

“How to follow a conversation without listening to the words”

Nick CAMPBELL

National Institute for Information and Communications Technology,
nick@nict.go.jp

ATR Spoken Language Communication Research Laboratories
nick@atr.jp

Machine intelligence,
ambient computing,
sensitive devices,
proactive machines

- “We need machines that better understand people, so that people can better understand machines.”
- “Sensitive machines can be made proactive, not needing to wait for instructions but being ready with an appropriate act.”

ヒューマンコミュニケーションの「場」が読めるロボットの研究開発 (041307003)

The “Robot’s Ears” Project

研究代表者

ニック・キャンベル

国際電気通信基礎技術研究所 音声言語コミュニケーション研究所 音声音響処理室学研究室 主幹研究員

Nick CAMPBELL

ATR Spoken Language Communication Research Laboratories

研究分担者

定延 利之十 千原 國宏 井村 誠幸

Sadanobu Toshiyuki, Chihara Kunihiro, Imura Masataka,

ダミアン ドゥシャン 岩橋 直人 中川 明子

Damien Douxchamps, Iwahashi Naoto, Nakagawa Akiko

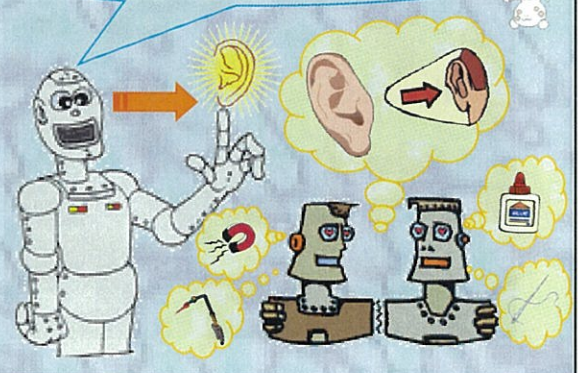
神戸大学 国際文化学部 国際文化学科 情報コミュニケーション論講座

竹奈良先端科学技術大学院大学 情報科学研究科 情報処理学専攻 情報処理学講座

国際電気通信基礎技術研究所 音声言語コミュニケーション研究所 音声音響処理室学研究室

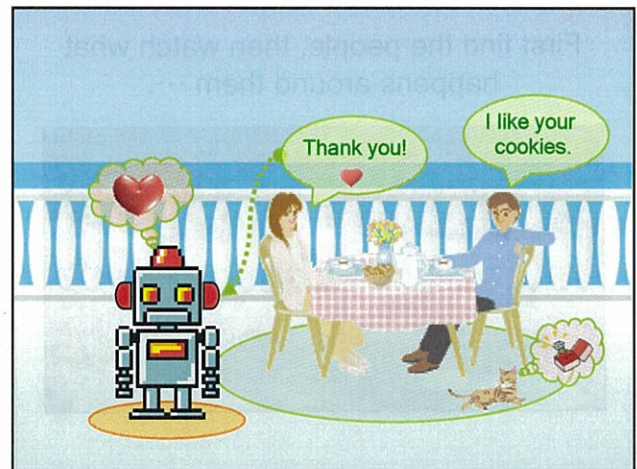
研究期間 平成16年度～平成18年度

How do you fix the ears on a robot?



You give them better eyes!

- By watching how people interact
- And listening to how they speak
- But not bothering much with what they say!





Following Meetings Data

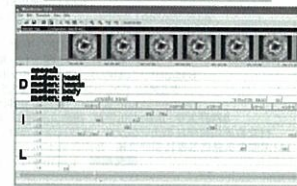
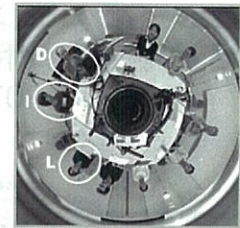
- 2-10 people sitting round a table
 - upper-body movement only
- who is participating how (to what degree)
 - talking / listening / thinking / waiting to speak
- basic primitives only
 - no speech recognition
 - no eye tracking
 - (simple unobtrusive devices)

the capture environment



Annotation:

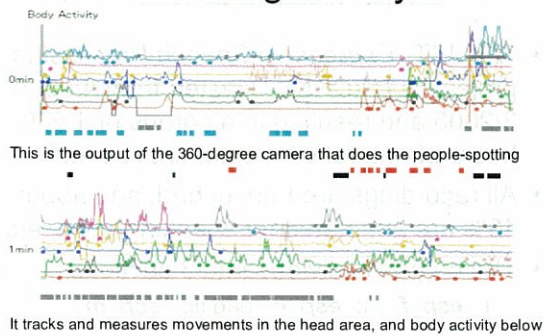
- speech
 - on/off
 - long/short
- motion
 - head
 - body
 - hands



