sik3 sp3ns of fiej 3 sno pit3s
faif 0ik' 0ik' s3l3ps 0v blu
fjiz ?3n meibi e sn3k fo h3z
b3u03 b3p w3 ?3ls3 nid ?e
sm3l pl3stik sn3k ?3n e bik
t3i f30g f30m 03 kits fj k33n
sk3p 03g 0fj3 3ntu jri: 33t'
b3ks ?3n w3 w3l go mit h3
w3nzdeis ?3d d3 tr33n st3if3]

[plis kol stala ask ha tu biftj
di: 0fjks w0 h3 f30m d3 st3u
sik3 sp3ns of fiej 3 sno pit3s
faif 0ik' 0ik' s3l3ps 0v blu
fjiz ?3n meibi e sn3k fo h3z
b3u03 b3p w3 ?3ls3 nid ?e
sm3l pl3stik sn3k ?3n e bik
t3i f30g f30m 03 kits fj k33n
sk3p 03g 0fj3 3ntu jri: 33t'
b3ks ?3n w3 w3l go mit h3
w3nzdeis ?3d d3 tr33n st3if3]

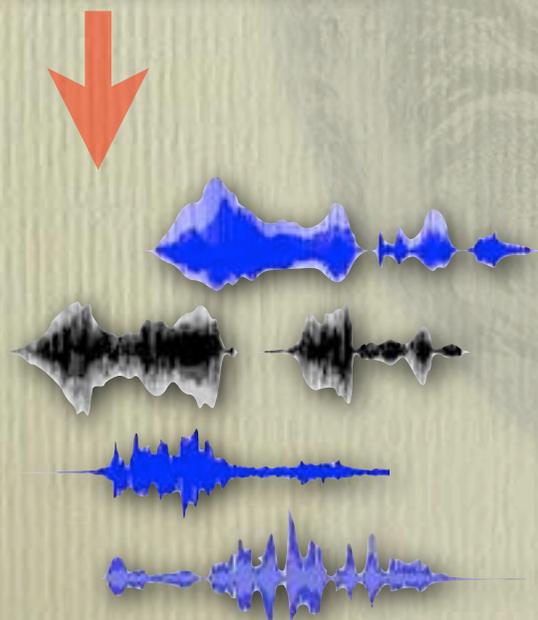
[plis kol stala ask ha tu biftj
di: 0fjks w0 h3 f30m d3 st3u
sik3 sp3ns of fiej 3 sno pit3s
faif 0ik' 0ik' s3l3ps 0v blu
fjiz ?3n meibi e sn3k fo h3z
b3u03 b3p w3 ?3ls3 nid ?e
sm3l pl3stik sn3k ?3n e bik
t3i f30g f30m 03 kits fj k33n
sk3p 03g 0fj3 3ntu jri: 33t'
b3ks ?3n w3 w3l go mit h3
w3nzdeis ?3d d3 tr33n st3if3]

Pron. clustering only based on SAA

N speakers



$$\{d_{mn}\} \approx \{p_{mn}\} ?$$



**Pron. Structure
Analysis**

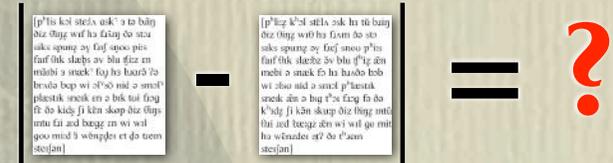


IPA-based reference pron. distance

Optimal alignment bet. two transcripts [Shen et al., '13]

- Dynamic Time Warping (DTW)

- DTW can minimize the accumulated distortion.

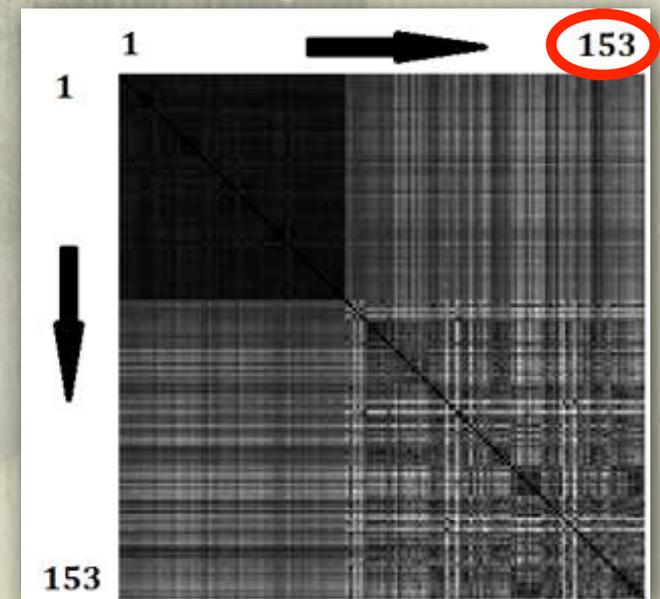


p_h	l	i	z	k	A	l	s	t	E	l	a
p_h	l	i	s	k	o	l_G	s	t	E	l	v

b	l	@	u
b	l	u	

?	o	l	s	o
A	s	o		

- Similar to edit-distance-based alignment of transcripts [Wieling et. al, '12]
- DTW requires a distance matrix of all the 153 IPA symbols used.
 - 20 productions for each by a phonetician
 - HMM is built for each symbol (SD-HMM)
 - HMM = Hidden Markov Model
 - Acoustic distance is obtained from each HMM (phone) pair.

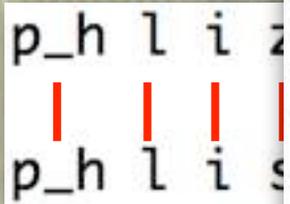


IPA-based reference pronunciation distance

Optimal alignment

Dynamic Time Warping

DTW can



Similar to

DTW requires

20 products

HMM is k

HMM

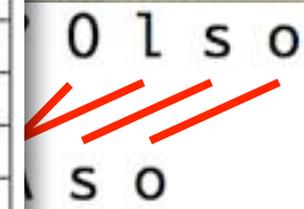
Acoustic
each HM

1. i	2. ĩ	3. i:	4. j	5. ĩ	6. ĩ
7. y	8. ɨ	9. ɪ	10. ɪ:	11. ɨ	12. ĩ
13. e	14. ě	15. ě	16. ɛ	17. ě	18. ě
19. æ	20. æ	21. æ:	22. æ̃	23. a	24. ā
25. ɨ	26. ɨ	27. ĩ	28. ʉ	29. ʉ	30. ɔ
31. ɜ	32. ɜ	33. ɛ	34. ě	35. ũ	36. ɵ
37. ɵ	38. ɔ	39. ɔ̃	40. ɔ̣	41. ɔ̃	42. ɔ̣
43. ʊ	44. ũ	45. ũ	46. u	47. ũ	48. u:
49. ü	50. ü	51. ũ:	52. ʊ	53. ʏ	54. o
55. ɵ	56. ɵ	57. ʌ	58. ʌ̃	59. ɔ	60. ɔ:
61. ɔ̃	62. ɔ̃	63. ʌ	64. ʌ:	65. ä	66. ä
67. p	68. p ^h	69. p̃	70. b	71. b̃	72. ḅ
73. ɸ	74. β	75. β̃	76. β̣	77. f	78. v
79. ɣ	80. ʋ	81. m	82. m̃	83. ṃ	84. n
85. ñ	86. ṇ	87. ŋ	88. ɲ	89. ŋ	90. N
91. t	92. t ^h	93. t̃	94. ṭ	95. t'	96. t̄
97. d	98. d̃	99. ḍ	100. ɖ	101. s	102. ʂ
103. s̃	104. z	105. z̃	106. ɹ	107. ɹ̃	108. ɹ̣
109. r	110. ɾ	111. ɾ̃	112. l	113. l̃	114. ḷ
115. θ	116. ð	117. ɸ	118. z	119. z̃	120. ʃ
121. ʒ	122. ɸ̣	123. j	124. j̃	125. k	126. k ^h
127. k̃	128. k'	129. k ^h	130. ḳ	131. g	132. g
133. ɡ̃	134. ɡ̣	135. x	136. ɣ	137. ɣ̃	138. ɰ
139. ʔ	140. h	141. h̃	142. w	143. ɰ	144. pɸ
145. tθ	146. dð	147. ts	148. dz	149. tɛ	150. dz
151. tʃ	152. dʒ	153. kx			

et al., '13]

