








# Index of /mlmi04/MLMI-Talk-010/ scannedSlides

<u>Name</u>	<u>Last modified</u>	<u>Size</u>	<u>Description</u>
 <a href="#">Parent Directory</a>	14-Oct-2004 17:35	-	
 <a href="#">p0010s001.jpg</a>	06-Oct-2004 11:06	332k	
 <a href="#">p0010s002.jpg</a>	06-Oct-2004 11:06	326k	
 <a href="#">p0010s003.jpg</a>	06-Oct-2004 11:06	378k	
 <a href="#">p0010s004.jpg</a>	06-Oct-2004 11:06	348k	
 <a href="#">p0010s005.jpg</a>	06-Oct-2004 11:06	404k	
 <a href="#">p0010s006.jpg</a>	06-Oct-2004 11:06	365k	
 <a href="#">p0010s007.jpg</a>	06-Oct-2004 11:06	384k	
 <a href="#">p0010s008.jpg</a>	06-Oct-2004 11:06	374k	
 <a href="#">p0010s009.jpg</a>	06-Oct-2004 11:06	312k	
 <a href="#">p0010s010.jpg</a>	06-Oct-2004 11:06	331k	
 <a href="#">p0010s011.jpg</a>	06-Oct-2004 11:06	322k	
 <a href="#">p0010s012.jpg</a>	06-Oct-2004 11:06	317k	
 <a href="#">p0010s013.jpg</a>	06-Oct-2004 11:06	301k	
 <a href="#">p0010s014.jpg</a>	06-Oct-2004 11:06	346k	
 <a href="#">p0010s015.jpg</a>	06-Oct-2004 11:06	321k	
 <a href="#">p0010s016.jpg</a>	06-Oct-2004 11:06	349k	
 <a href="#">p0010s017.jpg</a>	06-Oct-2004 11:06	318k	
 <a href="#">p0010s018.jpg</a>	06-Oct-2004 11:06	344k	
 <a href="#">p0010s019.jpg</a>	06-Oct-2004 11:06	305k	

	<a href="#">p0010s020.jpg</a>	06-Oct-2004 11:06	295k
	<a href="#">p0010s021.jpg</a>	06-Oct-2004 11:06	212k
	<a href="#">p0010s022.jpg</a>	06-Oct-2004 11:06	237k
	<a href="#">p0010s023.jpg</a>	06-Oct-2004 11:06	285k
	<a href="#">p0010s024.jpg</a>	06-Oct-2004 11:06	348k
	<a href="#">p0010s025.jpg</a>	06-Oct-2004 11:06	165k
	<a href="#">thumb/</a>	06-Oct-2004 11:40	-

---

Apache/1.3.31 Server at mmm.idiap.ch Port 80

# Summary

Conversational  
Telephone  
Speech  
(CTS)



Meetings



We present an overview of our effort to port the SRI CTS system to the Meeting task using:

- Adaptation (language and acoustic models)
- Preprocessing (noise reduction & array processing)
- Postprocessing (cross-talk elimination)

# The 2004 ICSI-SRI-UW Meeting Recognition System

International Computer Science Institute  
Berkeley, CA, USA



Speech Technology & Research Laboratory  
SRI International, Menlo Park, CA, USA



Signal, Speech & Language Processing Laboratory  
University of Washington, Seattle, WA, USA



# Cast of Characters

- ICSI** *Chuck Wooters* (data czar, segmentation, postprocessing)  
*Nikki Mirghafori* (text normalization, acoustic adaptation)  
*Tuomo Pirinen* (array processing)  
*David Gelbart* (noise filtering)  
*Barbara Peskin* (advising)
- SRI** *Andreas Stolcke* (acoustic modeling, system architecture)  
*Martin Graciarena* (POF experiments)  
*Jing Zheng* (structured MAP adaptation)  
*Ramana Gadde* (help with segmentation)  
*Wen Wang* (Fisher text processing, SuperARV LM)
- UW** *Ivan Bulyko* (language models)  
*Scott Otterson* (segmentation features)  
*Mari Ostendorf* (advising)

# “Map”

- Meeting Evals- History
- System architecture
- Baseline results
- Adaptation
- Preprocessing
- Postprocessing
- Summary and conclusions

# The “Meeting” Task

Conducted by the U.S. National Institute of Standards and Technology (NIST).

First NIST eval on Meeting Data held in 2002 (RT-02).

Second evaluation held this in the Spring of this year (RT-04).

