

# Detection of Laughter-in-Interaction in Multichannel Close-Talk Microphone Recordings of Meetings

Kornel Laskowski and Tanja Schultz

Cognitive Systems Lab, Universität Karlsruhe  
Language Technologies Institute, Carnegie Mellon University

09 September, 2008

# Why Detect Laughter in Meetings?

- evidence suggests that it is the most frequently occurring and most robust behavior which external observers associate with perceived emotion
  - marked valence
  - marked activation
- ① automatic emotion recognition in meetings
  - enable indexing, search and summary, mediated by para-propositional content
  - also necessary for autonomous machine participation
- ② detection and tracking of humour and seriousness

# Why Detect Laughter in Meetings?

- evidence suggests that it is the most frequently occurring and most robust behavior which external observers associate with perceived emotion
  - marked valence
  - marked activation
- ① automatic emotion recognition in meetings
  - enable indexing, search and summary, mediated by para-propositional content
  - also necessary for autonomous machine participation
- ② detection and tracking of humour and seriousness

# Why Detect Laughter in Meetings?

- evidence suggests that it is the most frequently occurring and most robust behavior which external observers associate with perceived emotion
  - marked valence
  - marked activation
- ① automatic emotion recognition in meetings
  - enable indexing, search and summary, mediated by para-propositional content
  - also necessary for autonomous machine participation
- ② detection and tracking of humour and seriousness

# Why Detect Laughter in Meetings?

- evidence suggests that it is the most frequently occurring and most robust behavior which external observers associate with perceived emotion
  - marked valence
  - marked activation
- ① automatic emotion recognition in meetings
  - enable indexing, search and summary, mediated by para-propositional content
  - also necessary for autonomous machine participation
- ② detection and tracking of humour and seriousness

# Why Detect Laughter in Meetings?

- evidence suggests that it is the most frequently occurring and most robust behavior which external observers associate with perceived emotion
  - marked valence
  - marked activation
- ① automatic emotion recognition in meetings
  - enable indexing, search and summary, mediated by para-propositional content
  - also necessary for autonomous machine participation
- ② detection and tracking of humour and seriousness

# Classifying Emotional Valence

- data: ISL Meeting Corpus (Burger et al, 2002)
- annotation: perceived valence (Laskowski & Burger, 2006)
- task: classify **segmented utterances** as exhibiting one of {negative, neutral, positive}

Classification	Accuracy, %
	EVAL
guessing with uniform prior	33.3
guessing with TRAINSET prior	≈67
guessing majority TRAINSET class	≈81
presence of $\mathcal{L}$ only	≈ <b>92</b>
prosody features	≈84
all features (except presence of $\mathcal{L}$ )	≈87

# Classifying Emotional Activation

- also known as emotional *arousal*
- data: ICSI Meeting Corpus (Janin et al, 2003)
- annotation: hotspots (Wrede & Shriberg, 2004; Wrede et al, 2005)
- task: classify **60-second intervals** as one of  
 $\{\text{involvementContaining}, \neg\text{involvementContaining}\}$

Classification	Accuracy, %		
	TRAIN	DEV	EVAL
guessing with uniform prior	50.0	50.0	50.0
guessing with TRAINSET prior	61.3	60.9	61.2
guessing majority TRAINSET class	73.7	72.9	73.7
features from $\mathcal{L}$ only	79.2	80.0	<b>80.6</b>
features from $\mathcal{L}$ and $\mathcal{S}$	84.3	82.7	<b>83.0</b>



# Goals of this Work

## Detection of Laughter-in-Interaction in Multichannel Close-Talk Microphone Recordings of Meetings

- ① propose a framework for detecting laughter from audio only
  - close-talk microphones on all participants
- ② attempt to detect **all** laughter
  - ① temporally isolated from the laugher's speech
  - ② occurring within dialog acts among verbal productions
- ③ attempt to detect **without prior knowledge**
  - no inactive channel exclusion
  - expect to encounter many false alarms
- ④ attribute laughter to specific participants

# Goals of this Work

## Detection of Laughter-in-Interaction in Multichannel Close-Talk Microphone Recordings of Meetings

- ① propose a framework for detecting laughter from audio only
  - close-talk microphones on all participants
- ② attempt to detect **all** laughter
  - ① temporally isolated from the laughter's speech
  - ② occurring within dialog acts among verbal productions
- ③ attempt to detect **without prior knowledge**
  - no inactive channel exclusion
  - expect to encounter many false alarms
- ④ attribute laughter to specific participants

