

## Audio, Video, and their Joint Processing with Applications to Teleconferencing

Cha Zhang and Zhengyou Zhang  
Communication and Collaboration Systems Group  
Microsoft Research  
7/11/2011



### Telepresence



Bell Labs Video Conferencing System (1967)  
HP Halo (2005)  
Polycom RPX (2006)  
Cisco Telepresence (2006)

And Tandberg, Telanetix, Teliris, ...  
2  
Microsoft Research

### Telepresence System Objectives

(Wainhouse Research 2007)

- High quality audio and video
- Simplicity
- High reliability
- Environmental excellence



## Telepresence System Objectives

(Wainhouse Research 2007)

- High quality audio and video
  - Audio: Clear, low noise, low latency, echo-free, spatial audio
  - Video: Clear, life size, low noise, low latency, eye contact, 3D
- Simplicity
- High reliability
- Environmental excellence

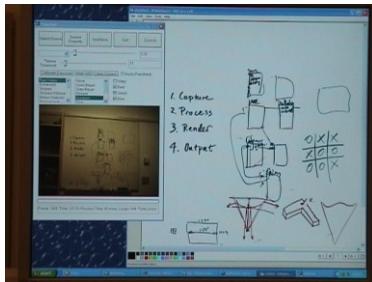
Microsoft<sup>4</sup>  
Research

## Related Systems at MSR

- Interactive whiteboard
- RingCam/RoundTable
- Interactive mobile teleconferencing
- Accelerated instant replay
- Personal telepresence station
- Viewport
- Embedded social proxies

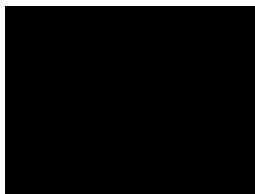
Microsoft<sup>5</sup>  
Research

## Interactive Whiteboard



Microsoft<sup>6</sup>  
Research

## RingCam and RoundTable



Microsoft<sup>7</sup>  
Research

---

---

---

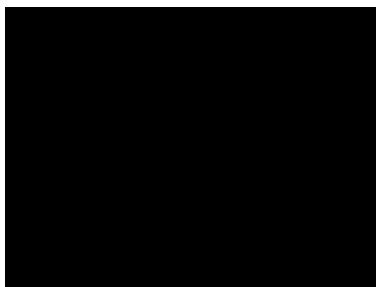
---

---

---

---

## Interactive Mobile Teleconferencing



Microsoft<sup>8</sup>  
Research

---

---

---

---

---

---

---

## Accelerated Instant Replay



Microsoft<sup>9</sup>  
Research

---

---

---

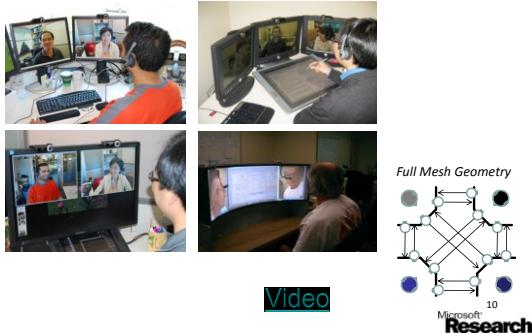
---

---

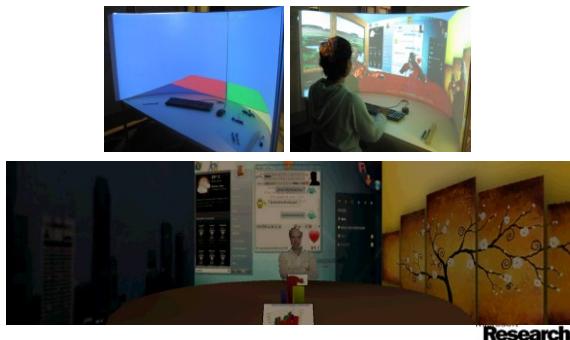
---

---

## Personal Telepresence Station



## Viewport



## Embodied Social Proxies



## Outline

- Introduction
- Sound source localization with compact mic arrays
- Spatial sound and multi-channel echo cancellation
- Video-based tracking
- Other video processing techniques
- Conclusions and future work

<sup>13</sup>  
Microsoft  
Research

## Outline

- Introduction
- Sound source localization with compact mic arrays
  - A maximum likelihood solution
  - Room modeling for enhancing range sensitivity
  - Joint audio/visual speaker localization
- Spatial sound and multi-channel echo cancellation
- Video-based tracking
- Other video processing techniques
- Conclusions and future work

<sup>14</sup>  
Microsoft  
Research

## Microsoft RoundTable



- Panoramic view from 5 high-resolution cameras
- Circular 6-microphone array for speech capturing and sound source localization
- Automated speaker switching

<sup>15</sup>  
Microsoft  
Research

## SSL Based on Time Difference of Arrival (TDOA)

- Cross-correlation

$$R_{ik}(\tau) = \int x_i(t)x_k(t-\tau)dt,$$

when  $\tau = \frac{d}{c}$ , one gets the max value. In freq. domain:

$$R_{ik}(\tau) = \int X_i(\omega)X_k^*(\omega)e^{j\omega\tau} d\omega,$$



In case of multiple mics, solve the following using hypothesis testing:

$$\begin{aligned} \mathcal{R}(s) &= \sum_{i=1}^P \sum_{k=1}^P \int X_i(\omega)X_k^*(\omega)e^{j\omega(\tau_i-\tau_k)} d\omega, \\ &= \int \left[ \sum_{i=1}^P X_i(\omega)e^{j\omega\tau_i} \right] \left[ \sum_{k=1}^P X_k(\omega)e^{j\omega\tau_k} \right]^* d\omega, \\ &= \int \left| \sum_{i=1}^P X_i(\omega)e^{j\omega\tau_i} \right|^2 d\omega, \end{aligned}$$

Microsoft Research

## SSL Based on TDOA (cont.)

- Generalized cross-correlation

$$\mathcal{R}(s) = \sum_{i=1}^P \sum_{k=1}^P \int \Psi_{ik}(\omega)X_i(\omega)X_k^*(\omega)e^{j\omega(\tau_i-\tau_k)} d\omega,$$

- SRP-PHAT

$$\Psi_{ik}(\omega) = \frac{1}{|X_i(\omega)X_k^*(\omega)|} = \frac{1}{|X_i(\omega)||X_k(\omega)|}$$

$$\mathcal{R}(s) = \int \left| \sum_{i=1}^P \frac{X_i(\omega)e^{j\omega\tau_i}}{|X_i(\omega)|} \right|^2 d\omega,$$

- ML based SRP (weight optimal for 2-mic only)

$$\Psi_{ij}(\omega) = \frac{|X_i(\omega)||X_j(\omega)|}{|N_i(\omega)|^2|X_j(\omega)|^2 + |N_j(\omega)|^2|X_i(\omega)|^2},$$

17  
Microsoft Research

## Motivation

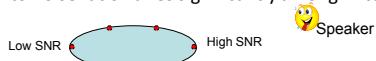
- RoundTable uses *directional* microphones

- All mics point outwards

- For sound capture, experimental results show:

Directional mics + mic selection > omni-directional mics + beamforming

- Signal to noise ratio varies significantly among mics



- Early work uses three closest mics for SSL, but performance is not great

18  
Microsoft Research

## Signal Model

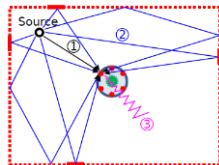
- **Model**  $\mathbf{X}(\omega) = \underbrace{S(\omega)\mathbf{G}(\omega)}_{\text{Received signal}} + \underbrace{S(\omega)\mathbf{H}(\omega)}_{\text{Environment response}} + \underbrace{\mathbf{N}(\omega)}_{\text{Additive noise}},$

Received signal:  $\mathbf{X}(\omega) = [X_1(\omega), \dots, X_P(\omega)]^T,$

Direct propagation gain:  $\mathbf{G}(\omega) = [g_1(\omega)e^{-j\omega\tau_1}, \dots, g_P(\omega)e^{-j\omega\tau_P}]^T,$

Environment response:  $\mathbf{H}(\omega) = [H_1(\omega), \dots, H_P(\omega)]^T,$

Additive noise:  $\mathbf{N}(\omega) = [N_1(\omega), \dots, N_P(\omega)]^T.$



19  
Microsoft Research

## Noise Modeling

- Define:

$$\mathbf{N}^c(\omega) = S(\omega)\mathbf{H}(\omega) + \mathbf{N}(\omega),$$

and assume it follows a zero-mean, independent between frequencies, joint Gaussian distribution

$$p(\mathbf{N}^c(\omega)) = \rho \exp \left\{ -\frac{1}{2} [\mathbf{N}^c(\omega)]^H \mathbf{Q}^{-1}(\omega) \mathbf{N}^c(\omega) \right\},$$

$$\begin{aligned} \mathbf{Q}(\omega) &= E\{\mathbf{N}^c(\omega)[\mathbf{N}^c(\omega)]^H\} \\ &= E\{\mathbf{N}(\omega)\mathbf{N}^H(\omega)\} + |S(\omega)|^2 E\{\mathbf{H}(\omega)\mathbf{H}^H(\omega)\} \end{aligned}$$

20  
Microsoft Research

## Estimate Noise Covariance Matrix

$$\begin{aligned} \mathbf{Q}(\omega) &= E\{\mathbf{N}^c(\omega)[\mathbf{N}^c(\omega)]^H\} \\ &= E\{\mathbf{N}(\omega)\mathbf{N}^H(\omega)\} + |S(\omega)|^2 E\{\mathbf{H}(\omega)\mathbf{H}^H(\omega)\} \end{aligned}$$

- First term:

$$E(N_i(\omega)N_j^*(\omega)) = \lim_{K \rightarrow \infty} \frac{1}{K} \sum_{k=1}^K N_{ik}(\omega)N_{jk}^*(\omega)$$

- Second term (Wang 97, Rui 04):

$$|S(\omega)|^2 E\{\mathbf{H}(\omega)\mathbf{H}^H(\omega)\} \approx \text{diag}(\lambda_1, \dots, \lambda_P)$$

$$\begin{aligned} \lambda_i &= E\{|H_i(\omega)|^2 |S(\omega)|^2\} \\ &\approx \gamma(|X_i(\omega)|^2 - E\{|N_i(\omega)|^2\}) \end{aligned}$$

$0 < \gamma < 1$  is an empirical parameter

21  
Microsoft Research

## Formulation and Solution

- Maximize the likelihood of the received signals

$$p(\mathbf{X}|S, \mathbf{G}, \mathbf{Q}) = \prod_{\omega} p(\mathbf{X}(\omega)|S(\omega), \mathbf{G}(\omega), \mathbf{Q}(\omega)),$$

- Results:

$$S(\omega) = \frac{\mathbf{G}^H(\omega)\mathbf{Q}^{-1}(\omega)\mathbf{G}(\omega)}{\mathbf{G}^H(\omega)\mathbf{Q}^{-1}(\omega)\mathbf{G}(\omega)}$$

and the ML based SSL algorithm shall maximize:

$$\int_{\omega} \frac{[\mathbf{G}^H(\omega)\mathbf{Q}^{-1}(\omega)\mathbf{X}(\omega)]^H \mathbf{G}^H(\omega)\mathbf{Q}^{-1}(\omega)\mathbf{X}(\omega)}{\mathbf{G}^H(\omega)\mathbf{Q}^{-1}(\omega)\mathbf{G}(\omega)} d\omega$$

through hypothesis testing.

22  
Microsoft Research

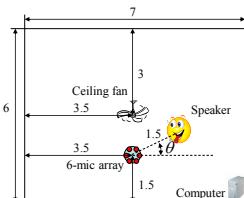
## Some High-Level Messages

- It's a very general ML framework for multi-microphone SSL
- It does not require known microphone gain pattern (gain is estimated by the algorithm)
- It can weight microphones differently based on their noise variances
- It can handle certain amount of reverberation
- SRP-PHAT is a special case of the framework when SNR is extremely high
- It is equivalent to using a minimum variance distortionless response (MVDR) beamformer for SSL
- For detailed derivations, see Zhang 07.

23  
Microsoft Research

## A Synthetic Environment

- An omni-directional mic array in a virtual room
  - Two simulated additive noise sources (ceiling fan and computer) at adjustable SNR
  - Adjustable reverberation level
  - Speaker is placed at different locations ( $\theta=0, 36^\circ, 72^\circ, \dots$ ).
- Average SSL error rate of all locations is reported.



24  
Microsoft Research

## Results

Input SNR	SRP-PHAT		Proposed ML SSL	
	< 2°	< 10°	< 2°	< 10°
Reverberation = 100 ms	25 dB	<b>93.28%</b>	<b>94.41%</b>	92.36% 93.82%
	20 dB	<b>83.35%</b>	84.92%	83.19% 84.96%
	15 dB	79.24%	81.15%	<b>80.30%</b> 82.35%
	10 dB	73.87%	76.71%	<b>76.29%</b> 78.75%
	5 dB	66.59%	70.87%	<b>71.94%</b> 75.30%
	0 dB	54.02%	60.78%	<b>64.19%</b> 69.57%

Input SNR	SRP-PHAT		Proposed ML SSL	
	< 2°	< 10°	< 2°	< 10°
Reverberation = 500 ms	25 dB	<b>61.53%</b>	<b>75.67%</b>	60.76% 74.57%
	20 dB	<b>60.14%</b>	<b>73.88%</b>	59.79% 73.23%
	15 dB	54.76%	<b>67.14%</b>	<b>54.78%</b> 66.97%
	10 dB	52.45%	64.66%	<b>53.40%</b> 65.22%
	5 dB	49.53%	61.96%	<b>50.80%</b> 62.86%
	0 dB	45.11%	57.91%	<b>47.47%</b> 60.14%

## Real-World Data

- 99 meetings captured in more than 50 rooms with various number of attendees
- 4 minutes each
- Speaker direction manually labeled (6706 frames in total)
- Input SNR ranges from 5 dB to 25 dB
- Report results on the percentage of frames that are within 6° (~20 pixels) and 14° (~40 pixels)

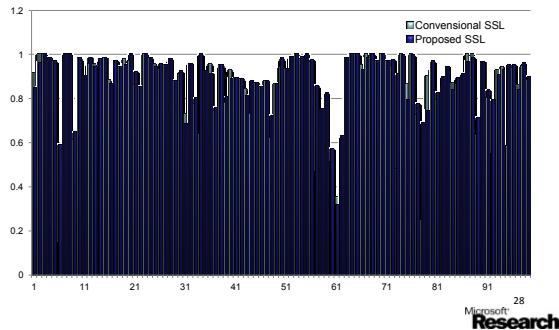
26  
Microsoft  
Research

## Results

Conventional Approach (closest 3 mics)		SRP-PHAT (6 mics)		Proposed method (6 mics)	
< 6°	< 14°	< 6°	< 14°	< 6°	< 14°
80.55%	86.85%	81.73%	88.13%	<b>83.49%</b>	<b>90.13%</b>
15.1% reduction in error		24.9% reduction in error			

27  
Microsoft  
Research

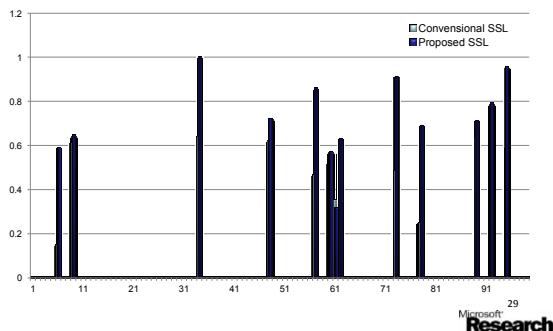
## Results (cont.)



Microsoft Research

28

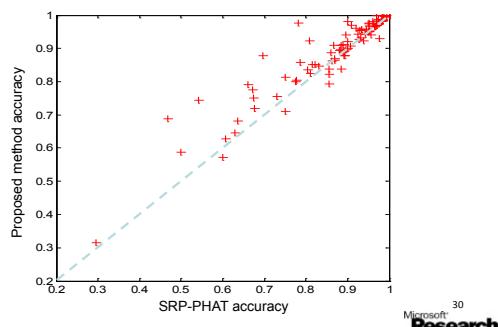
## Results (cont.)



Microsoft Research

29

## Results (cont.)



Microsoft Research

30

## References

- M. Brandstein and H. Silverman, "A practical methodology for speech localization with microphone arrays," *Computer, Speech, and Language*, vol. 11, no. 2, pp. 91–126, 1997.
- T. Gustafsson, B. Rao, and M. Trivedi, "Source localization in reverberant environments: performance bounds and ml estimation," in *Proc. of ICASSP*, 2001.
- K. Harmanci, J. Tabrikian, and J. L. Krolik, "Relationships between adaptive minimum variance beamforming and optimal source localization," *IEEE Trans. on Signal Processing*, vol. 48, no. 1, pp. 1–12, 2000.
- C. Knapp and G. Carter, "The generalized correlation method for estimation of time delay," *IEEE Trans. on Acoustics, Speech, and Signal Processing*, vol. ASSP-24, no. 4, pp. 320–327, 1976.
- Y. Rui and D. Florencio, "Time delay estimation in the presence of correlated noise and reverberation," in *Proc. of ICASSP*, 2004.
- Y. Rui, D. Florencio, W. Lam, and J. Su, "Sound source localization for circular arrays of directional microphones," in *Proc. of ICASSP*, 2005.
- H. Wang and P. Chu, "Voice source localization for automatic camera pointing system in videoconferencing," in *Proc. of IEEE ICASSP*, 1997.
- C. Zhang, D. Florencio, D. Ba, and Z. Zhang, "Maximum likelihood sound source localization and beamforming for directional microphone arrays in distributed meetings," *IEEE Trans. on Multimedia*, pp. 538–548, Vol. 10, No. 3, Apr. 2008.
- C. Zhang, Z. Zhang, and D. Florencio, "Maximum likelihood sound source localization for multiple directional microphones," in *Proc. ICASSP*, 2007.

