

The difference between speech and voice

— Computational description and control of
sentiment information embedded in speech —

5 March 2019, at Chulalongkorn univ.

Yoshinori Sagisaka

Waseda University
Applied Math Dept. / GITI /
Language and Speech Science Res. Lab.

Waseda University

1

Speech and Voice

What is Speech ?

What is the difference between
Speech and Voice ?

Voice \supset Speech

Speech is Voice containing
linguistic information

What is **linguistic information** ?

Waseda University

2

Contents of this talk

Description & control for communicative speech

The analysis of N (uhm) to say “Good morning!”

N !! 🎵 NN !? N?? N ! NN~?

Sentiment description of speech by color

Can you see and say the color of vowels ?

/a/ /i/ /u/ /e/ /o/

Waseda University

3

Contents of this talk

Description & control for communicative speech

The analysis of N (uhm) to say “Good morning!”

N !! 🎵 NN !? N?? N ! NN~?

Sentiment description of speech by color

Can you see and say the color of vowels ?

/a/ /i/ /u/ /e/ /o/

Waseda University

4

Studies for synthesis and extraction of communicative information at my lab.

F0 generation for adv+adj speech (Sp. Prosody '04)

- Correlation between comm. F0 and lexicon

Specification of I/O for communicative prosody

(ICASSP'05, INTERSPEECH'05, Sp. Prosody '06)

- Input expression by multi-dim. impressions
- Prosody control scheme using lexical attributes

Extraction of speech impression (SNLP'07)

Universality of lexicon-prosody mapping (SNLP'09)

F0 estimation using lexical attributes (SNLP'13)

Waseda University

5

Contents of this talk

Description & control for communicative speech

- Prosodic differences in communication
- Sentiment description using impressions
- Communicative F0 control using impressions

Sentiment description of speech by color

Can you see and say the color of vowels ?

/a/ /i/ /u/ /e/ /o/

Waseda University

6

Prosodic differences in communication

Communicative speech has ! 🎵, 😊 and ●

Greetings Good morning ! 🎵
おはようございます ! 🎵
สวัสดี ! 🎵

Thanks Thank you very much 😊
どうもありがとうございます 😊
ขอบคุณมาก 😊

Apologies I am very sorry ●
どうもすみません ●
ฉันเสียใจมาก ●

Waseda University

7

Need for Input/Output specification for communicative speech synthesis

To synthesize communicative speech

Thank you very much 😊

どうもありがとうございます 😊

ขอบคุณมาก 😊

We need extra-info 😊 for a synthesis system

• Linguistic content (text) + ?

Beyond text-to-speech synthesis

• Using what input ?

• How to control by what factors ?

Waseda University

8

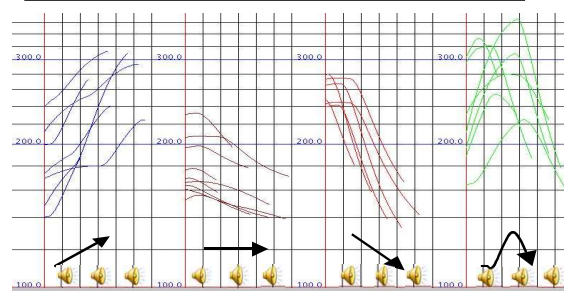
Information conveyed by prosody

utterance	information	expressions
A) N !! 🎵	Praise	(This is nice !)
B) NN !?	Doubt	(Really ?)
A) N ??	Incomprehension	(Why not ?)
A) N !	Attention	(Look at this !)
B) NN ~ 😊	Agreement	(I see.)

Waseda University

9

Four typical F0 patterns of “n” observed in real daily conversations



Waseda University

10

Prototypical F0 patterns by height & dynamics

Dynamics pattern	rise	flat	fall	rise&fall
height				
high	Really? ♪ Then, what happen?? ♪ ♪ Liar!! ♪ Really? Is that true? I didn't know that.	And? And? ♪ ♪ Then what? ♪ I agree!! Okay. I don't know... Is that okay?	Real? I didn't know that!! ♪ Of course!! ♪ That is nice ♪ That is fine. Okay... ↓	Never mind! ♪ That is okay.
low	Really ↓	I am not sure... I don't quite agree...	I didn't know that... ↓ Yes... but... ↓ ↓	No, I don't like that ↓