# Goals of this Work

## Detection of **Laughter-in-Interaction** in Multichannel Close-Talk Microphone Recordings of Meetings

- ➊ propose a framework for detecting laughter from audio only
  - close-talk microphones on all participants
- ➋ attempt to detect **all** laughter
  - ➊ temporally isolated from the laugher's speech
  - ➋ occurring within dialog acts among verbal productions
- ➌ attempt to detect **without prior knowledge**
  - no inactive channel exclusion
  - expect to encounter many false alarms
- ➍ attribute laughter to specific participants

# Goals of this Work

## Detection of Laughter-in-Interaction in Multichannel Close-Talk Microphone Recordings of Meetings

- ➊ propose a framework for detecting laughter from audio only
  - close-talk microphones on all participants
- ➋ attempt to detect **all** laughter
  - ➊ temporally isolated from the laugher's speech
  - ➋ occurring within dialog acts among verbal productions
- ➌ attempt to detect **without prior knowledge**
  - no inactive channel exclusion
  - expect to encounter many false alarms
- ➍ attribute laughter to specific participants

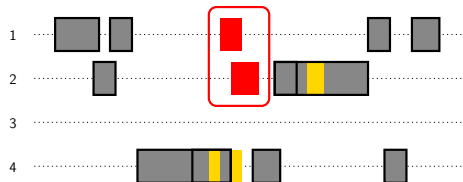
# Goals of this Work

## Detection of Laughter-in-Interaction in Multichannel Close-Talk Microphone Recordings of Meetings

- ① propose a framework for detecting laughter from audio only
  - close-talk microphones on all participants
- ② attempt to detect **all** laughter
  - ① temporally isolated from the laughter's speech
  - ② occurring within dialog acts among verbal productions
- ③ attempt to detect **without prior knowledge**
  - no inactive channel exclusion
  - expect to encounter many false alarms
- ④ attribute laughter to specific participants



# Detecting All Laughter from All Audio



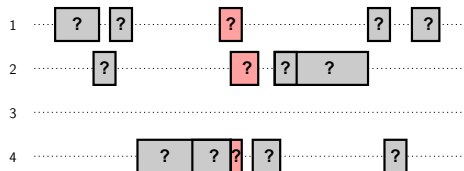
Past work has focused on:

- a subset of laughter (improving recall)
  - isolated laughter
    - loud, clear, unambiguous laughter
- and/or a subset of audio (improving precision)
  - segmented intervals
  - 100ms





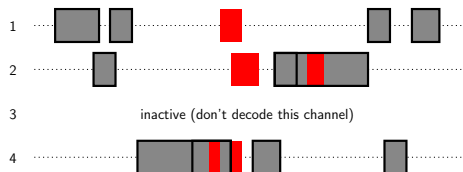
# Detecting All Laughter from All Audio



Past work has focused on:

- a subset of laughter (improving recall)
  - isolated laughter
  - loud, clear, unambiguous laughter
- and/or a subset of audio (improving precision)
  - segmented intervals
  - only active channels

# Detecting All Laughter from All Audio



Past work has focused on:

- a subset of laughter (improving recall)
  - isolated laughter
  - loud, clear, unambiguous laughter
- and/or a subset of audio (improving precision)
  - segmented intervals
  - only active channels

# Brief Comparison with Related Work

Aspect	$\mathcal{L}/\mathcal{S}$ class.		$\mathcal{L}/\neg\mathcal{L}$ segm.			this work
	[1]	[2]	[3]	[4]	[5]	
close-talk microphones	✓	✓	✓	✓		✓
farfield microphones					✓	
single channel at-a-time	✓	✓	✓	✓		
multi-channel at-a-time					✓	✓
participant attribution	✓	✓	✓	✓		✓
only group laughter					✓	
only isolated laughter	✓		✓	✓		
only clear laughter		✓				
rely on pre-segmentation	✓	✓	?			
rely on channel exclusion			?	✓		

[1] (Truong & van Leeuwen, 2005); [2] (Truong & van Leeuwen, 2007a); [3] (Truong & van Leeuwen, 2007b); [4] (Knox & Mirghafori, 2007); [5] (Kennedy & Ellis, 2004).