	Recording Conditions	
	Close mics	Distant mics
RT-02	Headset Lapel	Single Distant Mic
RT-04	Individual Headset Mics (IHM) (no lapel)	Multiple Distant Mics (MDM)
		Single Distant Mic (SDM)

# The “Meeting” Task

Conducted by the U.S. National Institute of Standards and Technology (NIST).

First NIST eval on Meeting Data held in 2002 (RT-02).

Second evaluation held this in the Spring of this year (RT-04).

	Recording Conditions	
	Close mics	Distant mics
RT-02	Headset Lapel	Single Distant Mic
RT-04	Individual Headset Mics (IHM) (no lapel)	Multiple Distant Mics (MDM)
		Single Distant Mic (SDM)

*primary*

*primary*

*required contrast*



# Dev and Eval Data

Development data same as RT-02 Eval with updated transcripts.  
Evaluation data consists of 2 meetings from each site. 1-2 hours each.  
-Recognize 11 minute segments only.

		CMU	ICSI	LDC	NIST
RT-04 Eval Data	Personal	Head	Head	Head	Head
	Tabletop	1	6	4	7
Dev Data (a.k.a. RT-02 eval)	Personal	Lapel	Head	Lapel	Head
	Tabletop	1	6	10	7

- CMU has only 1 table-top mic, others have 4-10.
- RT-02 eval (with updated transcripts) set served as devtest.

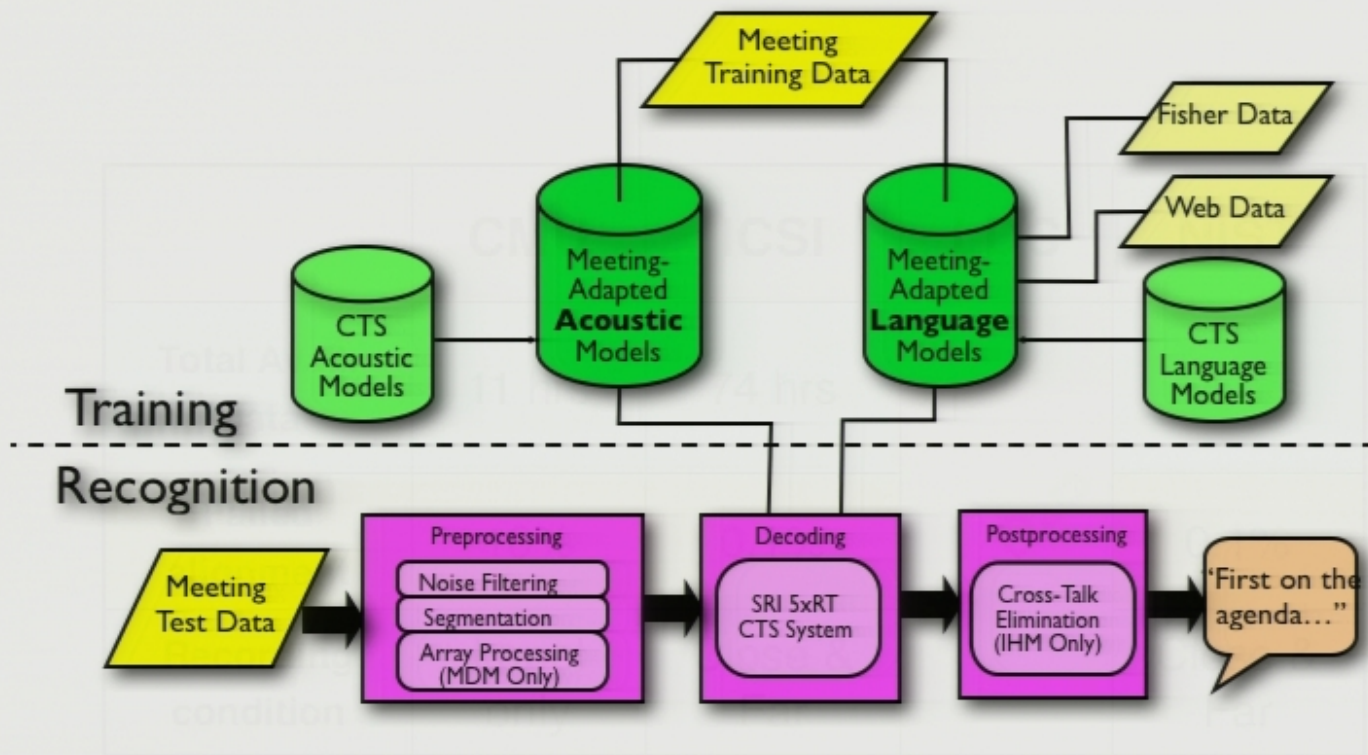
# Meeting Training Data

Data available for building the RT-04 Meeting recognition system.

