Speech Analytics Evaluation Project (OpenSAT)

Fred Byers NIST/Information Technology Laboratory Gaithersburg, MD



#PSCR2020



DISCLAIMER

Certain commercial entities, equipment, or materials may be identified in this document in order to describe an experimental procedure or concept adequately.

Such identification is not intended to imply recommendation or endorsement by the National Institute of Standards and Technology, nor is it intended to imply that the entities, materials, or equipment are necessarily the best available for the purpose.

* Please note, unless mentioned in reference to a NIST Publication, all information and data presented is preliminary/in-progress and subject to change

NIST

PSCR2020



6.44

Agenda

- Section 1
 - OpenSAT overview
 - Analytics (Tasks)
 - Data
 - Evaluation process
- Section 2
 - Equations and 2019 Test Results
 - Moving forward

OpenSAT Evaluation Series

NIST Open **Speech Analytics Evaluations**

Open

Anyone can participate

SAT Speech Analytic Technologies

Evaluation

A process lasting several months

Series

The Evaluation is reoccurring

OpenSAT Speech Analytics Testing



Goal

Advance the performance of evolving applications that rely on speech analytic technologies.

Purpose

To facilitate data-driven research to evaluate and improve the performance of speech analytics for public safety applications.



Speech Analytic Challenges in Public Safety Communications



Acoustics

Loud background sounds

Varying volume levels

Multiple types of background sounds

Speakers in the background

Crowd noise

Type of mic



Speech

Increased vocal effort (Lombard effect)

Urgency

Stress



How OpenSAT Evaluations Improve Applications

Researchers participate in OpenSAT



Multiple research teams simultaneously and independently

- work with the same data and the same metrics
- focus on common speech analytic technologies (tasks)
- train and develop these tasks for specific data
- improve the core technologies that applications depend on



Tapping the researcher competitive spirit

- NIST provides researchers the infrastructure, tools and testing model
- researchers compare system performance from the evaluation data
- state-of-the-art solutions are brought to the industry by researchers
- industry can compare their technology performance against others

The Core Technologies (Tasks) Evaluated

Currently, Three Tasks are Evaluated in OpenSAT



¹ Speech Activity Detection



Automatically detect the beginning and ending of speaking segments by timestamps. 2 Automatic Speech Recognition



```
going inside now
```

Automatically transcribe verbatim all audible speech to readable text.

³ Keyword Search

***	** alarm	* * * * *	* ***
***	*****	****	****
***	need back	kup **	* * * *

Automatically detect all spoken instances of a predetermined list of words and phrases.

To Evaluate You Need Relevant Data

Test data relevant to public safety

Relevant and high-quality data improves the evaluation of evolving technologies

- Hands-free voice interface
- Voice, video, GPS situational awareness system
- Automated transcription for reporting



The Linguistics Data Consortium (LDC) at the University of Pennsylvania (UPenn)

The Linguistic Data Consortium (LDC) recorded and transcribed dozens of fire rescue gameplaying communications.

Audio recordings were collected in controlled conditions.

The collection of recordings and transcriptions were organized to create specialized data sets for distribution to researchers.

This collection was funded by DHS.

LDC Catalog	Corpus Name
	Researchers and NIST
LDC2019E37	SAFE-T Corpus Speech Recording Audio Training Data R1 V1.1
LDC2019E38	SAFE-T Speech Recording Development Audio and Transcripts
LDC2019E36	SAFE-T Speech Recording Training Data Transcripts V1.1
LDC2019E50	OpenSAT19 Public Safety Communications Simulation Evaluation Data V1.1
LDC2019E53	OpenSAT19 Public Safety Communications Simulation Development Data V1.1
	NIST only
LDC2019R03	SAFE-T Corpus - Speech Recording Metadata for Dev Selection
LDC2019R10	SAFE-T Corpus Speech Recording Audio R1 V1.1
LDC2019R09	SAFE-T Corpus Speech Recording Metadata for Eval Selection
LDC2019R17	SAFE-T Speech Recording Evaluation Data Transcripts V1.1
LDC2019R19	SAFE-T Speech Recording Evaluation Audio for NIST V1.2
LDC2019R22	OpenSAT19 PSC Simulation Evaluation Data for NIST Preview
LDC2019R25	SAFE-T Corpus Speech Recording Audio R2 V2.0
LDC2019R31	SAFE-T Speech Recording Development Evaluation Superset Audio
LDC2019R32	SAFE-T Speech Recording Training Data Transcripts R2
LDC2019R33	SAFE-T Speech Recording Development Evaluation Superset Transcripts
LDC2019R34	SAFE-T Corpus Speech Recording Audio Training Data R2

Audio Collection Process

Audio Collected

Communications are recorded during the game.



