How can speech technologies support learners to improve their skills of speaking, listening, conversation, and more?

Nobuaki MINEMATSU Graduate School of Engineering, The University of Tokyo



Biography of Nobuaki MINEMATSU



Nobuaki MINEMATSU earned the doctor of Engineering in 1995 from UTokyo and since 2012, he has been a professor there. From 2002 to 2003, he was a visiting researcher at KTH, Sweden. He has a wide interest in speech communication covering speech science and speech engineering, especially he has an expert knowledge on Computer-Aided Language Learning (CALL). When he was a high-school student, he wanted to be a teacher of English, and when he was a university student, he was an amateur actor on English stages. He has published more than 450 journal and conference papers and received paper awards from RISP, JSAI, ICIST, O-COCOSDA, IEICE and an encouragement award from PSJ. He gave tutorial and invited talks on CALL at conferences such as APSIPA2011, INTERSPEECH2012, O-COCOSDA2014, and CASTEL/J2017. He was a distinguished lecturer of APSIPA from 2015 to 2016. He served as secretary of Speech Prosody 2004 and INTERSPEECH2010, co-organizer of SLaTE2010, and program chair of O-COCOSDA2018. He is the general chair of Speech Prosody 2020.



Outline of the presentation

CALL for speaking (reading aloud), listening, conversation, and more Computer-Aided Language Learning with speech technologies

with speech synthesis technologies



with speech analysis technologies





with speech recognition technologies

with new speech technologies being developed in our new project



Outline of the presentation

CALL for speaking (reading aloud), listening, conversation, and more Computer-Aided Language Learning with speech technologies

with speech synthesis technologies



with speech analysis technologies





with speech recognition technologies

with new speech technologies being developed in our new project



Japanese prosody — lexical level —

Figure Fi

♀ From the 1st mora to the 2nd mora, the pitch level generally goes up (L → H).
♀ The pitch level goes down somewhere in the word, never goes up again.
♀ Word accents are classified based on the falling position of pitch (→ accent type)
♀ Accent nucleus = the previous mora (syllable) before the pitch downfall.



<u> </u>	о Ф	900	90	9
さんがつ	ひこーき	かんごふ	いもーと	おはフ
頭高型	日間	事型	尾高型	平板
initial high	middle	e high	tail high	unacce
	起仇	犬式		平板
	acce	nted		unacce
1型 type 1	2型 type 2	3型 type 3	4型 type 4	0 캡 type
-4型	-3型	-2型	-1型	
And the Public Design in the Real Public Design in the Public Design in	And the second sec			



Notorious word accent sandhi (change) in Japanese

Version It is true that each word has its own accent type. **Examples of accent changes when speaking** \bigcirc A noun + another = a compound noun Verb conjugation ◎ あるく → あるきます, あるいて, あるい \bigcirc A bunsetsu + another = an accentual phrase \bigcirc http:// http://

Solution However, it's also true that accent control is done on a phrase level when speaking. Solution of the terminal of terminal o

IS HILL COMPANY AND DECY	0000	0000	0000	0000	0
	さんがつ	ひこーき	かんごふ	いもーと	お
	頭高型	中語	高型	尾高型	3
	initial high	middle	e high	tail high	una
	2	起位	犬式		2
		acce	nted		una
た、あるかない	1型 type 1	2型 type 2	3型 type 3	4型 type 4	t
raço	-4型	-3型	-2型	-1型	

かれは + たべる → かれはたべる

Lexical accent control of Japanese is **SOOOO MYSTERIOUS!!**





A fact that is not rare.

An email from a Canadian user of our system, OJAD

Dear OJAD,



I live in Canada and have been studying Japanese for 10 years - 4 of which were at university. I also worked in Japan for two years. And I had NEVER heard of Pitch Accent - I knew from my Japanese friends that there were different ways to say "hashi" and "ame", but I couldn't hear the difference, and I didn't know that pitch accent actually is a feature that ALL Japanese words have. I finally came across the topic of pitch accent when I started searching for pronunciation books while in Japan, and found a Japanese book that talked about 高低アクセント. I was intrigued and started looking online for more information.