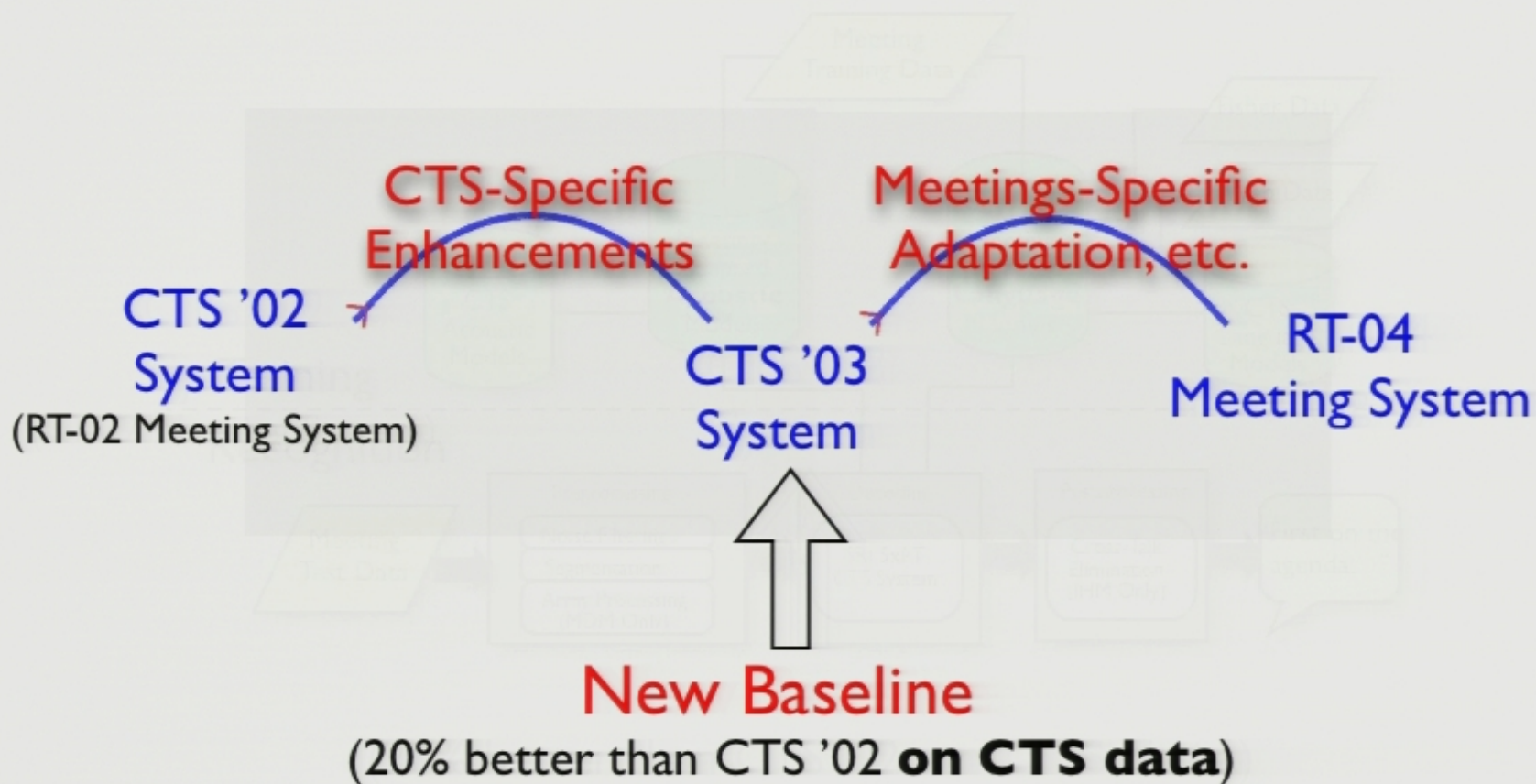
	CMU	ICSI	LDC	NIST
<b>Total Avail. Data</b>	11 hrs	74 hrs		14 hrs
<b>Failed alignment</b>	10%	0.1%	<i>No Data</i>	0.1%
<b>Recording condition</b>	Lapel only	Close & Far		Close & Far