# Outline of this Talk

1. Introduction (about to be over)
2. Data
3. Multiparticipant 3-state Vocal Activity Detector
4. Experiments
5. Analysis
6. Conclusions (& Unqualified Recommendations)

# ICSI Meeting Corpus

- the complete corpus (Janin et al, 2003)
  - 75 naturally occurring meetings
  - longitudinal CTM recordings of several work groups
  - 3-9 instrumented participants per meeting
- we use a subset of 67 meetings
  - types: Bed (15), Bmr (29), Bro (23)
  - 23 unique participants
  - 3 participants attend both Bmr and Bro
  - 1 participant attends both Bmr and Bed
- in particular, as elsewhere,
  - TRAINSET: 26 Bmr meetings
  - TESTSET: 3 Bmr meetings

# ICSI Meeting Corpus

- the complete corpus (Janin et al, 2003)
  - 75 naturally occurring meetings
  - longitudinal CTM recordings of several work groups
  - 3-9 instrumented participants per meeting
- we use a subset of 67 meetings
  - types: Bed (15), Bmr (29), Bro (23)
  - 23 unique participants
  - 3 participants attend both Bmr and Bro
  - 1 participant attends both Bmr and Bed
- in particular, as elsewhere,
  - TRAINSET: 26 Bmr meetings
  - TESTSET: 3 Bmr meetings

# ICSI Meeting Corpus

- the complete corpus (Janin et al, 2003)
  - 75 naturally occurring meetings
  - longitudinal CTM recordings of several work groups
  - 3-9 instrumented participants per meeting
- we use a subset of 67 meetings
  - types: Bed (15), Bmr (29), Bro (23)
  - 23 unique participants
  - 3 participants attend both Bmr and Bro
  - 1 participant attends both Bmr and Bed
- in particular, as elsewhere,
  - TRAINSET: 26 Bmr meetings
  - TESTSET: 3 Bmr meetings

# Reference Segmentation

- speech,  $\mathcal{S}$ 
  - forced alignment of words and word fragments
  - available in the ICSI MRDA Corpus (Shriberg et al, 2004)
  - bridge inter-lexeme gaps shorter than 300 ms
  - as in NIST Rich Transcription Meeting Recognition evaluations
- laughter,  $\mathcal{L}$ 
  - produced semi-automatically (Laskowski & Burger, 2007d)
  - $\geq 99\%$  of laughter markup, as originally transcribed
  - bouts include terminal “recovery” in-/exhalation, if present
  - augmented with voicing classification,  $\mathcal{L} \equiv \mathcal{L}_V \cup \mathcal{L}_U$
- “laughed speech” (Nwokah et al, 1999),  $\mathcal{S} \cap \mathcal{L}$ 
  - here, mapped to laughter  $\mathcal{L}$
  - each participant can be producing  $\mathcal{L}$ ,  $\mathcal{S}$ , or neither



# Reference Segmentation

- speech,  $\mathcal{S}$ 
  - forced alignment of words and word fragments
  - available in the ICSI MRDA Corpus (Shriberg et al, 2004)
  - bridge inter-lexeme gaps shorter than 300 ms
  - as in NIST Rich Transcription Meeting Recognition evaluations
- laughter,  $\mathcal{L}$ 
  - produced semi-automatically (Laskowski & Burger, 2007d)
  - $\geq 99\%$  of laughter markup, as originally transcribed
  - bouts include terminal “recovery” in-/exhalation, if present
  - augmented with voicing classification,  $\mathcal{L} \equiv \mathcal{L}_V \cup \mathcal{L}_U$
- “laughed speech” (Nwokah et al, 1999),  $\mathcal{S} \cap \mathcal{L}$ 
  - here, mapped to laughter  $\mathcal{L}$
  - each participant can be producing  $\mathcal{L}$ ,  $\mathcal{S}$ , or neither

# Reference Segmentation

- speech,  $\mathcal{S}$ 
  - forced alignment of words and word fragments
  - available in the ICSI MRDA Corpus (Shriberg et al, 2004)
  - bridge inter-lexeme gaps shorter than 300 ms
  - as in NIST Rich Transcription Meeting Recognition evaluations
- laughter,  $\mathcal{L}$ 
  - produced semi-automatically (Laskowski & Burger, 2007d)
  - $\geq 99\%$  of laughter markup, as originally transcribed
  - bouts include terminal “recovery” in-/exhalation, if present
  - augmented with voicing classification,  $\mathcal{L} \equiv \mathcal{L}_V \cup \mathcal{L}_U$
- “laughed speech” (Nwokah et al, 1999),  $\mathcal{S} \cap \mathcal{L}$ 
  - here, mapped to laughter  $\mathcal{L}$
  - each participant can be producing  $\mathcal{L}$ ,  $\mathcal{S}$ , or neither

