Speaker Recognition for Forensic Applications

Joseph P. Campbell

j.campbell@ieee.org



16 June 2014



This work was sponsored under Air Force contract FA8721-05-C-0002. Opinions, interpretations, conclusions, and recommendations are those of the authors and are not necessarily endorsed by the United States Government.





- In forensic or investigative speaker comparison, speech utterances are compared by humans and/or machines for use in court or investigation
 - High-stakes application affecting people's lives demands highest scientific standards
- Unfortunately, methods used in practice vary widely and not always for the better*†
- Methods and practices grounded in science are critical for proper application and nonapplication[‡] of speaker comparison to a variety of international investigative and forensic applications
- Provide a critical analysis of current techniques employed and lessons learned
- Crucial to improve communication between automatic speaker recognition researchers, legal scholars, and forensic practitioners internationally
 - Legal, policy, and societal questions such as allowing speaker comparisons in court
 - Requirements for expert witnesses
 - Requirements for specific automatic or human-based methods to be considered scientific
- You can help!



‡ Schwartz, R., et al., When to Punt on Speaker Comparison?, 162nd Meeting ASA, 2011.



- Background
- Approaches
- Activities
- Request
- Future
- Conclusion



- Forensics: seeks to establish facts of interest using science and technology in the context of the law or in a court of law
- Investigation: systematic inquiry, examination, study, and survey of facts, circumstances, situations, incidents, and scenarios to render a conclusion





Variations of Speaker Comparison*

	Evidential forensic speaker comparison	Investigatory forensic speaker comparison	Speaker comparison within investigatory voice biometrics (AFIS/ASIS-style)
Presentation in court?	Yes	Ν	0
Number of comparisons	Single comparison (or a relatively small set of comparisons within a complex case)		Large or very large number of comparisons
Methods	Auditory + acoustic; HASR etc. (e.g., see Gold & French†)		Fully automatic;* i.e., the investigator makes the database search without listening to the voices in the database
Reports	In a way accepted by the court, usually in the form of some kind of probability statement (not a categorical yes/no) (see Gold & French†)	Either in a way accepted by the court (although it is not intended for a court) or in a simplified form, which might also include a yes/no statement	In the form of a hit list of one or more speakers from the database according to criteria specified by the user (e.g., specifying the size of the hit list; using certain costs for false identification or false rejection; specifying a threshold for a not- in-the-database decision).

Odyssey Keynote: FSR - 5 Joe Campbell, 16 June 2014 * Michael Jessen, handout, Forensic Phonetics course, Summer School in Forensic Linguistics, 2013. † E. Gold, P. French, *International practices in forensic speaker comparison*, IJSLL, 2011.



Forensic Speaker Recognition Examples

- Atlanta Centennial Park Bombing (1996)
 - "There is a bomb in Centennial Park. You have thirty minutes." – 13-second 911 call
 - Are the caller and the suspect in custody the same person?
- Trayvon Martin (2012)
 - Zimmerman claims justified shooting
 - Orlando Sentinel hires "voice experts"
 - "Who was crying for help?"

Trayvon Martin shooting: It's not George Zimmerman crying for help on 911 recording, 2 experts say

5:38 p.m. EST, March 31, 2012 | By Jeff Weiner, Orlando Sentinel

As the Trayvon Martin controversy splinters into a debate about self-defense, a central question remains: Who was heard crying for help on a 911 call in the moments before the teen was shot?

A leading expert in the field of forensic voice identification sought to answer that question by analyzing the recordings for the Orlando Sentinel.







Forensic Speaker Recognition Is it Really That Difficult?

Trial	Truth	Human	Automatic
1	Т	FALSE	TRUE
2	F	FALSE	FALSE
3	F	FALSE	FALSE
4	F	FALSE	FALSE
5	Т	TRUE	TRUE
6	F	FALSE	FALSE
7	Т	FALSE	TRUE
8	F	TRUE	FALSE
9	F	FALSE	FALSE
10	Т	TRUE	TRUE
11	F	TRUE	TRUE
12	F	FALSE	FALSE
13	F	FALSE	FALSE
14	Т	TRUE	TRUE
15	Т	TRUE	TRUE

- NIST Human Assisted Speaker Recognition (HASR)
- Conventional NIST SRE uses too many trials (comparisons) for human processing
- Select a subset of trials for HASR
 - Find most confusable trials using baseline automatic system, then
 - Select most confusable trials by professional, not expert, listeners
- HASR protocol allows listening
- This is difficult, but real forensic data can be more difficult



Incorrect Responses



Challenges in Speaker Recognition* for Humans and Machines [1]

- NIST Speaker Recognition Evaluations (SRE & HASR), Netherlands Forensic Institute and Organization for Applied Scientific Research (NFI-TNO), etc. have addressed significant challenges, e.g.,
- Channel mismatch



Distance to microphone

 Forensic and investigative speaker recognition has additional challenges...