# System Architecture

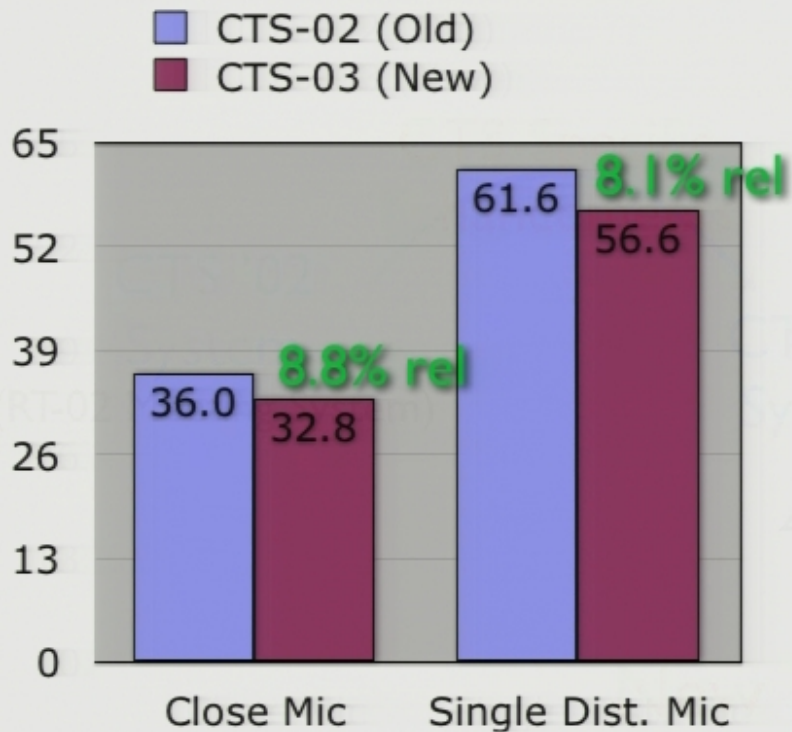
Data available for building the RT-04 Meeting recognition system.



# Baseline System Development



# Old Baseline vs. New Baseline on Meeting Data

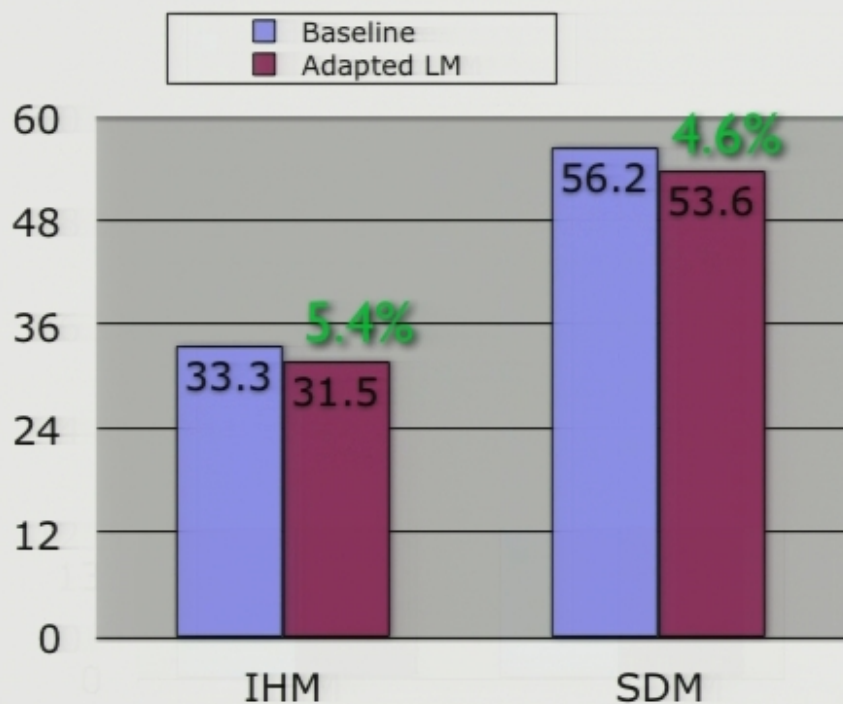


Not as much gain as on CTS (20%), but a nice improvement.

Error rate on RT-02 eval meeting data.