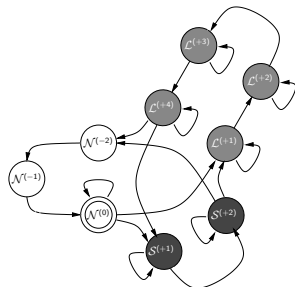
# Reference Segmentation

- speech,  $\mathcal{S}$ 
  - forced alignment of words and word fragments
  - available in the ICSI MRDA Corpus (Shriberg et al, 2004)
  - bridge inter-lexeme gaps shorter than 300 ms
  - as in NIST Rich Transcription Meeting Recognition evaluations
- laughter,  $\mathcal{L}$ 
  - produced semi-automatically (Laskowski & Burger, 2007d)
  - $\geq 99\%$  of laughter markup, as originally transcribed
  - bouts include terminal “recovery” in-/exhalation, if present
  - augmented with voicing classification,  $\mathcal{L} \equiv \mathcal{L}_V \cup \mathcal{L}_U$
- “laughed speech” (Nwokah et al, 1999),  $\mathcal{S} \cap \mathcal{L}$ 
  - here, mapped to laughter  $\mathcal{L}$
  - each participant can be producing  $\mathcal{L}$ ,  $\mathcal{S}$ , or neither

# Multiparticipant 3-state Vocal Activity Detector

- hidden Markov model
- pruned Viterbi (beam) decoding
- topology
  - single participant state subspace
  - multiparticipant state space, pruning
- multiparticipant transition probability model
- standard MFCC features, plus crosstalk suppression features
- multiparticipant emission probability model

# Single Participant (SP) State Subspace



- each participant can be
  - speaking,  $\mathcal{S}$
  - laughing,  $\mathcal{L}$
  - silent,  $\mathcal{N}$

- frame step  $\Delta T = 0.1$  s

- explicit minimum duration constraints

$$\begin{aligned} \mathbf{T}_{min} &\equiv (T_{min}^{\mathcal{S}}, T_{min}^{\mathcal{L}}, T_{min}^{\mathcal{N}}) \\ &= \Delta T \cdot (N_{min}^{\mathcal{S}}, N_{min}^{\mathcal{L}}, N_{min}^{\mathcal{N}}) \end{aligned}$$

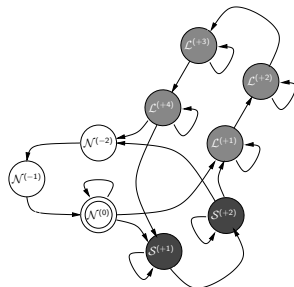
- number of states in 1-participant subspace

$$N = N_{min}^{\mathcal{S}} + N_{min}^{\mathcal{L}} + N_{min}^{\mathcal{N}}$$

- in example shown,  $N = 9$



# Single Participant (SP) State Subspace



- each participant can be
  - speaking,  $\mathcal{S}$
  - laughing,  $\mathcal{L}$
  - silent,  $\mathcal{N}$

- frame step  $\Delta T = 0.1$  s
- explicit minimum duration constraints

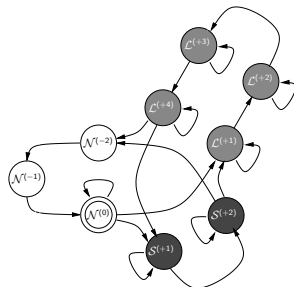
$$\begin{aligned} \mathbf{T}_{min} &\equiv (T_{min}^{\mathcal{S}}, T_{min}^{\mathcal{L}}, T_{min}^{\mathcal{N}}) \\ &= \Delta T \cdot (N_{min}^{\mathcal{S}}, N_{min}^{\mathcal{L}}, N_{min}^{\mathcal{N}}) \end{aligned}$$

- number of states in 1-participant subspace

$$N = N_{min}^{\mathcal{S}} + N_{min}^{\mathcal{L}} + N_{min}^{\mathcal{N}}$$

- in example shown,  $N = 9$

# Single Participant (SP) State Subspace



- each participant can be
  - speaking,  $\mathcal{S}$
  - laughing,  $\mathcal{L}$
  - silent,  $\mathcal{N}$