<sup>31</sup>  
Microsoft  
Research

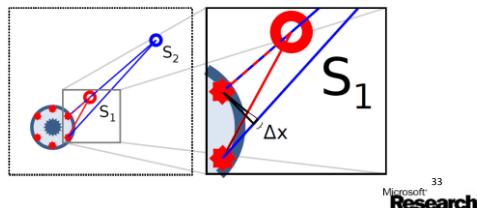
## Outline

- Introduction
- Sound source localization with compact mic arrays
  - A maximum likelihood solution
  - Room modeling for enhancing range sensitivity
  - Joint audio/visual speaker localization
- Spatial sound and multi-channel echo cancellation
- Video-based tracking
- Other video processing techniques
- Conclusions and future work

<sup>32</sup>  
Microsoft  
Research

## Motivation

- Compact mic array has poor resolution in range and elevation for SSL
- Reverberation degrades SSL performance



<sup>33</sup>  
Microsoft  
Research

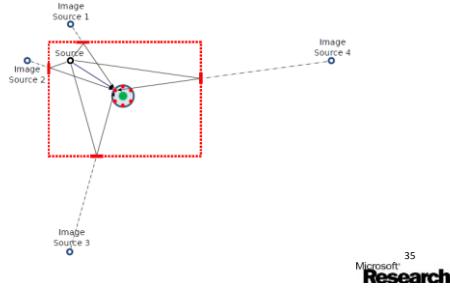
## Modeling the Reverberation

- Energy approximation
  - Wang 97, Rui 04, Zhang 07
- Measure/estimate full room impulse response (Weinstein 94)
  - Very challenging
  - RIR change rapidly and significantly with the position and orientation of the source

34  
Microsoft Research

## Our Approach: Room Modeling

- The image model (Allen 79)



## A New Signal Model

$$X_i(\omega) = \underbrace{\alpha_i(\omega) e^{-j\omega\tau_i} S(\omega)}_{①} + \underbrace{H_i(\omega) S(\omega) + N_i(\omega)}_{②}$$

$$X_i(\omega) = \underbrace{\sum_{r=0}^R \alpha_i^{(r)}(\omega) e^{-j\omega\tau_i^{(r)}} S(\omega)}_{①} + \underbrace{H_i(\omega) S(\omega) + N_i(\omega)}_{②}$$

36  
Microsoft Research

## Extend the ML Solution

- Consider:  $X_i(\omega) = \underbrace{\sum_{r=0}^R \alpha_i^{(r)}(\omega) e^{-j\omega\tau_i^{(r)}}}_{\text{①}} S(\omega) + \underbrace{H_i(\omega) S(\omega)}_{\text{②}} + \underbrace{N_i(\omega)}_{\text{③}}$

- Let:  $G_i(\omega) = \sum_{r=0}^R \alpha_i^{(r)}(\omega) e^{-j\omega\tau_i^{(r)}} = g_i(\omega) e^{-j\varphi_i(\omega)}$

$$g_i(\omega) = \left| \sum_{r=0}^R \alpha_i^{(r)}(\omega) e^{-j\omega\tau_i^{(r)}} \right| e^{-j\varphi_i(\omega)} = \frac{\sum_{r=0}^R \alpha_i^{(r)}(\omega) e^{-j\omega\tau_i^{(r)}}}{\left| \sum_{r=0}^R \alpha_i^{(r)}(\omega) e^{-j\omega\tau_i^{(r)}} \right|}$$

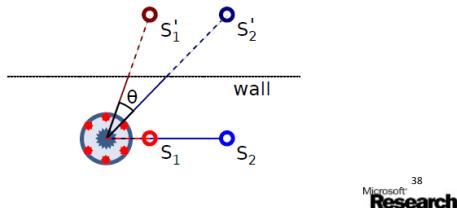
- Approximate:  $e^{-j\varphi_i(\omega)} \approx \frac{\sum_{r=0}^R \frac{\rho_i^{(r)}}{r_i^{(r)}} e^{-j\omega\tau_i^{(r)}}}{\left| \sum_{r=0}^R \frac{\rho_i^{(r)}}{r_i^{(r)}} e^{-j\omega\tau_i^{(r)}} \right|}$  known

- Same solution as before

37  
Microsoft Research

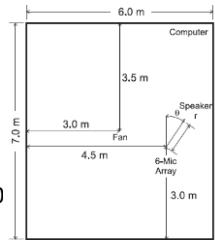
## Range Discrimination

- The reflected sound source can help range discrimination
- Reverberation is now our friend!



## Exp. Results – Synthetic Data

- Synthetic room
- Speaker at 1.3m to the array, tested at 0° elevation, and azimuth  $\theta = 0^\circ, 36^\circ, 72^\circ, \dots$
- Reverberation  $T_{60}$  at 250 ms and 500 ms



39  
Microsoft Research

## Exp. Results – Synthetic Data

Closest Mic. SNR	VAD Voice Frames	ML-SSL			R-ML-SSL		
		$\Delta\theta > 2^\circ$	$\Delta\phi > 1^\circ$	$\Delta r > .15 m$	$\Delta\theta > 2^\circ$	$\Delta\phi > 1^\circ$	$\Delta r > .15 m$
25 dB	48	0.4 %	92.0 %	84.6 %	0.4 %	0.2 %	0.2 %
20 dB	47	0.2 %	90.6 %	81.6 %	0.2 %	0.0 %	0.0 %
15 dB	44	1.1 %	88.2 %	80.9 %	1.1 %	0.9 %	0.4 %
10 dB	38	3.4 %	85.6 %	81.6 %	3.4 %	2.9 %	2.1 %
5 dB	29	15.0 %	86.0 %	90.1 %	15.0 %	14.3 %	6.9 %
0 dB	16	24.8 %	85.1 %	94.3 %	24.8 %	22.2 %	15.5 %

Table I

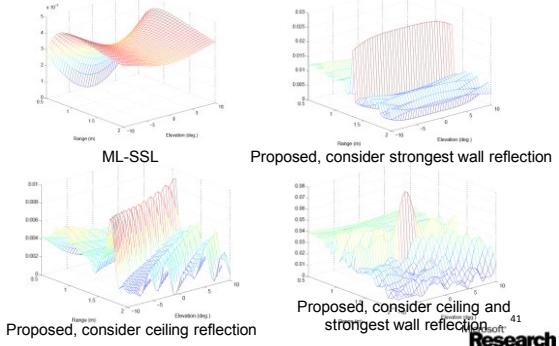
ERROR RATES FOR SYNTHETIC DATA,  $T_{60} = 250 \text{ ms}$ .

Closest Mic. SNR	VAD Voice Frames	ML-SSL			R-ML-SSL		
		$\Delta\theta > 2^\circ$	$\Delta\phi > 1^\circ$	$\Delta r > .15 m$	$\Delta\theta > 2^\circ$	$\Delta\phi > 1^\circ$	$\Delta r > .15 m$
25 dB	47	1.3 %	94.5 %	87.5 %	1.3 %	1.9 %	1.2 %
20 dB	45	1.1 %	93.2 %	86.3 %	1.1 %	1.7 %	1.5 %
15 dB	41	1.4 %	93.8 %	86.8 %	1.4 %	1.1 %	0.9 %
10 dB	38	2.6 %	94.1 %	88.5 %	2.6 %	2.6 %	1.6 %
5 dB	33	6.5 %	92.3 %	89.0 %	6.5 %	5.4 %	2.2 %
0 dB	23	21.4 %	91.0 %	93.0 %	21.4 %	19.3 %	12.4 %

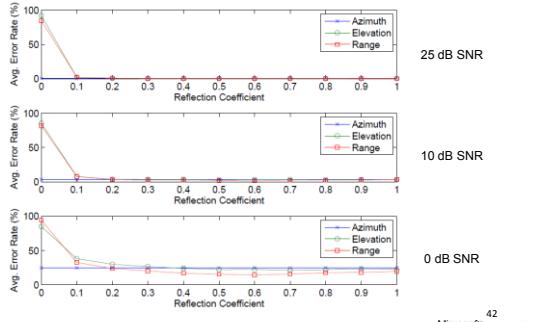
Table II

ERROR RATES FOR SYNTHETIC DATA,  $T_{60} = 500 \text{ ms}$ .Microsoft  
Research 40

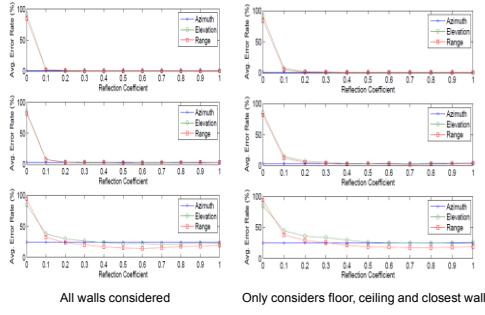
## Exp. Results – Synthetic Data

Microsoft  
Research 41

## Sensitivity to Wall Reflection Coef.

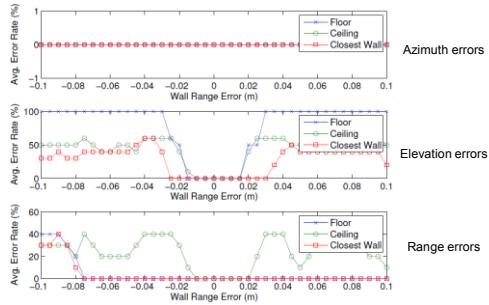
Microsoft  
Research 42

## Sensitivity to Number of Walls



Microsoft  
Research 43

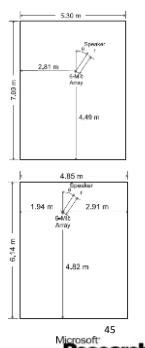
## Sensitivity to Wall Distance Error



Microsoft  
Research 44

## Exp. Results – Real Room Data

- Room 1
  - Distance to wall estimated using RoundTable device
  - 3 Closest reflectors used
- Room 2
  - Distance to wall measured by an ultrasonic range finder (1cm res.)
  - Reflection coef. set to 0.3



Microsoft  
Research 45

## Exp. Results – Real Room Data

Speaker Position	# of Frames	ML-SSL		R-ML-SSL			
		$\Delta\theta > 10^\circ$	$\Delta\phi > 5^\circ$	$\Delta r > 0.15$	$\Delta\theta > 10^\circ$	$\Delta\phi > 5^\circ$	$\Delta r > 0.15$
$\theta \approx 0^\circ, \phi \approx 2.5^\circ, r \approx 1.70 m$	11	0%	0%	55%	0%	0%	0%
$\theta \approx 60^\circ, \phi \approx 23.5^\circ, r \approx 1.47 m$	21	0%	5%	95%	0%	9%	0%
$\theta \approx 120^\circ, \phi \approx 2.5^\circ, r \approx 1.37 m$	9	0%	0%	33%	0%	11%	11%

Room 1

Speaker Position	# of Frames	ML-SSL		R-ML-SSL			
		$\Delta\theta > 10^\circ$	$\Delta\phi > 5^\circ$	$\Delta r > 0.15$	$\Delta\theta > 10^\circ$	$\Delta\phi > 5^\circ$	$\Delta r > 0.15$
$\theta \approx 60^\circ, \phi \approx 26^\circ, r \approx 1.85 m$	7	0%	0%	96%	0%	5%	1%
$\theta \approx 120^\circ, \phi \approx 240^\circ, r \approx 2.40 m$	71	0%	0%	87%	0%	3%	14%
$\theta \approx 120^\circ, \phi \approx 322^\circ, r \approx 1.60 m$	72	3%	3%	96%	3%	3%	3%
$\theta \approx 120^\circ, \phi \approx 222^\circ, r \approx 2.15 m$	68	0%	3%	90%	0%	2%	3%
$\theta \approx 180^\circ, \phi \approx 15^\circ, r \approx 3.10 m$	82	1%	10%	99%	1%	2%	27%
$\theta \approx 210^\circ, \phi \approx 12^\circ, r \approx 4.30 m$	76	0%	24%	90%	0%	1%	17%
$\theta \approx 210^\circ, \phi \approx 117^\circ, r \approx 1.45 m$	75	12%	12%	96%	12%	6%	33%
$\theta \approx 300^\circ, \phi \approx 40^\circ, r \approx 1.40 m$	84	7%	7%	98%	7%	6%	34%

Room 2

46  
Microsoft  
Research

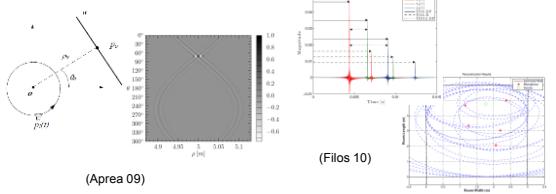
## How to Estimate Room Model

- Visual methods
  - Range sensors
  - Landmark based (Kimber 09)
- Acoustic methods
  - Acoustic imaging (Moebus 07, O'Donovan 08)
  - Image model (Allen 79)
    - Single mic (Aprea 09, Antonacci 10, Dokmanic 11)
    - Mic arrays (Filos 10, Canclini 11)

47  
Microsoft  
Research

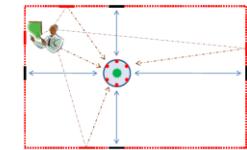
## Image Model based Methods

- Two steps:
  - Measure room impulse responses (RIRs)
  - Infer room model



48  
Microsoft  
Research

## Representation



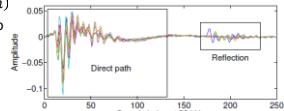
$$R = \{(a_i, \underbrace{d_i, \theta_i, \phi_i}_{\substack{\text{3D position of wall } i \\ \text{w.r.t. array center}}})\}_{i=1}^6$$

↗  
Reflection coef.  
of wall  $i$ .

49  
Microsoft Research

## RIR Modeling

- Signal observed at mic  $m$ :
    - $y_m(n) = h_m(n) * s(n) + u_m(n)$
    - $h_m(n)$ : RIR from array center to mic  $m$
  - Only consider early reflections



$h_m^{(dp)}(n)$ : direct path IR from array center to mic  $m$

$h_m^{(r_i, \theta_i, \phi_i)}(n)$ : IR at mic  $m$  from perfect wall at  $(r_i, \theta_i, \phi_i)$  (WIR)

*R* : number of modeled walls

Then

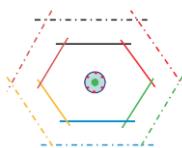
$$h_m(n) \approx h_m^{(dp)}(n) + \sum_{i=1}^R \rho^{(i)} h_m^{(\tau_i, \theta_i, \phi_i)}(n) + v_m(n)$$

50  
Microsoft Research

## Infer the Wall Locations

- In vector form:
$$\mathbf{h} \approx \mathbf{h}^{(dp)}(n) + \sum_{i=1}^R \rho^{(i)} \mathbf{h}^{(r_i, \theta_i, \phi_i)} + \mathbf{v}$$

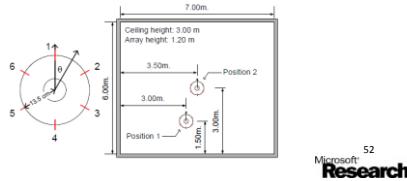
  - Solve:
$$\min_{\mathbf{a}} \|\mathbf{h}_{room} - \mathbf{H}\mathbf{a}\|_2^2 + \lambda \|\mathbf{a}\|_1$$
  - Note we obtain a sparse solution for  $\mathbf{a}$  due to  $L_1$  optimization.
  - Solved using LASSO (Tibshirani 96).



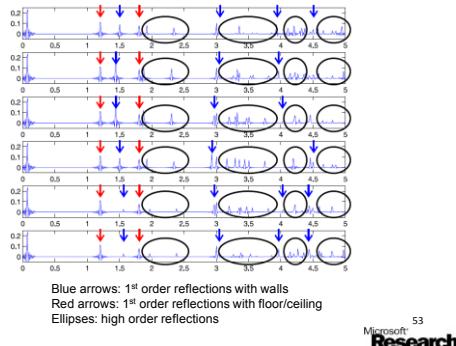
51  
Microsoft Research

## Exp. Results – Synthetic Data

- 6x7x3m room, 300ms reverberation
- Reflection coef. of 0.77
- SNR ~20dB
- Position 2 is more challenging



## Exp. Results – Synthetic Data



## Exp. Results – Synthetic Data

Ground Truth			Estimates		
$r$ (m)	$\theta$ ( $^{\circ}$ )	$\phi$ ( $^{\circ}$ )	$r$ (m)	$\theta$ ( $^{\circ}$ )	$\phi$ ( $^{\circ}$ )
1.200	0	-90	1.200	0	90
1.800	0	90	1.800	0	90
4.500	0	0	4.502	0	0
4.000	90	0	4.000	90	0
1.500	180	0	1.500	180	0
3.000	270	0	3.005	270	0

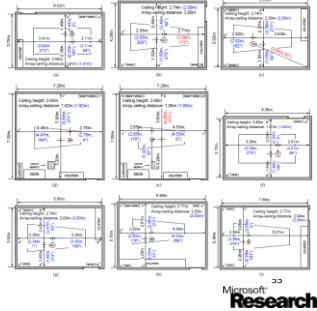
Position 1

Ground Truth			Estimates		
$r$ (m)	$\theta$ ( $^{\circ}$ )	$\phi$ ( $^{\circ}$ )	$r$ (m)	$\theta$ ( $^{\circ}$ )	$\phi$ ( $^{\circ}$ )
1.200	0.0	-90	1.200	0.0	90
1.800	0.0	90	1.800	0.0	90
3.000	0.0	0	3.001	0.0	0
3.500	90.0	0	3.501	90.0	0
3.000	180.0	0	3.001	180.0	0
3.500	270.0	0	3.501	270.0	0

Position 2

## Exp. Results – Real Room Data

- 9 conference rooms
- RoundTable device
  - WIR measured in anechoic chamber
  - Device on the table, no floor



## Exp. Results – Real Room Data

Room	Ground Truth			Estimates			Refs
	$r$ (m)	$\theta$ ( $^{\circ}$ )	$\phi$ ( $^{\circ}$ )	$r$ (m)	$\theta$ ( $^{\circ}$ )	$\phi$ ( $^{\circ}$ )	
A	1.30	0	90	1.81	0	90	—
	1.40	0	0	1.41	3	90	111
	2.11	90	0	2.11	94	0	111
	2.31	90	0	2.31	114	0	100
B	2.91	270	0	3.92	273	0	100
	2.00	0	90	2.00	0	90	—
	2.40	0	0	2.93	358	0	101
	2.71	180	0	2.48	180	0	100
C	1.83	270	0	1.83	270	0	101
	2.00	0	90	2.00	0	90	—
	3.63	0	0	3.94	30	0	100
	1.55	90	0	1.55	90	0	101
D	1.80	0	0	1.80	0	90	102
	2.37	270	0	2.38	271	0	101
	1.82	0	90	1.82	0	90	—
	2.75	0	0	2.75	4	0	101
E	1.40	90	0	1.40	90	0	100
	2.40	90	0	2.40	94	0	100
	2.75	0	0	2.75	0	0	100
	3.89	270	0	3.88	271	0	101
F	1.81	0	90	1.82	0	90	—
	2.23	0	0	2.24	0	0	101
	2.01	90	0	2.01	89	0	001
	1.80	0	0	1.80	149	0	101
G	3.35	270	0	3.36	273	0	010
	2.00	0	90	2.00	0	90	—
	2.36	90	0	2.34	1	0	111
	3.29	180	0	3.30	188	0	100
H	1.45	270	0	1.46	274	0	111
	2.00	0	90	2.02	0	90	—
	2.34	90	0	2.33	89	0	010
	2.28	90	0	2.29	91	0	110
I	2.26	180	0	2.62	178	0	121
	2.70	0	0	2.80	0	0	121
	2.04	0	90	2.04	0	90	—
	3.08	0	0	3.09	3	0	101
	5.27	90	0	5.27	90	0	—
	3.10	0	0	2.26	181	0	101
	2.39	270	0	2.38	274	0	102

Most walls are identified correctly, with a typical resolution of 1 cm.  
Some distant walls are not correctly detected, though they won't impact SSL performance.