[p'ɪŋ k'ɔl stɪk hɪ tɪ bʌŋ
dɪz ɪŋg wɔθ hɪ fɪm ðə st
sɪks spɪŋg zɪ ðeɪ snəʊ p'ɪs
fɪn'fɪk stɪkz zɪ bɪt p'ɪz zɪ
mɛtɪ ə snæk fɪ hɪ hʌndə kɔb
wɪ əlsoʊ nɪd ə snɪk p'ɪstɪk
snɪk ʌn ə hɪz p'ɪs frɛŋ fɪ ðə
k'ɔdʒ fɪ kɪn skɪp dɪz ɪŋg mɪt
hɪz nɔt bɪtʒ dɪz wɪ wɔt qɪ mɪ
hɪ wɪnɪdɪz sɪ? ðə c'ɔsn
stɛɪʃn]

= ?



Mieling et. al., '12]

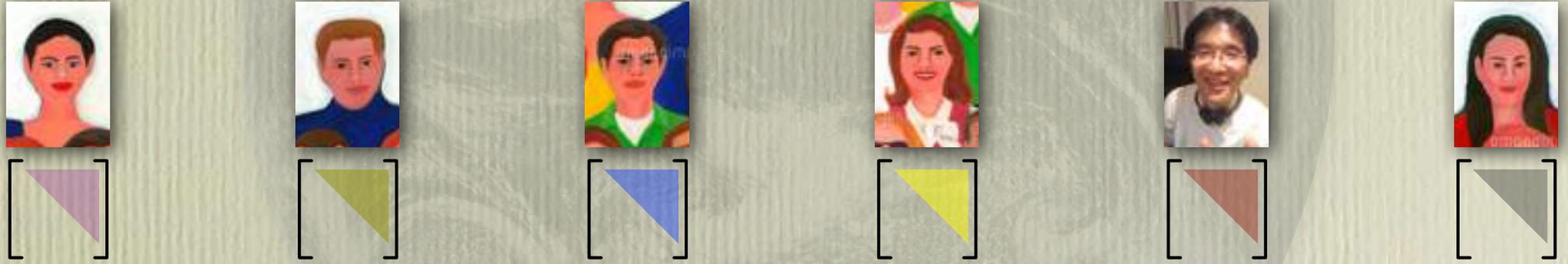
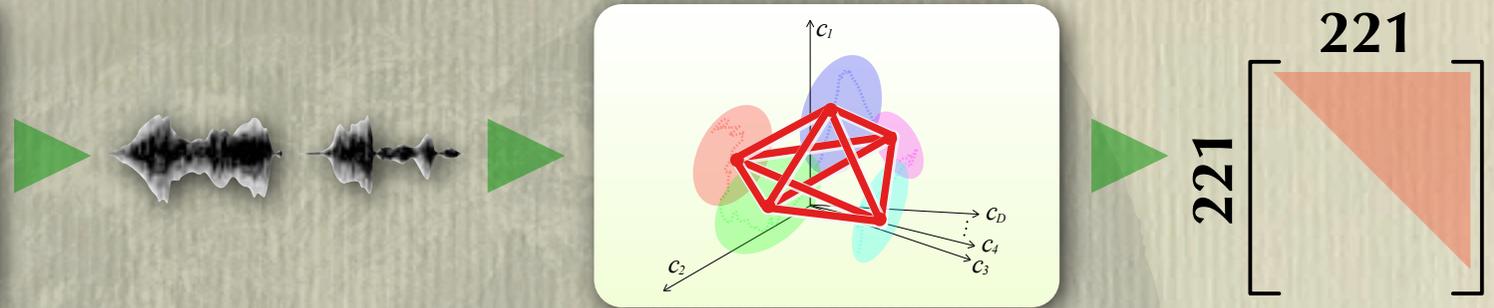
symbols used.

153

Pron. distance calculation using structure

A common paragraph to pron. structure

Please call Stella. Ask her to bring these things with her from the store: Six spoons of fresh snow peas, five thick slabs of blue cheese, and maybe a snack





Use of IPA transcripts to prepare reference distances

- DTW-based calculation of the reference distance bet. transcripts



Prediction of the ref. distances using pron. structures

- SVR-based supervised prediction using structures as input features



Use of **phonemic** transcripts to calculate distances

- Corresponds to calculate pron. distances somewhat coarsely.



/ pəli:z kəl^v stɪlə æsk hɜ: t^wə bɪɪŋ / #symbols = 153

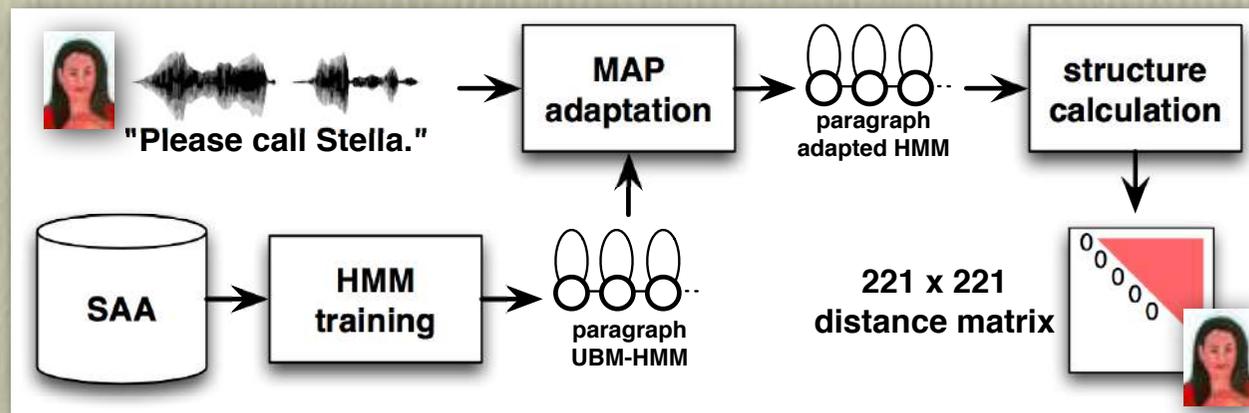
 [p ah l iy z k ao l s t ih l ah ae s k #symbols = 39
 hh ah r t ow b r ih ng]