Odyssey Keynote: FSR - 8 Joe Campbell, 16 June 2014

8

6

4

2

0

*Campbell, et al., "Forensic Speaker Recognition," *IEEE Signal Processing* **LINC** *Magazine, Special Issue on Digital Forensics,* v26, i2, Mar 2009, p 95-103.

LINCOLN LABORATORY MASSACHUSETTS INSTITUTE OF TECHNOLOGY



Challenges in Speaker Recognition for Humans and Machines [2]

To be

addressed

- Talkers
 - Unfamiliar to examiner
 - Familiar conversants
 - Multiple talkers
- Stresses
 - Emotional, Loading, Physical
- Styles
 - Conversational, read, orated, loud, yell,... accommodation
 - Plotting, deceptive, disguise
- States
 - Mentally ill, medicated
- Situational mismatch
 - crime vs interview voice

- Language
 - Foreign to examiner
 - Mismatched samples (dialect, accent too)
- Speech samples/segments
 - Few, short, noisy, distorted, noncontemporaneous
 - Few regions of interest*
- Combinations of above factors!
- Mismatch galore!
 - Between samples, models, background, hyperparams



Challenges in Speaker Recognition for Humans and Machines [3]

To be

addressed

- Presentation
 - Scoring
 - Decision?
 - Opinions
 - Court vs Investigation
 - Priors?
- Calibration
- Warnings to users
- The Court's questions
 Negotiable?
- Assigning voice samples to speaker models dilemma
 - Human Automatic systems



- Quantify degradations
 - Negative factors
- Daubert test/factors
- Machine & Human
 Fusion
- Machine vs. Human
- When to punt (not accept a case)*
 - For Machine?
 - For Human?
 - For Human & Machine?



Forensic Speaker Recognition Real Case Data

Situational mismatch



Complex situation



Investigative example



- Running case, triple homicide
 - Suspect ran from scene with a victim's cell phone talking via Bluetooth to a friend to get him
 - Known recordings are calls in jail from suspect to his friends
- Stressed overlapping talkers
 - Dangerous situation
 - 911 call
- Threat call
 - Prompt action?
- It's always something, every case!



- Judge considers the admissibility of scientific evidence
 - Judges are generally not scientists
- US Federal Court and ~half US State Courts under FRE 702
 FRE 702 Testimony by Expert Witnesses to assist the trier of fact
- Daubert* pretrial hearing to assess whether a scientific theory/technique in question
 - 1. Has been or can be tested
 - 2. Has been subjected to peer review and publication
 - 3. Has a known or potential error rate
 - 4. Has existing <u>standards</u> controlling its use that are maintained
 - 5. Has been generally accepted by the scientific community
- Other US States use the Frye test (#5) or case-by-case rules
- Judge Nelson's Order!
- Influencing Canada and UK



- Judge considers the admissibility of scientific evidence
 - Judges are generally not scientists
- US Federal Court and ~half US State Courts under FRE 702
 - FRE 702 Testimony by
- Daubert* pretrial hearing theory/technique in que
 - 1. Has been or can be <u>te</u>
 - 2. Has been subjected to
 - 3. Has a known or poten
 - 4. Has existing standard
 - 5. Has been generally ac
- Other US States use th
- Judge Nelson's Order! Judge Debra S. Nelson, "Order excluding the opinion testimony of Mr. Owen and Dr. Reich," Florida vs. Zimmerman, Circuit Court for 18th Judicial Circuit, 22 June 2013.
- Influencing Canada and UK



HLN WEEKEND EXPRESS

"the Court accepts the opinions of [defense witnesses] that reliable comparison of normal speech to the screams in the 911 call is not possible."