Now that I'm back home and continuing my self-study of Japanese, I have discovered pitch accent, and it makes SOOOO much more sense to me! I always wondered how I could sound "more Japanese" and get rid of my painfully obvious foreign accent. And I knew that Japanese people sounded different than me, but I didn't know what it was or how to train myself to copy it. Thanks to finding online resources like your website, I am enjoying learning about Japanese and am now able to hear pitch changes in Japanese. I also feel like my pronunciation is improving. And of all the resources I've found online, your website is definitely my favorite! *^v^*

I just wanted to say THANK YOU for creating your wonderful website!





Japanese prosody — phrase and sentence level —

An interesting example of comparison between Chinese and Japanese Pitch changes acoustically observed in a Chinese utterance (weather forecast)





Solution Pitch changes acoustically observed in a Japanese utterance (weather forecast)







Two kinds of language teachers

Those teaching to human learners and those to machine learners More-than-50-year history of teaching Japanese prosody to machine learners

おはようございます









M. Suzuki, et al., "Accent sandhi estimation of Tokyo dialect of Japanese using conditional random fields," Trans. IEICE, E100-D, 4, 655-661, 2017
N. Minematsu, et al., "Development and evaluation of online infrastructure to aid teaching and learning of Japanese prosody," Trans. IEICE, E100-D, 4, 662-669, 2017



1.5-min promotion video for Suzuki-kun of OJAD

Suzuki-kun = prosodic reading tutor of Tokyo Japanese in OJAD ♀ "The first and only teaching material to explain prosodic control of TJ for any given text."





1.5-min promotion video for Suzuki-kun of OJAD Suzuki-kun = prosodic reading tutor of Tokyo Japanese in OJAD ♀ "The first and only teaching material to explain prosodic control of TJ for any given text." ato Bahasa Jepang ke-5







1.5-min promotion video for Suzuki-kun of OJAD Suzuki-kun = prosodic reading tutor of Tokyo Japanese in OJAD "The first and only teaching material to explain prosodic control of TJ for any given text."

140 tutorial workshops in 40 countries



Outline of the presentation

CALL for speaking (reading aloud), listening, conversation, and more Computer-Aided Language Learning with speech technologies

with speech synthesis technologies



with speech analysis technologies





with speech recognition technologies

with new speech technologies being developed in our new project



Outline of the presentation

CALL for speaking (reading aloud), listening, conversation, and more Computer-Aided Language Learning with speech technologies



with speech synthesis technologies



with speech analysis technologies





with speech recognition technologies

with new speech technologies being developed in our new project



Real-world conditions of learners' listening

Seconditions of general listening material for learners Monologues and dialogues in a clean (no noise) condition. Subset of the second Solution But noise level or acoustic distortion involved is generally very mild. Real-world conditions that learners will be faced with Announcements in a train or bus background noises Utterances from very tall or small speakers Animation voices designed with a voice changer characters Speeches heard in a very big hall reverberation (echo) **Q** Utterances transmitted through a radio channel channel distortion **Necessity of robust listening** Acoustics and speech quality can be easily degraded with various factors. Solution Natives can listen but learners may have severe difficulty in listening to those utterances.

age and gender (vocal tract length)



A honest confession from a young Japanese pilot

命がけのリスニング

何か新しいことにチャレンジしてみようってことで、軽い気持ちで飛行機の練習をはじめたのが半年前。 泥沼にはまりつつも、なんとかもう少しで取れそうなところまで来た。

まだ終わったわけではないのだけど、一番大変だったのは無線交信。 予想以上にパイロットはしゃべる仕事だということを痛感。

特に、私が練習している空域は、北米でも有数の混雑地帯で、無線交信の量がはんぱじゃない。 飛行中は、常に誰かがしゃべっているという感じ。そして、管制塔の指示がわからなかったら、最悪の場 合、衝突の可能性もあるわけで、まさに命がけ。

それなのに、

- ・無線なのでノイズが大きい
- ・管制官が早口で訛っている
- ・パイロットも訛っている
- ・エンジンの音が大きい
- ・無線機がボロくて、たまに聞こえなくなる
- ・エリアによっては妨害電波が出ているみたい
- ・コクピットの中はただでさえ緊張して、頭の中が白くなる



desperate efforts needed for listening to air traffic controllers



ヽ)中が白くなる





A training method for robust listening

High Variability Phonetic Training (HVPT) Solution Section Content and C Speakers, speaking style, gender, age, accents, background noises, etc Often used by language teachers with good knowledge of phonetics Many papers or reports of the effectiveness of HVPT Lively+1993, Masuda+2012, Wong+2014, Hwang+2015, Solution For the second state of the second st **Technically-enhanced HVPT** Speech analysis-resynthesis technologies

