

## Abstract: What is this challenge about?

A public evaluation on the performance of systems dealing with two important subtasks in Computational Auditory Scene Analysis: the classification of acoustic scenes and the detection of acoustic events. Main challenge aims:

- Help the research community move a step forward in better defining the specific tasks.
- Offer a comprehensive evaluation framework and methodology for the tasks.
- Provide incentive for researchers to pursue research in this field by making available a well-structured and fully documented dataset.
- Help to shed light on controversies that currently exist in the tasks.

## Challenge results

### Task SC: Acoustic Scene Classification

A train/test classification task: 10 classes, 10 x 30-sec audio recordings per class. Five-fold crossvalidation.

Participants	Code	Method	Lang
Chum et al.	CHR	Various features at 2 frame sizes, classified either: (a) per-frame SVM + majority voting; (b) HMM	Matlab
Elizalde	ELF	See poster: "An I-Vector Based Approach ..."	Matlab
Geiger et al.	GSR	See poster: "Large-Scale Audio Feature Extraction and SVM ..."	Weka/HTK
Krijnders and ten Holt	KH	"A Tone-Fit Feature Representation ..."	Python
Li et al.	LTT	Wavelets, MFCCs and others, classified in 5-second windows by treebagger, majority voting	Matlab
Nam et al.	NHL	Feature learning by sparse RBM, then event detection and max-pooling, classified by SVM	Matlab
Nogueira et al.	NRI	MFCCs + MFCC temporal modulations + event density estimation + binaural modelling features, feature selection, classified by SVM	Matlab
Olivetti	OE	Normalised compression distance (Vorbis), Euclidean embedding, classified by Random Forest	Python
Patil and Elhilali	PE	See poster: "Multiresolution Auditory Representations ..."	Matlab
Rakotomamonjy and Gasso	RG	See poster: "Histogram of Gradients of Time-Frequency Representations ..."	Matlab
Roma et al.	RNH	See poster: "Recurrence Quantification Analysis Features ..."	Matlab
Baseline		MFCCs, classified with a bag-of-frames approach	Python

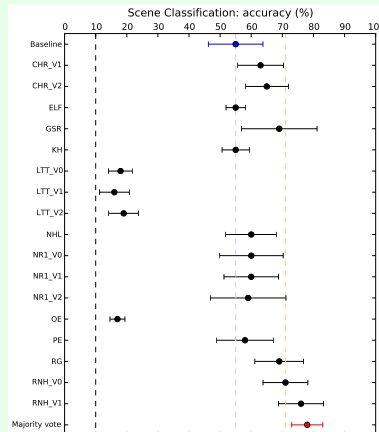


Figure: Classification accuracies (mean and 95% CI). Dashed lines indicate: chance performance (black); baseline (light blue); mean human accuracy (orange). "Majority vote" is a meta-classifier using the majority decision over all submissions.

	bus	busystreet	office	openairmarket	park	quietstreet	restaurant	supermarket	tube	tubestation
bus	81	3	4	1	4	6	2			
busystreet	69	14	2	1	2	1	3	3	5	
office	55	13	9	12	4	3	1	3		
openairmarket	1	2	59	13	9	12	3	2		
park	1	8	3	51	29	3	2	1	1	
quietstreet	5	4	3	29	43	9	5		1	
restaurant	1	1	16	5	53	21	2	3		
supermarket	6	5	6	6	4	7	10	42	7	7
tube	7	1	1	2	2	5	3	44	28	
tubestation	5	16	1	4	1	2	3	8	19	41

Table: Confusion matrix, aggregated over all submitted algorithms.

### Task OL: Acoustic Event Detection (Live Monophonic Audio)

Task: identify the start, end and class of events in an office scene. Monophonic sequences, real recordings. 16 classes.