Waseda University

11

Word description for perceptual impression

Speech 12 prototypical single utterances “n”
F0 average height (3) × F0 dynamics (4)

Listeners 5 Adults (2 male, 3 female)

Word description for impression

Possibly followed phrases or words
Imagined speaker's attitudes

Waseda University

12

Perceptual impressions

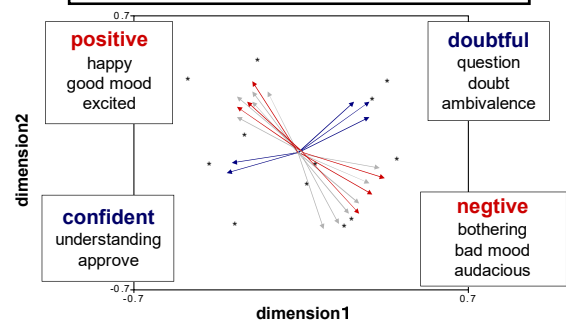
26 descriptions for perceptual impressions

- | | | |
|------------------|--------------------|---------------|
| 1: doubt | 10: dark | 19: gentle |
| 2: ambivalence | 11: cheerful | 20: audacious |
| 3: deny | 12: weak | 21: exciting |
| 4: question | 13: interested | 22: bothering |
| 5: objection | 14: not interested | 23: delight |
| 6: approve | 15: good mood | 24: anger |
| 7: agreement | 16: bad mood | 25: happy |
| 8: understanding | 17: light | 26: annoying |
| 9: bright | 18: heavy | |

Waseda University

13

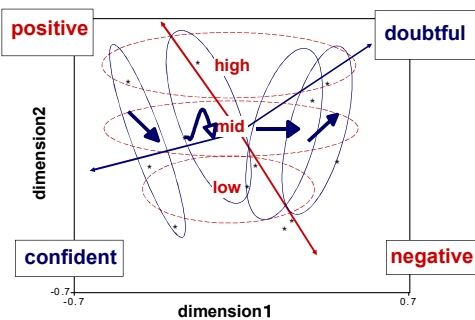
Word vectors in dim1&dim2-plane



Waseda University

14

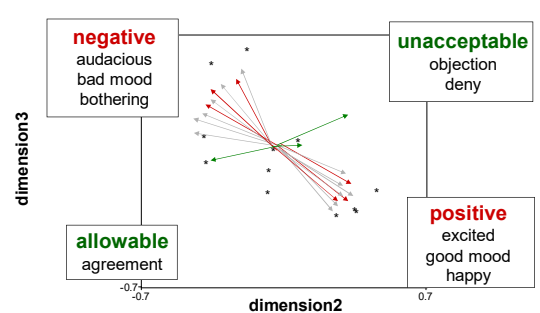
F0 height and dynamics in dim1&dim2-plane



Waseda University

15

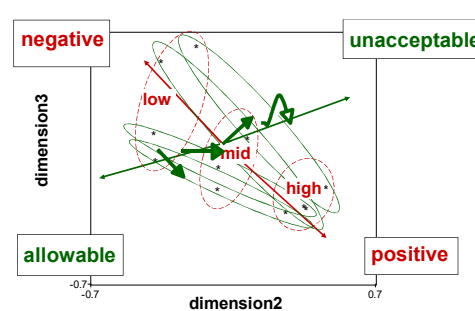
Word vectors in dim2&dim3-plane



Waseda University

16

F0 height and dynamics in dim2&dim3-plane

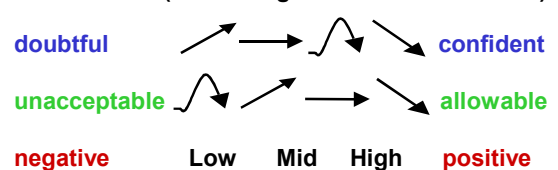


Waseda University

17

Three dimensional expressions for perceptual impressions

F0 characteristics corresponds to three dimensional perceptual impressions (Greenberg et al. IEEE ICASSP'05)



Waseda University

18

Computation of communicative prosody using input lexical attributes

speaker utterance possible expressions

- | | | |
|----|--------|----------------|
| A) | N !! 🎵 | This is nice ! |
| B) | NN !? | Really ? |
| A) | N?? | Why not ? |
| A) | N ! | Look at this ! |
| B) | NN ~ 😊 | I see. |