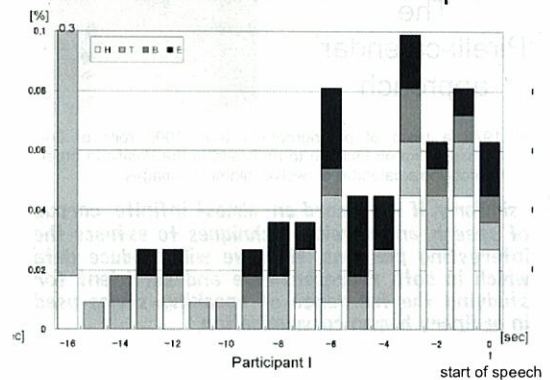
First find the people, then watch what happens around them ...



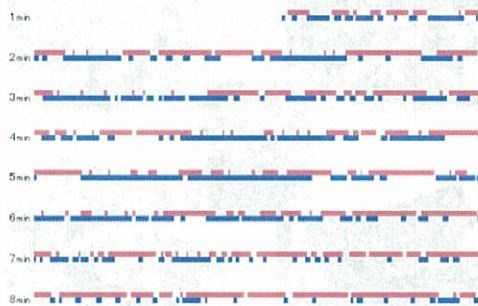
Tracking activity



Movement before start of speech



Short and long utterances



Tracking speech

10073	うん	467	ズー	228	うん	134	へー
9692	お	455	ズー	227	ん	134	はいはいはいはい
8607	はい	450	ん	226	へー	134	そうです
4216	laugh	446	うー	226	ハハハ	133	お
3487	うー	396	おー	225	ん	133	おそうめんですか
2906	ええ	395	あー	200	そうすね	130	おそうめんですか
1702	はい	393	はいはいはい	199	はー	129	はー
1573	うー	387	あー	195	ハ	129	い
1348	ん	372	ん	192	おの	127	はー
1139	ん	369	ん	190	ええ	125	ハハハハハ
1098	おの	369	だから	188	あー	119	はいはい
1084	あー	368	あー	187	ん	119	はー
981	はい	366	あー	180	ん	114	ハハ
942	おの	345	おの	180	あー	113	は
941	ん	337	ん	173	ん	113	で
910	そう	335	え	172	アハハハ	113	で
749	ん	311	で	168	はい	112	はあ
714	あー	305	ズー	164	うー	110	フフ
701	あ	274	うん、うん、うん	161	はー	110	おの
630	あー	266	ハハハハ	160	お	110	え
613	あ、はい	266	で	159	そうです、ね	109	ん
592	うん、うん	266	えー	151	あー	108	はあー
555	あー	258	で	143	だから	106	そうです、ねえ
500	ん	248	う	139	アハハハハ	105	ん
469	え	242	へー	137	そう、そう、そう	104	いや

「表現豊かな発話音声の コンピュータ処理システム」

"Expressive
Speech
Processing"

the JST CREST-ESP Project

科学技術振興事業団 第131号：
「高度メディア社会の生活情報技術」

The JST/CREST ESP project

Learning how people talk ...

- by listening to LOTS of examples
- recorded in everyday conversations
- over a period of 5 years
- all transcribed and annotated
- and acoustically mapped

The 'Pirelli-calendar' approach



in 1970 a team of photographers took 1000 rolls of 36-exposure film on location to an island in the Pacific in order to produce a calendar of twelve (glamour) images.

- > **similarly, if we record an 'almost infinite' corpus of speech, and develop techniques to extract the interesting portions, then we will produce data which is both representative and sufficient for studying the full range of speaking-styles used in ordinary human communication.**