# Language Model Adaptation

- Started from CTS LM
- Added components for meeting domain
- Extended the Vocabulary



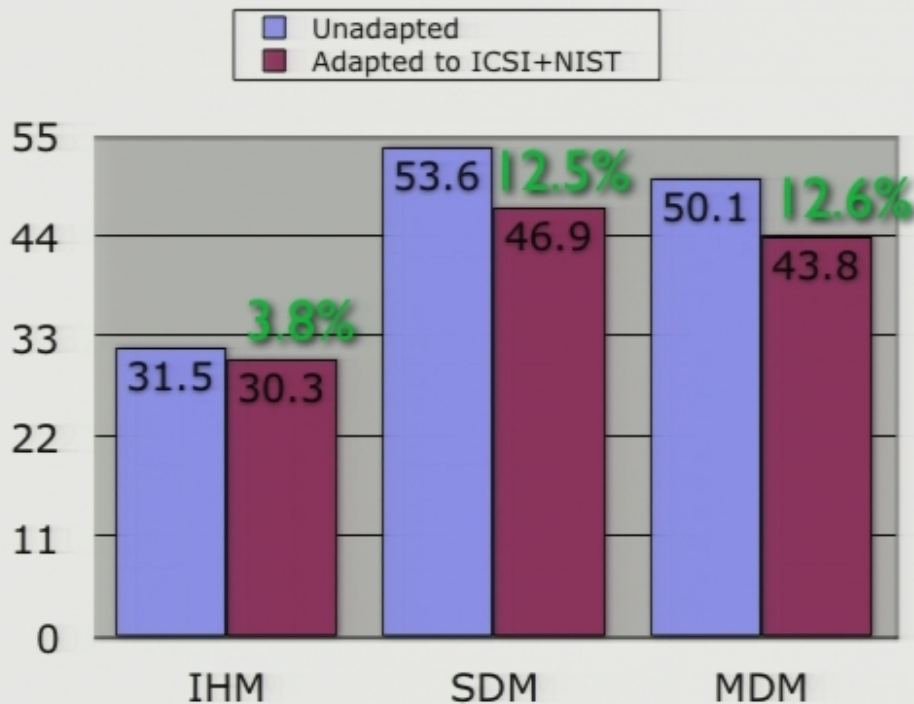
Roughly 5% improvement on both IHM and SDM.

Error rate on RT-04 dev data

Error rate on RT-02 eval meeting data.

# Acoustic Model Adaptation

- MAP-adapt acoustic models for each meeting source (using data from one or several sources)



Best to adapt to *all* distant mic channels

Captures acoustic variability due to both speaker and mic location (noise levels, reverberation)

Adapting to only central mic gives about 45-75% of possible gain.

Error rate on RT-04 dev data

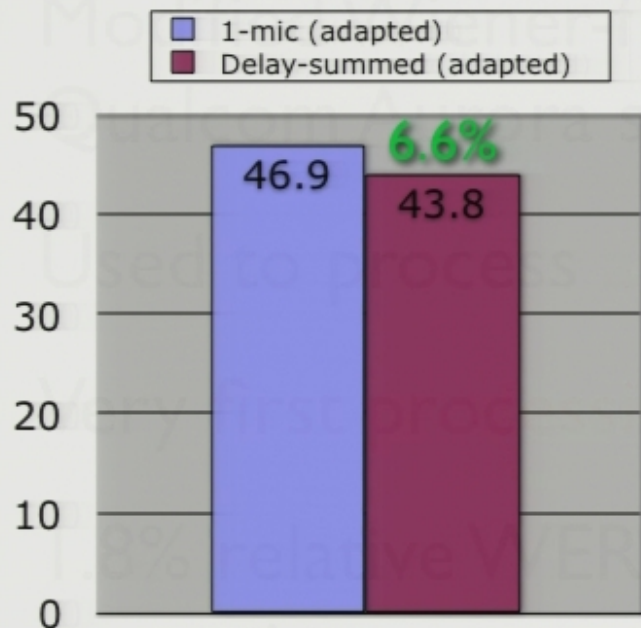
# Noise Filtering

- Modified Wiener-filtering developed for ICSI-OGI-  
Qualcom Aurora system
- Used to process **all distant mic channels**
- Very **first processing step**, on full meeting length
- 1.8% relative WER reduction with baseline  
recognizer



# Array Processing

- Idea: add all distant mic signals after shifting (shift to compensate for delays)
- Speech is enhanced
- Noise is phase delayed and partially cancelled in summing



Baseline used “central” or “best” single mic.

Acoustic adaptation: still better (5% relative) to adapt to all distant channels, NOT to delay-summed signal.