Communicative prosody control using lexical attributes
Shao, Greenberg and Sagisaka (SNLP'13)

Waseda University

19

Prosody generation for text-to-speech

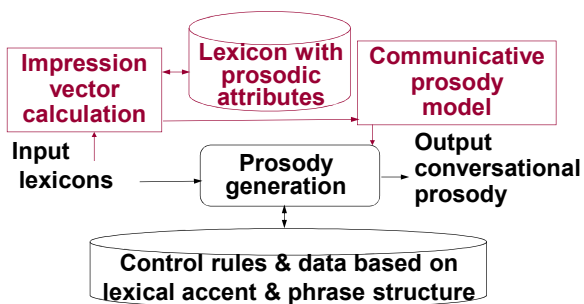
Input lexicons → Prosody generation → Output read prosody

Control rules & data based on lexical accent & phrase structure

Waseda University

20

Communicative prosody generation



Waseda University

21

Example of impression vector computation

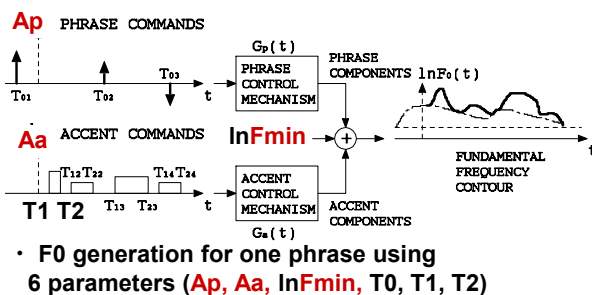
Input	これ	本当に	きれいだ	なの
	C'est	vraiment	beau	?
	this is	really	beautiful	?
Impression	(degree)	positive	doubtful	
value		v1 = +3	v2 = +1	v3 = -4

1st dim. doubtful-confident v3 = -4
 2nd dim. unacceptable-allowable unassigned 0
 3rd dim. negative - positive v1 × v2 = 3
 Normalized impression vector = (-0.8, 0, 0.6)

Waseda University

22

F0 generation by command-response model

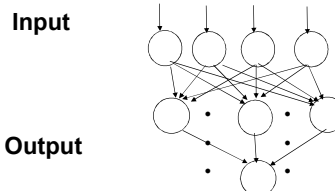


Waseda University

23

Communicative F0 control parameter generation from an impression using a NN

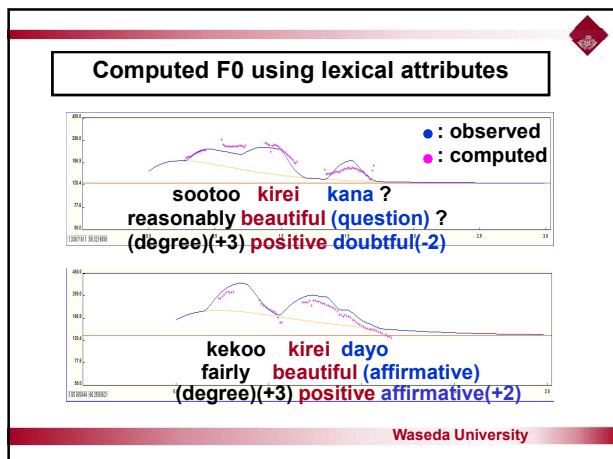
Impression vector → Reading F0 generation parameters (Ap Aa Fmin)



Communicative F0 generation parameters (Ap Aa Fmin)

Waseda University

24



25

Contents of this talk

Description & control for communicative speech

The analysis of N (uhm) to say “Good morning!”
N !! ♪ NN !? N?? N ! NN~?

Sentiment description of speech by color

Can you see and say the color of vowels ?
/a/ /i/ /u/ /e/ /o/

Waseda University

26

Contents of this talk

Description & control for communicative speech

The analysis of N (uhm) to say “Good morning!”
N !! ♪ NN !? N?? N ! NN~?