The JST/CREST 'ESP' corpus

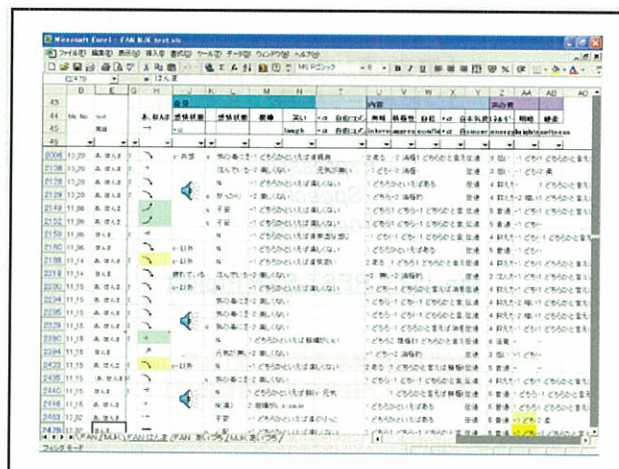
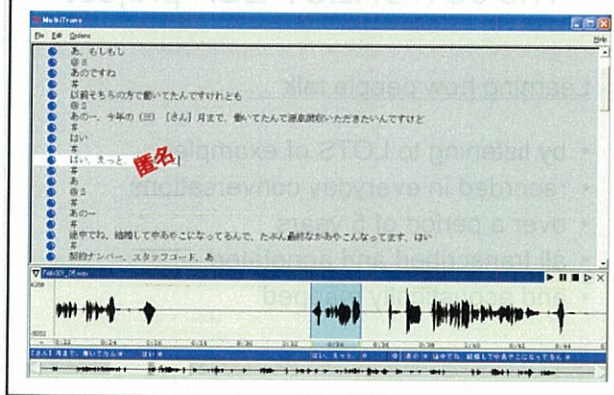
- The JST "Expressive Speech Processing" project (ATR/CREST) lasted from 4/2000 to 3/2005 and resulted in a corpus of 1,500 hours of natural conversational speech
- All recordings are transcribed, and about 10% are annotated for speaking-style, etc.
- The corpus is divided into 3 sections :
i: **esp_f**, ii: **esp_c**, and iii: **esp_m**

Sections of the ESP corpus

- esp_f**
 - one female speaker, head-mounted mic, 600 hours of daily spoken interactions, emotion/speech-act/etc ...
- esp_c**
 - 10 adult speakers, 5m 5f, 2 chinese, 2 english,
 - 30-minute telephone conversations x 10 weeks
 - all conversations in japanese, free content
- esp_m**
 - multi-speaker, head-mounted microphones, variety of interaction settings (like esp_f but many more voices)



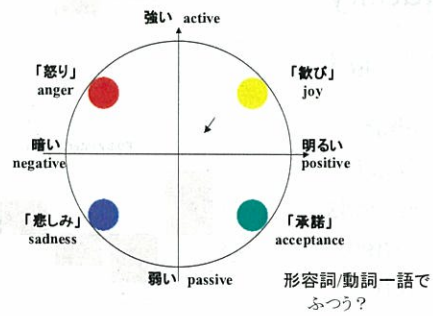
Transcription



What's the difference?

- The words are (almost) the same
 - Yet the whole 'feeling' is different
- Shows speaker - listener relationships
- Transcription of the speech may be the same
 - but these small differences are very important!
- And 'GRUNTS' are very common!

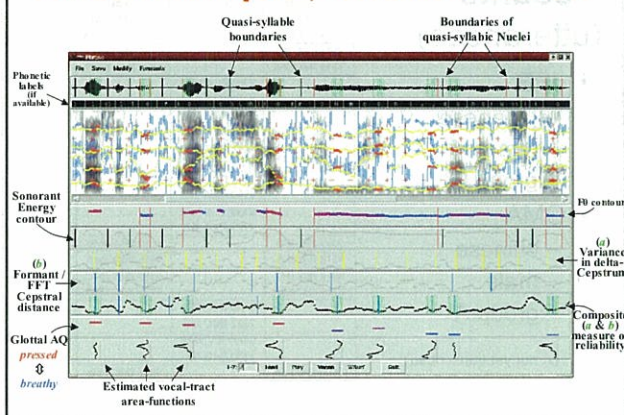
"Emotion" Labelling



Discourse Act Labelling

a	あいさつ	greeting	o	反論	argue
b	会話終了	closing	p	提案、申し出	suggest, offer
c	自己紹介	introduce-self	q	気づき	notice
d	話題紹介	introduce-topic	s	つなぎ	connector
e	情報提供	give-information	r	依頼、命令	request-action
f	意見、希望	give-opinion	t	文句	complain
g	応答肯定	affirm	u	褒める	flatter
h	応答否定	negate	w	独り言	talking-to-self
i	受け入れ	accept	x	言い詰まり	disfluency
j	拒絶	reject	y	演技	acting
k	了解、理解、納得	acknowledge	z	繰り返し	repeat
l	割り込み、相づち	interject	r*	要求	request (a~z)
m	感謝	thank	v*	確認を与える	verify (a~z)
n	謝罪	apologize	w*	よく分からない場合	