♀ can convert a single utterance into acoustically various versions with its message unchanged. Solution Usability or validity of training with artificially converted audio samples Not only for dictation tasks but also for comprehension tasks

H. Zhang, et al., "Computer-aided high variability phonetic training to improve robustness of learners' listening comprehension," Proc. ICPhS, 2019



Examples of speech conversion

Variously converted speech can be obtained easily. "February 14th is a day for people who are falling in love." **Original** VTL x 1.5 (giant), VTL / 1.5 (fairy) **Reverb** a big cathedral **Noise** babble noise (voice noise) 2G mobile phone, air traffic control (ATC) **Channel Generation** With quantitative control of degree of distortion A small girl is praying in a cathedral, surrounded by chatty tourists and her pray is recorded and transmitted via a 2G mobile phone network.





Specific types of distortion with little troubles to native listeners but big troubles to non-native listeners should be good material for robust listening training?







Very harsh EIKEN grade 2 listening test

4-choice questions after listening to monologues or dialogues \bigcirc Male \rightarrow giant or giant pilot (ATC) \bigcirc Female \rightarrow fairy or fairy pilot (ATC)





Accuracy of Japanese college students and native speakers

TOEIC	original	G/F	ATC	G/F + ATC
400-600	58.3			
600-800	78.2			
800-990	81.5			
Native				

Question: What is one thing the girl says?

- 1 She is not good at sports.
- 2 She will not go to college.
- 3 She needs more time to study.
- She wants to practice basketball more. 4



Very harsh EIKEN grade 2 listening test

4-choice questions after listening to monologues or dialogues \bigcirc Male \rightarrow giant or giant pilot (ATC) \bigcirc Female \rightarrow fairy or fairy pilot (ATC)





Accuracy of Japanese college students and native speakers

TOEIC	original	G/F	ATC	G/F + ATC
400-600	58.3	50.0	30.6	32.8
600-800	78.2	62.0	35.1	23.4
800-990	81.5	79.6	45.4	25.0
Native				

Question: What is one thing the girl says?

- 1 She is not good at sports.
- 2 She will not go to college.
- 3 She needs more time to study.
- She wants to practice basketball more. 4



Very harsh EIKEN grade 2 listening test

4-choice questions after listening to monologues or dialogues \bigcirc Male \rightarrow giant or giant pilot (ATC) \bigcirc Female \rightarrow fairy or fairy pilot (ATC)





Accuracy of Japanese college students and native speakers

TOEIC	original	G/F	ATC	G/F + ATC
400-600	58.3	50.0	30.6	32.8
600-800	78.2	62.0	35.1	23.4
800-990	81.5	79.6	45.4	25.0
Native	100	100	100	93.6

Question: What is one thing the girl says?

- 1 She is not good at sports.
- 2 She will not go to college.
- 3 She needs more time to study.
- She wants to practice basketball more. 4



Harsh listening exam -> harsh listening drills

Hash listening training drills were developed using ATC distortions.

Harsh listening test of EIKEN grade 2 **Pre**: Solution Fraining: Listening drills with varying degrees of ATC distortions only

Harsh listening test of EIKEN grade 2 (= Pre) **Post:**

July

Pre

1. Listening robustness is improved against ATC distortion? 2. Listening robustness is transferred to other kinds of distortion?