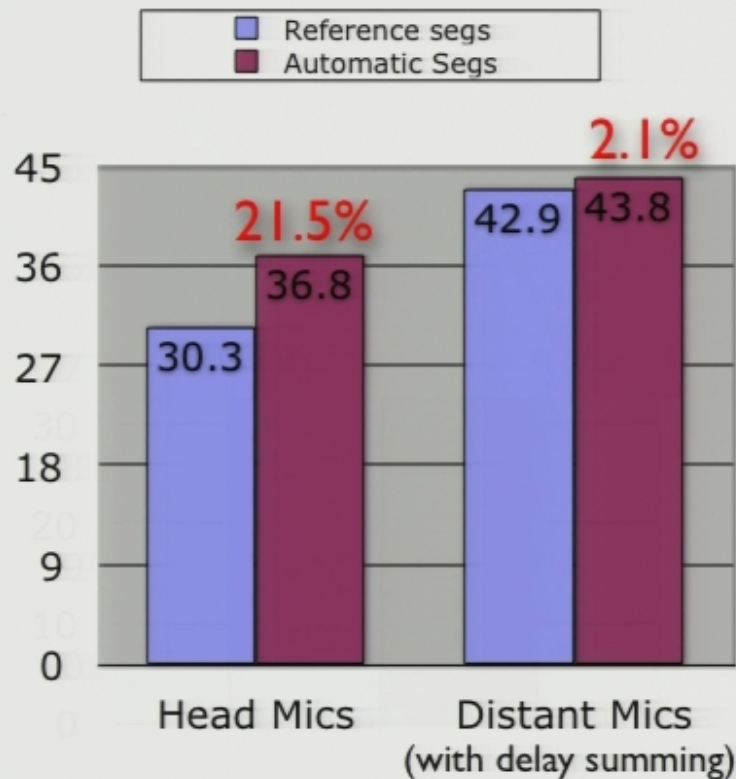
Run after segmentation

Works best AFTER noise filtering

Error rate on RT-04 dev data

# Segmentation

- Segmenting meeting into isolated utterances.
- How much are we hurt by the automatic segmentors?



Biggest problem: cross-talk in IHM

Tried to address with postprocessing

Also developing cross-channel features (UW)

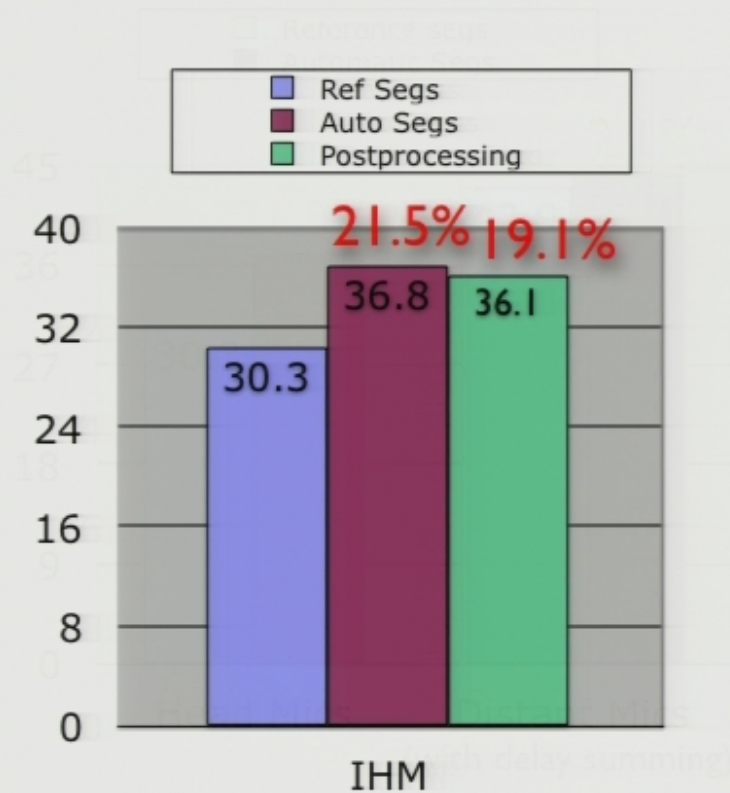
Run after segmentation

Works best AFTER noise filtering

## Error rate on RT-04 dev data

# Cross-talk Suppression

- Standard technique: delete words with low posteriors
- More sophisticated: look for words overlapping more than 50%, and delete those with low posteriors