Acoustic Analysis / Visualisation tool



An example of how speaking-style can change

- 13,604 conversational utterances
- 1 female Japanese speaker (age 32-35)
- listener/speech-act/emotion labels

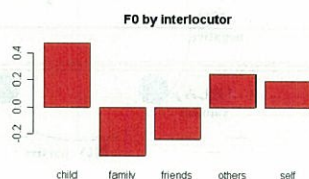
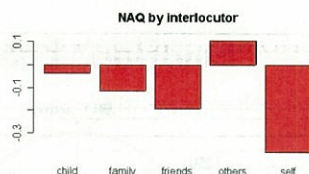
Interlocutor:

Child	Family	Friends	Others	Self
139	3623	9044	632	116

Voice quality

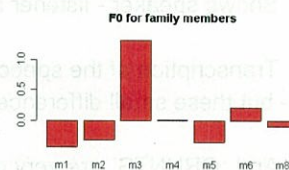
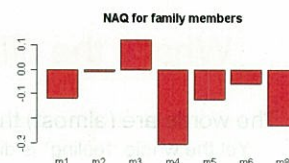
Talking to:

- child
- family
- friends
- others
- self



NAQ & F0 by family

- m1 - mother
- m2 - father
- m3 - baby girl
- m4 - husband
- m5 - big sister
- m6 - nephew
- m8 - aunt



Conversational Speech

- Together, the listener and the speaker signal both
 - (a) discourse & social relationships
 - (b) propositional content.
- Non-verbal utterances are commonly used to signal paralinguistic information.
- These two forms can be distinguished as
 - I-type (*information*) utterances
 - A-type (*affect*) utterances

esp_c: paired conversations

- 10 speakers, 5 male, 5 female
 - 2 Chinese, 2 English-native-speakers
 - all conversations in Japanese
- 3 groups (A,B, and C)
 - A - talking with foreigners
 - B - talking with strangers
 - C - talking with family members
 - all talking with each other ... for 10 sessions

groupings for natural dialogues

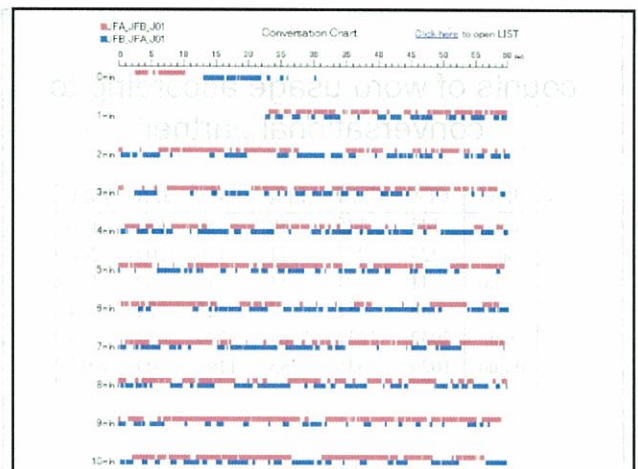
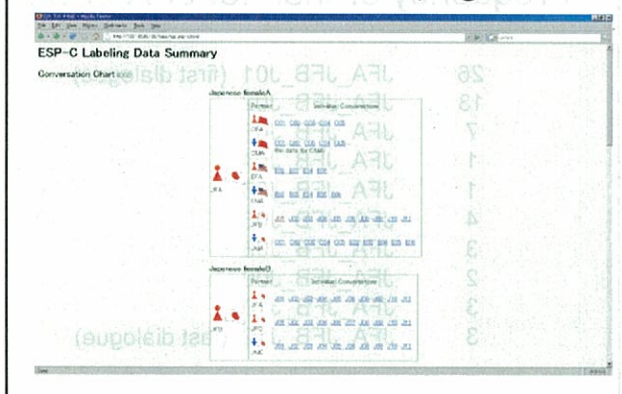
female		male	
(cfa efa		cma ema)	(foreign)
/		\	Group A
jfa	-	jma	
jfb		jmb	Group B
jfc	-	jmc	
			Group C
(fam)		(fam)	(intimate)

counts (utterances)