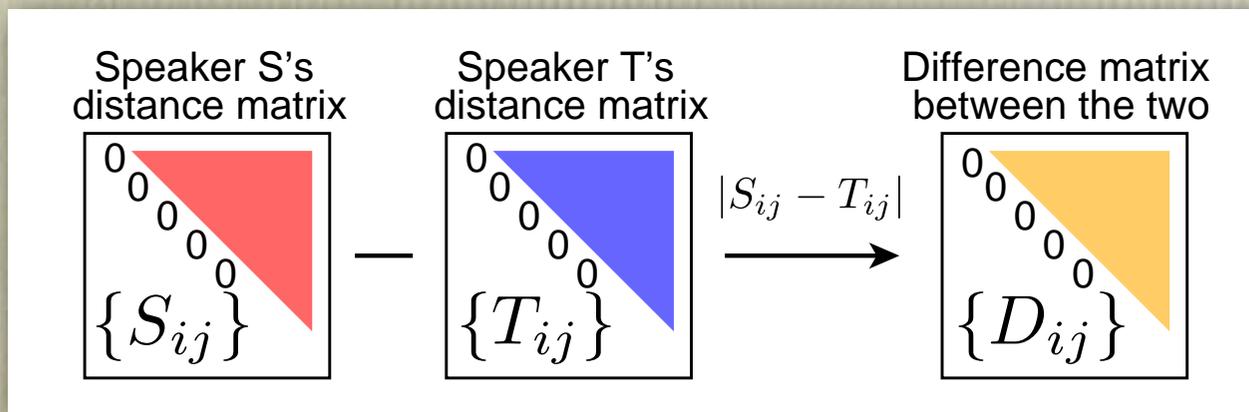
SVR-based prediction of IPA distances [Kasahara'14]



Pronunciation structure extraction from an SAA sample



Differential features from two pronunciation structures



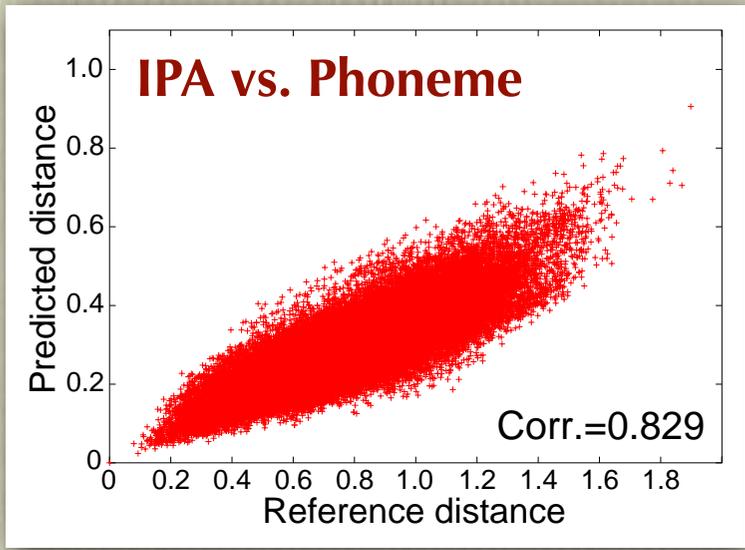
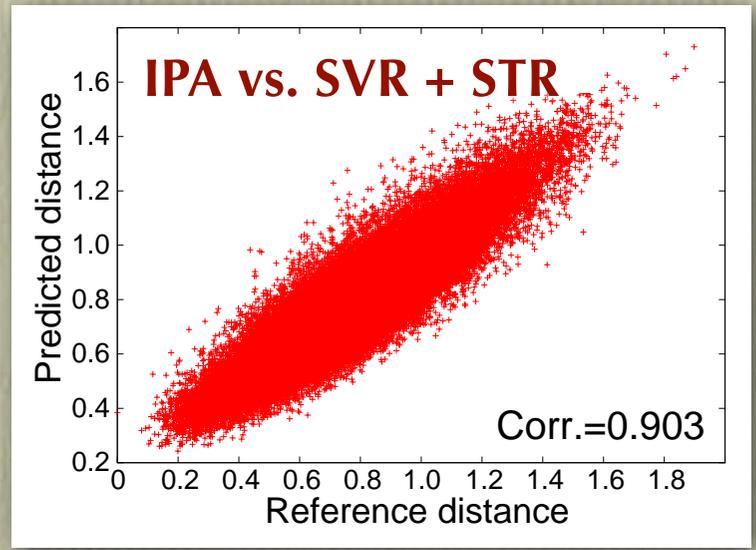


Three kinds of distances

- IPA-based distances
- SVR-based predicted distances
- Phoneme-based distances



Assessment in terms of correlation

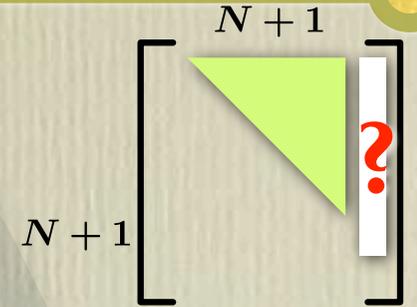


- IPA phone : minimum unit of sounds perceived by phoneticians
- Phoneme : minimum unit of sounds perceived by general listeners
- Our method can estimate distances better than general listeners.**

Application of **speaker-pair-open** prediction

TED talks browser from your viewpoint

- If TED talkers provide their SAA readings....
- If these readings are transcribed by phoneticians....



Visualization of pronunciation diversity [Kawase et al., '14]



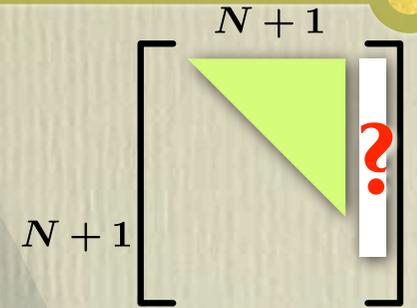
TED



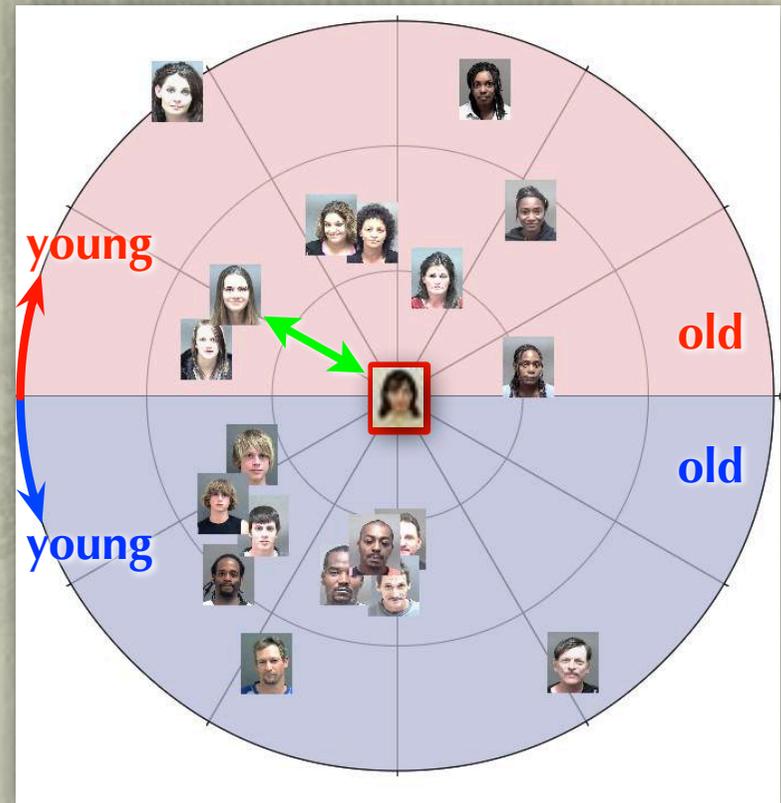
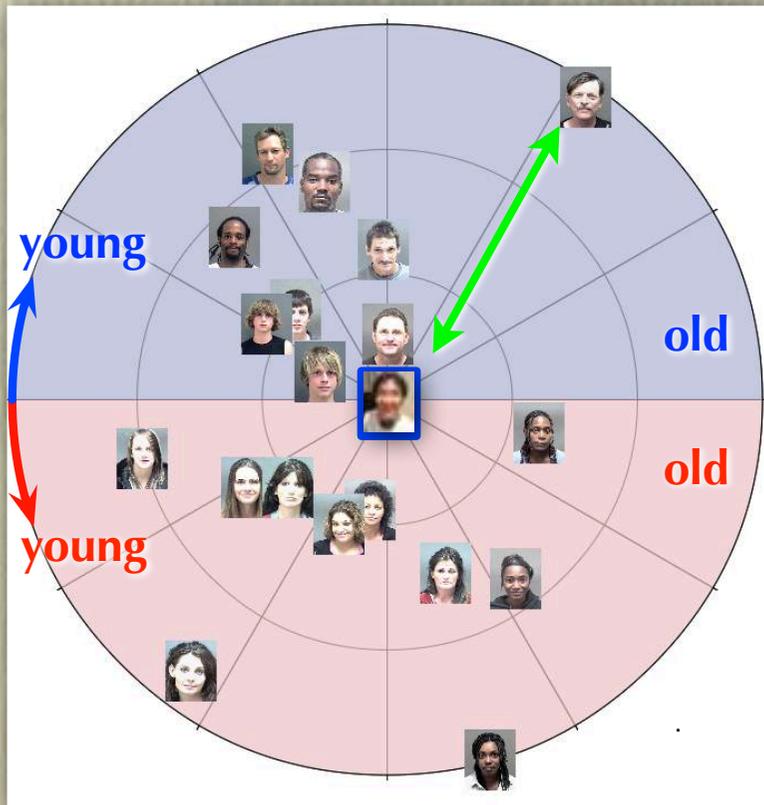
Application of **speaker-pair-open** prediction