- frame step  $\Delta T = 0.1$  s
- explicit minimum duration constraints

$$\begin{aligned} \mathbf{T}_{min} &\equiv (T_{min}^{\mathcal{S}}, T_{min}^{\mathcal{L}}, T_{min}^{\mathcal{N}}) \\ &= \Delta T \cdot (N_{min}^{\mathcal{S}}, N_{min}^{\mathcal{L}}, N_{min}^{\mathcal{N}}) \end{aligned}$$

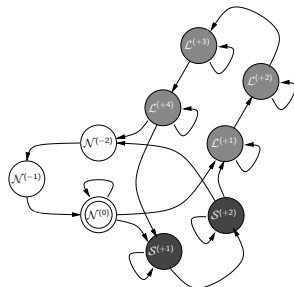
- number of states in 1-participant subspace

$$N = N_{min}^{\mathcal{S}} + N_{min}^{\mathcal{L}} + N_{min}^{\mathcal{N}}$$

- in example shown,  $N = 9$



# Single Participant (SP) State Subspace



- each participant can be
  - speaking,  $\mathcal{S}$
  - laughing,  $\mathcal{L}$
  - silent,  $\mathcal{N}$

- frame step  $\Delta T = 0.1$  s
- explicit minimum duration constraints

$$\begin{aligned} \mathbf{T}_{min} &\equiv (T_{min}^{\mathcal{S}}, T_{min}^{\mathcal{L}}, T_{min}^{\mathcal{N}}) \\ &= \Delta T \cdot (N_{min}^{\mathcal{S}}, N_{min}^{\mathcal{L}}, N_{min}^{\mathcal{N}}) \end{aligned}$$

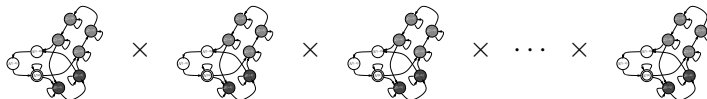
- number of states in 1-participant subspace

$$N = N_{min}^{\mathcal{S}} + N_{min}^{\mathcal{L}} + N_{min}^{\mathcal{N}}$$

- in example shown,  $N = 9$

# Multiparticipant (MP) State Space

- for a conversation of  $K$  participants,
- form the Cartesian product of  $K$  factors:

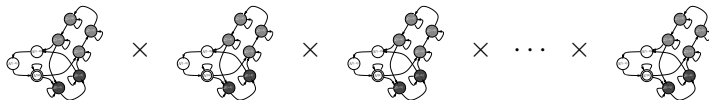


- each MP state:  $K$ -vector of  $N$  SP states
- total number of MP states in topology:  $N^K$
- impose maximum simultaneous vocalization constraints

$$\mathbf{K}_{max} = (K_{max}^S, K_{max}^L, K_{max}^N)$$

# Multiparticipant (MP) State Space

- for a conversation of  $K$  participants,
- form the Cartesian product of  $K$  factors:



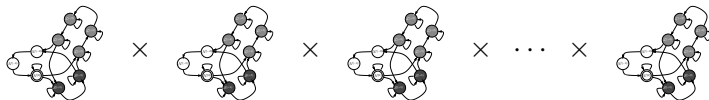
- each MP state:  $K$ -vector of  $N$  SP states
- total number of MP states in topology:  $N^K$
- impose maximum simultaneous vocalization constraints

$$\mathbf{K}_{max} = (K_{max}^S, K_{max}^L, K_{max}^N)$$

- ie.  $K_{max}^L$  max. # participants laughing at the same time

# Multiparticipant (MP) State Space

- for a conversation of  $K$  participants,
- form the Cartesian product of  $K$  factors:



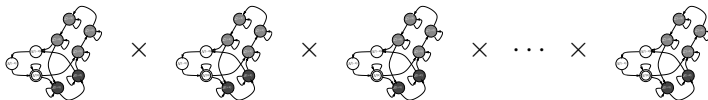
- each MP state:  $K$ -vector of  $N$  SP states
- total number of MP states in topology:  $N^K$
- impose maximum simultaneous vocalization constraints

$$\mathbf{K}_{max} = (K_{max}^S, K_{max}^L, K_{max}^{\neg N})$$

• i.e.  $K_{max}^L$  max. # participants laughing at the same time

# Multiparticipant (MP) State Space

- for a conversation of  $K$  participants,
- form the Cartesian product of  $K$  factors:



- each MP state:  $K$ -vector of  $N$  SP states
- total number of MP states in topology:  $N^K$
- impose maximum simultaneous vocalization constraints

$$\mathbf{K}_{max} = (K_{max}^S, K_{max}^L, K_{max}^{\neg N})$$

- ie.  $K_{max}^L$ : max. # participants laughing at the same time

# Transition Probability Model

- Example,  $K = 5$ :

- at time  $t$ ,  $\mathbf{q}_t = \mathbf{S}_i = [\mathcal{S}^{(2)}, \mathcal{N}^{(0)}, \mathcal{N}^{(0)}, \mathcal{N}^{(0)}]$
- at time  $t + 1$ ,  $\mathbf{q}_{t+1} = \mathbf{S}_j = [\mathcal{N}^{(-2)}, \mathcal{N}^{(0)}, \mathcal{S}^{(1)}, \mathcal{L}^{(1)}]$
- what is  $a_{ij} = P(\mathbf{q}_{t+1} = \mathbf{S}_j | \mathbf{q}_t = \mathbf{S}_i)$  ?

- 1  $a_{ij} = 0$  if the SP transition from  $\mathbf{S}_i$  to  $\mathbf{S}_j$  for any participant is not licensed by the SP topology
- 2 otherwise, ML estimate using ngram counts from best flat-start Viterbi path over training corpus
- 3 NOTE: each participant's index  $k$  in  $\mathbf{S}$  is arbitrary
  - for all  $K$ -symbol permutations/rotations  $\mathbf{R}$
  - want  $P(\mathbf{S}_j | \mathbf{S}_i) \equiv P(\mathbf{R} \cdot \mathbf{S}_j | \mathbf{R} \cdot \mathbf{S}_i)$
  - during model training & querying, rotate each  $\mathbf{q}_t$  into a fixed ordering of the  $N$  single-participant states

# Transition Probability Model

- Example,  $K = 5$ :

- at time  $t$ ,  $\mathbf{q}_t = \mathbf{S}_i = [\mathcal{S}^{(2)}, \mathcal{N}^{(0)}, \mathcal{N}^{(0)}, \mathcal{N}^{(0)}]$
- at time  $t + 1$ ,  $\mathbf{q}_{t+1} = \mathbf{S}_j = [\mathcal{N}^{(-2)}, \mathcal{N}^{(0)}, \mathcal{S}^{(1)}, \mathcal{L}^{(1)}]$
- what is  $a_{ij} = P(\mathbf{q}_{t+1} = \mathbf{S}_j | \mathbf{q}_t = \mathbf{S}_i)$  ?

- 1  $a_{ij} = 0$  if the SP transition from  $\mathbf{S}_i$  to  $\mathbf{S}_j$  for any participant is not licensed by the SP topology
- 2 otherwise, ML estimate using ngram counts from best flat-start Viterbi path over training corpus
- 3 NOTE: each participant's index  $k$  in  $\mathbf{S}$  is arbitrary
  - for all  $K$ -symbol permutations/rotations  $\mathbf{R}$
  - want  $P(\mathbf{S}_j | \mathbf{S}_i) \equiv P(\mathbf{R} \cdot \mathbf{S}_j | \mathbf{R} \cdot \mathbf{S}_i)$
  - during model training & querying, rotate each  $\mathbf{q}_t$  into a fixed ordering of the  $N$  single-participant states

# Transition Probability Model

- Example,  $K = 5$ :

- at time  $t$ ,  $\mathbf{q}_t = \mathbf{S}_i = [\mathcal{S}^{(2)}, \mathcal{N}^{(0)}, \mathcal{N}^{(0)}, \mathcal{N}^{(0)}]$
- at time  $t + 1$ ,  $\mathbf{q}_{t+1} = \mathbf{S}_j = [\mathcal{N}^{(-2)}, \mathcal{N}^{(0)}, \mathcal{S}^{(1)}, \mathcal{L}^{(1)}]$
- what is  $a_{ij} = P(\mathbf{q}_{t+1} = \mathbf{S}_j | \mathbf{q}_t = \mathbf{S}_i)$  ?