Microsoft Research

## References

- J. Allen and D. Berkley, "Image method for efficiently simulating small room acoustics," *J. Acoust. Soc. Am.*, vol. 65, no. 4, pp. 943-950, 1979.
- F. Antonacci, A. Sarti, and S. Tubaro, "Geometric reconstruction of the environment from its response to multiple acoustic emissions," in *Proc. of ICASSP*, 2010.
- D. Aprea, F. Antonacci, A. Sarti, and S. Tubaro, "Acoustic reconstruction of the geometry of an environment through acquisition of a controlled emission," in *Proc. of EUSIPCO*, 2009.
- D. Ba, F. Ribeiro, C. Zhang and D. Florencio, "L1 Regularized Room Modeling with Compact Microphone Arrays," in *Proc. of ICASSP*, 2010.
- A. Canclini, P. Annibale, F. Antonacci, A. Sarti, R. Rabenstein, and S. Tubaro, "From direction of arrival estimates to localization of plane reflectors in a two dimensional geometry," in *Proc. of ICASSP*, 2011.
- I. Dokmanic, Y.M. Lu, and M. Vetterli, "Can one hear the shape of a room: the 2-D polygonal case," in *Proc. of ICASSP*, 2011.
- J. Filos, E.A.P. Habets, and P.A. Naylor, "A Two-Step Approach to Blindly Infer Room Geometries," in *Proc. of the International Workshop on Acoustic Echo and Noise Control (IWAENC)*, 2010.
- D. Kimber, C. Chen, E.G. Rieffel, J. Shingu, and J. Vaughan, "Marking up a world: visual markup for creating and manipulating virtual models," in *Proc. of IMERSCOM*, 2009.
- M. Moebus and A.M. Zoubir, "Three-Dimensional Ultrasound Imaging in Air using a 2D Array on a Fixed Platform," in *Proc. of ICASSP*, 2007.
- A. O'Donovan, R. Duraiswami, and D. Zotkin, "Imaging concert hall acoustics using visual and audio cameras," in *Proc. of ICASSP*, 2008.

57  
Microsoft Research

## References

---



---



---



---



---



---



---



---



---



---



---



---

- F. Ribeiro, C. Zhang, D. Florencio, and D. Ba, "Using Reverberation to Improve Range and Elevation Discrimination for Small Array Sound Source Localization," *IEEE Transactions on Audio, Speech and Language Processing*, pp. 1781-1792, Vol. 18, No. 7, Sep. 2010.
- Y. Rui and D. Florencio, "Time delay estimation in the presence of correlated noise and reverberation," in Proc. of ICASSP, 2004.
- R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society, Series B (Methodological)*, vol. 58, no. 1, pp. 267-288, 1996.
- E. Weinstein, A. Oppenheim, M. Feder, and J. Buck, "Iterative and sequential algorithms for multisensor signal enhancement," *IEEE Trans. on Signal Processing*, vol. 42, no. 4, pp. 846-859, 1994.
- H. Wang and P. Chu, "Voice source localization for automatic camera pointing system in videoconferencing," in Proc. of IEEE ICASSP, 1997.
- C. Zhang, D. Florencio, D. Ba and Z. Zhang, "Maximum likelihood sound source localization and beamforming for directional microphone arrays in distributed meetings," *IEEE Trans. on Multimedia*, pp. 538-548, Vol. 10, No. 3, Apr. 2008.

<sup>58</sup>  
Microsoft  
Research

---



---



---



---



---



---



---



---



---



---



---

## Outline

---



---



---



---



---



---



---



---



---



---



---



---

- Introduction
- Sound source localization with compact mic arrays
  - A maximum likelihood solution
  - Room modeling for enhancing range sensitivity
  - Joint audio/visual speaker localization
- Spatial sound and multi-channel echo cancellation
- Video-based tracking
- Other video processing techniques
- Conclusions and future work

<sup>59</sup>  
Microsoft  
Research

---



---



---



---



---



---



---



---



---



---

## Motivation

---



---



---



---



---



---



---



---



---



---



---

- SSL only is still not accurate enough
  - Reverberation
  - Noise
- Can we use the panoramic video for help?



<sup>60</sup>  
Microsoft  
Research

---



---



---



---



---



---



---



---



---



---

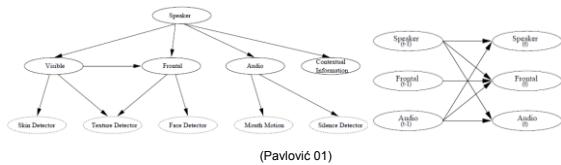
## Previous Work

- SSL only speaker pointing (Polycom iPower 900)
  - Sensitive to reflections and noises
- SSL + face detection decision level fusion (DLF)
  - Very popular (Kapralos 02, Yoshimi 02)
  - RingCam used this approach (Cutler 02)



## Previous Work

- Graphic models for audio/visual fusion (Pavlović 01, Beal 02)
  - Intuitive way to fuse multimodal data
  - Inference could be time-consuming



62  
Microsoft Research

## Previous Work

- Audio visual synchrony (Cutler00, Hershey00, Nock03)
  - Use mutual information between AV signals
  - Search for high correlation between AV
- Speaker tracking (Vermaak 01, Chen 04)
  - Fuse audio/video information using particle filters
  - Still need detection first



Light pixel represents high mutual information between AV (Nock03)

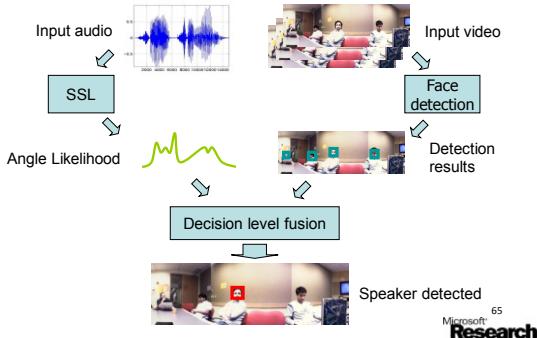
63  
Microsoft Research

# Why Is The Problem Hard?

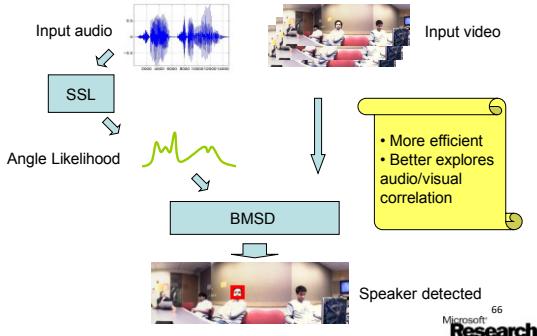
- People don't always look at the camera
  - There can be many meeting attendees
  - Color calibration is hard, hence skin color-based face detection doesn't work
  - Extremely low computation budget (100 MIPS)
  - Image resolution is low, face can be as small as 10x10 pixel



## Decision Level Fusion

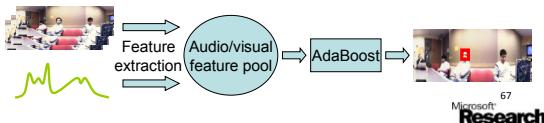


BMSD

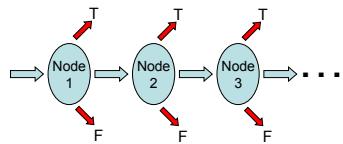


## AdaBoost

- Adaptively focus on the difficult examples
  - Boost the weak learner to a strong classifier by linear combination
- $$H(x) = \sum \alpha_t * h_t(x)$$
- Greedily pick the best feature according to the current sample weight distribution
  - Can be used for automatic feature selection



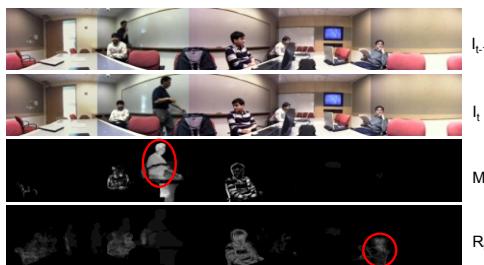
## Attention Cascade



- Each node can contain one or multiple features
- Most negative examples are rejected at early stages
- Extremely efficient

68  
Microsoft Research

## Video Feature Extraction



$$M_t = \min(|I_t - I_{t-1}|, |I_t - I_{t-2}|)$$

$$R_t = \alpha M_t + (1 - \alpha) R_{t-1}$$

69  
Microsoft Research

## Video Feature Extraction (cont.)

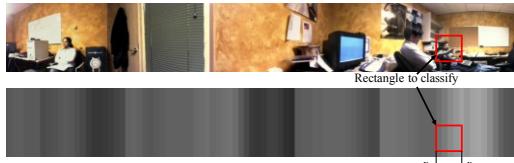


A diagram showing two black rectangular boxes side-by-side, representing the left and right eye regions.

Haar features (Viola 01)  
on all three images.

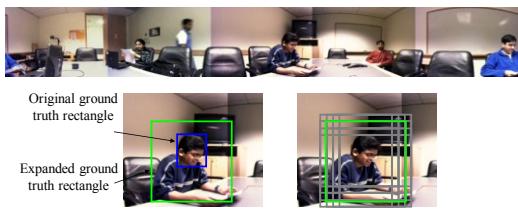
70  
Microsoft Research

## Audio Feature Extraction



71  
Microsoft Research

## Training and Testing



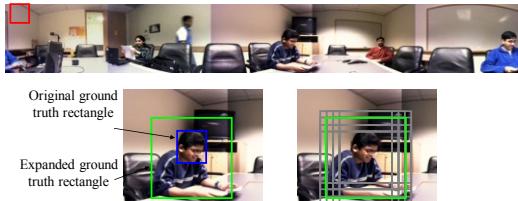
72

## Training and Testing



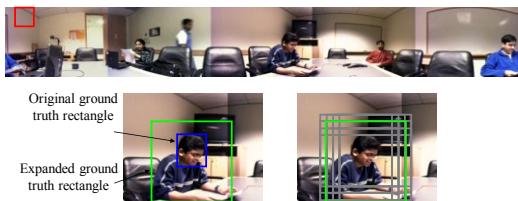
73  
Microsoft Research

## Training and Testing



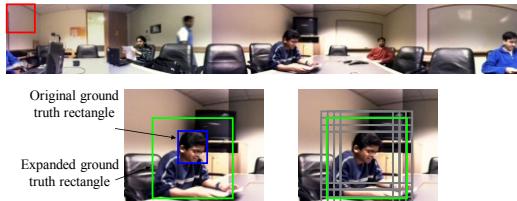
74  
Microsoft Research

## Training and Testing



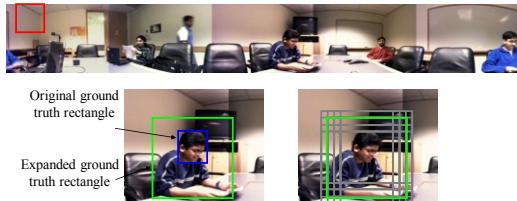
75  
Microsoft Research

## Training and Testing



76  
Microsoft Research

## Training and Testing



77  
Microsoft Research

## Merge of Detected Windows

- Multiple detections could exist for a frame
  - Output a single speaker location
  - Two methods to merge the windows
    - Projection and merge (PAM)
      - Project to azimuth axis, accumulate histogram
      - Find peak of histogram
      - Average nearby windows
    - Top N merge (TNM)
      - Select top N windows with highest confidence
      - Apply PAM above

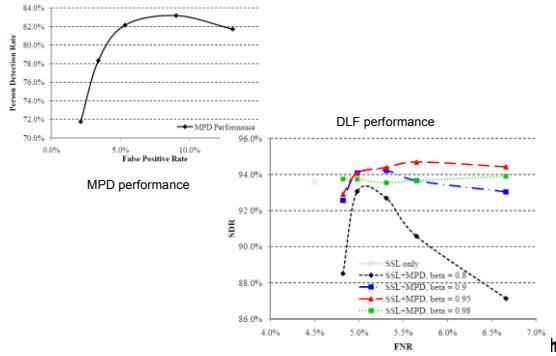
78

## Experiments

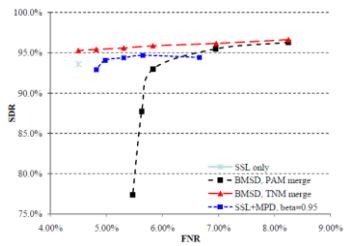
- Training: 93 meetings, 3204 labeled frames in 10 different conference rooms
- Testing: 82 meetings, 1502 labeled frames in 10 meeting rooms (different from training)
- Speaker detection in each frame independently
- Compare SSL-only, DLF (SSL+MPD) and BMSD

79  
Microsoft Research

### Decision Level Fusion

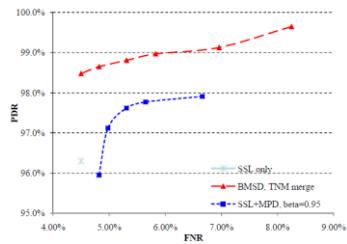


### Speaker Detection Performance



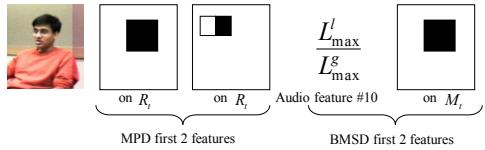
81  
Microsoft Research

## Person Detection Performance



<sup>82</sup>  
Microsoft  
Research

## Features Selected by AdaBoost



- BMSD much faster than DLF
  - Audio features help pruning
  - Audio features are extremely fast to compute

<sup>83</sup>  
Microsoft  
Research

## Examples



<sup>84</sup>  
Microsoft  
Research

## References

- M. Beal, H. Attias, and N. Jojic, "Audio-video sensor fusion with probabilistic graphical models," in Proc. of ECCV, 2002.
  - Y. Chen and R. Yu, "Real-time speaker tracking using particle filter sensor fusion," Proceedings of the IEEE, vol. 92, no. 3, pp. 485–494, 2004.
  - R. Cutler and L. Davis, "Look who's talking: speaker detection using video and audio correlation," in Proc. of IEEE ICME, 2000.
  - R. Cutler, Y. Rui, A. Gupta, J. Cadiz, J. Tashev, L. He, A. Colburn, Z. Zhang, Z. Liu, and S. Silverberg, "Distributed meetings: a meeting capture and broadcasting system," in Proc. ACM Conf. on Multimedia, 2002.
  - J. Hershey and J. Movellan, "Audio vision: using audio-visual synchrony to locate sounds," in Advances in Neural Information Processing Systems, 2000.
  - B. Kapralos, M. Jenkin, and E. Milios, "Audio-video localization of multiple speakers in a video teleconferencing setting," York University, Canada, Tech. Rep., 2002.
  - H. Nock, G. Kyriagis, and C. Neti, "Speaker localization using audio-visual synchrony: an empirical study," in Proc. of CVPR, 2003.
  - V. Pavlovic, A. Garcia, J. Rehg, and T. Huang, "Multimodal speaker detection using error feedback dynamic Bayesian networks," in Proc. of IEEE CVPR, 2001.
  - J. Vermaak, M. Gangnet, A. Black, and P. Perez, "Sequential Monte Carlo fusion of sound and vision for speaker tracking," in Proc. of IEEE ICCV, 2001.
  - B. Yoshimi and G. Pingali, "A multimodal speaker detection and tracking system for teleconferencing," in Proc. ACM Conf. on Multimedia, 2002.
  - C. Zhang, P. Yin, Y. Rui, R. Cutler, P. Viola, X. Sun, N. Pinto and Z. Zhang, "Boosting-Based Multimodal Speaker Detection for Distributed Meeting Videos," IEEE Trans. on Multimedia, pp. 1541-1552, Vol. 10, No. 8, Dec. 2008. 85

Microsoft Research

## Outline

- Introduction
  - Sound source localization with compact mic arrays
  - Spatial sound and multi-channel echo cancellation
    - Head phone based spatial sound
    - Loudspeaker spatial sound
    - Multi-channel echo cancellation
  - Video-based tracking
  - Other video processing techniques
  - Conclusions and future work

86  
Microsoft Research

## Motivation

- Current audio conferencing systems
    - Monaural → adequate for 1-to-1
    - Poor when #people > 2
  - Why poor?
    - All the voice streams are intermixed into a single one
    - Huge cognitive load: *Do 2 things simultaneously*
      - Associate voice signals to the speaker
      - Comprehend what is being discussed

87  
Microsoft Research




---

---

---

---

---

---

---

## Benefits of Spatial Audio

- Human's cocktail party effect
  - Selective attention
  - Only spend effort on comprehension
- Brain rejects incoherent signals at two ears
  - Reverberation & noise are disregarded (not for monol)
- Benefits (Baldis 01)
  - Memory
  - Focal Assurance
  - Perceived Comprehension
  - Listener's Preference

89  
Microsoft  
Research

---

---

---

---

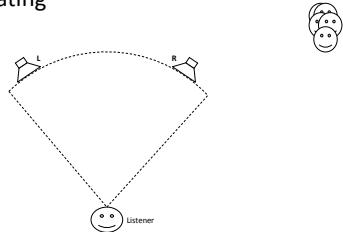
---

---

---

## Multiparty Spatial Audio Conferencing

- Virtual seating



90  
Microsoft  
Research

---

---

---

---

---

---

---

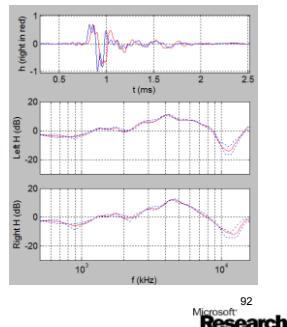
## Headphone based Spatial Sound

- Situations requiring privacy
    - Open offices, public spaces, home office, etc
    - On-the-go: mobile devices
  - Modern hardware
    - high-quality stereo ear-canal headphones
    - Compact Bluetooth headsets (monaural)
  - Controlled environment than loudspeakers
  - Easy job for AEC!