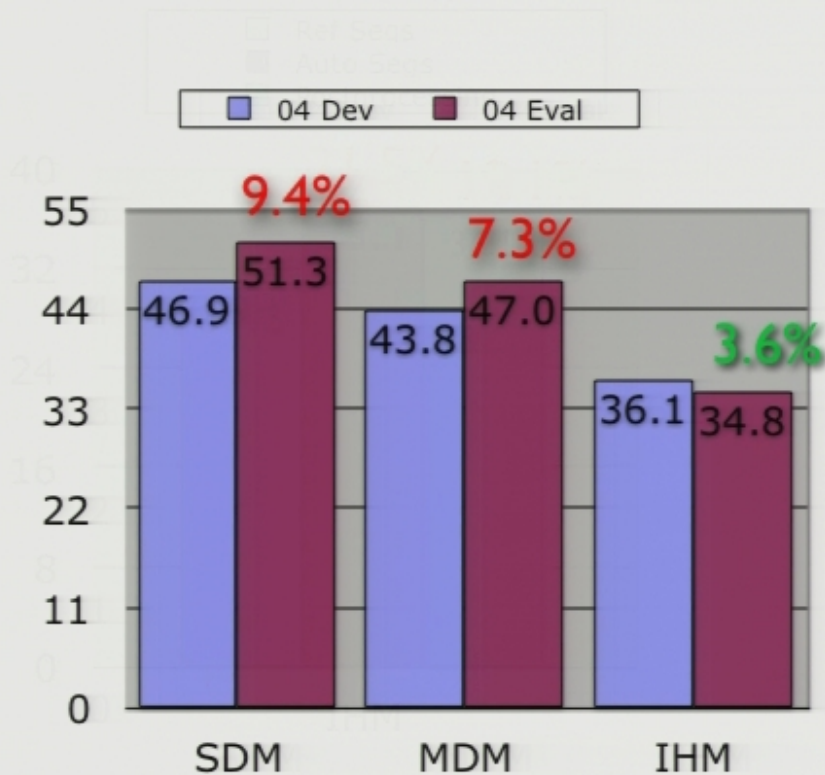
Biggest problem: cross-talk in IHM  
Still lots of room for improvement  
(features that look across channels?  
or confidence measures?)  
channel features (UVV)

Post eval diagnostics: helpful for lapel,  
not headmounted mics

Error rate on RT-04 dev data

# RT-04 Dev vs Eval

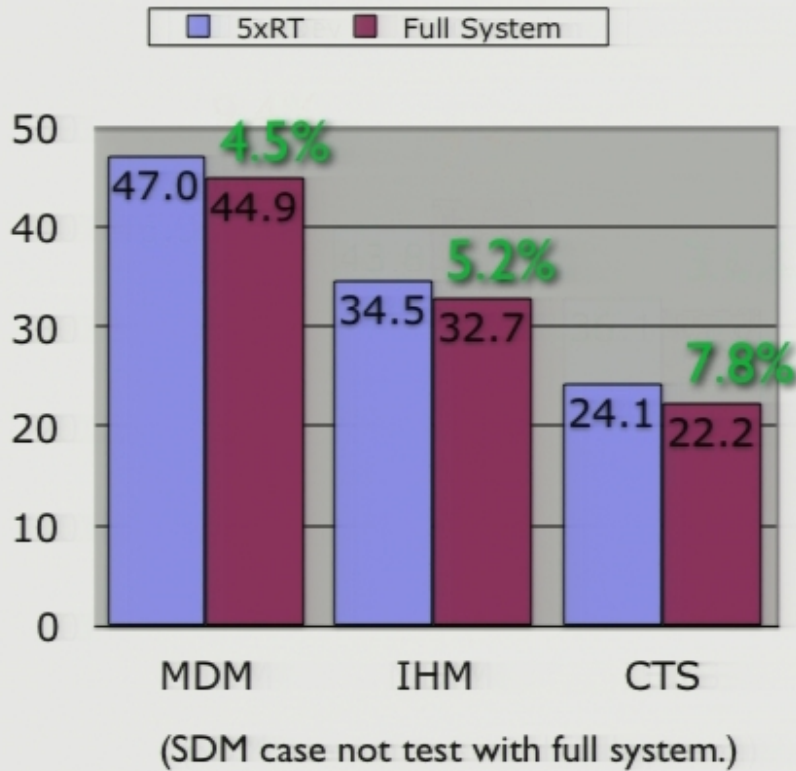
- How does the development data compare to the evaluation data?
  - More sophisticated, look for words overlapping more than 50%, and delete those with low posteriors



The more difficult sources (CMU and LDC) constituted a larger portion of the 04 Eval data.

# Full CTS System Results

- Full CTS eval system (20xRT) to compare to the evaluation data?
- Performance on RT-04 Eval data



Relative improvement on Meeting data is less than on CTS.