Judge Debra Nelson





- Background
- Approaches
 - Human example
 - Machine example
 - Human and Machine
- Activities
- Request
- Future
- Conclusion



Forensic Speaker Recognition Structured Listening

SDAAT: Super Depotic Applysis and Appet				
CDAAT XGam Month int				
▼ 20060223_142001_5431.1.wav				
32124				
-32124				
:.word <s></s>	hello yes			
b:.phn # hh	ax 1 ov y eh s			
time 0:00.8 0:00.9 0:01.0 0:01	.1 0:01.2 0:01.3 0:01.4 0:01.5 0:01.6 0:01.7 0:01.8 0:01.9 0:02.0 0:02.1 0:02.2			
😅 Zoom All Zoom In Zoom Out Zoom Sel Snap to Ite	em Select Item Transcribe			
Annotations				
Phenomena Type:	substitution			
	A: [eh]			
Rule:	B: [ih]			
	CONTEXT:			
Degree:	1 = I heard neither A or B 0 = I Heard A 1 = I heard A 2 = I heard B Voice Quality:			
Comments:	I heard something allrerent			
ROI Annotation Status:	1 = Indeterminant 0 = Incomplete 1 = Partially Complete 2 = Complete			
Progress: 2/1554				

- Transcribe speech into words
- Set of rules for American English transformations (other lx in process)
- Apply transformation rules to each Region of Interest (ROI) by listening to the speech to score it and produce likelihood information
- Examples: $/ih/ \rightarrow /eh/$ substitution, e.g., pin \rightarrow pen

 $/I/ \rightarrow$ (reduced) reduction, e.g., almond \rightarrow ahmond

• Final output is a report; e.g., "there is support for the hypothesis that the samples come from the same speaker" (with an explanation)



- Challenges
 - Long time to make decision
 - Complete detailed analysis can take over a week
 - Skilled analyst required
 - Native in language of samples
 - Might lack required amount of data for each evaluation
 - Enough data for Regions of Interest in structured listening
 - Reliability?
 - Process can be subjective
 - Performance of various methods not well quantified
- Combine with automatic method?



- Background
- Approaches
 - Human example
 - Machine example
 - Human and Machine
- Activities
- Request
- Future
- Conclusion



Automatic Speaker Recognition System Architecture



- Pre-processing
 - Input: Raw speech signal
 - Output: Salient features about the speaker
- Classifier
 - Input: Known and Questioned features
 - Output: Likelihood score

- Calibration
 - Input: Likelihood score
 - Output:
 - Likelihood score
 - Match probability
 - Decision
- Fusion (optional)
 - Multiple classifier inputs
 - Combined output

\bigotimes

Anatomy of a Speaker Comparison System



Ex. i-vector Speaker Recognition System





- Produce consistent results interpretable by humans
 - Across conditions
- Mitigate mismatches between known, questioned, and multiple training data sets
- Calibrator is trained for this purpose





- Background
- Approaches
 - Human example
 - Machine example
 - Human and Machine
- Activities
- Request
- Future
- Conclusion



- How to combine human and automatic speaker recognition?†
 - Separate processes followed by score combination?
 - Weigh each process dynamically?
 - Mitigate observation bias?
 - Deal with variations of subjects, humans, machines, and samples?
 - Consistency and repeatability?
 - Logically consistent results?
- Best Practices are needed to address these questions
- Evaluations of processes are needed
 - NFI-TNO Forensic Speaker Recognition Evaluation, 2003
 - NIST's Human Assisted Speaker Recognition (HASR)
 - Next NIST SRE more like forensic domain samples?
 - Others?

Progress toward Daubert factors?

Odyssey Keynote: FSR - 23 Joe Campbell, 16 June 2014



- Background
- Approaches
- Activities
 - US: SWG-Speaker
- Request
- Future
- Conclusion



US Advances in the Science and Practice of Speaker Recognition

– Vocabulary

- NRC/NAS, Strengthening Forensic Science in the United States: A Path Forward, 2009
- Investigatory Voice Working Group
 - Use Case Committee Report, 2009
 - Collection Standards Committee Report, 2009
- Schwartz, et al., When to Punt on Speaker Comparison?, 162nd ASA, San Diego, 2011
- Standard Operating Procedure for Forensic Speaker Recognition, 2013
- ANSI/NIST-ITL Type-11 Record Standard, 2013
- Scientific Working Group for Forensic and Investigative Speaker Recognition (SWG-Speaker), 2013
 - Research, Dev, Test, Eval Science in Law
 - Best Practices





- The Best Practices Committee seeks to improve forensic science through establishing best practices recommendations
- Develop best practices for
 - Collection protocols
 - Speech materials
 - <u>Audio recording</u> and <u>collecting related data</u> used for speaker recognition
 - Transmission of this audio and related data
 - ANSI/NIST-ITL Type-11 Record and companion Records
 - Proper <u>application</u> of speaker recognition technologies
 - Examination and reporting guidelines
 - Bias concerns
- Training, examiner certification, and laboratory accreditation advice

Ex. "Does the Questioned Voice Recording Share the Same Source as a Known Voice Recording?"?



Records

Type-1: Mandatory record submitted with each transaction, "transaction header information"

Type-2: Transaction related data; e.g., subject's name and other biographic information, reason for booking, any charges, etc.