TED talks browser from your viewpoint

- If TED talkers provide their SAA readings....
- If these readings are transcribed by phoneticians....



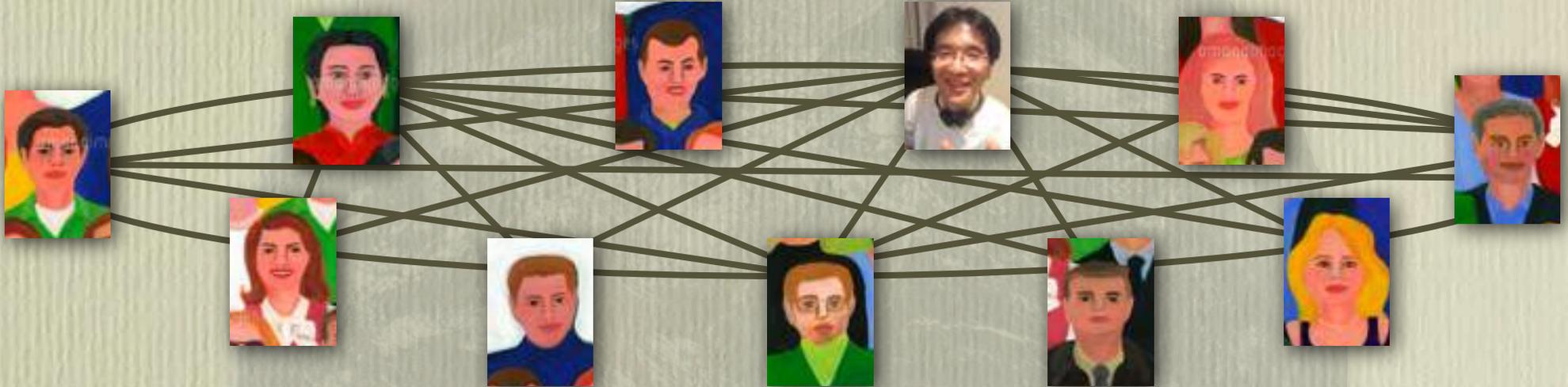
Visualization of pronunciation diversity [Kawase et al., '14]



Y. Kawase, et al., “Visualization of pronunciation diversity of World Englishes from a speaker’s self-centered viewpoint”

Possible application of **spk-open** prediction

Individual-based and really global map of WE pron.