91  
Microsoft Research

## Background

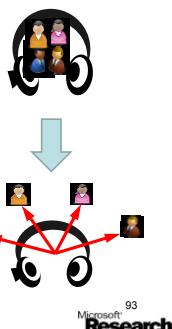
- HRTF/HRIR (Algazi 01)
    - Linear response of each ear to
      - Each azimuth & elevation
      - Captured in an anechoic room
    - Contributing factors
      - Head shadow
      - Pinna/concha/torso
      - Air absorption
    - Gives sense of directivity
      - Individualization is hard



92  
Microsoft Research

## Background

- Headphone “virtualization” (Zlzer 02)
    - Input: regular audio samples
    - Output: simulating a listening environment
      - Localization & externalization
  - Model the signal path (Schroeder 62, Stautner 82)
    - transmitter, medium & listener
    - Room response
      - Most important factor for externalization
      - Classical solution: various forms of artificial “reverb”
        - e.g. Schroeder’s original comb filters



Microsoft  
Research

## Challenges & Solution

- Unique requirements for conferencing
  - True room response is subtle
  - HRTF calibration is difficult
  - Sensitivity to human voice
- Higher realism is the solution
  - Capture “combined” transfer function
  - One or few models suffices
  - Localized parameterization possible
  - More computation but acceptable today

94  
Microsoft Research

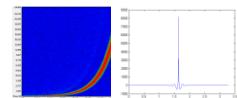
## Offline Measurement

- Overall transfer function

$$Y(\omega) = H_1 H_2 H_3 X(\omega) + N$$

- Probe signal

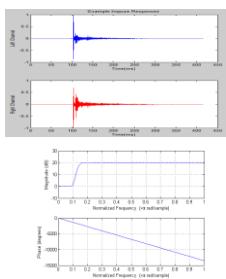
- Exponential chirp
  - $x(t) = \sin k_e e^{k_s t}$



95  
Microsoft Research

## Post-processing

- Sounds bad if use as is
- Alignment & truncation
- Equalization
  - Match the spectrum of the original and processed signal
- Acoustic ratio
  - Between the direct component & the rest
  - Parameterization & personalization



96  
Microsoft Research

## Online Synthesis

- Partitioning the filters in time
  - Short filter  $h^s$ 
    - Location dependent, up to 10ms
    - Source, HRTF & early reflections
  - Long filter  $h^l$ 
    - Location independent, up to 100+ ms
    - Reverberation: omni-directional

- Lower complexity implementation

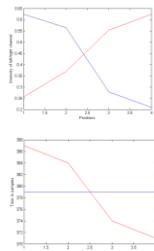
$$y_i(t;i) = x_i(t) * h_i^s(t;i) + h^l(t) * \sum_j x_j(t)$$

$$y_l(t) = \sum_i y_i(t;i)$$

97  
Microsoft Research

## Animation

- Important for added realism
  - Natural head motion while talking
  - With video, help resolve HRTF mismatch
- Classical limitation of non-model based approach
  - Non-parameterized  $h_i^s(t;i) \rightarrow h_i^s(t;\phi)$
  - Too many to measure
- Efficient implementation
  - Extract and interpolate three cues: AR, IID, ITD
  - Locally modify transfer function  $h_i^s(t;\phi) = \Delta\alpha h_i^s(t - \Delta\theta, i)$
  - TBD: smooth, large range animation



98  
Microsoft Research

## Demo

- Please put on headphone



99  
Microsoft Research

## References

---



---



---



---



---



---



---



---



---

- V. R. Algazi, R. O. Duda, D. M. Thompson, and C. Avendano, "The CIPIC HRTF database," in *Proc. IEEE Workshop on Applications of Signal Processing to Audio*, 2001.
- J. Baldiss, "Effects of spatial audio on memory, compression, and preference during desktop conferences," in *Proceedings of the CHI Conference on Human Factors in Computing Systems*, 2001, pp. 166-173.
- W. Chen and Z. Zhang, "Highly realistic audio spatialization for multiparty conferencing using headphones," in *Proc. of MMSP*, 2009.
- M. R. Schroeder, "Natural-sounding artificial reverberation," *J. Audio Eng. Soc.*, vol. 10, no. 3, pp. 219-233, 1962.
- J. Staunton and M. Puckette, "Designing multichannel reverberators," *Computer Music J.*, vol. 6, no. 1, pp. 52-65, 1982.
- U. Zlter, Ed., *DAFX: Digital audio effects*. John Wiley & Sons, LTD, 2002, ch. 3.

<sup>100</sup>  
Microsoft  
Research

---



---



---



---



---



---



---



---



---

## Outline

---



---



---



---



---



---



---



---



---

- Introduction
- Sound source localization with compact mic arrays
- Spatial sound and multi-channel echo cancellation
  - Head phone based spatial sound
  - Loudspeaker spatial sound
  - Multi-channel echo cancellation
- Video-based tracking
- Other video processing techniques
- Conclusions and future work

<sup>101</sup>  
Microsoft  
Research

---



---



---



---



---



---



---



---



---

## Motivation

---



---



---



---



---



---



---



---



---

- People do not like to wear headphones in a remote meeting
- Need a loud speaker solution
  - Challenges
    - User's head moves around
    - Requires multi-channel echo cancellation

<sup>102</sup>  
Microsoft  
Research

---



---



---



---



---



---



---



---



---

## Amplitude Panning

- Well known method (Pullki 97)

*Traditional Spatial*

- Delay and Gain Modulation

- Delay

$$\Delta_R = D - D \cos(\lambda(\Phi - \Theta))$$

$$\Delta_L = D - D \cos(\lambda(\Phi + \Theta))$$

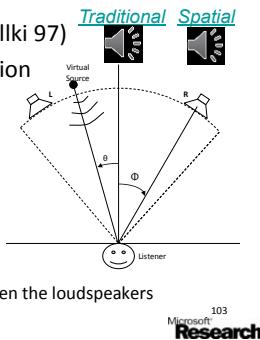
$$D=0.45\text{ms} \quad 1 \leq \lambda \leq \pi/(2\Phi)$$

- Gain

$$G_R = \cos(\lambda(\Phi - \Theta)/2)$$

$$G_L = \cos(\lambda(\Phi + \Theta)/2)$$

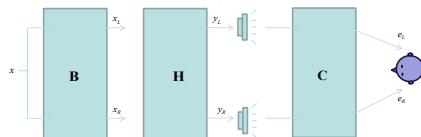
- Only create spatial sound between the loudspeakers



103  
Microsoft Research

## HRTF Based Schemes

- HRTF based method through crosstalk cancellation
  - Binaural sound reproduction over **loudspeaker** with crosstalk canceller



104  
Microsoft Research

- Organization
  - B : Binaural synthesis
  - H : crosstalk canceller : inversion of C
  - C : transmission path or acoustic channel

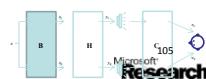
## Binaural Synthesis

- Locating sound images around the listener by filtering the monaural input signal with the HRTF for a given angle of incidence.

$$x = \begin{bmatrix} x_L \\ x_R \end{bmatrix} = \begin{bmatrix} B_L \\ B_R \end{bmatrix} x = Bx$$

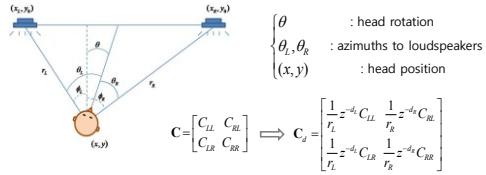
$x$  : input monaural signal  
 $\begin{bmatrix} x_L \text{ and } x_R \\ B_L \text{ and } B_R \end{bmatrix}$  : binaural signals for left and right ear  
 $B$  : HRTF for left and right ear

- Convolution with pre-measured HRTF

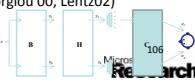


105  
Microsoft Research

## The Dynamic Acoustic Channel

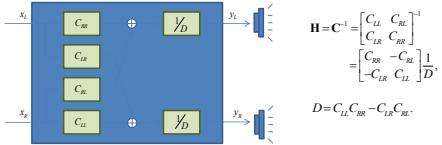


- Adjusting azimuth for HRTFs according to tracking information (Gardner 97)
- Involving time delay and level attenuation
- Tracking using electromagnetic trackers (Georgiou 00, Lentz02)
- Tracking using webcams (Kyriakakis 98, López 99, Lentz 06, Kim 08)



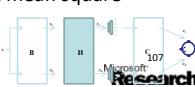
## Dynamic Crosstalk Cancellation

- Inverse the acoustic transfer matrix  $\mathbf{C}$

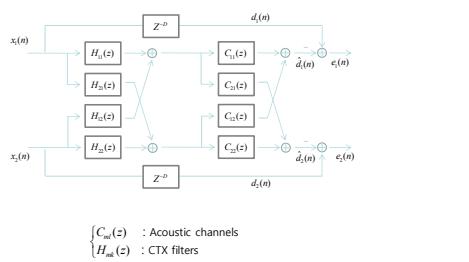


- Direct inverse is unstable

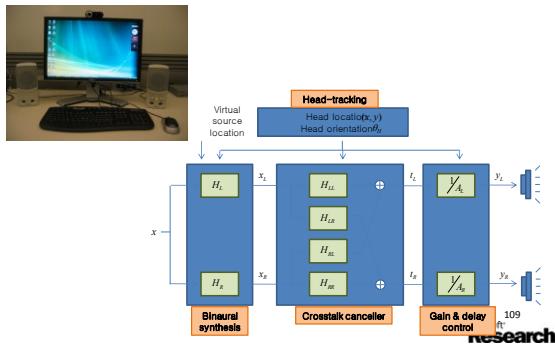
- $\mathbf{H}$  is adaptively obtained by least mean square (LMS) method



## Crosstalk Cancellation with LMS

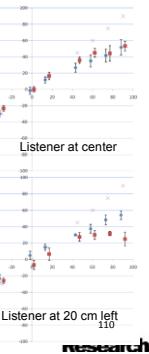
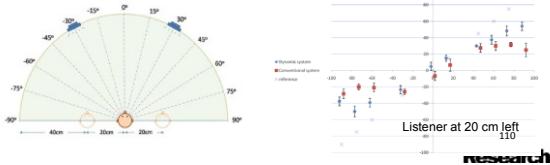


## Our Complete System



## Listening Test

- Pink noise sound source
- Subjects asked to identify virtual sound direction, which is compared with the ground truth



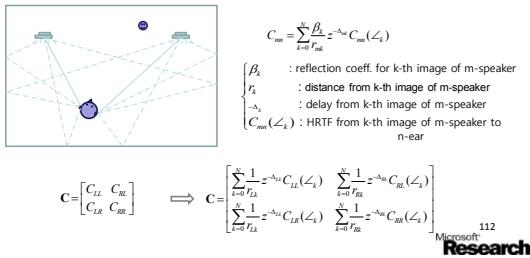
## Adding a Room Model

- Crosstalk canceller suffers from room reverberation. HRTFs only represent the *direct path* of acoustic channel.
- Room model?
  - Simplified room model often sufficient (e.g., SSL)
  - Acoustic measurement possible;

## Room Model Based Acoustic Channel Model

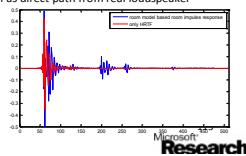
- New acoustic channel model

– Summation of the responses for the individual rays from 6 walls as well as direct path from real loudspeaker.



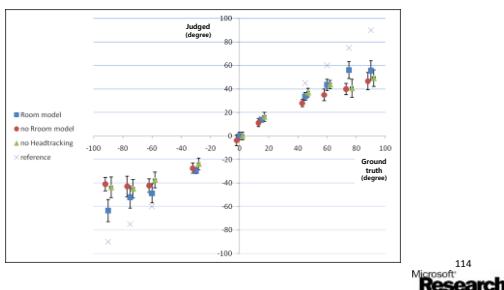
## Room Model Based Crosstalk Canceller

- Procedure
  - Room model estimation
    - Obtaining reflection coefficients and configuration of walls
  - Calculating each ray individually
    - Taking account of attenuation by reflection and distance of the path
  - Determining HRTFs
    - Corresponding to direction of the incoming ray.
  - Summation of the responses
    - For the individual rays from 6 walls as well as direct path from real loudspeaker
- Example of transmission channel based on room model
  - $C_{RR}$
  - Including first reflections from 6 walls



## User Study

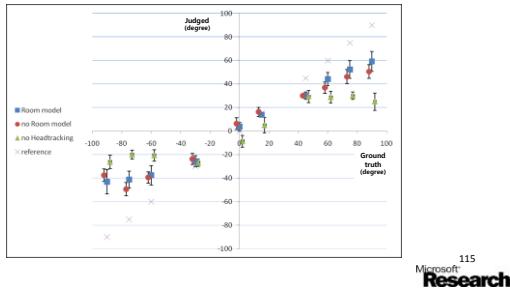
- Listener at center in reverberant room



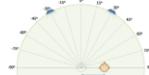
## User Study



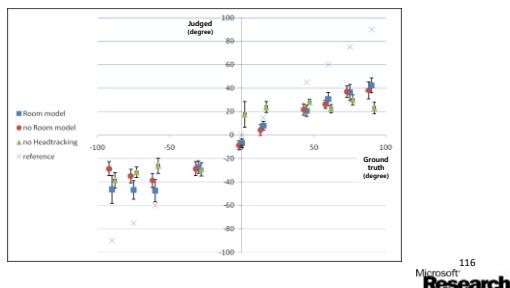
- Listener at 20cm left in reverberant room



## User Study



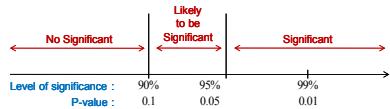
- Listener at 20cm right in reverberant room



## User Study

- Student t-test

- Assessing whether the means of two groups are statistically different from each other
- Comparing absolute value of difference between ground-truth and judged azimuth
  - $|reference_{i,n} - Judged_{i,n}|$ , where i and n are azimuth index and subject index

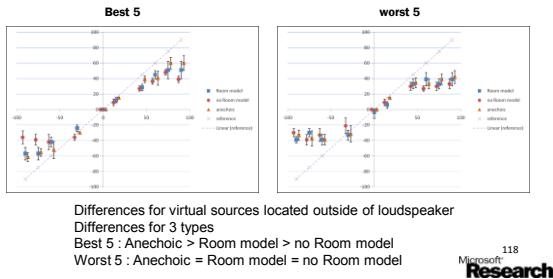


- T-test Results

- Room-model vs. no Room-model  $p = 0.0007893 < 0.05$   
→ There is significant difference between the two method
- Head-tracking vs. no Head-tracking  $p = 0.0000000483052 < 0.05$   
→ There is significant difference between the two method.

## User Study

- Results comparison between 5 good listeners and 5 bad listeners (user at center, no head tracking)



118  
Microsoft Research

## References

- W. Gardner, "3-D audio using loudspeakers," Ph.D. thesis, Massachusetts Institute of Technology, 1997.
- P. Georgiou, A. Mouchtaris, I. Roumeliotis, and C. Kyriakakis, "Immersive Sound Rendering Using Laser-Based Tracking", *Proc. 109th Convention of the Audio Eng. Soc.*, Paper 5227, 2000.
- S. Kim, D. Kong, and S. Jang, "Adaptive Virtual Surround Sound Rendering System for an Arbitrary Listening Position", *J. Audio Eng. Soc.*, Vol. 56, No. 4, 2008.
- C. Kyriakakis and T. Holman, "Video-based head tracking for improvements in multichannel loudspeaker audio", *105th Audio Engineering Society*, San Francisco, CA, 1998.
- T. Lentz, G. Behler, "Dynamic Crosstalk Cancellation for Binaural Synthesis in Virtual Reality Environments", *J. Audio Eng. Soc.*, Vol. 54, Issue 4, pp. 283-294, 2006.
- T. Lentz, O. Schmitz, "Realisation of an adaptive cross-talk cancellation system for a moving listener," *21st AES Conference on Architectural Acoustics and Sound Reinforcement*, 2002.
- J. Loper and A. Gonzalez, "3-D Audio With Dynamic Tracking for Multimedia Environments," 2nd COST-G6 Workshop on Digital Audio Effects, 1999.
- V. Pullki, "Virtual Sound Source Positioning Using Vector Base Amplitude Panning," *J. Audio Eng. Soc.*, Vol. 45, pp. 456-466, 1997.
- M. Song, C. Zhang, D. Florencio and H.-G. Kang, "Personal 3D Audio System with Loudspeakers," *International Workshop on Hot Topics in 3D*, in conjunction with ICME, 2010.
- M. Song, C. Zhang, D. Florencio and H.-G. Kang, "Enhanced Binaural Loudspeaker Audio System with Room Modeling," in *Proc. of MMSp*, 2010.

119  
Microsoft Research

## Outline

- Introduction
- Sound source localization with compact mic arrays
- Spatial sound and multi-channel echo cancellation
  - Head phone based spatial sound
  - Loudspeaker spatial sound
  - Multi-channel echo cancellation
- Video-based tracking
- Other video processing techniques
- Conclusions and future work

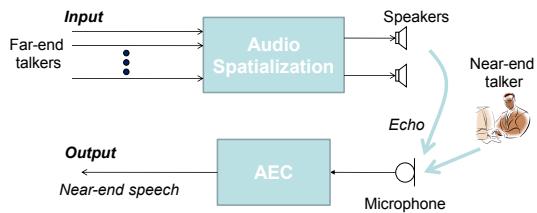
120  
Microsoft Research

## Challenges and Previous Work

- Stereo ACE suffers from the mis-convergence problem (Sonfhi 95)
- In order to de-correlate the speaker signals
  - Adding nonlinearity (Benesty 98)
  - Adding noise (Sonfhi95, Gilloire 98)
- Adaptive filtering
  - NLMS (Yensen 01)
  - Kalman filtering (Enzner 06)
  - Constrained Kalman filtering (Zhang 08)

<sup>121</sup>  
Microsoft  
Research

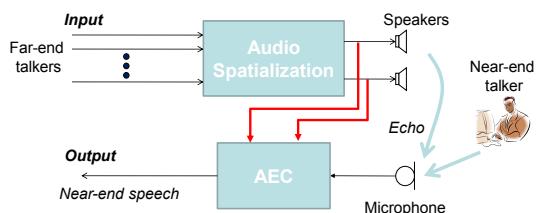
## Multichannel AEC



- Question: Which reference signals to use?