Sentiment description of speech by color

- Problems in the description of speech & color
- Introduction of sentiment analysis of multi-media
- Sentiment correlations between speech & color
- Communicative speech description using color

Waseda University

27

Multi-dimensional description for perceptual information

Description using multiply chosen 26 words

doubtful – confident
question, doubt, ambivalence, understanding, approve

unacceptable – allowable
deny, objection, agreement

negative – positive
dark, weakly, not interested, bad mood, heavy, bothering, audacious, anger, annoying, cheerful, delight, gentle, good mood, excited, happy, light, interested, bright

Waseda University

28

Problems in the description of speech & color

Speech description using impressions
Effective, however,
Fundamentally difficult using discrete symbols

Color has the same description problem using language form
e.g. burgundy red, emerald green, navy blue
→ Limitation using language form

Mapping between speech and color
A new possibility for multi-modal info processing

Waseda University

29

Sentiment analysis studies on correlations between different media

Correlation between **images** and **words**
Bouba / kiki effect (W. Köhler 1929)
Sport onomatopoeia (Fujino 2008)

Correlation between **sound** and **color**
Non-verbal mapping (Nagata 2003)
Analysis of synesthesia, an ability to hear color

Studies on **color** association from **speech** input
Color imaging from vowels (Wrembel 2008)
Mapping from speech to color (Watanabe 2014)

Waseda University

30

Sentiment analysis studies on correlations between different media

Correlation between images and words
Bouba / kiki effect (W. Köhler 1929)
 Sport onomatopoeia (Fujino 2008)

Correlation between sound and color
 Non-verbal mapping (Nagata 2003)
 Analysis of synesthesia, an ability to hear color

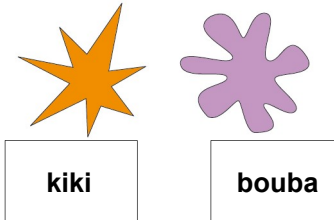
Studies on color association from speech input
 Color imaging from vowels (Wrembel 2008)
 Mapping from speech to color (Watanabe 2014)

Waseda University

31

Bouba / kiki effect

Which figure corresponds to bouba & kiki ?
Wolfgang Köhler 1929



Waseda University

32

Sentiment analysis studies on correlations between different media

Correlation between images and words
 Bouba / kiki effect (W. Köhler 1929)
Sport onomatopoeia (Fujino 2008)

Correlation between sound and color
 Non-verbal mapping (Nagata 2003)
 Analysis of synesthesia, an ability to hear color

Studies on color association from speech input
 Color imaging from vowels (Wrembel 2008)
 Mapping from speech to color (Watanabe 2014)

Waseda University

33

Sport onomatopoeia

Different F0 patterns of tsukuri, kuzushi and kake in Judo



Fujino 2008 (Shoogakukan)

Waseda University

34

Sentiment analysis studies on correlations between different media

Correlation between images and words
 Bouba / kiki effect (W. Köhler 1929)
 Sport onomatopoeia (Fujino 2008)

Correlation between sound and color
Non-verbal mapping (Nagata 2003)
 Analysis of synesthesia, an ability to hear color

Studies on color association from speech input
 Color imaging from vowels (Wrembel 2008)
 Mapping from speech to color (Watanabe 2014)

Waseda University

35

Non-verbal mapping between sound & color

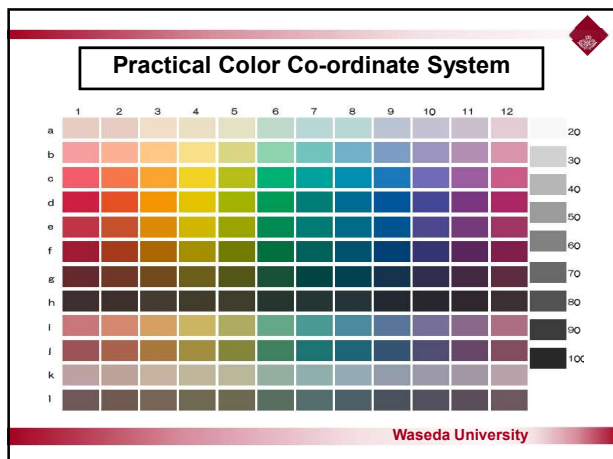
Non-verbal mapping **sound to color (Nagata 2003)**
 Compared synesthesia and ordinary people

Sound stimuli consisting of
 Key : C major to B major, C minor to B minor
 Timbre : Pure tone and various types of harmonics
 Height : Rising scale of D major

Color selection
 153 Hues from
 PCCS(Practical Color Co-ordinate System)