Can be used for WE communication facilitator

- Easy access to speakers with pronunciation similar to yours



at a restaurant



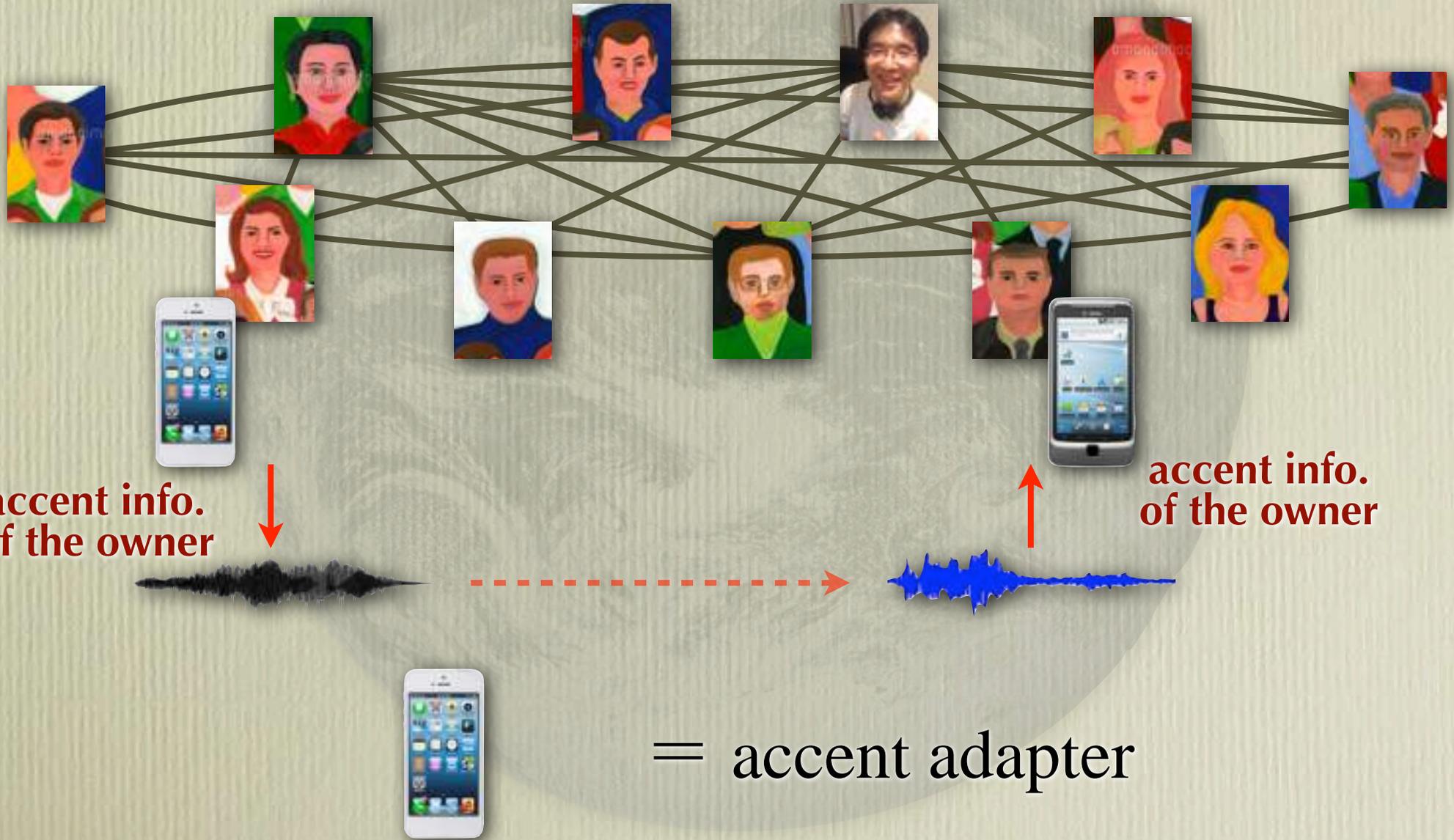
at a hotel



at a ticket office

Possible application of **spk-open** prediction

Adaptation of others' accents to a listener's own accent



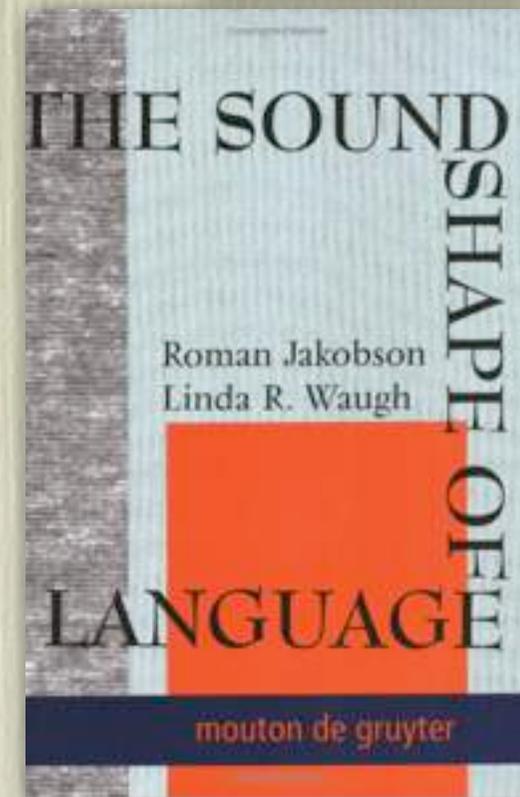
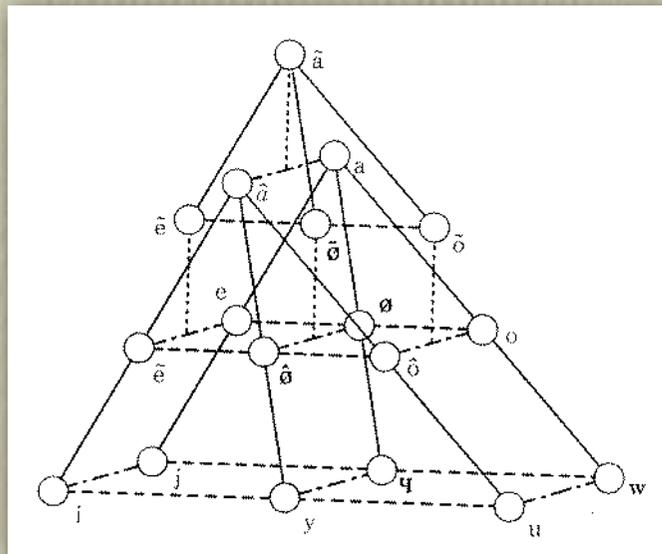
古典的音韻論に見られる主張

Roman Jakobson (1896-1982)

● The sound shape of language (1949)

Physiologically identical sounds may possess different values in conformity with the whole sound system, i.e. with their relations to the other sounds.

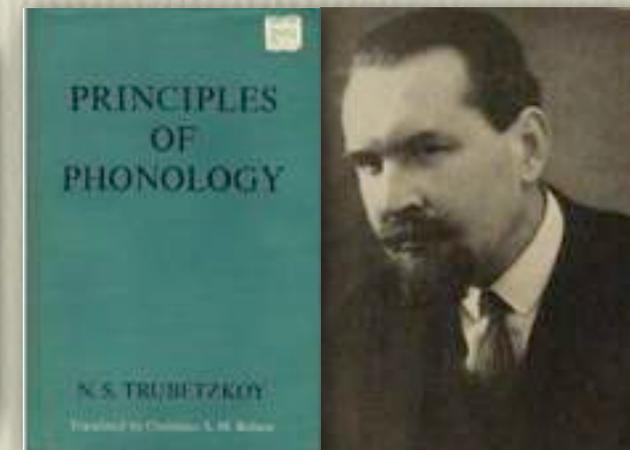
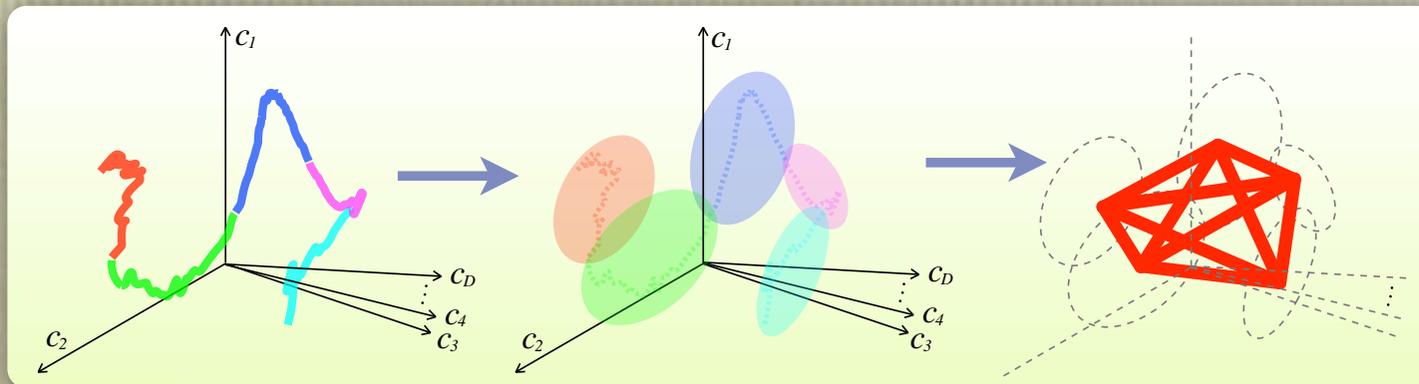
We have to put aside the accidental properties of individual sounds and substitute a general expression that is the common denominator of these variables.



古典的音韻論に見られる主張

Nikolay Trubetskoj (1890-1938)

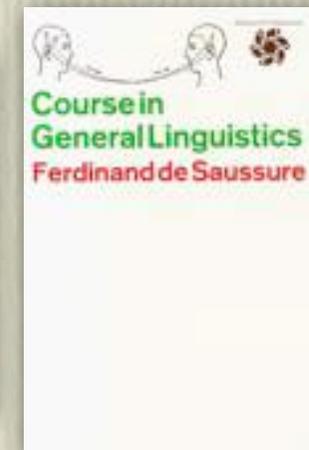
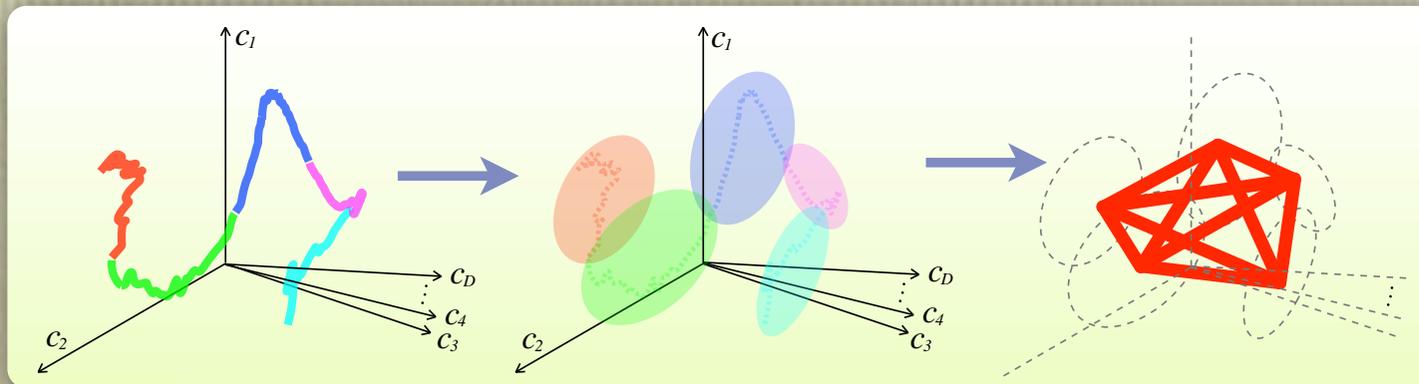
- The Principle of Phonology (1939)
- The phonemes should not be considered as building blocks out of which individual words are assembled. Each word is a phonic entity, a Gestalt, and is also recognized as such by the hearer.
- As a Gestalt, each word contains something more than sum of its constituents (phonemes), namely, the principle of unity that holds the phoneme sequence together and lends individuality to a word.