Participants	Code	Method	Lang
Chauhan et al.	CPS	Feature extraction - Segmentation - Likelihood ratio test classification	Matlab
Diment et al.	DHV	MFCCs (features) - HMMs (detection)	Matlab
Gemmeke et al.	GVV	See poster: "An Exemplar-Based NMF Approach ..."	Matlab
Niessen et al.	NVM	See poster: "Hierarchical Modeling Using Automated Sub-Clustering ..."	Matlab
Nogueira et al.	NR2	See poster: "Automatic event classification using front end single channel noise reduction, ..."	Matlab
Schröder et al.	SCS	See poster: "On the Use of Spectro-Temporal Features ..."	Matlab
Vuegen et al.	VVK	See poster: "An MFCC-GMM Approach ..."	Matlab
Baseline		NMF with pre-extracted bases (detection)	Matlab

System	Evaluation Method									
	Event-Based				Class-Wise Event-Based				Frame-Based	
	F (%)	F <sub>offset</sub> (%)	AEER	AEER <sub>offset</sub>	F (%)	F <sub>offset</sub> (%)	AEER	AEER <sub>offset</sub>	F (%)	AEER
CPS	2.23	1.65	2.285	2.301	0.65	0.49	1.872	1.891	3.82	2.116
DHV	26.67	22.43	2.519	2.676	30.72	25.29	2.182	2.370	26.0	3.128
GVV	15.52	13.46	1.779	1.831	13.21	12.03	1.556	<b>1.606</b>	31.94	1.084
NVM_1	32.57	24.95	1.864	2.095	29.37	21.80	1.639	1.899	40.85	1.115
NVM_2	34.16	26.28	1.852	2.095	33.05	24.88	1.602	1.877	42.76	1.102
NVM_3	34.51	27.01	1.827	2.052	33.52	24.65	1.575	1.846	45.50	1.212
NVM_4	30.47	24.68	1.906	2.083	28.17	21.62	1.650	1.849	42.86	1.360
NR2	19.21	15.26	3.076	3.244	21.54	17.64	2.857	3.010	34.66	1.885
SCS_1	39.47	36.74	1.669	1.749	36.33	34.20	1.579	1.677	53.02	1.167
SCS_2	<b>45.17</b>	<b>41.06</b>	<b>1.601</b>	<b>1.727</b>	<b>41.51</b>	<b>38.32</b>	<b>1.511</b>	1.646	<b>61.52</b>	1.016
VVK	30.77	25.40	2.054	2.224	24.55	20.36	1.762	1.949	43.42	<b>1.001</b>
Baseline	7.38	1.58	5.900	6.318	9.00	1.86	5.960	6.462	10.72	2.590

Table: Evaluation measures: the F-measure and the acoustic event error rate (AEER). Each is calculated three ways: per event, per event (normalised for class prevalence), and per audio frame.

### Task OS: Acoustic Event Detection (Synthetic Polyphonic Audio)

Task: identify the start, end and class of events in an office scene. Synthetic mixtures of recorded events, with controlled polyphony. 16 classes.

Participants	Code	Method	Lang
Diment et al.	DHV	MFCCs (features) - HMMs (detection)	Matlab
Gemmeke et al.	GVV	See poster: "An Exemplar-Based NMF Approach ..."	Matlab
Vuegen et al.	VVK	See poster: "An MFCC-GMM Approach ..."	Matlab
Baseline		NMF with pre-extracted bases (detection)	Matlab

System	Event-Based				Class-Wise Event-Based				Frame-Based	
	F (%)	F <sub>offset</sub> (%)	AEER	AEER <sub>offset</sub>	F (%)	F <sub>offset</sub> (%)	AEER	AEER <sub>offset</sub>	F (%)	AEER
DHV	<b>8.45</b>	6.18	4.741	4.860	<b>9.73</b>	<b>7.58</b>	4.028	4.147	<b>13.08</b>	8.426
GVV	7.69	<b>7.33</b>	1.913	1.920	6.69	6.51	1.584	1.591	10.30	<b>1.553</b>
VVK	5.80	5.28	<b>1.885</b>	<b>1.895</b>	5.10	4.77	<b>1.436</b>	<b>1.445</b>	5.77	2.106
Baseline	4.98	0.24	6.507	6.895	6.69	0.18	5.389	5.782	6.88	3.047

Table: Evaluation measures (for details see OL task above).