Waseda University

36



37

Sentiment analysis studies on correlations between different media

- Correlation between images and words
 - Bouba / kiki effect (W. Köhler 1929)
 - Sport onomatopoeia (Fujino 2008)
- Correlation between sound and color
 - Non-verbal mapping (Nagata 2003)
 - Analysis of synesthesia, an ability to hear color
- Studies on color association from speech input
 - Color imaging from vowels (Wrembel 2008)
 - Mapping from speech to color (Watanabe 2014)

Waseda University

38

Color imaging from vowels

After listening vowel sounds
Imagined color category was selected

Selected tendency

Vowel	Hue
front /i/	→ yellow
mid /e/	→ green
back /u/ /o/	→ brown blue black
open /a/	→ red

(M. Wrembel et al, "Sounds like a rainbow" ISCA WS 2008)

Waseda University

39

Sentiment analysis studies on correlations between different media

- Correlation between images and words
 - Bouba / kiki effect (W. Köhler 1929)
 - Sport onomatopoeia (Fujino 2008)
- Correlation between sound and color
 - Non-verbal mapping (Nagata 2003)
 - Analysis of synesthesia, an ability to hear color
- Studies on color association from speech input
 - Color imaging from vowels (Wrembel 2008)
 - Mapping from speech to color (Watanabe 2014)

Waseda University

40

Mapping from speech to color

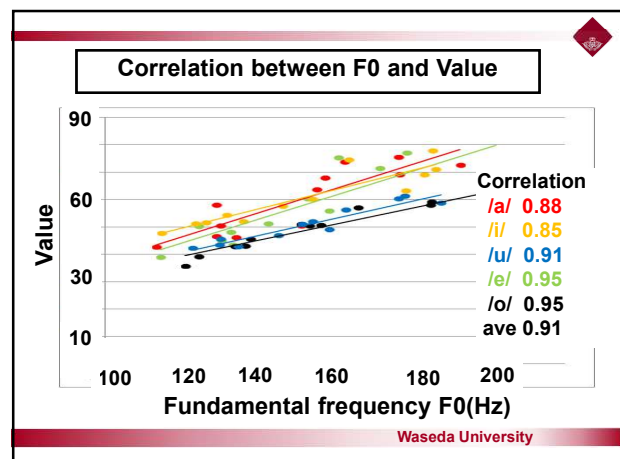
After listening vowel sounds (Watanabe 2014)
Imagined color category was selected

Speech stimuli
Communicative speech consisting of 5 Japanese vowels /a/, /i/, /u/, /e/ and /o/ with 12 F0 patterns (3 heights × 4 dynamics)

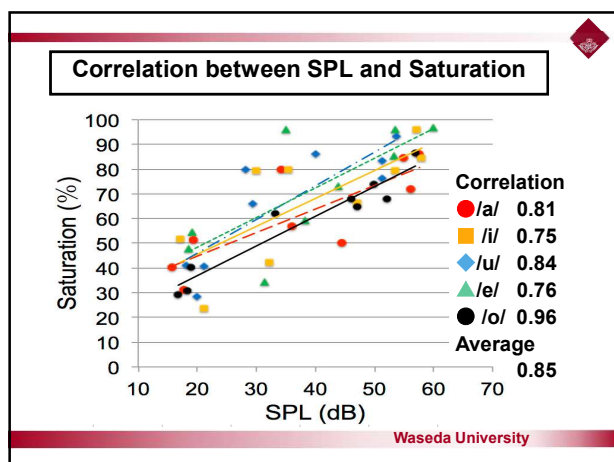
Color selection
153 Hues from PCCS(Practical Color Co-ordinate System)

