

New Resources for Recognition of Confusable Linguistic Varieties: The LRE11 Corpus

Stephanie Strassel, Kevin Walker, Karen Jones, Dave Graff, Christopher Cieri

> Linguistic Data Consortium University of Pennsylvania, USA





- Data Requirements
- Language Selection
- Broadcast Collection
- Telephone Collection
- Segment Selection
- Auditing
- Auditor Agreement
- Data Distribution and Corpus Summary
- Conclusions



Data Requirements for LRE 2011

- Distribution of previous LRE data to new participants
 - Previous test sets
 - LRE 2009 training data, including large broadcast news corpus
- New resources for LRE 2011
 - As in LRE 2009, includes both conversational telephone speech (CTS) and broadcast narrowband speech (BNBS)
 - Both genres for most but not all languages
 - Arabic varieties limited to broadcast-only (MSA) or telephone-only (Iraqi, Levantine, Maghrebi)
 - Target 24 languages/dialects, some of which may be mutually intelligible to some extent by humans
 - 400 segments per language
 - At least 2 unique sources per language
 - Broadcast source is provider-program (so CNN Larry King is different source than CNN Headline News)



- Reviewed information sources like Ethnologue
- Compiled list of candidates plus confusability index score
 - 0 Not likely to be confusable with another candidate language*
 - 1 Possibly confusable with another candidate language; languages are related and may be confused by (some) systems if not by (most) humans
 - 2 Likely confusable with another candidate language; at least some evidence that (some) humans may find the varieties mutually intelligible to some extent
- Candidate set of 38 languages whittled down to 24 with NIST and sponsor input, and considering
 - Availability of broadcast sources
 - Availability of claques/auditors

*Throughout we use *language* as shorthand for a linguistic variety that may be referred to by different sources as a language or dialect



Potential Confusability for LRE 2011 Languages

	ISO 639-3 or	Confusablity	Language(s) of Possible		ISO 639-3 or	Confusablity	Language(s) of
Language	code	Score	Confusion	Language	code	Score	Confusion
Arabic Iraqi	acm	2	other Arabic	Mandarin	cmn	0	
Arabic Levantine	alv	2	other Arabic	Pashto	pus	0	
Arabic Maghrebi	arm	2	other Arabic	Polish	pol	1	other Slavic
Arabic MSA	ara	2	other Arabic	Punjabi, Western	pnb	1	other Indic
Bengali	ben	1	other Indic	Russian	rus	1	other Slavic
Czech	ces	1	slk	Slovak	slk	1	ces
Dari	prs	2	fas	Spanish	spa	0	
English (American)	eng	1	emi	Tamil	tam	0	
English (Indian)	emi	1	eng	Thai	tha	1	lao
Farsi/Persian	fas	2	prs	Turkish	tur	0	
Hindi	hin	2	urd	Ukrainian	ukr	1	other Slavic
Lao	lao	2	tha	Urdu	urd	2	hin



Broadcast Collection

Multiple broadcast sources

- Existing, unexposed VOA1 data
- New and unexposed archival data from local satellite collections in Philadelphia, Tunis and Hong Kong
- New collection from cable, satellite and off-the air sources via portable collection platform installed in New Delhi
- New collection from streaming web radio sources

Variety of formats

- Satellite data is MPEG1 Audio Layer II (.mp2)
 - MPEG ADTS, layer II, v1, 128 kbps, 48 kHz, Stereo
 - MPEG ADTS, layer II, v1, 160 kbps, 48 kHz, Stereo
 - MPEG ADTS, layer II, v1, 192 kbps, 48 kHz, Stereo
 - MPEG ADTS, layer II, v1, 64 kbps, 44.1 kHz, Monaural
 - MPEG ADTS, layer II, v1, 64 kbps, 48 kHz, Stereo
- All streaming sources mp3, 128kbps bitrate, 44.1kHz sample rate



Telephone Collection

- Claque-based collection model
 - Claque is a native speaker informant
 - Eases recruitment burden
 - Claques later serve as auditors
- 2-5 claques recruited per language
- Each claque makes a single call to each of 15-30 individuals in their existing social network
 - Callee hears pre-recorded message and provides consent prior to call being recorded
 - Steps taken to ensure different claques' callees did not overlap
 - Claque call sides excluded from corpus
 - Require at least some calls within US to avoid bi-uniqueness of channel/language conditions

Calls collected on LDC's CTS platform in 8kHz, 8-bit plaw



- Full recordings passed through SAD system to distinguish speech vs. silence, music, other non-speech
- For CTS data we extract 2 segments per call, 30-35 seconds each
- For BNBS, additional bandwidth filter prior to selection
 - From the intersection of speech and bandwidth filters, continuous regions of 33+ seconds selected
 - For regions > 33 sec, single 33-sec segment chosen from center
 - No selection of multiple segments from single stretch of speech
 - When necessary to get a sufficient number of auditable segments for a given language, shorter continuous segments (down to a minimum of 10 sec) were selected
 - No concatenation of separate, short BN segments