# Summary

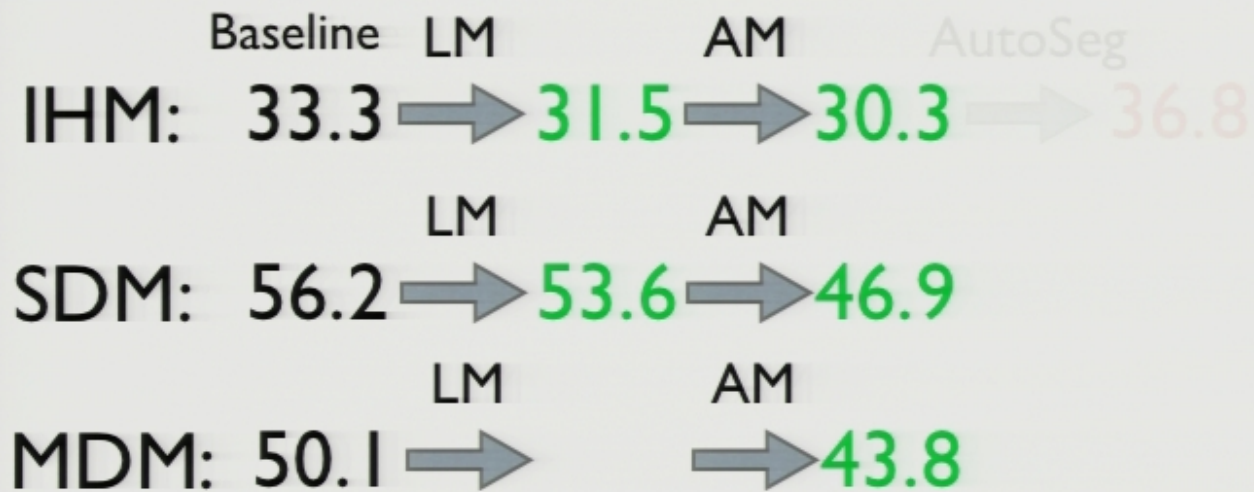
Baseline

IHM: 33.3

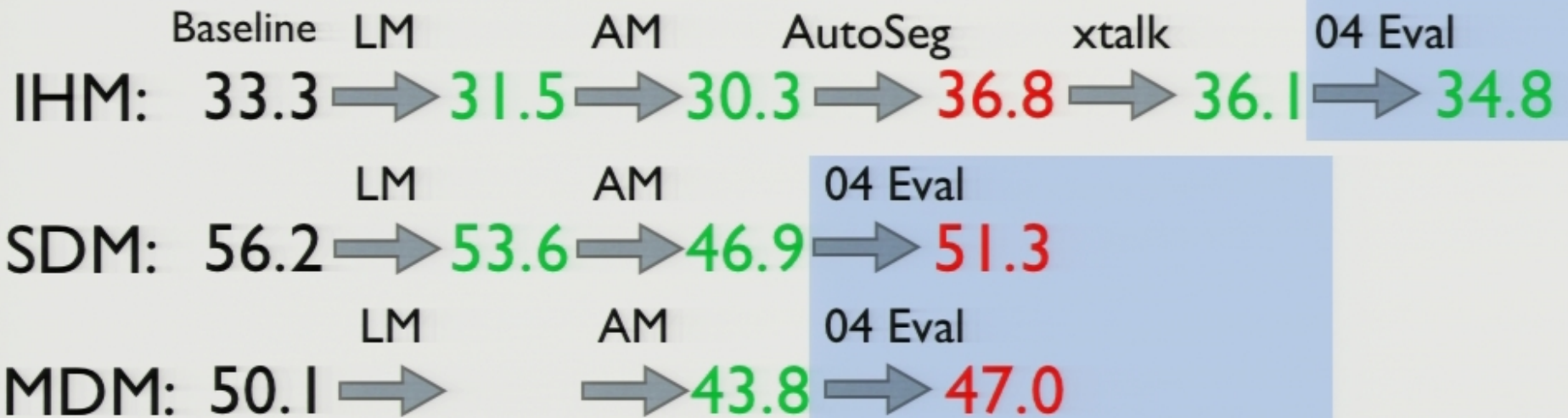
SDM: 56.2

MDM: 50.1

# Summary



# Summary





# Conclusions

## (Observations)

### **CTS system was relatively easy to adapt for the Meeting task**

- But: CTS improvements don't all carry over to Meetings
- Lots of challenges pre- and post-recognition
- AM adaptation very helpful (somewhat for individual, a lot for distant mics)
- Delay-sum array processing a big win
- Challenge for individual mics: speech detection (cross-talk)
- Challenge for distant mics: speaker tracking

File Edit Movie Favorites Window Help