<sup>122</sup>  
Microsoft  
Research

## Approach 1: Use Speaker Signals



- Our solution: Kalman Filtering (general RLS)
- Possible problem: Correlation between speaker signals

<sup>123</sup>  
Microsoft  
Research

## Problem Statement: Speaker based

- Speaker Signals:

$$\mathbf{x}_k^S := [\mathbf{x}_{k,0}^S \mathbf{x}_{k,1}^S \dots \mathbf{x}_{k,L-1}^S] \quad \mathbf{X}_k^S := [\mathbf{x}_k^S \mathbf{x}_{k+1}^S \dots \mathbf{x}_{k+L-1}^S]$$

- Speaker's room response:  $L$ -tap filter

$$\mathbf{h}_k^S := [\mathbf{h}_{k,0}^S \mathbf{h}_{k,1}^S \dots \mathbf{h}_{k,L-1}^S] \quad \mathbf{H}_k^S := [\mathbf{h}_k^S \mathbf{h}_{k+1}^S \dots \mathbf{h}_{k+L-1}^S]$$

- Microphone input: Echo

$$\mathbf{v}_k := [\mathbf{v}_{k,0}^T \mathbf{v}_{k,1}^T \dots \mathbf{v}_{k,L-1}^T]$$

where

$$\mathbb{E}[\mathbf{v}_k] = \mathbf{0} \quad \mathbb{E}[\mathbf{v}_k \mathbf{v}_k^T] = \mathbf{R}_k$$

124  
Microsoft Research

## Multichannel AEC: Speaker-based

- Kalman Filter (Generalized RLS)

$$\mathbf{K}_k := \lambda^{-1} \mathbf{P}_{k-1} \mathbf{X}_k / (\mathbf{Q}_k + \lambda^{-1} \mathbf{X}_k^T \mathbf{P}_{k-1} \mathbf{X}_k)$$

$$\mathbf{P}_k := \mathbf{P}_{k-1} - \mathbf{K}_k^T \mathbf{P}_{k-1} \mathbf{X}_k$$

$$\mathbf{H}_k := \mathbf{H}_{k-1} + \mathbf{K}_k \mathbf{e}_k$$

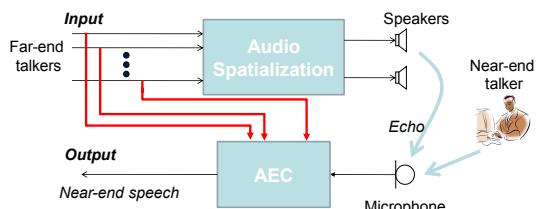
$$\mathbf{P}_{k-1} := \lambda^{-1} (\mathbf{I} - \mathbf{K}_k \mathbf{X}_k^T) \mathbf{P}_{k-1}$$

where  $\lambda$  is the aging factor and  $\mathbf{P}_k$  changes over time

- When  $\lambda^2 = 1$ , Kalman Filter  $\rightarrow$  Traditional RLS
- We model ambient noise.

125  
Microsoft Research

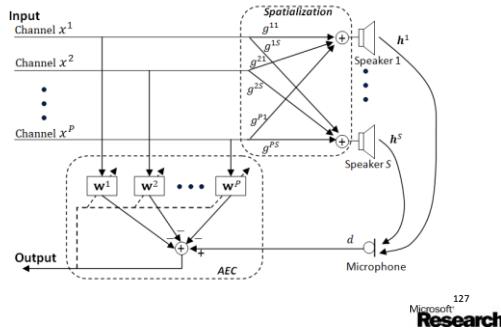
## Approach 2: Use Far-End Channels



- Cancel each individual far-end speech
- Our solution:** Constrained Kalman filtering

126  
Microsoft Research

## Multichannel AEC: CKF Diagram



## Problem Statement

- Remote channels:  $\{x^i | i = 1, \dots, P\}$
- Spatialization on  $S$  speakers:  $Y^s = \sum_{i=1}^P G^{is} X^i$
- Speaker's room response:  $L$ -tap filter  

$$\mathbf{H}_t^s = [H_t^s, H_{t-1}^s, \dots, H_{t-L+1}^s]^T$$
- Microphone input: Echo  

$$D_t = \sum_{i=1}^P \sum_{s=1}^S G^{is} (H_t^s X_t^i + H_{t-1}^s X_{t-1}^i + \dots + H_{t-L+1}^s X_{t-L+1}^i) = \sum_{i=1}^P \sum_{s=1}^S G^{is} \mathbf{H}_t^s X_t^i$$

128  
Microsoft Research

## Problem Statement (cont'd)

- Determine the echo cancellers:
  - one per remote channel  $i$ :  $L$ -tap filter  

$$\mathbf{W}_t^i = [W_t^i, W_{t-1}^i, \dots, W_{t-L+1}^i]^T$$
  - such that echo is cancelled, i.e.,  

$$D_t - \sum_{i=1}^P \mathbf{W}_t^i X_t^i = 0$$
- Constraint:  $\mathbf{W}_i^i$ 's are not mutually independent

$$D_t = \sum_{i=1}^P \sum_{s=1}^S G^{is} \mathbf{H}_t^s X_t^i \quad \longrightarrow \quad \text{Red box}$$

129  
Microsoft Research

## Constrained Kalman Filtering

- State Vector: Echo cancellers + Speaker RIR filters

$$\mathbf{S}_t = [\mathbf{w}_t^T, \dots, \mathbf{w}_t^T, \mathbf{H}_t^T, \dots, \mathbf{H}_t^T]^T$$

- System equation:  $\boxed{\mathbf{D}_t = \mathbf{A}_t^T \mathbf{S}_t + \mathbf{v}_t}$

- Observation equation:  $\mathbf{D}_t = \mathbf{A}_t^T \mathbf{S}_t + \mathbf{v}_t$  with  $\boxed{\mathbf{y}_t = [\mathbf{w}_t^T, \dots, \mathbf{w}_t^T, \mathbf{H}_t^T, \dots, \mathbf{H}_t^T]^T}$

- Constraint:  $\mathbf{C} \mathbf{S}_t = \mathbf{0}$  with  $\mathbf{C} = \begin{bmatrix} -1 & G^{11} \mathbf{1} & \cdots & G^{1S} \mathbf{1} \\ \vdots & \vdots & \ddots & \vdots \\ -1 & G^{P1} \mathbf{1} & \cdots & G^{PS} \mathbf{1} \end{bmatrix}$

- New observation equation: observation + constraint  
with  $\mathbf{Y}_t = [\mathbf{D}_t, 0, \dots, 0]^T$   $\mathbf{B}_t = \begin{bmatrix} \mathbf{A}_t^T \\ \mathbf{C} \end{bmatrix}$   $\mathbf{v}_t = \begin{bmatrix} \mathbf{v}_t \\ \mathbf{u}_t \end{bmatrix}$

130  
Microsoft  
Research

## Constrained Kalman Filtering (cont'd)

- Assumptions

$$E[\mathbf{n}_t] = \mathbf{0} \quad E[\mathbf{n}_t \mathbf{n}_t^T] = Q_t$$

$$E[\mathbf{v}_t] = \mathbf{0} \quad E[\mathbf{v}_t \mathbf{v}_t^T] = R_t = \begin{bmatrix} \sigma_v^2 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$$

Tuning parameter  
to control how hard  
the constraint be satisfied

- Equations

$$\mathbf{S}_t^- = \mathbf{S}_{t-1}$$

$$\mathbf{P}_t^- = \mathbf{P}_{t-1} + Q_t$$

$$\mathbf{K}_t = \mathbf{P}_t^- \mathbf{B}_t^H (\mathbf{B}_t \mathbf{P}_t^- \mathbf{B}_t^H + R_t)^{-1}$$

$$\mathbf{S}_t = \mathbf{S}_t^- + \mathbf{K}_t (\mathbf{Y}_t - \mathbf{B}_t \mathbf{S}_t^-)$$

$$\mathbf{P}_t = (\mathbf{I} - \mathbf{K}_t \mathbf{B}_t) \mathbf{P}_t^-$$

131  
Microsoft  
Research

## Benefits of Constrained KF

- The constraint is taken care of automatically, and can be imposed *with varying degrees*.
- All channels are taken into account simultaneously.  
→ Overlapping far-end talking is not an issue
- The AEC for each channel is updated continuously because of the constraint, even if it is inactive.  
→ AEC's are always up to date
- Ambient noise can be time varying.  
→ Use a separate noise tracker

132  
Microsoft  
Research

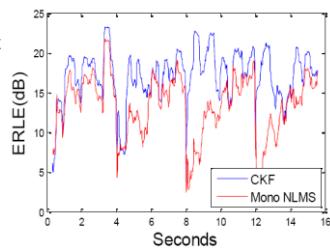
## Comparison with Prior Art

- T.N. Yensen, R.A. Goubran, and I. Lambadaris, "Synthetic Stereo Acoustic Echo Cancellation Structure for Multiple Participant VoIP Conferences", IEEE Transactions on Speech and Audio Processing, Vol. 9, No. 2, pp. 168-174, Feb. 2001.
- Same: One canceller per remote channel
- Differences:
  - Constrained vs. independent cancellers
    - Additional canceller is initiated before being active
    - A canceller is updated even if it is not active
  - Frequency vs. time domain
  - KF (RLS) vs. NLMS

133  
Microsoft Research

## Experimental Results

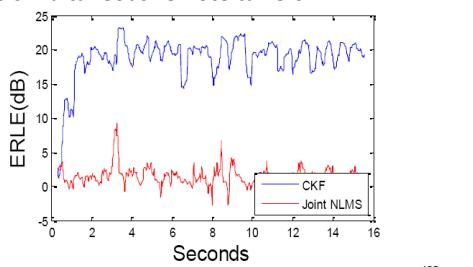
- Simulation setup:
  - 4 remote talkers at  $[-30^\circ, 30^\circ, 0^\circ, -45^\circ]$
  - Each talks for 4s
  - Noise: -20dB
  - Fixed RIR
- Comparison
  - Constrained KF
  - Multiple mono NLMS



134  
Microsoft Research

## Experiment: Overlapping Talkers

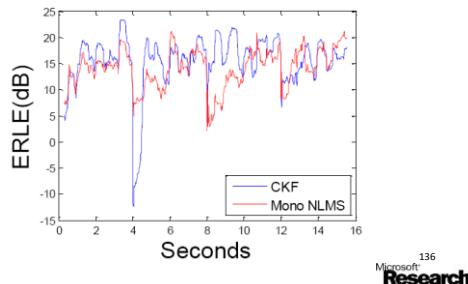
- Two simultaneous remote talkers



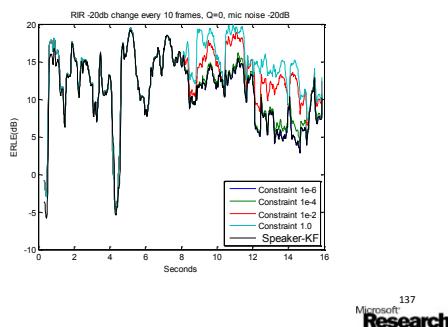
135  
Microsoft Research

## Experiment: Changing RIR

- 30dB change in RIR every 0.5 seconds



## Different Constraints



## Experiment with real data

- Original recording with near-end talker



- AEC with multiple mono NLMSs



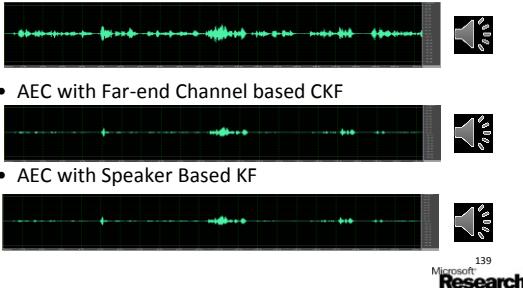
- AEC with CKF



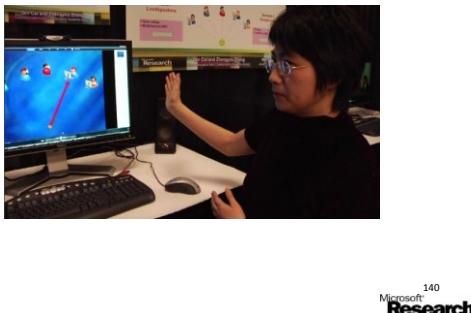
138 Microsoft Research

## Experiment with real data (2)

- Original Mic Input with Near-end Speech



## Demo Video



## References

- J. Baldini, "Effects of spatial audio on memory, compression, and preference during desktop conferences," in *Proc. of the CHI Conference on Human Factors in Computing Systems*, 2001.
  - J. Benesty, D. Morgan, and M. Sonihi, "A better understanding and an improved solution to the problem of stereophonic acoustic echo cancellation," *IEEE Trans. on Speech and Audio Processing*, vol. 6, pp. 156-165, 1998.
  - G. Enzner, and P. Vary, "Frequency-domain adaptive Kalman filter for acoustic echo control in hands-free telephones," *Signal Processing*, vol. 86, pp. 1140-1156, Oct, 2006.
  - A. Gilloire, and V. Turbin, "Using auditory properties to improve the behavior of stereophonic acoustic cancellers," in *Proc. of ICASSP*, 1998.
  - M. Sonihi, D. Morgan, and J. Hall, "Stereophonic acoustic echo cancellation – An overview of the fundamental problem," *IEEE Speech Processing Letter*, vol. 2, pp. 148-151, 1995.
  - T. N. Yensen, R. A. Goubran, and I. Lambadaris, "Synthetic Stereo Acoustic Echo Cancellation Structure for Multiple Participant VoIP Conference," *IEEE Trans. on Speech and Audio Processing*, vol. 9, no. 2, pp. 168-174, Feb, 2001.
  - Z. Zhang, Q. Cai and J. Stokes, "Multichannel Acoustic Echo Cancellation in Multiparty Spatial Audio Conferencing with Constrained Kalman Filtering," in *Proc. of IWAEWA*, 2008.

## Outline

- Introduction
- Sound source localization with compact mic arrays
- Spatial sound and multi-channel echo cancellation
- Video-based tracking
  - Multi-camera head pose tracking
  - Depth camera based facial expression tracking
  - Depth camera based Skeleton tracking
- Other video processing techniques
- Conclusions and future work

142  
Microsoft Research

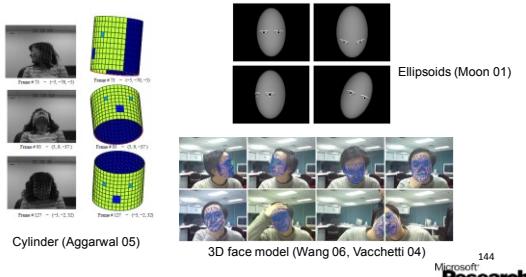
## Motivation

- Head Pose is very important for tele-conferencing
  - Face Relighting
  - Gaze Correction
  - Adaptive Displays
  - Video driven avatar animation, etc.

143  
Microsoft Research

## Single Camera Head Pose Tracking

- Model-based (see Lepetit 05 for a survey)



144  
Microsoft Research

## Motivation

- Single Camera Tracking has limitation on
  - Narrow Range of Tracked Poses
  - Occlusions
- Multiple Camera Tracking
  - Existing Setup in Video Conferencing
  - Multiple Views of the Face

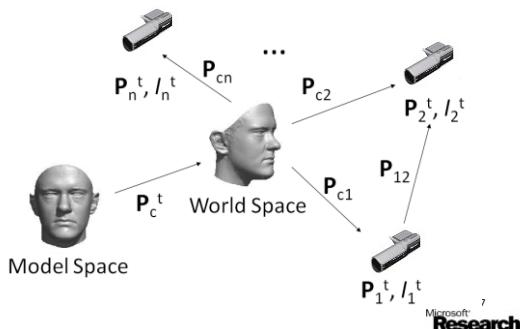
145  
Microsoft  
Research

## Challenges

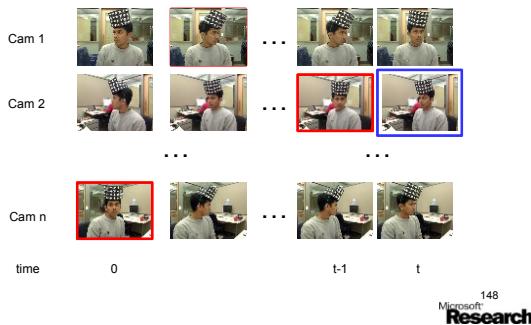
- Our Multi-Cam Tracking System
  - No Pre-calibration of Camera Placements
  - Real Time Processing
  - Drift Free
  - Wider Range of Tracked Poses
  - Continuous Tracking w/ Occlusion
  - Recover after Tracking is Lost

146  
Microsoft  
Research

## Model-based Approach



## Tracking by References



148  
Microsoft Research

## Generic Tracking Scheme

- Suitable for both single and multiple cameras
- References span on both Temporal and Spatial domains

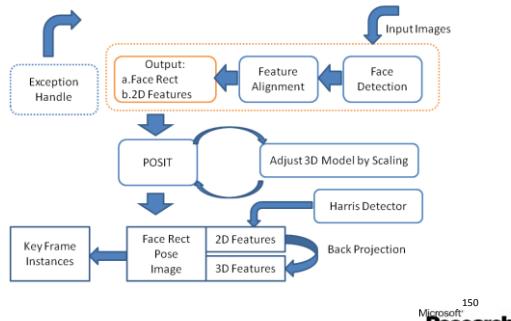
$$v(t) := \sum_{\text{reference } r} w_r v_r(t) \\ + \sum_{r, k} w_{rk} \rho \left( \eta_{rk}^t - \psi(\lambda_r^t \mu_r^t \nu_k^t \eta_k^t) \right)$$

- Camera Relative Geometry converges over time

$$\begin{bmatrix} l_{ij}^0 & l_{ij}^{old} & l_{ij}^2 \\ l_{ij}^0 & l_{ij}^{old} & l_{ij}^2 \end{bmatrix}^T \Delta_{ij}^{-1} \begin{bmatrix} l_{ij}^0 & l_{ij}^{old} \end{bmatrix}$$

149  
Microsoft Research

## Tracking Initialization



150  
Microsoft Research

## Tracking Algorithm

- **Input:** Image  $I^k$ , previous frames  $I^{k-1}$  & key frames  $I^j$
- **Output:** Current Pose  $\hat{M}^k$ .
- Select Pose of Reference frames  $\hat{M}^c$  on each camera
- Initialize Relative Pose  $\hat{M}_{ci}^k$  between Ref Camera  $c$  and other camera  $i$ .
- For each camera  $i$ , do
  - For each pose of ref. frame  $\hat{M}^c = X^c \cup Y^c$ , do
  - Find 2D feature matches between  $I^k$  and  $I^c$ , then back-project in  $I^k$  to get 3D points in the model  $U$ .
  - Discard outliers by applying RANSAC
  - end
- Apply POSIT for all inlier 3D-2D matches
- End

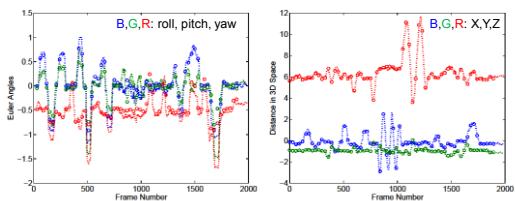
151  
Microsoft  
Research

## Results (3 Cameras)



152  
Microsoft  
Research

## Comp. w/ Ground Truth I (No Occlusion)

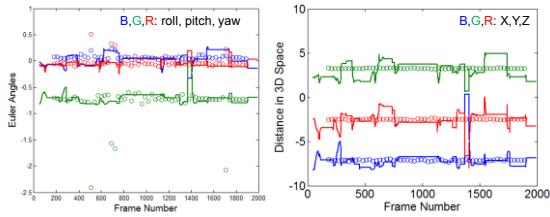


..... Tracked Pose from our Algorithm.  
o Ground Truth from Calibration.