- 1  $a_{ij} = 0$  if the SP transition from  $\mathbf{S}_i$  to  $\mathbf{S}_j$  for any participant is not licensed by the SP topology
- 2 otherwise, ML estimate using ngram counts from best flat-start Viterbi path over training corpus
- 3 NOTE: each participant's index  $k$  in  $\mathbf{S}$  is arbitrary
  - for all  $K$ -symbol permutations/rotations  $\mathbf{R}$
  - want  $P(\mathbf{S}_j | \mathbf{S}_i) \equiv P(\mathbf{R} \cdot \mathbf{S}_j | \mathbf{R} \cdot \mathbf{S}_i)$
  - during model training & querying, rotate each  $\mathbf{q}_t$  into a fixed ordering of the  $N$  single-participant states



# Transition Probability Model

- Example,  $K = 5$ :
  - at time  $t$ ,  $\mathbf{q}_t = \mathbf{S}_i = [\mathcal{S}^{(2)}, \mathcal{N}^{(0)}, \mathcal{N}^{(0)}, \mathcal{N}^{(0)}]$
  - at time  $t + 1$ ,  $\mathbf{q}_{t+1} = \mathbf{S}_j = [\mathcal{N}^{(-2)}, \mathcal{N}^{(0)}, \mathcal{S}^{(1)}, \mathcal{L}^{(1)}]$
  - what is  $a_{ij} = P(\mathbf{q}_{t+1} = \mathbf{S}_j | \mathbf{q}_t = \mathbf{S}_i)$  ?
- ①  $a_{ij} = 0$  if the SP transition from  $\mathbf{S}_i$  to  $\mathbf{S}_j$  for any participant is not licensed by the SP topology
- ② otherwise, ML estimate using ngram counts from best flat-start Viterbi path over training corpus
- ③ NOTE: each participant's index  $k$  in  $\mathbf{S}$  is arbitrary
  - for all  $K$ -symbol permutations/rotations  $\mathbf{R}$
  - want  $P(\mathbf{S}_j | \mathbf{S}_i) \equiv P(\mathbf{R} \cdot \mathbf{S}_j | \mathbf{R} \cdot \mathbf{S}_i)$
  - during model training & querying, rotate each  $\mathbf{q}_t$  into a fixed ordering of the  $N$  single-participant states

# Observables

- each of  $K$  participants is wearing a close-talk mic (CTM)
- extract 41 features from every CTM channel
  - log energy + MFCCs
  - $\Delta$ s and  $\Delta\Delta$ s
  - min and max normalized log energy differences (NLEDs)  
(Baskys & Stolcke, 2006)
- $41 \cdot K$  features per frame
- may vary from meeting to meeting (as  $K$  does)

# Observables

- each of  $K$  participants is wearing a close-talk mic (CTM)
- extract 41 features from every CTM channel
  - log energy + MFCCs
  - $\Delta$ s and  $\Delta\Delta$ s
  - min and max normalized log energy differences (NLEDs)  
(Boakye & Stolcke, 2006)
- 41 ·  $K$  features per frame
- may vary from meeting to meeting (as  $K$  does)

# Observables

- each of  $K$  participants is wearing a close-talk mic (CTM)
- extract 41 features from every CTM channel
  - log energy + MFCCs
  - $\Delta$ s and  $\Delta\Delta$ s
  - min and max normalized log energy differences (NLEDs)  
(Boakye & Stolcke, 2006)
- $41 \cdot K$  features per frame
- may vary from meeting to meeting (as  $K$  does)

# Observables

- each of  $K$  participants is wearing a close-talk mic (CTM)
- extract 41 features from every CTM channel
  - log energy + MFCCs
  - $\Delta$ s and  $\Delta\Delta$ s
  - min and max normalized log energy differences (NLEDs)  
(Boakye & Stolcke, 2006)
- $41 \cdot K$  features per frame
- may vary from meeting to meeting (as  $K$  does)

# Emission Probability Model

- variable feature length vector  $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_K]$
- train a single-channel GMM (64 components)
  - for  $\mathcal{S}$  and  $\mathcal{L}$
  - for  $\mathcal{N}_{all}$  and  $\mathcal{N}_{nearfield}$
- then approximate the joint MP emission with

$$P(\mathbf{X} | \mathbf{S}_i) = \prod_{k=1}^K P(\mathbf{X}[k] | \mathbf{S}_i[k]) \quad (1)$$