- JFB
 - japanese
 - female
 - group B
 - CMA
 - chinese
 - male
 - group A
- (x: not yet transcribed)

both female	mixed	both male
6425 EFA JFA		9348 EMA JMA
7359 JFA EFA		7433 JMA EMA
8827 CFA JFA		x (cma jma)
9145 JFA CFA		7530 JMA CMA
	9236 EFA JMA	
	8499 JMA EFA	
	7557 JMA CFA	
	x (cfa jma)	
	8237 JFA CMA	
	x (cma jfa)	
	8416 JFA EMA	
	8560 EMA JFA	
	10068 JFA JMA	
	7701 JMA JFA	
9069 JFA JFB		8614 JMA JMB
9378 JFB JFA		9465 JMB JMA
8044 JFB JFC		6983 JMB JMC
8234 JFC JFB		7735 JMC JMB
	7686 JFB JMC	
	7222 JMC JFB	
	10005 JFC JMB	
	7980 JMB JFC	
13000 JFC Fam		9961 JMC Fam

the top-level web pages



the structure of spoken language

- one of the main goals of this work is to model the structure of spoken dialogues
- to understand why they appear 'broken'
- to model the different types of information (not just linguistic) that are carried by the frequent (non-verbal) sounds and by differences in the structure of speech

a sample conversation

start	end	speaker	PLAY	transcription
0:02	0:04	JFA	start	十二月十二、日本曜日
0:05	0:05	JFA	start	四、時
0:06	0:07	JFA	start	四、五、分
0:08	0:10	JFA	start	待機中、開始します
0:13	0:14	JFB	start	十二月十二
0:14	0:15	JFB	start	午後四、時四、十五
0:16	0:18	JFB	start	四、十五分
0:18	0:18	JFB	start	あ、ええ
0:18	0:22	JFB	start	え、一、待機中、開始します
0:22	0:23	JFB	start	はい、はい
0:24	0:25	JFB	start	あ、はい
0:25	0:26	JFB	start	はい
0:26	0:29	JFB	start	はい
1:22	1:23	JFB	start	はい、はい
1:23	1:24	JFA	start	あ、はい、はい
1:24	1:24	JFB	start	あ、はい
1:24	1:25	JFB	start	はい、はい

the raw data

- transcribed speech
- start & end times
- filenames also show speaker, partner, & conversation number
- one utterance / line
- noises and non-speech sounds also transcribed

```

CFA JFA C01 200.369 0.491 #
CFA JFA C01 200.660 0.808 laugh
CFA JFA C01 201.668 0.869 あ、はい
CFA JFA C01 202.537 1.099 変わりました
CFA JFA C01 203.636 1.868 laugh
CFA JFA C01 205.504 0.670 うん
CFA JFA C01 206.174 0.744 #
CFA JFA C01 206.918 0.917 はい
CFA JFA C01 207.835 2.691 #
CFA JFA C01 210.526 0.602 はい
CFA JFA C01 211.128 2.791 #
CFA JFA C01 213.919 0.749 #
CFA JFA C01 214.668 2.685 そう、です、結構、し
う、四、有難、なりました
CFA JFA C01 217.353 0.785 はい
CFA JFA C01 218.138 0.561 #
CFA JFA C01 218.699 0.731 はい
CFA JFA C01 219.430 1.384 #
CFA JFA C01 220.814 1.088 行、って、ます
CFA JFA C01 221.902 0.738 #
CFA JFA C01 222.640 0.784 はい
CFA JFA C01 223.424 1.107 #
CFA JFA C01 224.531 1.356 あ、はい、です
CFA JFA C01 225.887 0.525 #
CFA JFA C01 226.412 0.600 はい
CFA JFA C01 227.012 2.795 #
CFA JFA C01 229.807 0.443 はい
CFA JFA C01 230.250 0.941 #
  
```

the hundred most common utterances

10073	うん	467	スー	228	ううん	134	へー
9692	あ、	455	スー	227	えっ	134	はいはいはいはい
8607	はい	450	んー	226	へー	134	そうです
4216	laugh	446	うーん	226	ハハハ	133	あ、
3487	うーん	396	えー	225	うん	133	あ、そうなんですか
2906	ええ	395	あー	200	そうですね	130	そうですね、
1702	はい	393	はいはいはい	199	はー	129	はい
1573	うーん	387	あーはい	193	ハ	129	い
1348	スー	372	ええ	192	その	127	はー
1139	ふん	369	ふーん	190	ええ	125	ハハハハハ
1098	あ、	369	だから	188	あ、あー	119	はいはい
1084	あ、	368	あー	187	ハ	119	はー
981	はい	366	あー	180	え、はい	114	ハハ
942	あ、	345	あ、あー	180	あ、あー	113	は
941	ふん	337	ふん	173	ん	113	てー
910	そう	335	え	172	アハハハ	113	はー
749	えー	311	でも	168	はい	112	は、あー
714	あー	305	スー	164	うーん	110	フフ
701	あ	274	うん、うん、うん	161	はー	110	その
630	あー	266	ハハハハ	160	あ、	110	もう
613	あ、はい	266	てー	159	そうですね	109	ふーん
592	うん、うん	266	えー	151	あー	108	あ、あー
555	あー	258	で	143	だから	106	そうですね
500	ん	248	う	139	アハハハハ	105	ん
469	ん	242	へー	137	そう、そう、	104	いや