153  
Microsoft  
Research

## Relative Pose Comp. w/ Ground Truth

II



<sup>154</sup>  
Microsoft  
Research

## Results (Occlusion)



<sup>156</sup>  
Microsoft  
Research

## Fair Comparison (no Occlusion)

# of Cams	Median Angle Error (degree)	Median Translation Error	Mean Angle Error (degree)	Mean Translation Error
1	5.9754	0.4138	19.0132	0.5504
2	5.4878	0.3481	11.8039	0.4885
3	4.4671	0.3404	11.0915	0.4634
3 (w/ pre camera calibration)	4.0827	0.3132	11.0577	0.4540

<sup>157</sup>  
Microsoft  
Research

## References

---



---



---



---



---



---



---



---



---

- G. Aggarwal, A. Veeraraghavan, and R. Chellappa, "3D Facial Pose Tracking in Uncalibrated Videos," *Lecture Notes in Computer Science*, 3776:515, 2005.
- Q. Cai, A. Sankaranarayanan, Q. Zhang, Z. Zhang, and Z. Liu, "Real Time Head Pose Tracking from Multiple Cameras with a Generic Model," in *IEEE Workshop on Analysis and Modeling of Faces and Gestures*, 2010.
- V. Lepetit and P. Fua, "Monocular model-based 3d tracking of rigid objects: A survey," *Foundations and Trends in Computer Graphics and Vision*, 1(1):1–89, October 2005.
- H. Moon, R. Chellappa, and A. Rosenfeld, "3d object tracking using shape-encoded particle propagation," in *IEEE International Conference on Computer Vision*, 2001.
- L. Vacchetti, V. Lepetit, and P. Fua, "Stable Real-Time 3D Tracking Using Online and Offline Information," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1385–1391, 2004.
- Q. Wang, W. Zhang, X. Tang, and H. Shum, "Real-Time Bayesian 3-D Pose Tracking," *IEEE Transactions on Circuits and Systems for Video Technology*, 16(12):1533, 2006.

<sup>158</sup>  
Microsoft  
Research

---



---



---



---



---



---



---



---

## Outline

---



---



---



---



---



---



---



---

- Introduction
- Sound source localization with compact mic arrays
- Spatial sound and multi-channel echo cancellation
- **Video-based tracking**
  - Multi-camera head pose tracking
  - Depth camera based facial expression tracking
  - Depth camera based Skeleton tracking
- Other video processing techniques
- Conclusions and future work

<sup>159</sup>  
Microsoft  
Research

---



---



---



---



---



---



---



---

## Deformable Face Tracking

---



---



---



---



---



---



---



---

- Many applications
  - Human computer interaction
  - Performance-driven facial animation
  - Face recognition
- Challenging
  - Limited number of features on the face
  - Dozens of parameters to estimate

<sup>160</sup>  
Microsoft  
Research

---



---



---



---



---



---



---



---

## Related Works

- Video based fitting
  - Active appearance models (AAM) (Xiao04)
  - 3D morphable models (Blanz99)
  - 3D deformable model + active shape model (ASM) (Vogler07)
  - Feature based (Zhang08)
- 3D scans based fitting
  - Iterative closest point (ICP) + 3D deformable model (Zhang04, Wang04, Weise09, Weise11)

<sup>161</sup>  
Microsoft  
Research

## Commodity Depth Cameras



<sup>162</sup>  
Microsoft  
Research

## Kinect Sensor

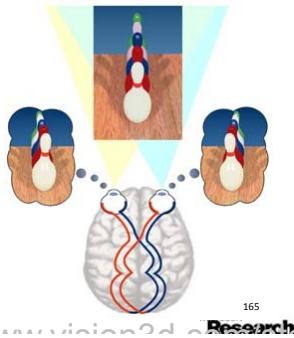


<sup>163</sup>  
Microsoft  
Research

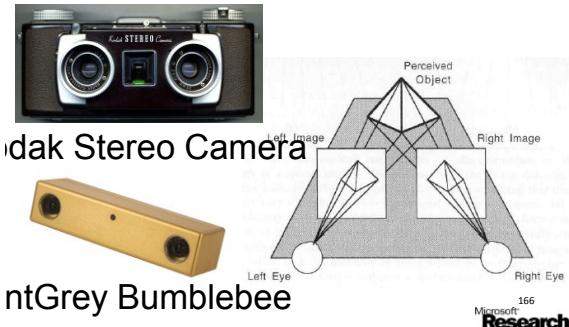


## Human Stereo Vision

**Difference**  
in your two eyes  
gives you the ability  
to perceive  
your surrounding  
environment  
in **3 Dimensions**



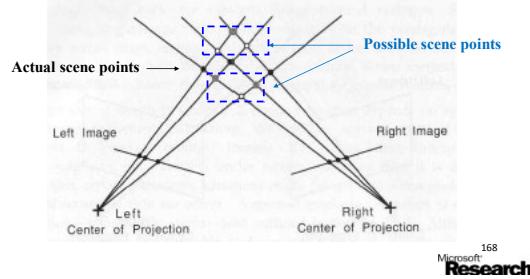
# Computer Stereo Vision



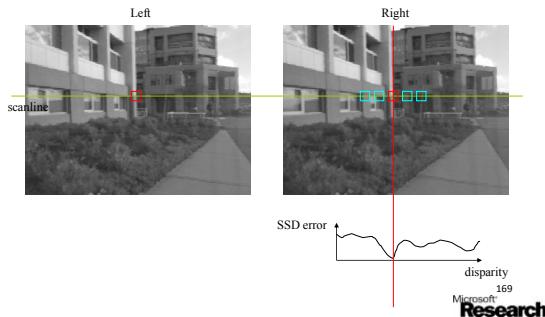
## Stereo Images



## Correspondence Problem

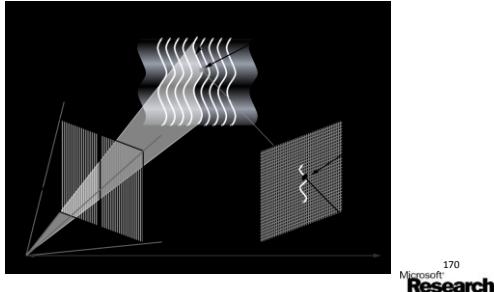


## Correspondence by Correlation



## How it works?

- Structured light 3D scanner




---



---



---



---



---



---



---



---



---

## How it works? Kinect Sensor

- Modified structured light 3D scanner

- IR projector
- IR camera
- Random pattern




---



---



---



---



---



---



---



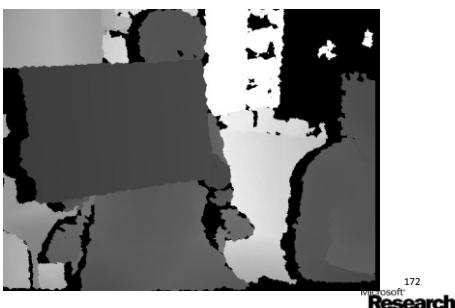
---



---

## Matching & Depth Map

- Correlation




---



---



---



---



---



---



---



---



---

## Overlay of Depth Map on IR Image




---

---

---

---

---

---

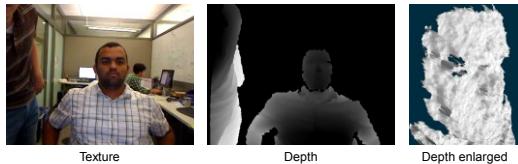
---

---

---

---

## Captured Image



174

Microsoft  
Research

---

---

---

---

---

---

---

---

---

---

## Challenge and Solution

- Depth information should help tracking
- Challenge: noisy depth input
- Solution:
  - Maximum likelihood solution with arbitrary noise covariance matrices
  - Regularization

175

Microsoft  
Research

---

---

---

---

---

---

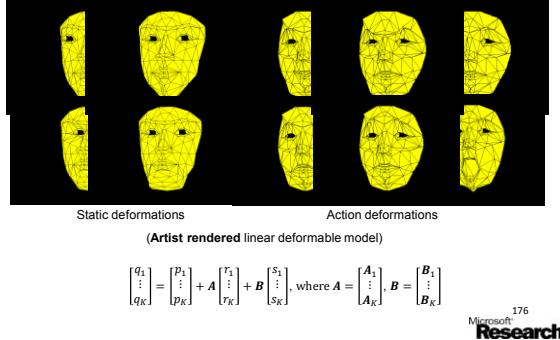
---

---

---

---

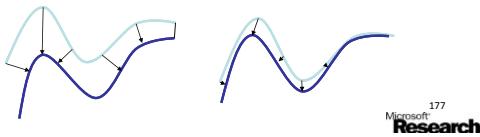
## Linear Deformable Model



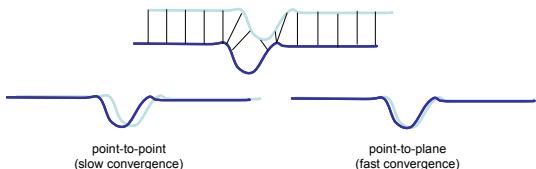
## Maximum Likelihood DMF

- Formulation,  $(q_k, g_k)$  correspondence pair:  

$$\begin{aligned} \mathbf{R}(\mathbf{p}_k + A_k \mathbf{r} + B_k \mathbf{s}) + \mathbf{t} &= \mathbf{g}_k + \mathbf{x}_k \\ \mathbf{x}_k &\sim N(\mathbf{0}, \Sigma_{x_k}) \end{aligned}$$
- Iterative closest point
  - Assume closest points correspond
  - Compute transformation
  - Iterate until convergence



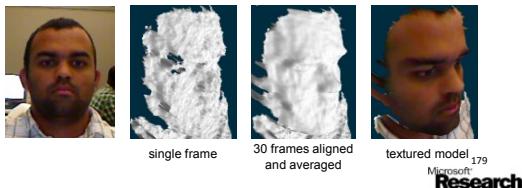
## Point-to-Point and Point-to-Plane Distance



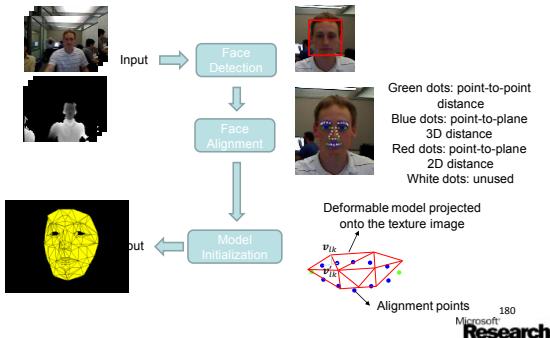
- In both cases, we have solution for arbitrary noise covariance matrix  $\Sigma_{x_k}$ 
  - See paper for details

## Model Initialization

- Initialization
  - Assume multiple neutral face frames available
  - Action deformations set to zero
  - Jointly solve shape deformations and rotation/translation of each frame



## Model Initialization

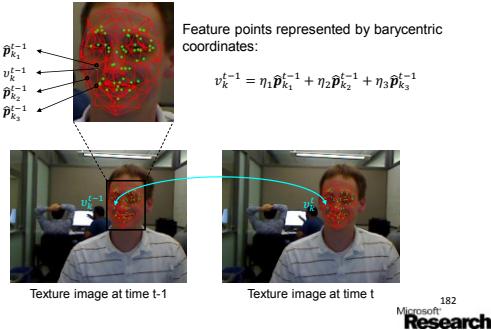


## Face Tracking

- Tracking
  - Shape deformations fixed
  - Based on feature point correspondence
  - Solve for action deformation, rotation and translation
  - Regularization
    - $L_2$  norm constraining the difference between neighboring frames' action deformations
    - $L_1$  norm constraining the number of non-zero action deformation parameters

Microsoft Research<sup>181</sup>

## Face Tracking



## Tracking Results



[Video](#)

<sup>183</sup>  
Microsoft  
Research

## Qualitative Results

Median tracking error in pixels

	ID+ $l_2$	ID+ $l_1$	ID+ $l_2+l_1$	NM+ $l_2$	NM+ $l_1$	NM+ $l_2+l_1$
Seq #1	3.56	2.88	2.78	2.85	2.69	2.66
Seq #2	4.48	3.78	3.71	4.30	3.64	3.55
Seq #3	3.98L	3.91	3.91	3.92L	3.91	3.50

ID: use identity covariance matrix for sensor noise

NM: use the proposed noise modeling scheme

$l_2$ : quadratic constraint between successive frames

$l_1$ : sparse constraint on the action transforms

L: lost tracking in the middle and never recover

<sup>184</sup>  
Microsoft  
Research

## References

- V. Blanz and T. Vetter, "A morphable model for the synthesis of 3D faces," in *Proc. of SIGGRAPH*, 1999.
- Q. Cai, D. Gallup, C. Zhang and Z. Zhang, "3D Deformable Face Tracking with a Commodity Depth Camera," in *Proc. of ECCV*, 2010.
- Y. Wang, X. Huang, C. Lee, S. Zhang, Z. Li, D. Samaras, D. Metaxas, A. Elgammal, P. Huang, "High resolution acquisition, learning and transfer of dynamic 3-D facial expressions," in *Proc. of EUROGRAPHICS*, 2004.
- T. Weise, S. Bouaziz, H. Li, M. Pauly, "Realtime performance-based facial animation," in *Proc. of SIGGRAPH*, 2011.
- T. Weise, H. Li, L. Gool, M. Pauly, "Face/off: Live facial puppetry," in *Symposium on Computer Animation*, 2009.
- C. Vogler, Z. Li, A. Kanaujia, S. Goldenstein, D. Metaxas, "The best of both worlds: Combining 3d deformable models with active shape models," in *Proc. of ICCV*, 2007.
- J. Xiao, S. Baker, I. Matthews, T. Kanade, "Real-time combined 2d+3d active appearance models," in *Proc. of CVPR*, 2004.
- L. Zhang, N. Snavely, B. Curless, S. Seitz, "Spacetime faces: high-resolution capture for modeling and animation," in *Proc. of SIGGRAPH*, 2004.
- W. Zhang, Q. Wang, X. Tang, "Real time feature based 3-D deformable face tracking," in *Proc. of ECCV*, 2008.

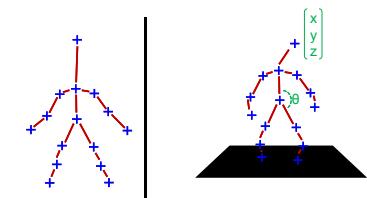
<sup>185</sup>  
Microsoft  
Research

## Outline

- Introduction
- Sound source localization with compact mic arrays
- Spatial sound and multi-channel echo cancellation
- Video-based tracking
  - Multi-camera head pose tracking
  - Depth camera based facial expression tracking
  - Depth camera based Skeleton tracking
- Other video processing techniques
- Conclusions and future work

<sup>186</sup>  
Microsoft  
Research

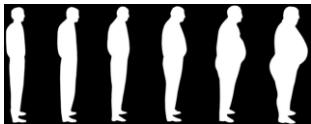
## Human Pose Estimation



Kinect tracks 20 body joints in real time.

<sup>187</sup>  
Microsoft  
Research

## Why Is It Hard?



188  
Microsoft  
Research

## Why Do We Want It?



And applications to be determined...

189  
Microsoft  
Research

## Motion Capture



[Vicon]



[Xsens]

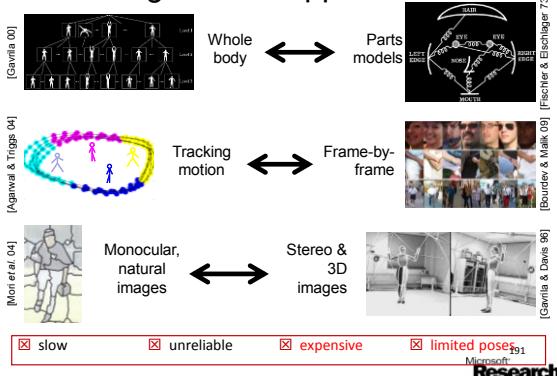
- very accurate
- high frame rate

- suit / sensors
- expensive

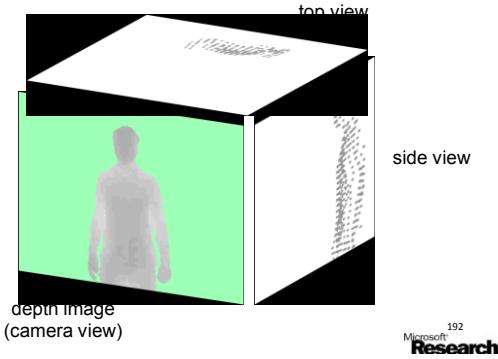
- large space
- calibration

190  
Microsoft  
Research

## RGB Image-based Approaches



## Depth Cameras



## RGB vs. Depth for Pose Estimation

### RGB

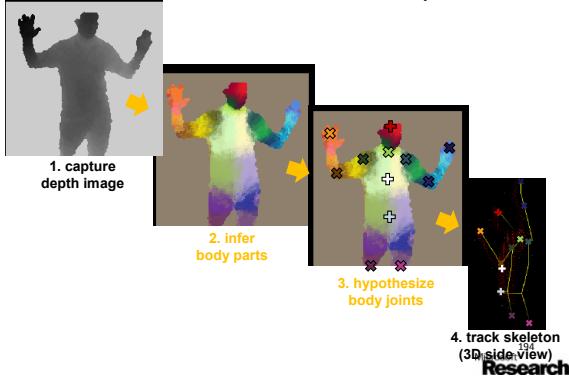
- Only works well lit
- Background clutter
- Scale unknown
- Clothing, skin colour



### DEPTH

- Works in low light
- Person 'pops' out from bg
- Scale known
- Uniform texture
- easy to simulate
- Shadows, missing pixels

## The Kinect Pose Estimation Pipeline




---



---



---



---



---



---



---

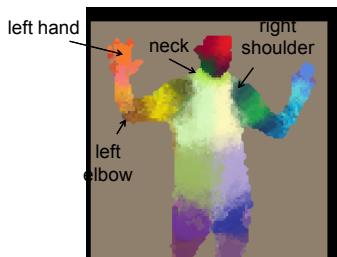


---



---

## Body Part Recognition



("left" = player left with camera acting as mirror)

195  
Microsoft Research

---



---



---



---



---



---



---



---



---

## Body Part Recognition

- No temporal information
  - Frame-by-frame
- Local pose estimate of parts
  - Independently treat
    - Each pixel & each body joint
  - Reduces training data and computation time
- Very fast
  - Depth image features
  - Decision forest classifier