Mid Dec End of Dec

3-week Training

Post



§ 18-day listening drills of different levels of ATC

🔮 音声ファイル

・音声ファイルをダウンロードすると、雑音レベルに合わせてlevel0~level3というフォルダが生成され、その下 に音声ファイルが保存されます。各音声に対する質問は問題セットをクリックして下さい。 ・雑音レベルを3段階用意していますが、level1でも難しすぎる場合はlevel1の難度を落としますので、その旨連 絡ください(level0.2, level0.5 のサンプルを上に掲載しています)。 連絡先は<u>こちら</u>。 ・音声ファイル (mp3) は10個ずつまとめて, zip ファイルとして提供しています。 スマホの場合, zip ファイル を解凍できないかもしれません。その場合は、以下のツールをダウンロードして zip ファイルを解凍して下さ v، iPhone : <u>iZip</u> Android: <u>file manager</u>

問題セット	音声ファイ	ル (0: 雑音	な
<u>DAY01 (11/27)</u>	<u>0</u>	<u>1</u>	<u>2</u>
<u>DAY02 (11/28)</u>	<u>0</u>	1	<u>2</u>
<u>DAY03 (11/29)</u>	<u>0</u>	<u>1</u>	<u>2</u>
<u>DAY04 (11/30)</u>	<u>0</u>	<u>1</u>	<u>2</u>
<u>DAY05 (12/1)</u>	<u>0</u>	<u>1</u>	<u>2</u>

July



し~3: 雑音MAX)



Mid Dec End of Dec



Pre → Drill → Post

Accuracy of pre-test and post-test

A half of the pre test examinees (55 students) undertook the post test.

Part	TOEIC	Ν	Orig.	GF	ATC	GF+ATC	Pa	rt	TOEIC	N	Orig.	GF	ATC	GF+AT	С
A	400-600	15	66.7	48.3	25.0	41.7	A		400–600	15	70.0	66.7	26.7	35.0	
	600-800	32	77.3	65.6	38.3	25.8			600-800	32	73.4	73.4	40.6	32.8	
	800–990	8	84.4	84.4	43.8	21.9			800–990	8	96.9	96.9	75.0	40.6	
В	400-600	15	50.0	43.3	28.3	23.3	E		400–600	15	66.7	48.3	38.3	23.3	
	600-800	32	65.6	48.4	39.1	30.5			600-800	32	61.7	51.6	42.2	35.2	
	800–990	8	78.1	62.5	37.5	28.1			800–990	8	87.5	84.4	62.5	31.3	

Error reduction rate

- ✓ Accuracy: 70% → 85 %
- \Rightarrow ERD = (30-15)/30 = 50 %
- A: monologue, B: dialogue

July

Part	TOEIC	Ν	Orig.	GF	ATC	GF+ATC
A	400–600	15	9.9	35.6	2.3	-11.5
	600-800	32	-17.2	22.7	3.7	9.4
	800–990	8	80.1	80.1	55.5	23.9
В	400–600	15	33.4	8.8	13.9	0
	600-800	32	-11.3	6.2	5.1	6.8
	800–990	8	42.9	58.4	40.0	4.5

Mid Dec End of Dec



Outline of the presentation

CALL for speaking (reading aloud), listening, conversation, and more Computer-Aided Language Learning with speech technologies



with speech synthesis technologies



with speech analysis technologies





with speech recognition technologies

with new speech technologies being developed in our new project



Outline of the presentation

CALL for speaking (reading aloud), listening, conversation, and more Computer-Aided Language Learning with speech technologies



with speech synthesis technologies



with speech analysis technologies





with speech recognition technologies

with new speech technologies being developed in our new project



Conversation is a multi-task speech activity.

Listening, understanding, and speaking running almost together









Conversation is a multi-task speech activity.

Listening, understanding, and speaking running almost together

Shadowing is a multi-task speech training. A special form of listen-and-repeat practice, with as short delay as possible







learner



Conversation is a multi-task speech activity.

Listening, understanding, and speaking running almost together

Shadowing is a multi-task speech training. A special form of listen-and-repeat practice, with as short delay as possible



native

ss=smoothness of shadow





with ASR

learner



Data collection and teachers' manual rating

Collection of samples from 125 university or college students

- \bigcirc 4 passages = 55 sentences
- Sour repetitions and 27,500 utterances all together
- Sentence selection for manual rating
 - 9 10 sentences were selected based on syntactic complexity and pronunciation difficulty.
 - Solution Fourth shadowings were rated manually by a unit of clause.
 - Strategies for rating
 - How correctly phonemes are produced (P).
 - We have correctly prosody is produced (S = Supra-segmental = Prosody).
 - Whether each word sounds as if it is produced after identifying that word (C = Correctness). \bigcirc 5-step scale (1–5) and the total score (P+S+C) varies from 3 to 15.
 - Raters 9
 - General General (AE+J) teachers of English



Spectrogram is converted to posteriogram

Phoneme posterior probabilities calculated by DNN



W DNN can be viewed as strong abstraction.