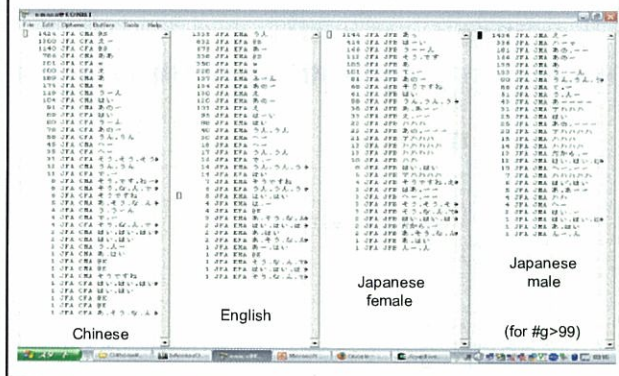
counts of word usage according to conversational partner

JFA:	CFA	CMA	EFA	EMA	JFB	JMA
a,a-	143	145	88	89	138	170
ano	224	277	221	176	209	266
demo	41	24	31	17	89	134
e-	48	51	37	25	74	94
hai	2932	2234	2181	3239	72	33
un,un	1029	546	585	1190	909	1037

Frequency of "hai" for JFA-JFB

26	JFA_JFB_J01 (first dialogue)
13	JFA_JFB_J02
7	JFA_JFB_J03
1	JFA_JFB_J04
1	JFA_JFB_J05
4	JFA_JFB_J06
3	JFA_JFB_J08
2	JFA_JFB_J09
3	JFA_JFB_J10
3	JFA_JFB_J12 (last dialogue)

Frequency of 'grunts' per partner

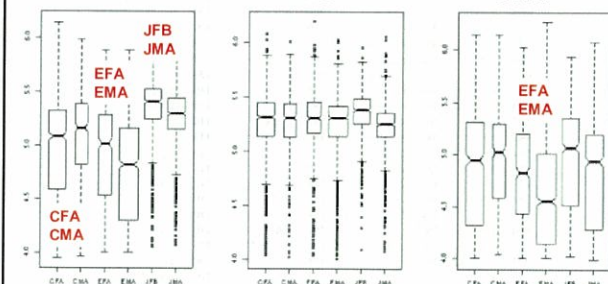


some acoustic measures for JFA

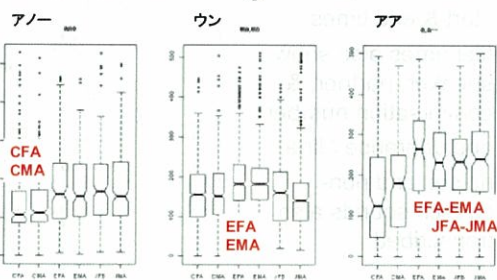
"a,a-"	CFA	CMA	EFA	EMA	JFB	JMA
f0r	125	181	266	232	234	241
f0m	201	214	220	192	206	198
pwr	28	29	29	28	31	31
pwm	38	39	36	35	42	41
"un,un"	CFA	CMA	EFA	EMA	JFB	JMA
f0r	154	152	182	181	161	141
f0m	172	175	162	145	198	174
pwr	28	29	27	26	29	27
pwm	37	40	36	35	42	39
"ano"	CFA	CMA	EFA	EMA	JFB	JMA
f0r	106	113	161	154	169	155
f0m	131	136	142	133	156	149
pwr	27	28	28	27	31	29
pwm	38	40	37	36	42	39

mean f0 values for JFA's うんうん

start mid end



pitch range variation in 3 utterances from speaker JFA differs according to interlocutor



utterance frequencies for JMA

a	a-	a-	a—	a.a-	a-hai
296	368	693	608	390	386
a.hai	a-n	ano	ano-	a!!	demo
577	368	337	494	927	272
e-	e-	ee	fun	fu-n	fu-n
665	254	2679	642	625	273
ha.ai	hai	ha-i	hai.hai.hai	n(ummi)	n-
978	7295	1657	378	265	456
n-	nanka	ne-	nee	@S	sou
410	273	367	284	3382	810
su-	su-	un	u-n	u-n	u-n
429	296	3717	2401	1243	333
un.un	@W	zu-	zu-		
351	3041	1348	467		