Waseda University

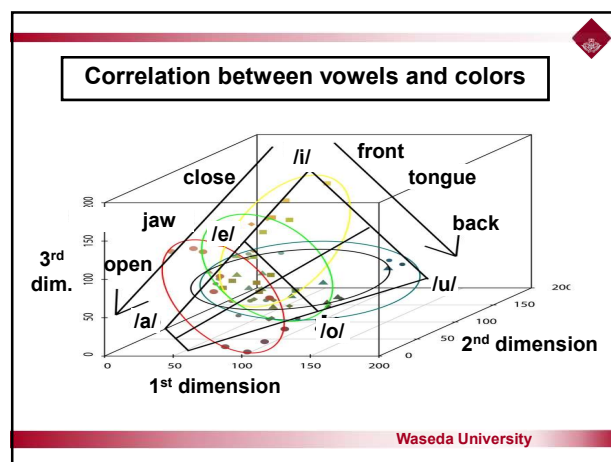
41



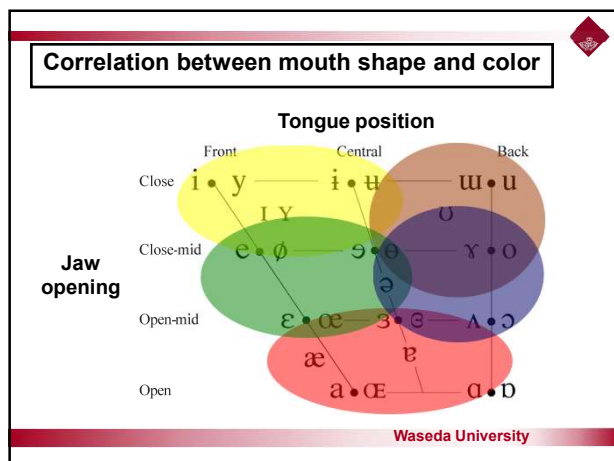
42



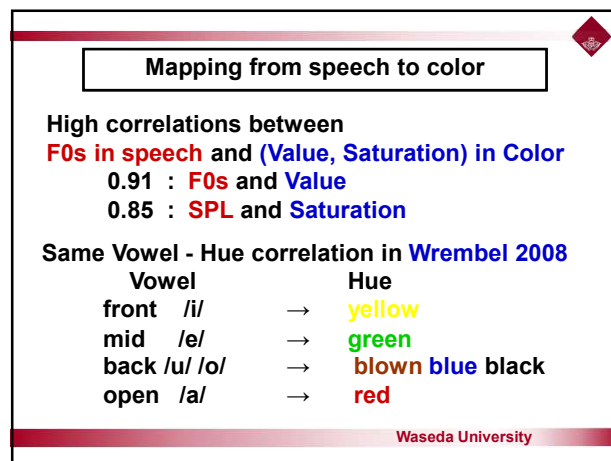
43



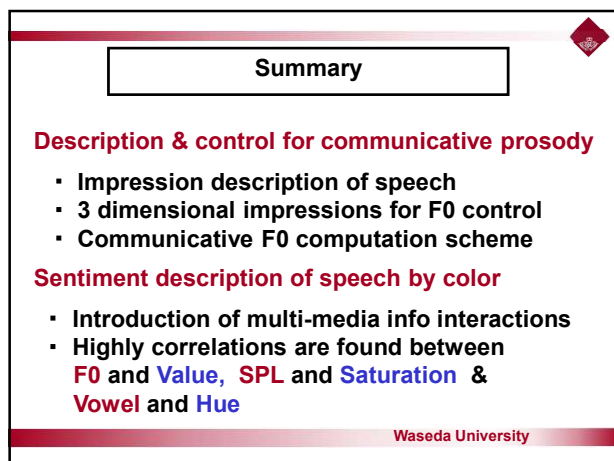
44



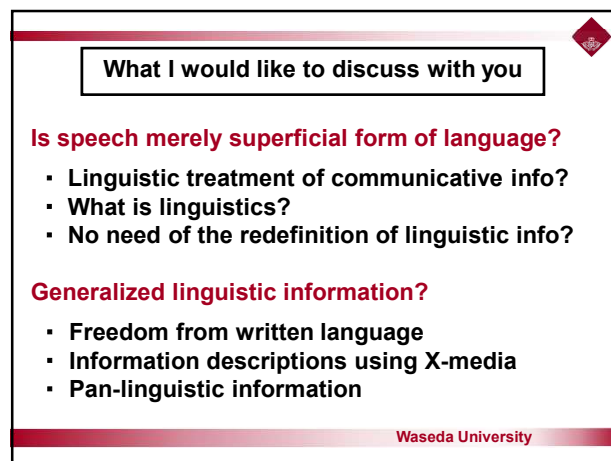
45



46



47



48

Acknowledgements

Thanks to all the current and previous
lab members

Special thanks to collaborators

Yoko Greenberg

Yoshitaka Fujino

Lu Shao

Kanako Watanabe

Waseda University

49

Waseda University

50