Background noises are fed through headphones in various noise types and varying volume levels. Participants were recruited to play a collaborative board game to rescue victims from a burning fire.

First responder background noises are fed through headphones of participants while playing the game.

The background noises and sense of urgency elicit varying levels of vocal effort intended to simulate real-world operational speech.

Selected audio files were transcribed to create development and evaluation data sets.

Data Sets Created



The recorded audio and transcripts are organized to create the SAFE-T Speech Corpus.

The SAFE-T Speech Corpus package contains the three types of data sets for researchers.

Apply Speech Analytic To the Data



Analytic Output Is Evaluated



Evaluation - Score



Scores, Ranking, Workshop, Knowledge Transfer



Scores & Ranking

Researchers see how their system's

performance compares to others.

Workshop

Researchers learn how the top ranked systems work from presentations and O8/

systems work from presentations and Q&A.

Technology Transfer

Researchers bring their knowledge to the

industry and system development efforts.

OpenSAT20 Web Site



OpenSAT Web Site Researcher Dashboard

(Screen Shot of Dashboard)

OpenSAT Series	Dashboard fbyers@nist.gov(Participant) -
Registration Workflow	Submission Management
Please follow the steps below. Green items can be visited at any time.	Please use the links below to manage submissions.
1. Create profile	Team2
2. Create/Join a site	SAD: (development, evaluation)
3. Create/Join a team	KWS: (development, evaluation)
4. Select Tasks	ASR: (development, evaluation)
5. Sign and upload license	License Agreements
6. Workflow completed!	Please use the links below to manage your license agreements.
	'OpenSAT 2020 Terms and Conditions' for Team2
1 Sites	'LDC Data License Agreement' for Team2
Owner of the following Site(s):	
Fred2	▲ Datasets
A	Please use the links below to view information about available datasets or to view download options.
L leams	SAD scoring and validation tool
Owner of the following team(s):	
Team2 (Fred2)	

OpenSAT Evaluation Plan Provides All The Details

NIST Open Speech Analytic Technologies 2020 Evaluation Plan (OpenSAT20)

1	Introductio	on2				
2	Planned Sc	chedule 2				
3	Data					
4	Tasks - Ov	verview and Performance Metrics				
5	Evaluation	Rules				
6	Evaluation	Protocol				
7	7 System Input, Output, and Performance Metrics					
App	oendix I	SAD7				
App	oendix II	KWS10				
App	oendix III	ASR				
App	oendix IV	SAD, KWS, and ASR - System Output Submission Packaging				
App	oendix V	SAD, KWS, and ASR - System Descriptions and Auxiliary Condition Reporting				

Section 2

Equations and 2019 Test Results

Speech Activity Detection Scoring

Speech Activity Detection (SAD) Output Evaluated

 Scoring Server

 Output →

 Speech segments -
Start-Stop
Missed
Out of sync
Incorrect

 $DCF\left(\theta\right)=0.75\times P_{\mathsf{FN}}\left(\theta\right)+0.25\times P_{\mathsf{FP}}\left(\theta\right)$

Detection Cost Function (DCF)

System goal is to detect the start and end of when a person speaks.

Scoring server calculates error probabilities from missed speech detection time and incorrect detection time.