F0 varies according to partner ...

mean f0 max f0 min f0

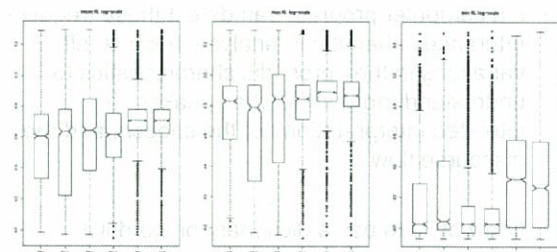


Figure 1: Plots of mean, maximum and minimum f0 values observed in the data of each of the interlocutors. The box-plots show median and interquartile values, with whiskers extending to 1.5 times the interquartile range. All F0 measurements are converted to their log values for ease of comparison.

as does signal amplitude (energy)

mean max min

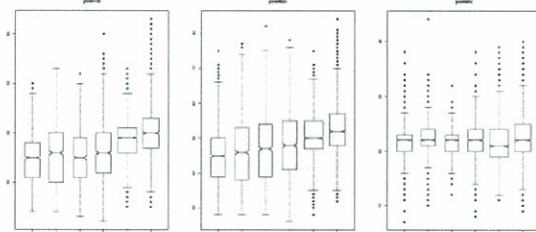


Figure 2: Plots of mean, maximum and minimum rms amplitude (speech signal power) values for each of the interlocutors.

... and also the ranges of pitch and power

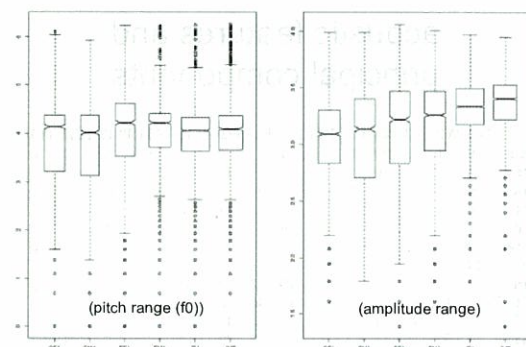
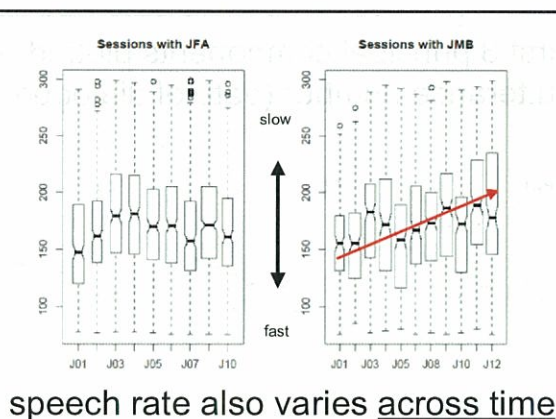


Figure 3: Ranges of fundamental frequency and power measurements for the affective utterances, plotted by interlocutor. As above, C,E,J stand for Chinese, English, Japanese respectively, and FM represent female and male interlocutors. Both f0 and power are plotted as log values.



speech rate also varies across time

Figure 4: Speaking rate changes over weekly sessions.

talk's focus : frequent utterances

- utterances that are frequently repeated can be used as *comparison points* ...
- small differences in voice quality and other prosodic parameters can carry significant information about the speaker state
- i.e., *even if we do not know the speaker well*, we can use common speech sounds as an indicator of changes in that persons affective state(s) and discourse intentions

Detecting frequent utterances

- if a computer program can detect these frequent utterances, then it can analyse these small variations in their prosodic characteristics to understand more about the speaker, the intended interpretations of the speech, and the discourse flow
- one way is to use a dictionary or word list (one-hundred words >50% of utterances!)

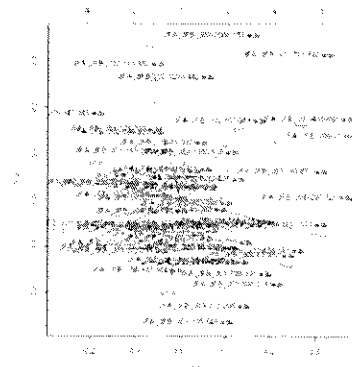
discussion

- Prosody varies according to the nature of relationships between speaker and partner
- By tracking frequent repeats, we can detect small changes in these relationships
- Small changes can be easily detected by comparison between frequent similar segments that are close in time