Broadcast Segment Selection



Although many speech segments are large enough to yield multiple 30 second sub-segments we do not further segment them in order to maximize the number of potential speakers in the corpus



Extracted segments converted to auditor format

- BNBS: 16 KHz, 16 bit
- CTS: 8 KHz single-channel
- Converted to pcm, ms-wav file format for browser compatibility
- Auditor presented with entire segment
 - BNBS: typically 33 sec long, but possibly as little as 10 sec
 - CTS: entire 30-35 second segment



Auditing Kit Construction

- Baseline: segments expected to be in the auditor's language
 - For BNBS, proportional selection of all available segments
 - For CTS, other claques' callee sides
- Up to 10% **distractor** segments from non-confusable language
 - Presented to the auditor at random, to keep them attentive
- Up to 10% **dual** segments, also assigned to other auditor(s)
- For languages with confusable/buddy languages, also include confusable segments comprising 10%, 25% or 100% of the baseline amount, as follows:
 - 10% for related/possibly confusable varieties (e.g. Polish/Slovak)
 - 25% for likely confusable varieties (e.g. Lao/Thai)
 - 100% for known confusable varieties (e.g. Hindi/Urdu)
- Given non-linear nature of the collection, actual kit makeup varies
 - So one kit may be predominately CTS, or have < 10% dual segments



 Preliminary online screening of potential auditors for language skills

- Questions about language background, education
- Listening test comprising 10 segments including target language, distractor and potentially confusable language segments
- Screening results also helped point out areas where auditor training was required, e.g. to clarify language labels
- Of ~130 who took screening, 84 passed and were hired and given additional training (typically in-person)
 - Telephone calling instructions
 - Bandwidth detection training via "Signal Quality Perception Test"
 - Train, then test on ability to distinguish wide- and narrow-band segments (described as "phone-like" quality and "studio-like" quality)
 - Auditors could revisit test anytime to refresh their memory





- Goal of auditing is to ensure that segments
 - Contain only speech
 - Are in the target language variety
 - Are narrowband
 - Contain only one speaker
 - Audio quality is acceptable
 - In contrast to previous LRE, no question about speaker uniqueness
- Auditors judge each segment via a web-based interface
 - Required to listen to entire segment
 - Instructed to use good-quality headphones
 - Formal auditing instructions explain how to answer each question



Auditing Interface

ara	Arabic, MSA
Aud	io Quality
	Is the segment all speech (no music or sound affects)?
	Is the segment an speech (no music of sound effects): $v_{yes} = v_{no}$
	is it all "telephone-like" in quality (not studio quality)? $v_{yes} v_{no}$
	clear and easy to understand
	© somewhat unclear
	very unclear, hard to understand
	Check all that apply: distortion noise drop-outs interference other
Com	nment (optional):
Lang	guage
	Is all of the speech in Arabic, MSA? ^O yes ^O no
	Click here if the content is offensive:
Com	nment (optional):
Spea	aker
Is al	If the speech from a single speaker? \bigcirc yes \bigcirc no
	What is the speaker's sex? \circ mole \circ formula \circ uncome
w	What is the speaker's dialect/accent?
	speaker uses the expected dialect/accent
	speaker uses a different dialect/accent
	not a native speaker
Com	ament (optional):
	Submit your answers



 The audited segments delivered for LRE11 were limited to just those where

- (a) we had only one auditor judgment on record, or
- (b) the two or more auditor judgments were in agreement
- When one of those was true, and the judgment indicated a usable segment (in the auditor's target language, and all speech), the segment was delivered to NIST
- Segments that showed discrepant auditor judgments or indeterminacy in manual language labeling were excluded from delivery
- Numbers reported here are from segments assigned during normal auditing process, not during a post-hoc consistency analysis task



Comparing multiple judgments

- where the expected language of the segment was the language of the auditors

- what is language-label agreement?

 e.g. two Bengali speakers judge clips that are purported to be Bengali (e.g. because of the collection source)



Within Language (1)

Arabic varieties

Auditor	Expected	Count	
Language	Language	Total	% Agreement
ara	ara	77	42.86%
arm	arm	41	85.37%
acm	acm	39	92.31%
alv	alv	54	98.15%

Farsi/Persian, Dari

Auditor	Expected	Count	
Language	Language	Total	% Agreement
fas	fas	291	98.63%
prs	prs	107	98.13%



Within Language (2)

Slavic varieties

Auditor	Expected	Count	
Language	Language	Total	% Agreement
ces	ces	48	100.00%
pol	pol	52	98.08%
slk	slk	72	98.61%
ukr	ukr	205	99.51%
rus	rus	44	100.00%

South Asian varieties

Auditor	Expected	Count	
Language	Language	Total	% Agreement
pnb	pnb	119	98.32%
ben	ben	49	97.96%
tam	tam	63	100.00%
hin	hin	148	89.19%
urd	urd	44	90.91%



Within Language (3)

Other languages

Auditor	Expected	Count	
Language	Language	Total	% Agreement
cmn	cmn	44	100.00%
spa	spa	147	97.28%
tur	tur	60	98.33%
tha	tha	83	93.98%

No dual annotation for Lao due to lack of auditors

?