 P_{FN} – Probability of a system missing speech P_{FP} – Probability of a system incorrectly detecting speech

 $P_{FN} = \frac{\text{total missed speech time}}{\text{total actual speech time}}$

 $P_{FP} = \frac{\text{total incorrect speech time}}{\text{total actual nonspeech time}}$

Speech Activity Detection Scores

Speech Activity Detection (SAD) Output Evaluated

Detection Cost Function (DCF)

System scores for both loud and quiet background noise levels

	Scoring Server		System	DCF (Loud BG)	DCF (Quiet BG)	DCF Average
Output		Carrier	1	0.0481	0.0624	0.0525
Output →	DCF	Score →	2	0.1093	0.1175	0.1099
	Calculated		3	0.1479	0.0854	0.1245
			4	0.1492	0.0911	0.1352
	$0.75 \dots D$ (0) $1.0.25$	·· D (0)		0.0		

 $DCF(\theta) = 0.75 \times P_{FN}(\theta) + 0.25 \times P_{FP}(\theta)$

0.0		1.0
Perfect	DCF	Bac

Speech Activity Detection Scores

Speech Activity Detection (SAD) Output Evaluated

Detection Cost Function (DCF)

Individual speakers' scores with loud background noise level

	Scoring Server		Speaker	FN (Prob. Miss)	FP (Prob. Incorrect)	DCF
Output >		Score >	1	0.8992	0.9751	0.9181
Output 7	DCF	Score →	2	0.0443	0.0345	0.0419
	Calculated		3	0.1763	0.0745	0.1508
			4	0.0229	0.0531	0.0305
$DCF(\theta) = 0.75 \times P_{FN}(\theta) + 0.25 \times P_{FP}(\theta)$			0.0		.0	

0.0		
Perfect	DCF	Bac

Automatic Speech Recognition Scoring

Automatic Speech Recognition (ASR) Output Evaluated



Word Error Rate (WER)

System goal is to transcribe all audible speech to readable text.

Scoring server detects and counts the number of missed words, inserted words, and mismatched words.

- N_{Del} Number of missed words
- N_{Ins} Number of inserted words (words added where there is no speaking)
- N_{subst} Number of non-matching words (e.g., misspelled, wrong word)

N_{Ref} – Total number of words in the transcript

Automated Speech Recognition Scores

Automatic Speech Recognition (ASR) Word Error Rate (WER) **Output Evaluated** System scores Scoring Server Loud and Quiet Backgrounds WER (Loud BG) WER (Quiet BG) Average WER System Output → Score \rightarrow WER 1 15.7% 15.0% 15.4% Calculated 2 36.1% 22.6% 29.3% WER = $\frac{(N_{Del} + N_{Ins} + N_{Subst})}{(N_{Del} + N_{Ins} + N_{Subst})}$ 0.0% ------ 100% N_{Ref} Perfect WER Bad

Automated Speech Recognition Scores

Automatic Speech Recognition (ASR) Output Evaluated

Word Error Rate (WER)

Speaker score examples for a system

Scoring Server		WER with Loud Background						
Outrast			Speaker	Total Words	Subst.	Missed	Inserted	WER
Output →	WER	Score →	798	150	10	11	4	16.7%
	Calculated		798	308	18	112	5	43.8%
			8801	131	7	16	1	18.3%
			8801	227	29	46	0	33.0%
$WER = \frac{(N_{Del} + N_{Ins} + N_{Subst})}{N_{Ref}}$				0.0% Perfect	 W	 ′ER	100% Bad	6

Keyword Search Scoring



System goal is to automatically detect all keywords in the audio.

Scoring server calculates error probabilities from the number of missed keywords and incorrect detections.

P_{FN} – Probability of a system missing a keyword

P_{FP} – Probability of a system detecting an incorrect word or an incorrect location

 $P_{FN} = \frac{\text{number of keywords missed (FN)}}{\text{total number of keywords}}$

$$P_{FP} = \frac{\text{number of keywords misplaced (FP)}}{\text{total duration (seconds)}} / \text{#target words}$$

Keyword Search Scores



Keyword Search Scores



Keyword Search Scores



Moving Forward

• 2020

- Complete OpenSAT20 Evaluation
- Compare OpenSAT20 and OpenSAT19 results (4th quarter)
- Obtain high-quality transcription of real-world operational audio

• 2021

- OpenSAT21
- Continue with the LDC data
- Include real-world operational data for testing
- Add the speaker diarization task (keeping track of speakers A, B, C, etc.)

Contact Us





NIST 100 Bureau Dr, MS 894 Gaithersburg, MD 20899



OpenSAT Web Site <u>https://sat.nist.gov</u> Group Web Site <u>https://www.nist.gov/itl/iad/mig/opensat</u>



NIST/ITL/ Information Access Division/Multimodal Information GroupFred Byersfrederick.byers@nist.govOffice: 301-975-2909

THANK YOU



#PSCR2020