Spectrogram is acoustic representation, including extra-linguistic features. See Posteriogram is phonetic/phonemic representation, suppressing those features.

Hidden Layers



DNN-based calculation of GOP

GOP = Goodness Of Pronunciation ğ

+phonemes intended by the model speaker

ime	Frame	Phoneme
t	1	а
	2	а
*	3	u
	•••	•••
	1232	sil





 $0.8 + 0.7 + 0.4 + \dots + 0.9$ GOP 0.63= 1232 **DNN-GOP**



Another method for utterance comparison

The two utterances are compared.

- Question of the Warping (DTW)
 - Alignment of two sequences of different length
- - Spectrum vectors are sensitive to age, gender, etc.



OTW-based comparison between the two



Correlations bet. human scores and machine scores

Sentence-based and speaker-based rating scores Solution of the second Regression model to predict human scores Solution States of DNN-GOP and some other features or scores are prepared for regression.

Table 2. Feature-based correlations with teachers' scores										
features	Р	S	С	P+S+C						
bGOP [16]	0.74	0.83	0.71	0.83						
pGOP	0.79	0.84	0.78	0.88						
vGOP	0.70	0.83	0.70	0.81						
cGOP	0.79	0.82	0.78	0.87						
v1GOP	0.63	0.78	0.64	0.75						
v2GOP	0.42	0.41	0.43	0.46						
v0GOP	0.71	0.75	0.78	0.78						
DNN-DTW	-0.66	-0.84	-0.69	-0.80						
RS	-0.34	-0.21	-0.29	-0.30						
WRR	0.79	0.81	0.71	0.84						

S. Kabashima, et al., "DNN-based scoring of language learners' proficiency using learners' shadowings and native listeners' responsive shadowings," Proc. Spoken Language Technology, 2018

Table 3. Model-based correlations in a speaker level										
models	Р	S	С	P+S+C						
bGOP [16]	0.74	0.83	0.71	0.83						
Lasso	0.84	0.89	0.76	0.90						
SVR	0.85	0.89	0.83	0.89						
Random Forest	0.77	0.84	0.79	0.86						
inter-rater	0.77	0.69	0.86	0.87						

Table 4. Model-based correlations in a sentence level

models	Ρ	S	С	P+S+C
Lasso	0.68	0.73	0.65	0.77
SVR	0.70	0.73	0.68	0.78
Random Forest	0.67	0.68	0.61	0.74
inter-rater	0.58	0.54	0.74	0.75



Outline of the presentation

CALL for speaking (reading aloud), listening, conversation, and more Computer-Aided Language Learning with speech technologies



with speech synthesis technologies



with speech analysis technologies





with speech recognition technologies

with new speech technologies being developed in our new project



Outline of the presentation

CALL for speaking (reading aloud), listening, conversation, and more Computer-Aided Language Learning with speech technologies



with speech synthesis technologies



with speech analysis technologies





with speech recognition technologies

with new speech technologies being developed in our new project



DNN-GOP and DNN-DTW

...

....

...

...

....

DNN

W DNN-GOP = comparison bet. an L2 utterance and native models



+phonemes intended by the model speaker

time	Frame	Phoneme	a l	Frame	sil	а	i	u	
	1	а	÷	1	-0.01	0.8	0.1	0.02	
	2	а		2	-0:0 1	0.7	0.1	0.1	
	3	u		3	-0.01	0.5	0,	0.4	
	1232	sil		1232	0.9	0	0.01	0	

Model Student

phoneme

Solution ONN-DTW = comparison bet. an L2 utterance and its native version

Native-likeness

DTW

D	NN-	-GO	Ρ

Sequence of posterior vectors

Sequence of ···· posterior vectors



Signal B







Americans encounter Japanese English for the fist time.

