# **Visual Analytics for Audio**

Mark Hasegawa-Johnson, Camille Goudeseune, Kai-Hsiang Lin, David Cohen, Xi Zhou, Xiaodan Zhuang, Kyung-Tae Kim, Hank Kaczmarski and Tom Huang

> **ILLINOIS NIPS Workshop on Visual Analytics December 11, 2009**

# **Visual Analytics for Audio**

- Visualizing Large Audio Testbeds
  - Multi-Day Timeline Audio
  - Milliphone = 1000 Microphones
- Signal Features
  - Multiscale FFT/ Multiscale Filterbanks
  - Auditory Salience
- Likelihood Features: Acoustic Event Detection
  - Minimum Bayes Risk Feature Selection
  - Tandem ANN/HMM Event Detection
  - UBM/MAP Event Verification
- Summary

# <u>Testbed #1: Portable Multi-Day</u> <u>Audio Timeliner</u>

- **Dramatis Personae:** Emergency first responders (EFRs)
- Analysis Object: One microphone, one month
- <u>Act 1, Scene 1:</u> EFRs arrive on scene, download surveillance audio to a handheld
- **Objective:** Event diagnosis, prognosis, & management



# <u>Testbed #2: 1000 Microphones =</u> <u>One Milliphone</u>

- **Dramatis Personae:** Command center data analysts
- <u>Analysis Object</u>: 1000 microphones, 24 hours
   <u>Act 2, Scene 1</u>: Analyst in a Virtual Reality Theater (the Beckman CUBE) seeks anomalies in a large dataset
   <u>Objective</u>: Find the anomalies

#### **Dataset #1: Meeting Room Audio** 30 annotators, 24 hours of data, 14 acoustic events



### Dataset #2: Willard Airport

# 24 hours of audio, `labeled' by commercial airplane takeoff & landing records (inadequate!)







### **Signal Features**

- Multi-Scale Features
  - Motivation: Processing in auditory cortex
  - Implementation: Multiscale FFT
- Auditory Salience
  - Motivation: Visual salience
  - Implementation: DoG in appropriate features
  - Real-world tests: Illinois labeling of AMI corpus

#### Motivation: Rate-Scale Features (Carlyon & Shamma, 2003)



#### <u>Multiscale FFT: Formulas</u> (Cohen et al., FODAVA 09-01)

 $X[2k] = F_{N/2} \{ x_0[n] + x_1[n] \}, \ 0 \le k < N/2$  $X[2k+1] = F_{N/2} \{ e^{-j2\pi n/N} (x_0[n] - x_1[n]) \}, \ 0 \le k < N/2$ 

#### $X[2k] = X_0[k] + X_1[k], \ 0 \le k < N/2$

 $X[2k+1] = X_0[k + \frac{1}{2}] - X_1[k + \frac{1}{2}], \ 0 \le k < N/2$ 

# Multiscale FFT: Butterflies

#### (Cohen et al., FODAVA 09-01)



### Digression: How to Display Audio Events?

Channel Model: Information in the Visual Channel

 $\psi(\vec{r}) = h(\vec{r}) * \phi(\vec{r}) + v(\vec{r})$ 

$$C = \int \left( 1 + \frac{|H(\vec{\vec{\Omega}})|^2 S_{\phi}(\vec{\Omega})}{S_v(\vec{\Omega})} \right) d\vec{\Omega}$$

To Maximize Information Transfer, Choose

$$\begin{array}{lll} \Delta \vec{r} &=& \mathrm{BW}\left(H(\vec{\Omega})\right) \\ &\approx& 2\mathrm{pixels} \\ \Delta \phi(\vec{\Omega}) &\propto& \sqrt{\frac{|H(\vec{\Omega})|^2}{S_v(\vec{\Omega})}} \\ &\approx& \frac{1}{128}\max\phi(\vec{r}) \end{array}$$