Type-11: Voice data and voice metadata for the subject in corresponding Type-2 ("voice data" can be marks for the subject in the original audio stored)

Type-20: Repository of original data if in digital format (original format, nonmanipulated, and unprocessed), which includes raw evidence (without redaction)

Type-xx: Other record types can be transferred that might not be used in the speaker recognition process; e.g., photo of subject, signed papers, etc.



Proposed Best Practices Process



- Validation methods involving humans?
 - Black Box Examiner Study? NIST Forensic SRE? Corpora? Funds?



- Background
- Approaches
- Activities
- Request
 - Future
 - Conclusion





- Develop candidate Best Practices
 - Submit to SWG-Speaker/OSAC for consideration
- Pursue Daubert factors
- Improve robustness
 - Core classifiers
 - Calibration
 - Limited in-domain training data
 - Combined processes
- Work with analysts/examiners to improve usability and performance
 - Human in the loop
- Rise to the challenges of forensic and investigative data
 - Handling speaker variability from stress and emotional state
- Participate in forensic/investigative-style evaluations



Standards, practices, evaluations, and data for forensic and investigative speaker recognition



OSAC Organization of Scientific Area Committees (NIST)

- Speaker Recognition Subcommittee (IT/Multimedia Sci Area Committee)
- SWG-Speaker transition to OSAC
- ANSI/NIST-ITL 1-2011 Update:2013, Data Format for the Interchange of...
 - Record <u>Type-11</u>: Forensic and investigatory voice record
- <u>IAFPA</u> International Association for Forensic Phonetics and Acoustics
 - Code of Practice, workshops, shared
 IJSLL journal with <u>IAFL</u> International
 Association of Forensic Linguists
- AES <u>Audio Forensics</u>: Techniques, Technologies, and Practice
- <u>FAS</u> ASA Forensic Acoustics Subcommittee (Speech Comm TC)

- <u>ENFSI</u> European Network of Forensic Science Institutes
 - <u>FSAAWG</u> Expert Working Group for Forensic Speech and Audio Analysis
 - Monopoly 2011 Methodological guidelines for semi-automatic and automatic speaker recognition for case assessment and interpretation
- <u>AGSE</u> Arbeitsgemeinschaft Sprechererkennung (WG of German speaking forensic speech and audio specialists from European Labs)
- <u>EAB</u> Workshop on Biometrics and Forensics
- <u>FSA</u> UK Home Office Forensic Speech and Audio Group
 - Bring forensic speech and audio under the regulation of ISO 17025



Standards, practices, evaluations, and data for forensic and investigative speaker recognition

- <u>ASSTA</u> Forensic Speech Science Committee Australasian Speech Science and Technology Assoc
- <u>EAFS</u> European Academy of Forensics Science
- <u>IAFS</u> International Association of Forensics Sciences
- <u>ICFIS</u> International Conference on Forensic Inference and Statistics
- <u>AAFS</u> American Academy of Forensic Sciences
- <u>NAS</u> National Academy of Sciences
- Evaluations and data
 - NIST SRE HASR (toward forensic)



- NFI FRITS (after 2003 NFI-TNO FSRE)
- Cambridge DyViS corpus

• ISCA

- <u>Interspeech</u>, ICSLP, Eurospeech special sessions on forensic speech
- <u>Odyssey</u> The Speaker and Language Recognition Workshop
- <u>SpLC</u> Speaker and Language Characterization SIG
- <u>AFCP</u> Association Francophone de la Communication Parlée SIG and <u>JEP</u>
- iberSPEECH
- Programs



SIIP Speaker Identification Integrated Project (EU FP7)

- <u>BBfor2</u> Bayesian Biometrics for Forensics project
- ICT COST Action IC1106 Integrating Biometrics and Forensics for the Digital Age



- Speaker recognition is successfully used today in a variety of applications, but must be applied responsibly with caution*
- Need to address factors in forensic domain that degrade recognition performance
 - Increase robustness and effectiveness
 - Improve calibration and efficiency
- Real case data can be extremely challenging
 - Opportunities for research!
- Many challenges to be addressed to satisfy Daubert test
 - Additional international challenges and opportunities
- Please contact me to share ideas!
- Sauna is a very Finnish way for decision making Pasi

Odyssey Keynote: FSR - 33 Joe Campbell, 16 June 2014 * Campbell, J., Shen, W., Schwartz, R., Bonastre, J.F., Matrouf, D., Forensic Speaker Recognition: A Need for Caution, IEEE Signal Processing Magazine, Vol. 26, Issue 2, p. 95-103, March 2009.



Questions?