How intelligible is JE to Americans with no exposure to JE?[Minematsu+'11] Ģ

Solution States and Americans (GA) read aloud sentences written by native speakers. All the recordings were judged as correct by the speakers themselves. Other 173 Americans listened to each utterance only once and repeat it. Topic and speaker always varied from utterance to utterance. Solution Frank Content of the second staff transcribed carefully the repetitions (20 repetitions on avg. / utterance).

> 200 Japanese **20 Americans**

Playing speech files selected from ERJ

Recording the response

Later, all the responses are transcribed.

800 JE + 600 AE utterances

intelligibility

Listening to each utterance only once

Repeating what the listener has heard.

17,416 JE + 12,859 AE transcriptions

173 American listeners



Data were collected at Indiana Univ. with support from Ordinate corp.



How correctly were JE repeated by Americans?

"The misquote was retracted with an apology."

- # the misquote was retracted with an apology
- # the misquote was retracted with an apology
- # the misquote was retracted with an apology #
- # the misquote was retracted with an apology #
- the misquote quote was retracted with an apology
- the misquote was [S>] attr(acted)- [
- the misquote was attra(cted)- was retracted with an apology
- the misquote was attract was retracted with an apology
- the misquote was attracted with an apology
- the misquote was attracted with an apology
- the misquote was retracted with an apology

- i don't know
- sammy's coat was instructed
- constructed
- distracted @
- was instructed with an apology
- @ by an apology
- something @ without apology
- @ was something
- instructed with an apology
- an apology
- someone was distracted with an apology
- is destructed apologize
- [Q]

goo.gl/jUAehX (O

- uh something was obstructed and needs an apology [N]
- distracted with an apology
- sammy's co:at was obstructed with apology
- @ is extracted with an apology
- @ was instructed with an apology
- the m(isquote)- mister was di- distracted with an apology
- attracted with # *polar
- the misquote was constructed with an apology #
- sammy was instructed with a apology [N]
- someone was instructed with an apology



How correctly were JE repeated by Americans?

stracted with an apology."

Sammy's coat was instructed???

the misquote w

1/1

- the misquote wa
- the misquote was attr
- the misquote was attrac
- the misquote was retrac
- the misquote was retract.
- the misquote was retract
- the misquote was retracte
- the misquote was retracted
 the misquote was retracted
 the misquote was retracted

 i don't know sammy's co constructed distracted @ was instructed with a an anology

The misquote was retracted with an apology.

ds an apology [N]

apology

distracted with an apology

d with an apology # apology D (W O)

the mis

samn som

goo.gl/iUAehX



How correctly were JE repeated by Americans?

"The misquote was retracted with an apology."

- # the misquote was retracted with an apology .
- # the misquote was retracted with an apology
- # the misquote was retracted with an apology #
- # the misquote was retracted with an apology # •
- the misquote quote was retracted with an apology •
- the misquote was [S>] attr(acted)-[.
- the misquote was attra(cted)- was retracted with an apology
- the misquote was attract was retracted with an apology .
- the misquote was attracted with an apology .
- the misquote was attracted with an apology .
- the misquote was retracted with an apology
- the misquote was retracted with an apology .
- the misquote was retracted with an apology .
- the misquote was retracted with an apology
- the misquote was retracted with an apology .
- the misquote was retracted with an apology .
- the misquote was retracted with an apology .
- the misquote was retracted with an apology .
- the misquote was retracted with an apology .
- the misquote was retracted with an apology
- the misquote was retracted with an apology
- the misquote was retracted with an apology

the misquote was a tracted with an apology

- i don't know
- sammy's coat was instructed .
- constructed .
- distracted @
- was instructed with an apology
- @ by an apology
- something @ without apology
- @ was something .
- instructed with an apology
- an apology •
- someone was distracted with an apology
- is destructed apologize
- [Q] •

goo.gl/jUAehX (O

- uh something was obstructed and needs an apology [N] •
- distracted with an apology .
- sammy's co:at was obstructed with apology .
- @ is extracted with an apology
- @ was instructed with an apology
- the m(isquote)- mister was di- distracted with an apology .
- attracted with # *polar
- the misquote was constructed with an apology #



Americans encounter Japanese English for the fist time.