#### <u>Salience</u>

### Salient objects/events **POP OUT** Computer vision: Difference of Gaussians

Surround-

Center

DoG



level 0 (= original image)

Smooth

down-sample

then

#### Audio Salience: Data Annotation First pass: subjects label "salient" vs. "non-salient" Second pass: subjects identify audio events

Cian - Underined File	Marine	Ontions Window Help	menu bar			
vide	eo viewer	Grid Text Subtitles Volume: 100 0	Audio Recognizer	volume control ta	ab	· · · ·
00		Rate:	0.01.00.750 485	1 00 v v	n († 2000) 19	20
	a lavigation			ection Mode 🔝 Loop Mode		
tiers	10.58.000 waveform	viewer <sup>1.000</sup> 00:00:59.500	00:01:00.000 00:	:01:00.500 00:01:01.000	00:01:01.500 00:0	1:02.000
	0:58.000 00:00:58.500	00:00.59.000 annotat	ion 1:00.000 00:	:01:00.500 00:01:01.000	00:01:01.500 00:0	1:02.000
door slam pj foot step				selection		
chair						
spoon and cup						
applause [0]						
laugher [9]	timeline v	iewer				
paper						

# Audio Salience: Signal Model



#### Audio Salience: Sample Results



### Audio Salience EER: ~76%





### **Likelihood Features**

- Acoustic Event Detection: Why is it hard?
  - Relevant information scattered across multiple scales, multiple (t,f) coordinates
  - Low SNR
- Minimum Bayes risk feature selection
- ANN/HMM hybrid AED
- UBM/MAP rescoring

### Non-Speech Audio Events

#### Motivation

"Activity detection and description is a key functionality of perceptually aware interfaces working in collaborative human communication environments... detection and classification of acoustic events may help to detect and describe human activity..." (CLEAR-AED Task Brief)

#### Difficulties

- Negative SNR (speech is "background noise")
- Unknown spectral structure
- Different spectral structure for each event type

### **Difficulty #1: Negative SNR**



# Difficulty #2: Class-Dependent Scale & Structure



#### Discriminative Feature Selection (Zhuang et al., ICASSP 2008)

- Problem: what acoustic features are relevant for detecting non-speech acoustic events?
- Input: (x<sub>i</sub> ∈ ℜ<sup>D</sup>) includes many acoustic features invented for speech processing (MFCC, PLP, energy, ZCR)
   Output: (f<sub>i</sub> ∈ ℜ<sup>d</sup>) selects the most useful features:

$$f_i = W x_i$$

where  $W^T = [w_1, \ldots, w_K]$ , and  $w_k$  is an indicator vector (only one non-zero element)

■ Hidden Markov Modeling: the label sequence
Y\* = [y<sub>1</sub><sup>\*</sup>,...,y<sub>N</sub><sup>\*</sup>], y<sub>i</sub> ∈ {keyjingle, footstep,...} is chosen by a hidden Markov model observing F = [f<sub>1</sub>,..., f<sub>N</sub>]:

 $Y^* = \arg \max p(F|Y)p(Y)$ 

#### Bayes Error Rate (Zhuang et al., ICASSP 2008)

#### Bayes Error Rate

Let  $w_k$  be an indicator vector (all zeros except for one element). The Bayes-optimal error rate of a classifier observing feature  $w_k^T x$  is

$$P(\text{error}) = \int \int P\left(y \neq \arg\max p(w_k^T x, y)\right) dy dx$$

Bayes Error Rate Approximated on a Database

$$\mathcal{F}(w_k) = \frac{1}{N} \sum_{i=1}^{N} \delta\left(y_i \neq \arg\max p(w_k^T x_i, y_i)\right)$$

#### Feature Selection Algorithms (Zhuang et al., ICASSP 2008)

#### Hard-Bayes-Error Feature Selection

For k = 1, ..., K, Choose the indicator vector  $w_k$  ( $w_k$  is all zeros except for one nonzero element) to minimize

$$\mathcal{F}(w_k) = \frac{1}{N} \sum_{i=1}^{N} \delta\left(y_i \neq \arg\max p(w_k^T x_i, y_i)\right)$$