196  
Microsoft Research

---



---



---



---



---



---



---

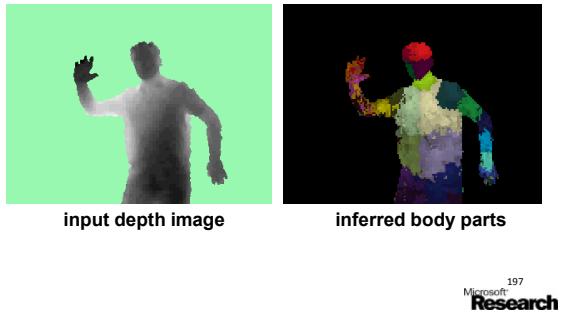


---



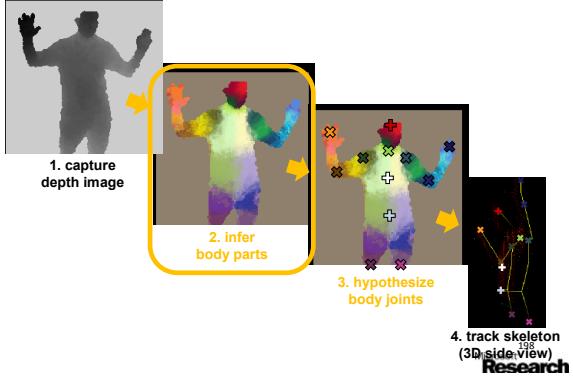
---

## Body Part Recognition



197  
Microsoft  
Research

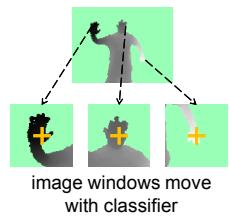
## The Kinect Pose Estimation Pipeline



Microsoft  
Research

## Classifying Pixels

- Compute  $P(c_i|w_i)$ 
  - pixels  $i = (x, y)$
  - body part  $c_i$
  - image window  $w_i$
- Discriminative approach
- Learn classifier  $P(c_i|w_i)$  from training data



199  
Microsoft  
Research

## Fast Depth Image Features

- Depth comparisons:

$$- f(i; \Delta) = d(i) - d(i') \\ \text{where } i' = i + \Delta$$

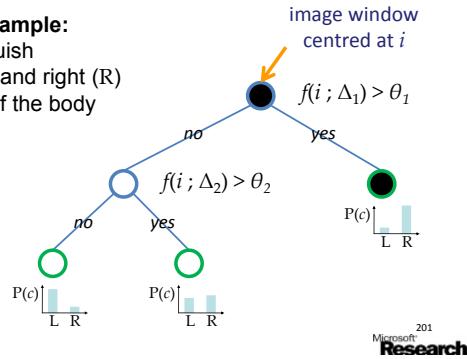
- Background pixels

-  $d$  = large constant



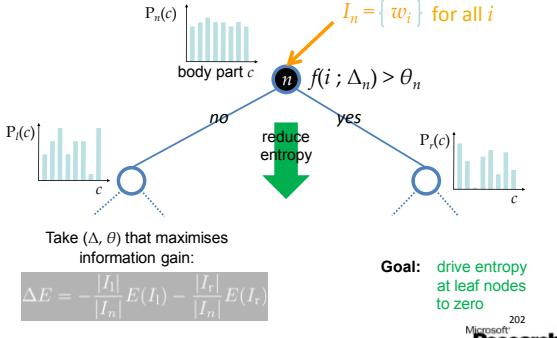
## Decision Tree Classification

**Toy example:**  
distinguish  
left (L) and right (R)  
sides of the body



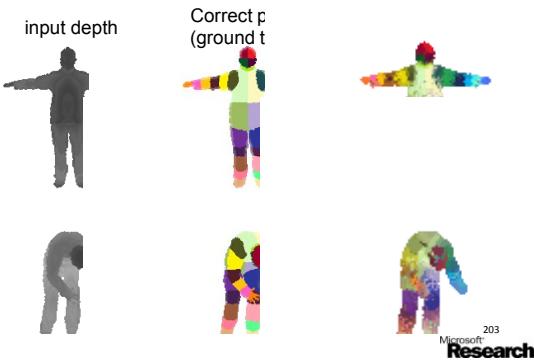
Microsoft Research 201

## Training Decision Trees (Breiman 84)



Microsoft Research 202

### Depth of Trees




---

---

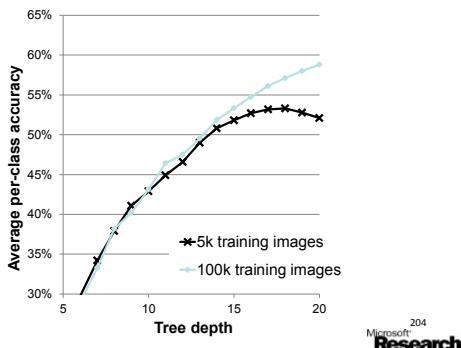
---

---

---

---

### Depth of Trees




---

---

---

---

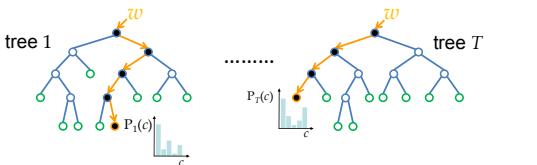
---

---

### Decision Forests

(Amit 97, Breiman 01, Geurts 06)

- Single trees tend to over-fit
- Train forest – ensemble of trees:



– different random subset of images

– average tree posteriors

$$P(c|w) = \sum_{t=1}^T P_t(c|w)$$

205 Microsoft Research

---

---

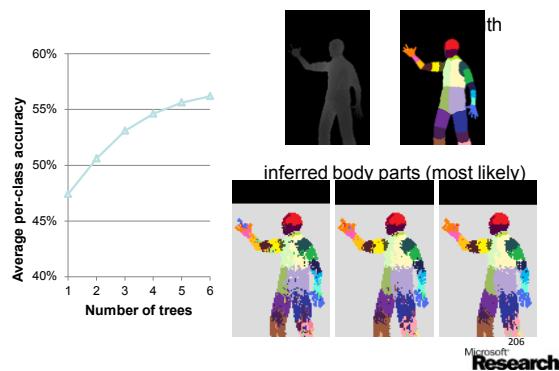
---

---

---

---

## Number of Trees




---

---

---

---

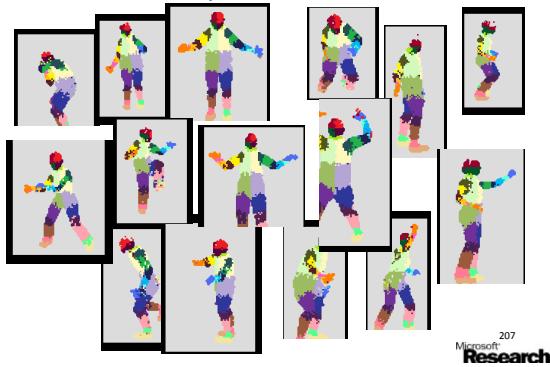
---

---

---

---

## Example Inferences




---

---

---

---

---

---

---

---

## Body Parts to Joint Hypotheses

- Depth image & probability mass
- Localize body parts in 3D
  - global centroid of prob. mass
  - local modes of density (mean shift)
- Map body parts to skeletal joints
  - many parts map directly to joints




---

---

---

---

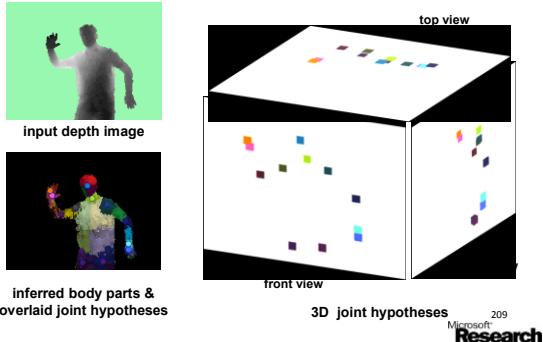
---

---

---

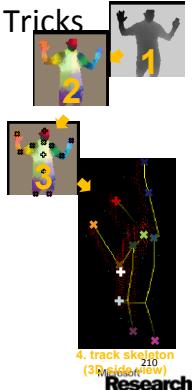
---

## 3D Joint Hypotheses



## ... and Some Other Tricks

- Exploit
  - 3D joint hypotheses
  - kinematic constraints
  - temporal coherence
  
- Predict
  - full skeleton
  - invisible joints
  - multi-player



## References

- Y. Amit and D. Geman, "Shape quantization and recognition with randomized trees," *Neural Computation*, pp. 1545–1588, Vol. 9, No. 7, 1997.
- L. Breiman, "Random forests," *Machine Learning*, pp. 5–32, Vol. 45, No. 1, 2001.
- L. Breiman, J. Friedman, C. Stone and R. Olshen, *Classification and Regression Trees*, Chapman and Hall/CRC, 1984.
- P. Geurts, D. Ernst and L. Wehenkel, "Extremely randomized trees," *Machine Learning*, pp. 3–42, Vol. 63, No. 1, 2006.
- J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, A. Blake, "Real-time human pose recognition in parts from single depth images," *Proc. of CVPR*, 2011.

## Outline

- Introduction
- Sound source localization with compact mic arrays
- Spatial sound and multi-channel echo cancellation
- Video-based tracking
- Other video processing techniques
  - Learning-based perceptual image quality enhancement
  - Monocular 3D and Foreground/background segmentation
  - Joint tracking and multiview video compression
- Conclusions and future work

<sup>212</sup>  
Microsoft  
Research

---



---



---



---



---



---



---



---

## Motivation

- Two factors affecting image quality during video conferencing:
  - Camera device
  - Lighting condition
- Observation
  - Camera devices are getting better
  - Lighting conditions will stay the same
- Goal:
  - Improving virtual lighting condition

<sup>213</sup>  
Microsoft  
Research

---



---



---



---



---



---



---



---

## Previous work

- Brightness and contrast improvement
  - Hardware:
    - Automatic gain control
    - Face tracking -> gain control parameters
  - Image processing techniques (exposure correction):
    - Bhukhanwala 94, Messina 03, Saitoh 99, Shi 04

<sup>214</sup>  
Microsoft  
Research

---



---



---



---



---



---

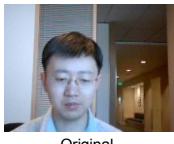


---

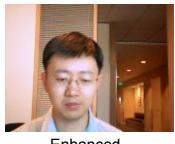


---

## Beyond exposure



Original



Enhanced

Good exposure

But pale, unpleasant

Warm, appealing

<sup>215</sup>  
Microsoft  
Research

---

---

---

---

---

---

---

## Color tone

- Stage lighting design
  - Color scheme is critical
    - The perceived look of the host / actors
    - Mood of the stage
- How does a color affect the look of a face
  - Subjective
  - No mathematical formula
  - Learning from experts

<sup>216</sup>  
Microsoft  
Research

---

---

---

---

---

---

---

## Learning-based color tone mapping

- Idea:
  - Learning color statistics from images taken by professional photographers (in good lighting conditions)
  - Color tone mapping for an input image
- Data collection:
  - 400 celebrity images



<sup>217</sup>  
Microsoft  
Research

---

---

---

---

---

---

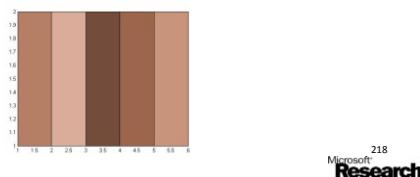
---

## Training

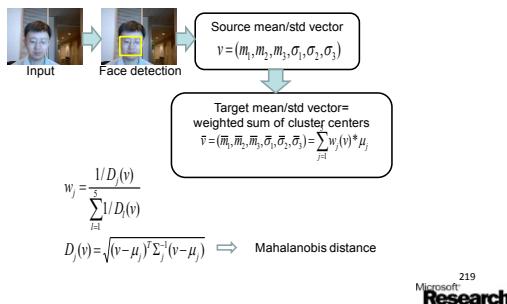
- Mixtures of Gaussians

- Feature:

- mean and variance on the face region  
 $v_i = (m_{i,r}, m_{i,g}, m_{i,b}, \sigma_{i,r}, \sigma_{i,g}, \sigma_{i,b}), i=1, \dots, 400$
- 5 mixture components  $(\mu_j, \Sigma_j), j=1, \dots, 5$



## Color tone mapping



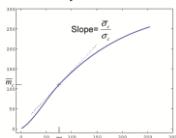
## Color tone mapping function

$y = f_c(x)$ : Color tone mapping function for channel c

Piecewise cubic spline with constraints:

$$f_c(m_c) = \bar{m}_c \quad f_c(0) = 0 \quad f_c(255) = 255$$

$$f'_c(m_c) = \frac{\bar{\sigma}_c}{\sigma_c} \quad f'_c(0) = 0.5 \frac{\bar{m}_c}{m_c} \quad f'_c(255) = 0.5 \frac{255 - \bar{m}_c}{255 - m_c}$$



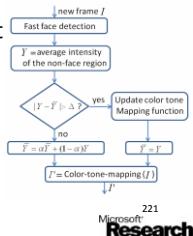
220  
Microsoft  
Research

## Application to video sequence

- Image intensities change over time:
  - Automatic gain control
  - Lighting change

- Global intensity change detect

- Updating color tone mapping function whenever a global intensity change is detected



221  
Microsoft Research

## Computation reduction

- Color table lookup:

- Each color tone mapping functions  $f_i(x)$  is stored as a color lookup table with 256 entries
- Color tone mapping for each pixel becomes 3 simple table lookup.

- CPU usage:

- 5% on a 3.2GHz PC with video resolution 320x240 and frame rate 30fps

222  
Microsoft Research

## Experiment results



223  
Microsoft Research

## User Study

- 16 video sequences
  - Offices with various lighting conditions
- Eight different webcams
  - Logitech, Creative, Veo, Dynex, Microsoft LiveCam, etc.
- User is asked to look at the side-by-side view for each sequence, and rate each one with a MOS score:
  - 1: very bad quality
  - 2: bad quality
  - 3: acceptable
  - 4: good quality
  - 5 very good quality

224  
Microsoft  
Research

## User study results

- 18 users responded
- Average improvement:
  - 2.55 ➔ 3.30

Sequence ID	Original Video	Enhanced Video
1	2.89	3.30
2	2.94	3.33
3	1.94	2.67
4	1.61	2.83
5	2.11	3.00
6	3.22	4.17
7	3.11	3.72
8	1.89	2.72
9	2.11	2.17
10	1.94	2.00
11	2.11	3.44
12	3.28	3.72
13	2.65	4.18
14	2.33	3.72
15	3.00	3.89
16	3.72	3.67
Average	2.88	3.30

225  
Microsoft  
Research

## Extend to Hardware Solution

- Observation:
  - Camera auto exposure responds to environment
  - Local Lighting for target (face) may not be sufficient
- Control exposure based on target
  - Extreme exposure leads to noise
  - Lack of color tone adjustment



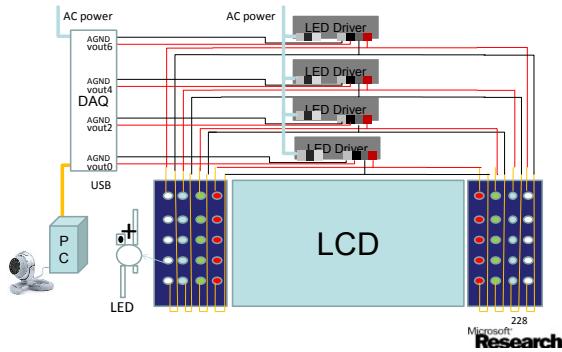
226  
Microsoft  
Research

## Our Solution

- Active lighting
  - Use computer controlled lights
  - Automatically determine the appropriate lighting to optimize face image quality

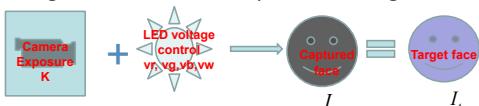


## Hardware Configuration



## Software: Optimization Problem

- Image is a function of exposure and light



- Goal: Find the right exposure  $k$  and voltages  $\mathbf{v}$  to minimize the difference between captured image and target image

$$\min_{k, \mathbf{v}} \|I(k, \mathbf{v}) - I_t\|$$

$$\mathbf{v} = [v_r, v_g, v_b, v_w]^T$$

229  
Microsoft Research

## Experiment Results



[Video](#)

<sup>230</sup>  
Microsoft  
Research

---

---

---

---

---

---

---

---

## References

- S. A. Bhukhanwala and T. V. Ramabadran, "Automated global enhancement of digitized photographs," *IEEE Trans. on Consumer Electronics*, 40(1), 1994.
- Z. Liu, C. Zhang and Z. Zhang, "Learning-Based Perceptual Image Quality Improvement for Video Conferencing," in *Proc. of ICME*, 2007.
- G. Messina, A. Castorina, S. Battiatto, and A. Bosco, "Image quality improvement by adaptive exposure correction techniques," In *Proc. of ICME*, 2003.
- F. Saitoh, "Image contrast enhancement using genetic algorithm," in *IEEE International Conference on SMC*, 1999.
- C. Shi, K. Yu, J. Li, and S. Li, "Automatic image quality improvement for videoconferencing," in *Proc. of ICASSP*, 2004.
- M. Sun, Z. Liu, J. Qiu, Z. Zhang, M. Sinclair, "Active lighting for Video Conferencing," *IEEE Trans. on CSE*, Vol. 19, No. 12, pp.1819-1829, Dec. 2009.

<sup>231</sup>  
Microsoft  
Research

---

---

---

---

---

---

---

---

## Outline

- Introduction
- Sound source localization with compact mic arrays
- Spatial sound and multi-channel echo cancellation
- Video-based tracking
- **Other video processing techniques**
  - Learning-based perceptual image quality enhancement
  - Monocular 3D and Foreground/background segmentation
  - Joint tracking and multiview video compression
- Conclusions and future work

<sup>232</sup>  
Microsoft  
Research

---

---

---

---

---

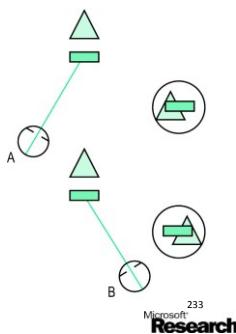
---

---

---

## What is Motion Parallax

- Apparent difference in the direction of movement or speed produced when the subject moves relative to his environment




---

---

---

---

---

---

---

## Track User Movement

- Motion tracking systems (Welch and Foxlin 02)
  - Mechanical sensing
  - Inertial sensing
  - Acoustic sensing
  - Magnetic sensing
  - Optical sensing
  - Radio frequency sensing
- For our applications, we use
  - Face tracking via webcams

<sup>234</sup>  
Microsoft  
Research

---

---

---

---

---

---

---

## Mono 3D

(Pastoor 99, Suenaga 05, Harrison 08)

- Pseudo-3D effects for teleconferencing
  - Box framing




---

---

---

---

---

---

---

## Mono 3D

- Pseudo-3D effects for teleconferencing

- Box framing

- Layered video



236  
Microsoft  
Research

---

---

---

---

---

---

---

## Video Segmentation

- Very challenging problem
- Existing scheme using:
  - Stereo (Kolmogorov 05)
  - Mono video (Criminisi 06, Sun 06, Yu 07)
  - Depth camera

237  
Microsoft  
Research

---

---

---

---

---

---

---

## Fg/Bg Segmentation with Depth Camera

- We explore the use of depth camera for live video foreground/background segmentation



ZCam system specification: Full RGB resolution: 640\*480;  
Operation range: 0.5-2.5 meter  
Depth data rate: 30fps (320°\*240)

238  
Microsoft  
Research

---

---

---

---

---

---

---

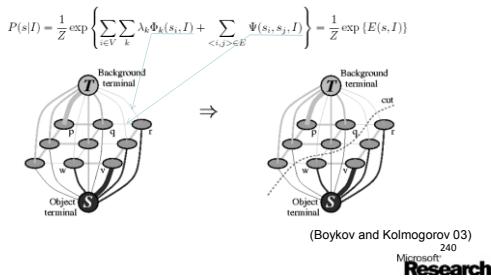
## Sample Video



239  
Microsoft  
Research

## Graph-Cut Based Segmentation

- We build a graph model for graph-cut segmentation.