How intelligible is JE to Americans with no exposure to JE?[Minematsu+'11] Ģ

Solution States and Americans (GA) read aloud sentences written by native speakers. All the recordings were judged as correct by the speakers themselves. Other 173 Americans listened to each utterance only once and repeat it. Topic and speaker always varied from utterance to utterance. Solution Frank Content of the second staff transcribed carefully the repetitions (20 repetitions on avg. / utterance).

> 200 Japanese **20 Americans**

Playing speech files selected from ERJ

Recording the response

Later, all the responses are transcribed.

800 JE + 600 AE utterances

intelligibility

Listening to each utterance only once

Repeating what the listener has heard.

173 American listeners



17,416 JE + 12,859 AE transcriptions

Data were collected at Indiana Univ. with support from Ordinate corp.



Intelligibility and comprehensibility

Intelligible/comprehensible enough pronunciations [Derwing+'09] **Q** Intelligibility

- Generative How many words in a given utterance can be identified correctly?
- Measured objectively by native listeners' transcription or oral repetition.
- Focuses on the results of listeners' recognition process. -> offline
- Comprehensibility
 - How easily, i.e., how smoothly, the content of a given utterance can be understood?
 - Measured objectively by monitoring brain activities or size of pupils
 - Focuses on how the recognition process is running. -> online











Shadowing = repeating without waiting and deep guessing





"Shadowing = repeating without waiting and deep guessing"

General form of shadowing



ss=smoothness of shadow native learner Objective measurement of SS (shadowability) pronunciations. **Proposed (inverse) form of shadowing**



learner

Shadowing = simultaneous reproduction of words intended by the learner, in a native pronunciation, not imitation of accented



learner



Collection of non-native and native Karaoke readings

¥ L1 = Vietnamese and L2 = Japanese Slow utterances are easy to shadow. Speaking rate control was introduced to recording. **Weight Four Setting With Speaking rate controlled** Model utterances from the CD of a textbook are used as reference. Solor of the text changes according to the speaking rate of the model utterances.





語呂合わせの記念日がたくさんあります >>>>> に愿んには / 『ごろあわせ』のきねんびがたくさんあります。

Record Stop Play 0:00 / 0:00

Task: 3 / 10 Progress: 1 / 15



Examples of readings and shadowings



A THE STAR THE THE SAN BUSINESS BUSINESS BUSINESS STORES

Utterances read by Vietnamese

Ward and state the second and the se

Shadowing utterances by Japanese listeners

Sand the Standing of Roads of Restant and the set

Texts and model utterances from CD

Read by Vietnamese

Read by Japanese

Utterances read by Japanese

Shadowed by Japanese listeners

Shadowing utterances by Japanese listeners



Conditions of the experiments

Non-native and native speakers and their Karaoke-style readings V Solution Native speakers of Japanese x 6 N \bigcirc 164 phrase utterances from the CD of the textbook \rightarrow 96 VJ + 68 NJ **With Native shadowers and their tasks** Solution Section Asked to reproduce in a native pronunciation what was heard as simultaneously as possible, not to imitate accented pronunciations. Solution Section Asked to rate comprehensibility of a given phrase utterance using a 7-degree scale.



"5"

1) shadowing 2) rating





Features extracted for correlation analysis



Ģ

164 phrase utterances of a model speaker

Karaoke-style recording

Experiments





Experiments

Features extracted for correlation analysis



164 phrase utterances of a model speaker

Karaoke-style recording

WRR

VJ-WRR

MS = Model Speech, RS = Responsive Shadowing

GOP (Goodness of Pronunciation) [Yue+2017] Accuracy of articulation, calculated as phoneme-based posterior probabilities WRR (Word-based Recognition Rate) **Delay of shadowing** Calculated as averaged phoneme boundary gap between VJ and RS



- Performance of the ASR system that is used as baseline system in the Japanese ASR community

Features extracted for correlation analysis

164 phrase utterances of a model speaker

Karaoke-style recording

WRR

VJ-WRR

Results of correlation analysis between the features and CS

VJ-GOP	VJ-WRR	
0.58	0.47	
RS-GOP	RS-WRR	VJ-RS-delay
0.74	0.53	-0.59

Y. Inoue et al., "A study of objective measurement of comprehensibility through native speakers' responsive shadowing of learners' utterances," Proc. INTERSPEECH, 2018