#### Soft-Bayes-Error Feature Selection

For k = 1, ..., K, Choose the indicator vector  $w_k$  ( $w_k$  is all zeros except for one nonzero element) to minimize

$$\mathcal{F}_{S}(w_{k}) = \frac{1}{N} \sum_{i=1}^{N} \operatorname{rank}\left(y_{i} \left| w_{k}^{T} x_{i} \right.\right)$$

# Acoustic Event Detection Accuracy

#### (Zhuang et al., ICASSP 2008)



- MFCC26DAZ = 26 Mel-frequency cepstral coefficients + deltas + acceleration
- DERIVE26DAZ = 26 Derived features + deltas + acceleration
- DERIVE78 = 78 Derived features

### Tandem ANN/HMM for AED



# **Rescore Using Supervectors**

Mixture Gaussian Model

$$p(\vec{x}|q) = \sum_{k} c_{qk} \mathcal{N}\left(\vec{x}; \vec{\mu}_{k}, \Sigma_{k}\right)$$

#### MAP Adaptation to the p'th Detected Segment

$$\vec{\mu}_k^{(p)} = \frac{\sum_t \gamma_k(t) \vec{x}_t + \nu \vec{\mu}_k}{\sum_t \gamma_k(t) + \nu}$$

Supervector Representation

$$\vec{s}_p = \begin{bmatrix} \Sigma_1^{-1/2} (\vec{\mu}_1^{(p)} - \vec{\mu}_1) \\ \vdots \\ \Sigma_K^{-1/2} (\vec{\mu}_K^{(p)} - \vec{\mu}_K) \end{bmatrix}$$

<u>Within Class Covariance</u> <u>Normalization (WCCN)</u>

• <u>Z-normalize</u> the supervector to reduce the effect of irrelevant variability using a robust regularized covariance matrix:

#### $S=(\gamma S+(1-\gamma)I)$

Z-normalization results in better linear separability

# AED Accuracy with Tandem & Supervectors

								Mark Brite My for							
	ap	cl	cm	со	ds	kj	kn	kt	la	pr	pw	st	Average		
MFCC	78.3	26.9	29.5	24.2	56.3	39.9	7.7	0.0	39.0	35.2	14.1	28.7	28.2		
FB	34.5	21.8	25.4	24.9	38.9	27.2	11.7	0.0	49.1	13.8	11.7	28.1	27.8		
Adaboost	44.4	25.5	31.3	31.2	57.3	33.2	13.5	1.9	51.3	36.7	17.6	36.8	34.0		
Adaboost+T	52.6	21.9	37.2	51.3	63.0	29.6	11.5	0.0	54.2	42.7	25.8	34.6	35.3		
Adaboost+S	44.4	25.0	33.7	31.2	56.6	33.2	20.9	35.5	51.3	36.7	19.2	41.3	37.5		
Adaboost+T+S	52.6	21.9	39.6	49.8	63.0	29.6	15.4	36.8	54.7	41.7	26.0	38.3	38.6		

Effectiveness of system components Metric: AED-ACC (%) = weighted multi-class F score

- 'T' = Tandem ANN/HMM recognizer
- 'S' = GMM-Supervector rescoring
- 'Adaboost' = minimum Bayes risk feature selection
- Adaboost+T+S was #1 overall in CLEAR 2007 competition

# **Summary & Conclusions**

- Audio events occur at scales ranging from 2ms to 20 minutes
  - Efficient computation: multiscale FFT
  - Classifier training: Minimum Bayes risk selection
- Audio events occur at negative SNR
  - Noise compensation: Tandem ANN/HMM
  - Noise compensation: Supervector WCCN
- Visualization is most effective with...
  - Signal features, match the salience of audio & image
  - Log likelihood features; task matching useful but not necessary