古典的音韻論に見られる主張

Ferdinand de Saussure (1857-1913)

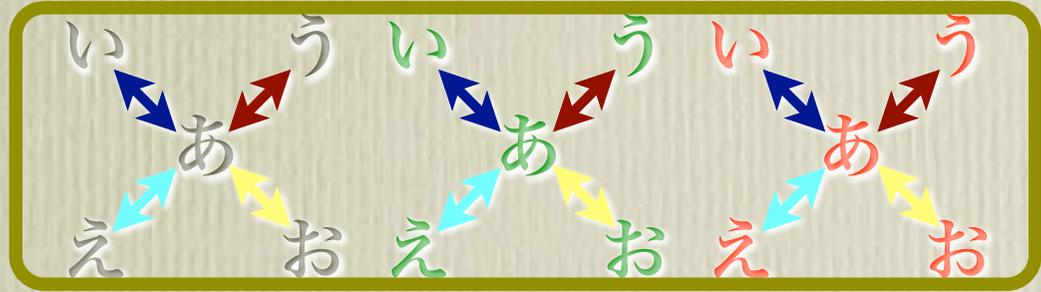
- Father of modern linguistics
- Course in General Linguistics (1916)
- What defines a linguistic element, conceptual or phonic, is the relation in which it stands to the other elements in the linguistic system.
- The important thing in the word is not the sound alone but the phonic differences that make it possible to distinguish this word from the others.
- Language is a system of only conceptual differences and phonic differences.



「あ」って何だろう？

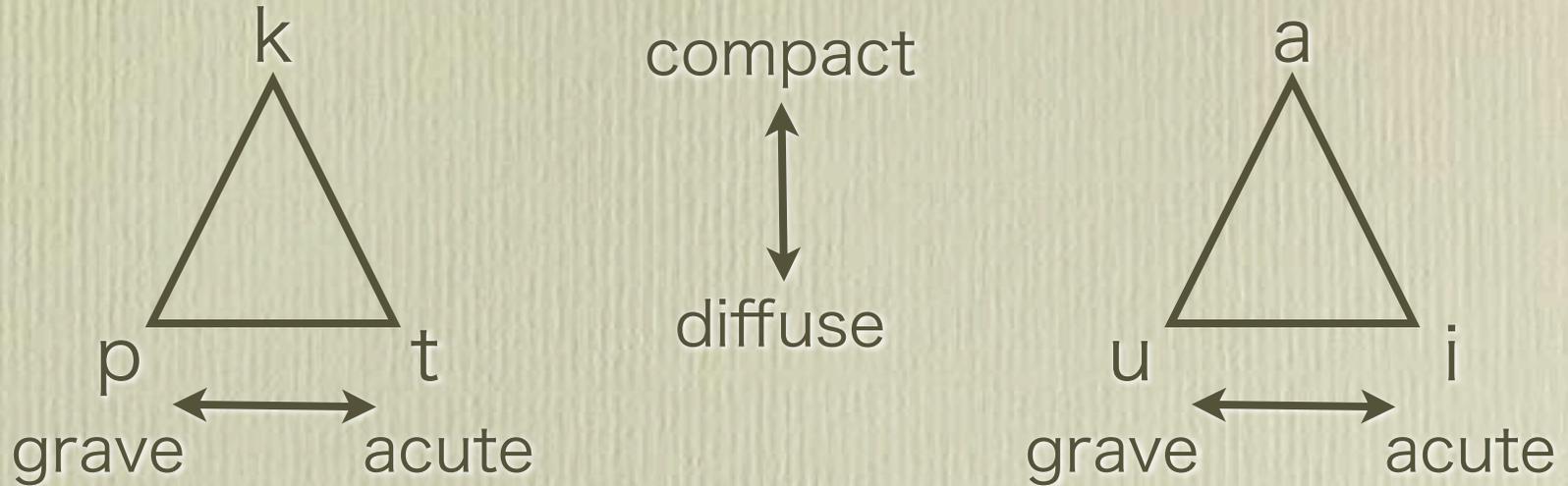
二つの見方, Element first or system first?

- An introduction to descriptive linguistics
 - Written by Gleason, H. A, 1961
 - A phoneme is a class of sounds that: (1) are phonetically similar and (2) show certain characteristic patterns of distribution in the language or dialect under consideration.
 - A phoneme is one element in the sound system of a language having a characteristic set of interrelations with each of the other elements in that system.
 - The phoneme cannot, therefore, be acoustically defined. The phoneme is a feature of language structure. That is, it is an abstraction from the psychological and acoustical patterns that enables a linguist to describe the observed repetitions of things that seem to function within the system as identical in spite of obvious differences. The phonemes of a language are a set of abstractions.



ヤコブソンの弁別素性

弁別素性を用いた音的差異・音体系の記述



弁別素性を用いた音素の記述

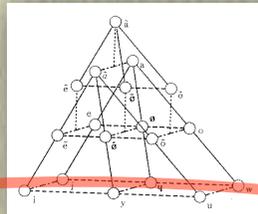
	p	b	m	f	v	θ	ð	t	d	n	s	z	l	r	ʃ	ʒ	tʃ	dʒ	j	ɹ	k	g	ŋ	w	ʔ	h
Back	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	+	+	+	+	+	+
High	-	-	-	-	-	-	-	-	-	-	-	-	-	-	+	+	+	+	+	+	+	+	+	+	-	-
Coronal	-	-	-	-	-	+	+	+	+	+	+	+	+	+	+	+	+	+	-	-	-	-	-	-	-	-
Anterior	+	+	+	+	+	+	+	+	+	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-
Labial	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	+	-	-
Continuant	-	-	-	+	+	+	+	-	-	-	+	+	+	-	+	+	-	-	+	+	-	-	-	+	-	+
Lateral	-	-	-	-	-	-	-	-	-	-	-	-	+	-	-	-	-	-	-	-	-	-	-	-	-	-
Nasal	-	-	+	-	-	-	-	-	-	+	-	-	-	-	-	-	-	-	-	-	-	-	+	-	-	-
Sonorant	-	-	+	-	-	-	-	-	-	+	-	-	+	+	-	-	-	-	+	+	-	-	+	+	-	-
Strident	-	-	-	+	+	-	-	-	-	-	+	+	-	-	+	+	+	+	-	-	-	-	-	-	-	-
Voiced	-	+	+	-	+	-	+	-	+	+	-	+	+	+	-	+	-	+	+	+	-	+	+	+	-	-