acoustic features and principal components

	PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8	PC9	PC10	PC11	PC12	PC13	PC14
fmean	-50	14	25	-20	11	-3	4	-7	8	3	14	27	71	0
fmax	-18	10	23	0	9	1	30	-26	18	18	36	-37	-16	0
fmin	-29	11	17	-16	12	-26	-27	15	7	-14	-62	-9	-26	0
fpct	-6	20	-13	-22	-10	46	-62	-25	11	18	12	-9	-2	0
fved	-8	-23	-30	-29	-7	6	36	-61	-15	9	-28	29	-7	0
pmean	-36	-26	-31	29	-16	-7	-7	-17	-28	-16	-59	30	0	0
pmax	-13	-12	0	42	-5	2	-27	1	-17	-30	8	56	-34	0
pmin	-20	-26	-37	7	-12	-34	-10	32	29	64	2	11	-1	0
ppct	-16	16	-15	-14	-51	30	45	46	23	-27	-6	8	-5	0
hli2	-13	-30	32	4	6	53	8	25	-44	41	-26	-5	1	0
hli3	9	-50	34	-12	-28	-9	-4	-9	20	-12	7	1	1	-67
hl	5	-57	11	-14	14	19	-7	-1	43	-21	8	0	0	60
a3	-8	0	-39	-1	63	40	-5	12	27	-10	0	-1	-1	-44
dn	5	18	22	54	-6	12	7	-27	49	14	-51	5	9	0

Principal Components Analysis



PCA in R (r-project statistical software):
prcomp(un, retx=T, center=T, scale.=T)

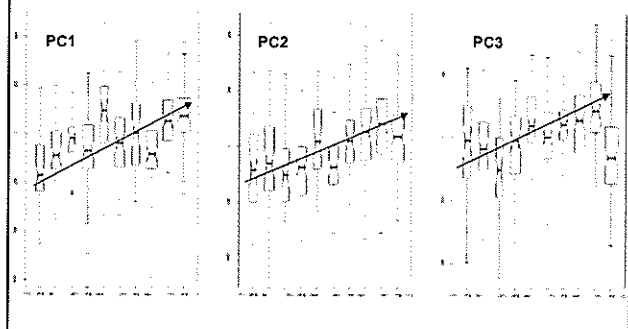
3 components: 50%, 7 components: 80%

Importance of components:

	PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8
Standard deviation	1.65	1.53	1.38	1.32	1.12	0.96	0.89	0.83
Proportion of Variance	0.19	0.16	0.13	0.12	0.08	0.06	0.05	0.04
Cumulative Proportion	0.19	0.36	<u>0.49</u>	0.62	0.71	0.78	<u>0.83</u>	0.88

	PC9	PC10	PC11	PC12	PC13	PC14
Standard deviation	0.74	0.71	0.61	0.29	0.23	0.0004
Proportion of Variance	0.03	0.03	0.026	0.006	0.004	0.0001
Cumulative Proportion	0.92	0.96	0.98	0.99	1.00	1.00

First 3 principal components plotted by utterance number (date of dialogue)



conclusion

- this paper has presented an initial analysis of the jst/crest esp_c corpus
- which is balanced according to speaker familiarity, sex, and ease of interaction
- showing that speakers modify their speaking habits according (perhaps?) to familiarity with the interlocutor
- and that the “ill-formed” nature of speech serves to carry important social information

... more conclusion

- ... and has presented a simple program for analysing the transcriptions of the conversations to
 - differentiate between frequent “wrappers” and their linguistic “fillers”
 - to provide a measure of the differences in “social prosody” in conversational interactions
- this is work in progress ...
 - so I hope to benefit from your helpful comments and suggestions

some initial findings : acoustic differences

- people change their speaking styles when talking with different people
 - this is plain common sense (nothing new!) but nice to be able to measure/quantify
- people also vary significantly their pitch ranges and voice quality settings, even for the same utterance

Latest devices



SONY RPU-C251 (desktop version)



... more conclusion



conclusion

- this paper has presented an initial analysis of the Japanese esp_c corpus
- which is balanced according to speaker familiarity, sex, and ease of interaction
- showing results according (perhaps) to familiarity with the interlocutor
- and that the "ill-formed" nature of speech serves to carry important social information

thank you for listening



some initial findings : acoustic differences

- people change their speaking styles when talking with different people
– this is plain common sense (nothing new)
but nice to be able to measure quantitatively
- people also vary significantly their pitch ranges and voice quality settings, even for the same utterance

SONY RPU-C251 (desktop version)