Auditor	Expected	Count	
Language	Language	Total	% Agreement
eng	eng	57	94.74%
eni	eni	123	98.37%



Dual Annotation Results

Includes all multiply judged segments by auditors of same language, regardless of expected segment language

Overall Language Agreement



Linguistic Data Consortium Dual Annotation Confusion Matrices

All Speech?

All Speech	Ν	Y	
N	309		
Y	451	1904	

All Narrowband?

Alinb	Ν	Υ	
Ν	67		
Y	118	2479	

Signal Quality Judgment

SigQual	diffic	easy	medium
difficult	19		
easy	51	2029	
medium	43	337	92

Single Speaker?

SnglSpkr	Ν	Y
N	206	
Y	220	2238

Speaker Sex

SpkrSex	Ν	Y	U	
N	1087			
Y	58	1383		
U	24	28	20	

Dialect Judgment

Dialect	marked	non-ntv	normal
marked	30		
non-ntv	11	4	
normal	165	20	1420

Includes all multiply judged segments by auditors of same language, regardless of expected language of segment



Reporting judgments where

- a segment was confirmed by an annotator to be in their language

- that language was the expected language

- independently judged by an annotator of another language to also be in that language

 E.g. a Hindi speaker verifies an expected Hindi segment to be Hindi, and an Urdu speaker judges the same segment to be Urdu





Arabic varieties

Expected	Auditor		
Language	Language	Count Total	% Confusion
alv	acm	108	4.63%
ara	acm	108	4.63%
ara	alv	121	19.01%
ara	arm	104	18.27%
arm	acm	111	0.90%
arm	alv	120	0.83%





Czech - Slovak

Expected	Auditor			
Language	Language	Count Total	% Confusion	
ces	slk	179	1.12%	
slk	ces	120	0.83%	

Thai - Lao

Expected Language	Auditor Language	Count Total	% Confusion
lao	tha	140	10.71%
tha	lao	73	6.85%

American English – Indian English

Expected	Auditor		
Language	Language	Count Total	% Confusion
eng	eni	160	28.75%
eni	eng	154	0.65%



Cross-Language

Hindi - Urdu

Expected	Auditor		
Language	Language	Count Total	% Confusion
hin	urd	496	25.40%
urd	hin	786	53.31%

Dari – Farsi/Persian

Expected	Auditor		
Language	Language	Count Total	% Confusion
fas	prs	307	0.33%
prs	fas	(18) <	77.78%





Data Distribution

- Data and audit results distributed to NIST in 6 incremental releases
- Packages contain
 - Full source audio recordings from which segments extracted, in original format
 - Auditor-versions of extracted segments
 - Audit results for segments that meet these criteria
 - Is the segment in the target language? (YES only)
 - Does the segment contain only speech? (YES only)
 - Is all the speech from one speaker? (YES or NO)
 - Does the entire segment sound like narrow-band signal? (YES or NO)
 - Segment metadata table

audid	numeric ID of audit submission in Ib_aud_ann table		noise	cmt free-text auditor comment on signal quality
segid	numeric ID of audited segment		spkr	cmt free-text auditor comment on speaker
Ingid	3-letter language ID as confirmed by auditor		Ing_cmt	free-text auditor comment on language
result	concatenation of responses to yes/no questions		ref	reference status
sex	speaker gender (M/F)		auditor	numeric ID of auditor
spkr	typ speaker's dialect category (native/non-native/etc)	•	src ·	path/name of source audio file
noise	amt auditor's judgment of noise level (easy/hard/etc)	•	duration	length in seconds of the audio segment
noise	typ auditor's list of noise conditions (distortion/etc)			





Linguistic Data Consortium



Conclusions

- Significant volumes of new telephone and broadcast data collection for 24 languages which include several confusable varieties
 - New collection strategies needed to support corpus requirements
- 84 auditors made 22,561 audit judgments yielding 9889 LRE segments
- Auditing kits constructed to support consistency analysis
 - Within-language agreement over 95% except for
 - Most Arabic varieties
 - Hindi, Urdu
 - Thai (but not Lao)
 - Cross-language confusion for
 - Some Arabic pairs, especially involving Modern Standard Arabic
 - Thai/Lao (asymmetry)
 - American English/Indian English (strong asymmetry)
 - Hindi/Urdu (asymmetry)
 - Farsi/Dari (asymmetry; small sample size)
- Corpus supported LRE 2011 evaluation and will be published in LDC catalog pending authorization by sponsors





- Thanks to Speech@FIT group in the Faculty of Information Technology at Brno University of Technology (BUT) in Czech Republic (Brno) for providing speech and bandwidth detection technologies
- Thanks to Alvin Martin and Craig Greenberg for their ongoing support
- Thanks to our collection partners in Tunis, New Delhi and Hong Kong