二つの構造主義とその変遷

1800 1900 2000



フェルディナン・ド・ソシュール
(1857-1913)
一般言語学



レヴィ・ストロース
(1905-)
構造人類学・神話学

構造主義思想

ローマン・ヤコブソン
(1896-1982)
構造音韻論

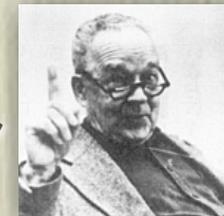
アメリカ構造言語学

ヨーロッパ構造言語学



エルンスト・マッハ
(1838-1916)
物理学・音響学・心理学

クリスチャン・フォン・エーレンフェルス
(1859-1932)
ゲシュタルト心理学



ジェームス・ギブソン
(1904-1979)
生態学的認識論・アフォーダンス



コネクショニスト

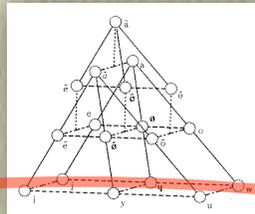
複雑系

二つの構造主義とその変遷

レヴィ・ストロース
(1905-)

構造人類学・神話学

構造主義思想



ローマン・ヤコブソン
(1896-1982)

構造音韻論

アメリカ構造言語学

ヨーロッパ構造言語学



フェルディナン・ド・ソシュール
(1857-1913)

一般言語学

1800

1900

2000

フェルディナン・ド・ソシュール(1857-1913)

「言語が含むのは、言語体系に先立って存在する観念でも音でもなく、ただこの体系から生じる観念的差異と音的差異だけである」

意味の世界が、各言語において、どのように領域分割しているのか？

音の世界が、各話者において、どのように領域分割しているのか？

要素の定義は、他者との差異を通して初めて可能になる。



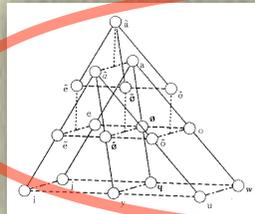
複雑系

二つの構造主義とその変遷

レヴィ・ストロース
(1905-)

構造主義思想

構造人類学・神話学



ローマン・ヤコブソン
(1896-1982)

アメリカ構造言語学

構造音韻論

ヨーロッパ構造言語学



フェルディナン・ド・ソシュール
(1857-1913)

一般言語学

1800

1900

2000



エルンスト・マッハ
(1838-1916)

物理学・音響学・心理学

クリスチャン・フォン・エーレンフェルス
(1859-1932)

ゲシュタルト心理学

ジェームス・ギブソン
(1904-1979)

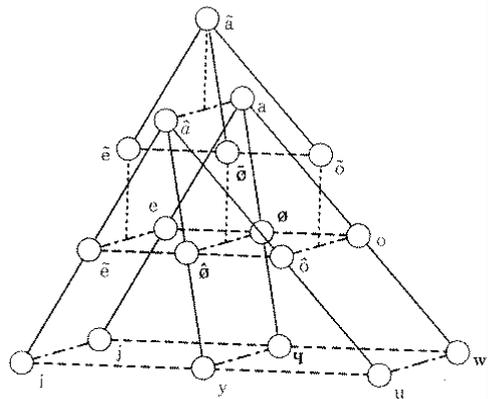
生態学的認識論・アフォーダンス



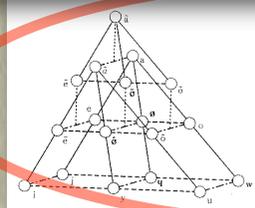
コネクショニスト



複雑系



	p	b	m	f	v	θ	ð	t	d	n	s	z	l	r	ʃ	ʒ	tʃ	dʒ	ɹ	ɻ	k	g	ŋ	w	ʔ	h
Back	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	+	+	+	+	+
High	-	-	-	-	-	-	-	-	-	-	-	-	-	-	+	+	+	+	+	+	+	+	+	+	-	-
Coronal	-	-	-	-	-	+	+	+	+	+	+	+	+	+	+	+	+	+	-	-	-	-	-	-	-	-
Anterior	+	+	+	+	+	+	+	+	+	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-
Labial	+	+	+	+	+	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	+	-	-
Continuant	-	-	-	+	+	+	+	-	-	-	+	+	+	-	+	+	-	-	+	+	-	-	-	+	-	+
Lateral	-	-	-	-	-	-	-	-	-	-	-	-	+	-	-	-	-	-	-	-	-	-	-	-	-	-
Nasal	-	-	+	-	-	-	-	-	-	+	-	-	-	-	-	-	-	-	-	-	-	-	-	+	-	-
Sonorant	-	-	+	-	-	-	-	-	-	+	-	-	+	+	-	-	-	-	+	+	-	-	+	+	-	-
Strident	-	-	-	+	+	-	-	-	-	-	-	+	+	-	-	+	+	+	+	-	-	-	-	-	-	-
Voiced	-	+	+	-	+	-	+	-	+	+	-	+	+	+	-	+	-	+	+	+	-	+	+	+	-	-



ローマン・ヤコブソン
(1896-1982)
構造音韻論

アメリカ構造言語学
ヨーロッパ構造言語学



フェルディナン・ド・ソシュール
(1857-1913)
一般言語学

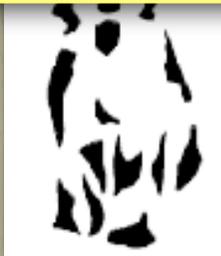


ローマン・ヤコブソン(1896-1913)

「二音素間の差異を幾つかの弁別素性を用いて表現」

音素群の幾何学構造

やがて音素そのものを素性の束として定義するが、これは勇み足か？



生態学的認識論・アン・オートタンス

コネクショニスト

複雑系

二つの構造主義とその変遷

1800



フェルディナン・ド・ソシュール
(1857-1913)
一般言語学

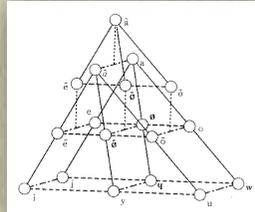


エルンスト・マッハ
(1838-1916)
物理学・音響学・心理学



クリスチャン・フォン・エーレンフェルス
(1859-1932)
ゲシュタルト心理学

1900



レヴィ・ストロース
(1905-)
構造人類学・神話学

構造主義思想

ローマン・ヤコブソン
(1896-1982)
アメリカ構造言語学
ヨーロッパ構造言語学

1900



ジェームス・ギブソン
(1904-1979)
生態学的認識論・アフォーダンス

コネクショニスト

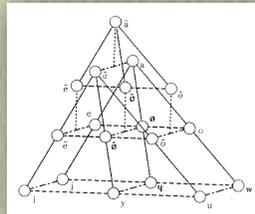
2000

複雑系

二つの構造主義とその変遷

レヴィ・ストロース
(1905-)
構造人類学・神話学

構造主義思想



ローマン・ヤコブソン
(1896-1982)
構造音韻論

アメリカ構造言語学

ヨーロッパ構造言語学



フェルディナン・ド・ソシュール
(1857-1913)
一般言語学

1800

1900

2000

レヴィ・ストロース(1905-)

「未開社会の親族構造の中に不変的な数学的構造を見いだした。また、抽象化を通して種々の神話の中に不変的な構造を見いだした」

構造=ある変換操作によって変化しない特性・関係

構造と変換は表裏一体であり、どの変換に着眼するかで構造の定義も異なる。



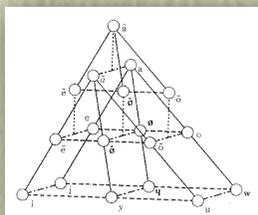
コネクション

複雑系

二つの構造主義とその変遷



フェルディナン・ド・ソシュール
(1857-1913)
一般言語学



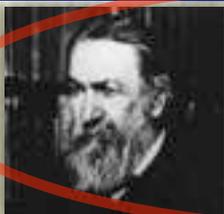
レヴィ・ストロース
(1905-)
構造人類学・神話学

構造主義思想

ローマン・ヤコブソン
(1896-1982)
構造音韻論

アメリカ構造言語学

ヨーロッパ構造言語学



エルンスト・マッハ
(1838-1916)
物理学・音響学・心理学



クリスチャン・フォン・エーレンフェルス
(1859-1932)
ゲシュタルト心理学



ジェームス・ギブソン
(1904-1979)
生態学的認識論・アフォーダンス

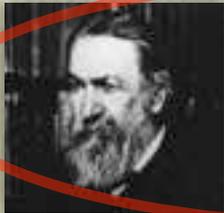
コネクショニスト

複雑系

二つの構造主義とその変遷

エルンスト・マッハ(1838-1916)