## The Cues for Segmentation

- ZCam provides color images, depth and depth confidence.



$$P(s|I) = \frac{1}{Z} \exp \left\{ \sum_{i \in V} \sum_k \lambda_k \Phi_k(s_i, I) + \sum_{\langle i, j \rangle \in E} \Psi(s_i, s_j, I) \right\} = \frac{1}{Z} \exp \{E(s, I)\}$$

241  
Microsoft  
Research

## Data Term

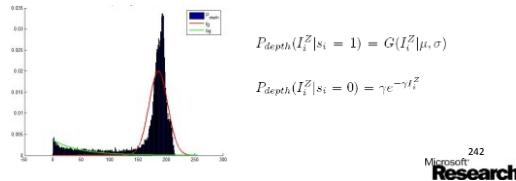
- Data term in the graph

$$\Phi_{color}(s_i, I) = -\ln \{P_{color}(I_i|s_i) \cdot P_{confidence}(s_i)\}$$

$$\Phi_{depth}(s_i, I) = -\ln \{P_{depth}(I_i|s_i) \cdot P_{confidence}(s_i)\}$$

The color likelihood is computed from the previous segmentation mask.

The depth likelihood is generated from depth histogram.



## Data Term

- Data term in the graph

$$\Phi_{depth}(s_i, I) \quad s_i = 1 \quad s_i = 0$$



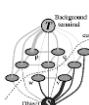
$$\Phi_{color}(s_i, I) = -\ln \{P_{color}(I_i|s_i) \cdot P_{confidence}(s_i)\}$$

$$\Phi_{depth}(s_i, I) = -\ln \{P_{depth}(I_i|s_i) \cdot P_{confidence}(s_i)\}$$

## Link Term

- Link term in the graph

$$\Psi(s_i, s_j, I) = [s_i \neq s_j] \cdot \frac{1}{dist(i, j)} \cdot e^{-\frac{|I_i - I_j|^2}{2\sigma^2}}$$



Horizontal link term



Vertical link term



Diagonal link term



## Results



Results showed with background replacement.

245  
Microsoft  
Research

## References

- Y. Boykov and V. Kolmogorov, "Computing geodesics and minimal surfaces via graph cuts", in *Proc. of ICCV*, 2003.
- A. Criminisi, G. Cross, A. Blake, and V. Kolmogorov, "Bilayer segmentation of live video," in *Proc. of CVPR*, 2006.
- C. Harrison and S. E. Hudson, "Pseudo-3d video conferencing with a generic webcam," in *Proc. of the 10th IEEE International Symposium on Multimedia*, 2008.
- V. Kolmogorov, A. Criminisi, A. Blake, G. Cross, and C. Rother, "Bilayer segmentation of binocular stereo video," in *Proc. of CVPR*, 2005.
- S. Pastoor, J. Liu, and S. Renault, "An experimental multimedia system allowing 3-d visualization and eye-controlled interaction without userworn devices," *IEEE Trans. on Multimedia*, vol. 1, no. 1, pp. 41–52, March 1999.
- T. Suenaga, Y. Matsumoto, and T. Ogawara, "3d display based on motion parallax using non-contact 3d measurement of head position," in *Proc. of the 17th Australasian conference on Computer-Human Interaction*, 2005.
- J. Sun, W. Zhang, X. Tang, and H. Y. Shum, "Background cut," in *Proc. of ECCV*, 2006.
- G. Welch and E. Foxlin, "Motion tracking: No silver bullet, but a respectable arsenal," *IEEE Computer Graphics and Applications*, vol. 22, no. 6, pp. 24–38, Nov.-Dec. 2002.
- T. Yu, C. Zhang, M. Cohen and Y. Rui, "Monocular Video Foreground/Background Segmentation by Tracking Spatial-Color Gaussian Mixture Models", *IEEE Workshop on Motion and Video Computing*, 2007.
- C. Zhang, Z. Yin and D. Florencio, "Improving Depth Perception with Motion Parallax and Its Application in Teleconferencing," *Proc. of MMSP*, 2009.

246  
Microsoft  
Research

## Outline

- Introduction
- Sound source localization with compact mic arrays
- Spatial sound and multi-channel echo cancellation
- Video-based tracking
- **Other video processing techniques**
  - Learning-based perceptual image quality enhancement
  - Monocular 3D and Foreground/background segmentation
  - Joint tracking and multiview video compression
- Conclusions and future work

247  
Microsoft  
Research

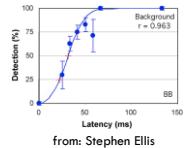
## Motion Parallax

- Improve 3D perception
  - More realistic 3D experience when combined with 3D displays
  - Works with legacy 2D displays

248  
Microsoft Research

## Motion Parallax Delay Sensitivity

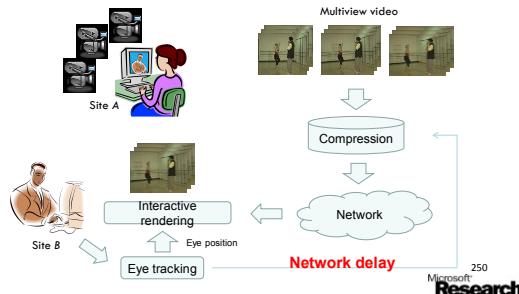
- Rendering video from exact viewer perspective is a must. This has impact on:
  - Realism
  - 3D clues
  - Comfort/sickness
- Motion parallax should be kept below **30 ms**.
- Motion Parallax based Rendering needs to be **LOCAL**.



249  
Microsoft Research

## Enable Motion Parallax for Multiview Video based 3D Teleconferencing

- Eye tracking + multiview video + view synthesis



250  
Microsoft Research

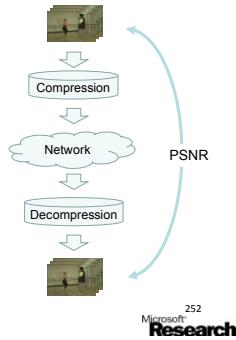
## Core Idea

- At the encoding site, predict the viewer's future viewing position
  - Encode 2 nearby frames (Kuruteppe 07)
  - Rate distortion optimization (Ramanathan 07)
- Encode each region of the views with different fidelity, *depending on the its use* (weight) on final (synthesized) image at the predicted viewing position.
- At the decoding site, decompress the multiview video, and render with the user's actual viewing position

251  
Microsoft  
Research

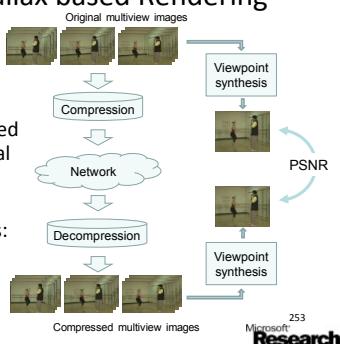
## Video Compression Error Criteria

- Traditionally:
  - Input: Video
  - Output: Video
  - Measure PSNR between original frames and the decompressed frames
  - What the user sees: the decompressed frames

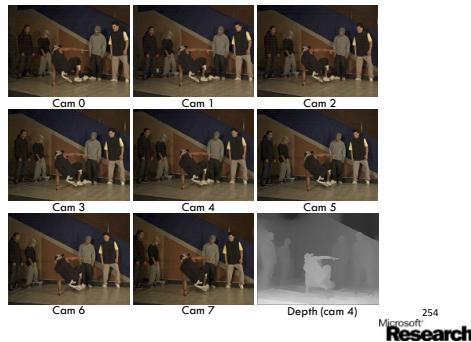


## Multiview Video Compression with Motion Parallax based Rendering

- Input: Video
- Output: Video
- Measure PSNR between synthesized views of the original and decompressed videos
- What the user sees: the **synthesized** frames

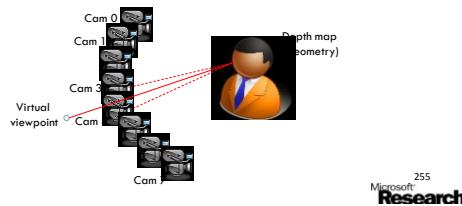


## An Example Multiview Frame

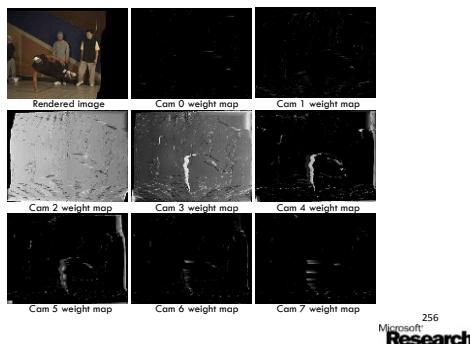


## Virtual View Synthesis

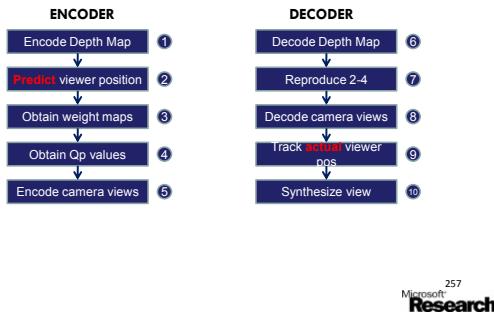
We project (each pixel of) the virtual view onto the depth map. Each point is then traced back to two nearest cameras for interpolation. Occlusion is explicitly computed. (Buehler 01)



## Weight Maps for Rendering

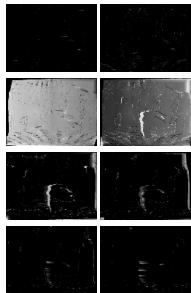


## Overall Diagram

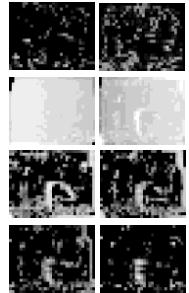


257  
Microsoft Research

## Weight Maps and Qp Maps



Weight maps



Qp maps

258  
Microsoft Research



Cam 2

Cam 3

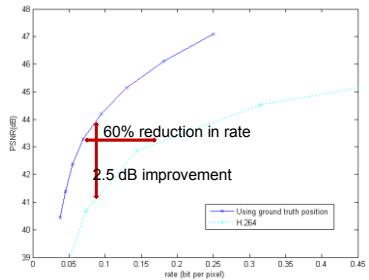
Cam 4

259  
Microsoft Research



260  
Microsoft  
Research

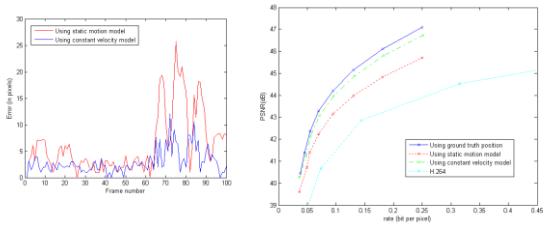
### Result with Perfect Viewer Positions



261  
Microsoft  
Research

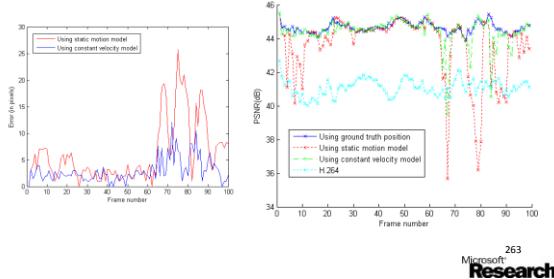
### What About Prediction Errors?

- Can cause significant performance drop



262  
Microsoft  
Research

## Frame-by-frame PSNR



263  
Microsoft  
Research

## Particle Filter (Condensation) based Tracking (Isard 98)

- A Posteriori Probability of the object status at time  $t$ :

$$p(\mathbf{x}_t | \mathbf{Z}_t) \propto p(\mathbf{z}_t | \mathbf{x}_t) p(\mathbf{x}_t | \mathbf{Z}_{t-1})$$

$\mathbf{x}_t$ : object status;  $\mathbf{Z}_t$ : input tracking video stream

where:

$$p(\mathbf{x}_t | \mathbf{Z}_{t-1}) = \int_{\mathbf{x}_{t-1}} p(\mathbf{x}_t | \mathbf{x}_{t-1}) p(\mathbf{x}_{t-1} | \mathbf{Z}_{t-1})$$

assuming object dynamics follows a temporal Markov chain.

264  
Microsoft  
Research

## Particle Filter (Condensation) based Tracking (Isard 98)

- A Posteriori Probability of the object status at time  $t$ :

$$p(\mathbf{x}_t | \mathbf{Z}_t) \propto p(\mathbf{z}_t | \mathbf{x}_t) p(\mathbf{x}_t | \mathbf{Z}_{t-1})$$

$\mathbf{x}_t$ : object status;  $\mathbf{Z}_t$ : input tracking video stream

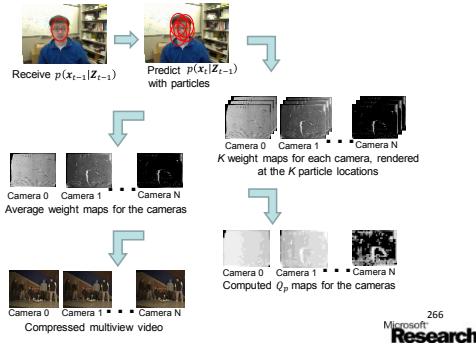
where: Best prediction at the encoder

$$p(\mathbf{x}_t | \mathbf{Z}_{t-1}) = \int_{\mathbf{x}_{t-1}} p(\mathbf{x}_t | \mathbf{x}_{t-1}) p(\mathbf{x}_{t-1} | \mathbf{Z}_{t-1})$$

assuming object dynamics follows a temporal Markov chain.

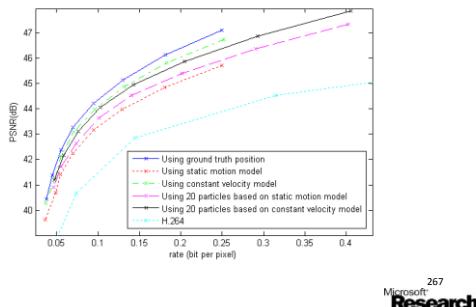
265  
Microsoft  
Research

## Joint Tracking and Compression



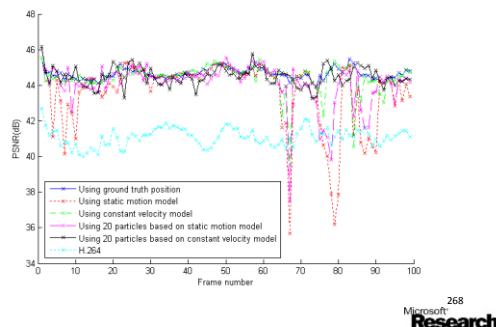
266  
Microsoft Research

## Results with Joint Tracking and Compression



267  
Microsoft Research

## Frame-by-frame PSNR



268  
Microsoft Research

## Frame-by-frame PSNR

Approach	Variance of PSNR (dB)
Using ground truth position	0.097
Using static motion model	3.819
Using constant velocity model	0.693
Using 20 particles based on static motion model	1.477
Using 20 particles based on constant velocity model	0.285
H.264	0.246

269  
Microsoft  
Research

## References

- C. Buehler, M. Bosse, L. McMillan, S. J. Gortler, and M. F. Cohen, "Unstructured lumigraph rendering," *ACM SIGGRAPH*, 2001.
- E. Kuruteppe, R. Civanlar and M. Tekalp, "Client-driven selective streaming of multiview video for interactive 3DTV", *IEEE Trans. On Circuits and Systems for Video Technology*, Vol. 17, No. 11, pp. 1558–1565, Nov. 2007.
- D. Florencio and C. Zhang, "Multiview video Compression and Streaming Based on Predicted Viewer Position," in *Proc. of ICASSP*, 2009.
- M. Isard and A. Blake, "CONDENSATION – conditional density propagation for visual tracking", *International Journal of Computer Vision*, vol. 29, no. 1, pp. 5–28, 1998.
- P. Ramanathan, M. Kalman and B. Girod, "Rate distortion optimized interactive light field streaming", *IEEE Trans. On Multimedia*, vol. 9, no. 4, pp. 813–825, June 2007.
- C. Zhang and D. Florencio, "Joint Tracking and Multiview Video Compression," in *Proc. of VCIP*, 2010.

270  
Microsoft  
Research

## Outline

- Introduction
- Sound source localization with compact mic arrays
- Spatial sound and multi-channel echo cancellation
- Real-time video processing
- Exploring depth sensors
- Conclusions and future work

271  
Microsoft  
Research

## Conclusions

- Compact mic array SSL and its further improvement with room models
- Audio/visual feature level fusion for speaker detection
- Headphone/loudspeaker audio spatialization
- Multi-channel echo cancellation with Kalman filtering
- Head pose/facial expression tracking
- Depth camera based skeleton tracking
- Learning based perceptual video quality enhancement
- Motion parallax in videoconferencing
- Joint tracking and multiview compression

<sup>272</sup>  
Microsoft  
Research

---

---

---

---

---

---

---

---

---

---

## Future Work

- Audio
  - Audio quality assessment
  - Personal HRTF measurement
  - Mic and speaker at unknown locations
  - Distributed mic arrays
  - ...
- Video
  - More reliable/accurate tracking, low computational cost
  - Understand human emotions
  - Depth image processing
  - 3D view synthesis, hybrid virtual/real environments
  - ...

<sup>273</sup>  
Microsoft  
Research

---

---

---

---

---

---

---

---

---

---

## Acknowledgement

- Our colleagues
  - Qin Cai, Wei-ge Chen, Ross Cutler, Philip Chou, Dinei Florencio, Rajesh Hedge, Zicheng Liu, Yong Rui, Jamie Shotton, Mike Sinclair, Jay Stokes, Jian Sun, Paul Viola, etc.
- Students
  - Demba Ba, Flávio Ribeiro, Myung-Suk Song, Zhaozheng Yin, Pei Yin, Aswin C. Sankaranarayanan, Qing Zhang, David Gallup, Mingxuan Sun, etc.
- Microsoft product groups
  - Microsoft Unified Communication Group, Microsoft Interactive Entertainment Division Kinect Team, etc.

<sup>274</sup>  
Microsoft  
Research

---

---

---

---

---

---

---

---

---

---

Thank you!

---

---

---

---

---

---

---

---