Experiments

Examples of readings and shadowings

Texts and model utterances from CD

Utterances read by Vietnamese

La the state the second of the second of the second

ARANTAR Startigoning Raca Star Rasa in Arthurista Star

Read by Japanese

Utterances read by Japanese

Shadowed by Japanese listeners

Shadowing utterances by Japanese listeners

Outline of the presentation

CALL for speaking (reading aloud), listening, conversation, and more Computer-Aided Language Learning with speech technologies

with speech synthesis technologies

with speech analysis technologies

with speech recognition technologies

with new speech technologies being developed in our new project

Out poster at AAAL2019

INTER-LEARNER SHADOWING WITH SPEECH TECHNOLOGIES ENABLES AUTOMATIC AND OBJECTIVE MEASUREMENT OF COMPREHENSIBILITY OF LEARNERS' UTTERANCES

Nobuaki Minematsu*, Yusuke Inoue*, Daisuke Saito*, Yutaka Yamauchi**, Kumi Kanamura*** (*UTokyo, **Soka Univ., ***Nagoya Univ. Economics)

RESEARCH QUESTION AND RELATED WORKS

Two different goals of pronunciation training

- @ Native-like vs. intelligible/comprehensible enough pronunciations
- ♀ Intelligibility and comprehensibility [Derwing+'09]
- ♀ I: How many words are correctly identified in given L2 utterances? \rightarrow Focuses on results of listeners' process of understanding \rightarrow offline
- ♀ C: How easily or smoothly given L2 utterances are understood?
- \rightarrow Focuses on how the understanding process is running while listening \rightarrow online

How to calculate intelligibility/comprehensibility of L2 speech?

Subjective assessment

- "How many words do you think you identified correctly?" 100%, 80%, 60%, 50%,,,,,
- @ "How easily did you understand the message?" Very easily, easily, rather easily,,,,,

Objective assessment

- ♀ L2 utterances were transcribed by native speakers after hearing them. Objective intelligibility = word-level correct transcription rate [Bernstein'03]
- Listeners' behaviors are observed using physiological sensors. Listening effort observed using EEG (Electroencephalogram) [Song+'18] **Cognitive load** observed as the size of pupils using an eye-tracker [Govender+'18] Both methods are too expensive to be used in classrooms.

Q: How to calculate comprehensibility in an inexpensive way?

- Objective assessment on intelligibility [Bernstein'03, Minematsu+'11]
- [©] Transcription or **oral repetition** is done **after** hearing L2 utterances.
- \bigcirc Transcription or repetition allows waiting and guessing. \rightarrow offline
- Objective assessment on comprehensibility using modified repetition [Inoue+'18]
- \bigcirc Repetition with almost no waiting or guessing \rightarrow shadowing!!!
- \bigcirc Native listeners' **reverse** shadowing of L2 utterances!!! \rightarrow **online**

Conventional form of shadowing

SS = Smoothness of Shadowing → GOP

♀ Learners' SS is automatically predicted. [Luo+'10][Yue+'17][Kabashima+'18]

Recording of L2 read sentences or phrases

- ♀ 6 intermediate and 6 advanced learners
- Read sentences were also collected from natives as reference

Native listeners' reverse shadowing [Inoue+'18]

"A simpler approach to calculate smoothness of shadowing" **Proposed (inverse) form of shadowing** Ş Solution ONN-GOP is calculated at every frame, that can be viewed as spoken annotation!! Annotations (labels) should be collected with simpler and more reliable techniques.

learner

native

A much simpler and reliable approach A native listener is asked to **shadow** an utterance given from a learner. Solution For the standard of the sext that was read by that learner. Read speech = most prepared speech, shadowed speech = least prepared (hastened) speech OTW between the two speeches will give us a sequence of smoothness of shadowing.

What will be possible with a huge amount of data?

L2 utterances with spoken annotations

Data collection and system development

INTER-LEARNER SHADOWING

SS = 60**SS** = "**B**+"

Logical feedback

- utterances

Conclusions

CALL for speaking (reading aloud), listening, conversation, and more Computer-Aided Language Learning with speech technologies

with speech synthesis technologies

with speech analysis technologies

with speech recognition technologies

with new speech technologies being developed in our new project

Conclusions