「物体の質量は，その物体の周りの全ての物体との関係で決る。他に何も無い空間の中では，ある物体の質量には，何の意味もない」
全ての自然現象は感性的諸要素の関数的相互依存関係に基づく複合体
全ての観測結果は人間の感覚量でしかない。空間・時間感覚も人間の感覚の一つ。
その意味において，ユークリッド空間よりも非ユークリッド空間の方が本質的。



エルンスト・マッハ
(1838-1916)
物理学・音響学・心理学



クリスチャン・フォン・エーレンフェルス
(1859-1932)
ゲシュタルト心理学



ジェームス・ギブソン
(1904-1979)
生態学的認識論・アフォーダンス

コネクショニスト

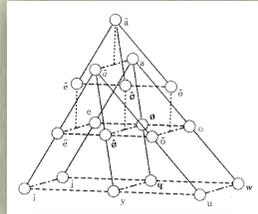
複雑系

二つの構造主義とその変遷

1800 1900 2000



フェルディナン・ド・ソシュール
(1857-1913)
一般言語学



レヴィ・ストロース
(1905-)
構造人類学・神話学

構造主義思想

ローマン・ヤコブソン
(1896-1982)
構造音韻論

アメリカ構造言語学

ヨーロッパ構造言語学



エルンスト・マッハ
(1838-1916)
物理学・音響学・心理学



クリスチャン・フォン・エーレンフェルス
(1859-1932)
ゲシュタルト心理学

ジェームス・ギブソン
(1904-1979)
生態学的認識論・アフォーダンス



コネクショニスト

複雑系

二つの構造主義とその変遷

クリスチャン・フォン・エーレンフェルス(1859-1932)

「知覚は、対象の個別刺激を統合して起こるものではなく、それ以前に全体的な枠組の中で認識が起こる」

ゲシュタルト＝全体＝「部分の単純な総和」以上のもの

移調性と不変性：ゲシュタルトにある種の変形を施しても、ゲシュタルトは不変
(メロディーと転調など)



エルンスト・マッハ
(1838-1916)
物理学・音響学・心理学

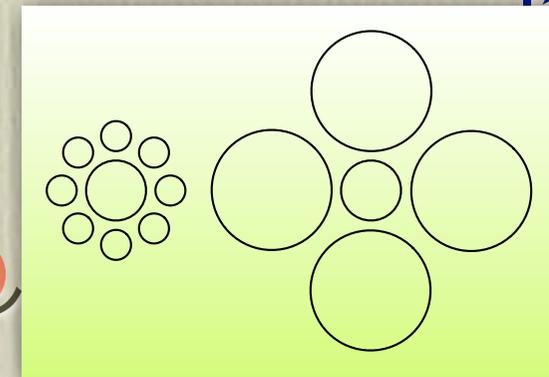


クリスチャン・フォン・エーレンフェルス
(1859-1932)

ゲシュタルト心理学

ジェームス・ギブソン
(1904-1979)

生態学的認識論・アフォーダンス



コネクショニスト

複雑系

二つの構造主義とその変遷

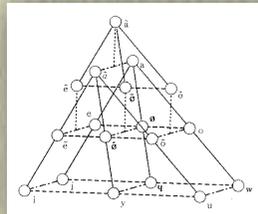
1800

1900

2000



フェルディナン・ド・ソシュール
(1857-1913)
一般言語学



レヴィ・ストロース
(1905-)
構造人類学・神話学

構造主義思想

ローマン・ヤコブソン
(1896-1982)
構造音韻論

アメリカ構造言語学

ヨーロッパ構造言語学



エルンスト・マッハ
(1838-1916)
物理学・音響学・心理学

クリスチャン・フォン・エーレンフェルス
(1859-1932)
ゲシュタルト心理学



ジェームス・ギブソン
(1904-1979)
生態学的認識論・アフォーダンス



コネクショニスト

複雑系

二つの構造主義とその変遷

ジェームス・ギブソン(1904-1979)

「視覚は『要素刺激感覚+そのまとめ上げ』ではなく、『対象の動き(変形)の中に見られる不変的特性』に視覚の本質がある」

変形と不変 ⇔ 移調性と不変性 ⇔ 変換と構造

面性の知覚：光が作る差異の構造に基づいて行なわれる（光の差異が成す不変項）。
情報は人間の頭の中で作り出すのではなく、環境そのものの中に存在している。



エルンスト・マッハ
(1838-1916)
物理学・音響学・心理学

クリスチャン・フォン・エーレンフェルス
(1859-1932)

ゲシュタルト心理学



ジェームス・ギブソン
(1904-1979)

生態学的認識論・アフォーダンス



コネクショニスト

複雑系

二つの構造主義とその変遷

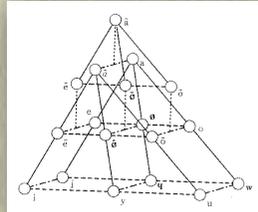
1800

1900

2000



フェルディナン・ド・ソシュール
(1857-1913)
一般言語学



レヴィ・ストロース
(1905-)
構造人類学・神話学

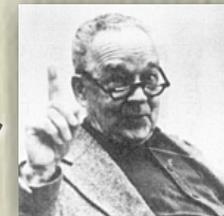
構造主義思想

ローマン・ヤコブソン
(1896-1982)
アメリカ構造言語学
ヨーロッパ構造言語学



エルンスト・マッハ
(1838-1916)
物理学・音響学・心理学

クリスチャン・フォン・エーレンフェルス
(1859-1932)
ゲシュタルト心理学



ジェームス・ギブソン
(1904-1979)
生態学的認識論・アフォーダンス



コネクショニスト

複雑系

二つの構造主義とその変遷

エルンスト・マッハ(1838-1916)

「私達は個々の要素についての研究から始めなければならないが、自然は要素と共に始まった訳ではない。私達が圧倒的な全体系から時々目をそらし、個別的な点に集中できるのは確かに幸福なことである。しかしさしあたり無視されていたことを、改めて補修し修正しながら研究することを怠ってはならない。」 (要素還元主義への警鐘)



エルンスト・マッハ
(1838-1916)
物理学・音響学・心理学



クリスチャン・フォン・エーレンフェルス
(1859-1932)

ゲシュタルト心理学

ジェームス・ギブソン
(1904-1979)

生態学的認識論・アフォーダンス



コネクショニスト

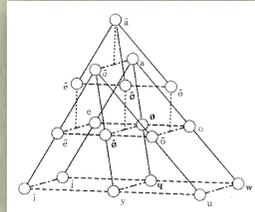
複雑系

要素感覚を集めるのではなく、要素間の関係・差異を統合することで
着眼する変換・変形操作に対して不変なる構造を見いだすこと。

レヴィ・ストロース
(1905-)

構造主義思想

構造人類学・神話学



ローマン・ヤコブソン
(1896-1982)

アメリカ構造言語学

構造音韻論

ヨーロッパ構造言語学



フェルディナン・ド・ソシュール
(1857-1913)

一般言語学

1800

1900

2000



エルンスト・マッハ
(1838-1916)

物理学・音響学・心理学

クリスチャン・フォン・エーレンフェルス
(1859-1932)

ゲシュタルト心理学



ジェームス・ギブソン
(1904-1979)

生態学的認識論・アフォーダンス



コネクショニスト

複雑系