# Emission Probability Model

- variable feature length vector  $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_K]$
- train a single-channel GMM (64 components)
  - for  $\mathcal{S}$  and  $\mathcal{L}$
  - for  $\mathcal{N}_{all}$  and  $\mathcal{N}_{nearfield}$
- then approximate the joint MP emission with

$$P(\mathbf{X} | \mathbf{S}_i) = \prod_{k=1}^K P(\mathbf{X}[k] | \mathbf{S}_i[k]) \quad (1)$$

# Emission Probability Model

- variable feature length vector  $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_K]$
- train a single-channel GMM (64 components)
  - for  $\mathcal{S}$  and  $\mathcal{L}$
  - for  $\mathcal{N}_{all}$  and  $\mathcal{N}_{nearfield}$
- then approximate the joint MP emission with

$$P(\mathbf{X} | \mathbf{S}_i) = \prod_{k=1}^K P(\mathbf{X}[k] | \mathbf{S}_i[k]) \quad (1)$$



# Described Experiments

- ① independent versus joint participant decoding
- ② sensitivity to minimum duration constraints
- ③ sensitivity to maximum overlap constraints
- ④ generalization to other (non-Bmr) meetings

**Evaluation:** recall (R), precision (P), and unweighted F

- goal here:  $\mathcal{L}$  versus  $\neg\mathcal{L} = \mathcal{S} \cup \mathcal{N}$
- sanity:  $\mathcal{S}$  versus  $\neg\mathcal{S} = \mathcal{L} \cup \mathcal{N}$
- sanity:  $\mathcal{V} = \mathcal{S} \cup \mathcal{L}$  versus  $\neg\mathcal{V} = \mathcal{N}$

# Described Experiments

- ① independent versus joint participant decoding
- ② sensitivity to minimum duration constraints
- ③ sensitivity to maximum overlap constraints
- ④ generalization to other (non-Bmr) meetings

**Evaluation:** recall (R), precision (P), and unweighted F

- goal here:  $\mathcal{L}$  versus  $\neg\mathcal{L} = \mathcal{S} \cup \mathcal{N}$
- sanity:  $\mathcal{S}$  versus  $\neg\mathcal{S} = \mathcal{L} \cup \mathcal{N}$
- sanity:  $\mathcal{V} = \mathcal{S} \cup \mathcal{L}$  versus  $\neg\mathcal{V} = \mathcal{N}$

# Described Experiments

- ① independent versus joint participant decoding
- ② sensitivity to minimum duration constraints
- ③ sensitivity to maximum overlap constraints
- ④ generalization to other (non-Bmr) meetings

**Evaluation:** recall (R), precision (P), and unweighted F

- goal here:  $\mathcal{L}$  versus  $\neg\mathcal{L} = \mathcal{S} \cup \mathcal{N}$
- sanity:  $\mathcal{S}$  versus  $\neg\mathcal{S} = \mathcal{L} \cup \mathcal{N}$
- sanity:  $\mathcal{V} = \mathcal{S} \cup \mathcal{L}$  versus  $\neg\mathcal{V} = \mathcal{N}$

# Single-participant vs Multiparticipant Decoding

- for decoding participants independently
  - $\mathcal{N}_{all}$  and  $\mathcal{N}_{farfield}$  both represent nearfield silence  $\mathcal{N}$
  - $\rightarrow$  3 competing models, rather than 4
- for decoding participant jointly, can use either 3 or 4 models

Decoding	$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
	F	R	P	F	R	P	F
indep, 3 AM	76.3	90.3	85.0	87.6	80.9	20.4	32.6
joint, 3 AM	78.8	89.7	86.0	87.8	59.2	20.6	30.6
► joint, 4 AM	<b>79.5</b>	83.6	90.0	86.7	55.2	<b>25.1</b>	<b>34.5</b>

- 1 joint decoding improves precision by reducing potential overlap
- 2 modeling farfield vocalization on CTMs significantly improves precision for  $\mathcal{S}$  and  $\mathcal{L}$

# Single-participant vs Multiparticipant Decoding

- for decoding participants independently
  - $\mathcal{N}_{all}$  and  $\mathcal{N}_{farfield}$  both represent nearfield silence  $\mathcal{N}$
  - $\rightarrow$  3 competing models, rather than 4
- for decoding participant jointly, can use either 3 or 4 models

Decoding	$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
	F	R	P	F	R	P	F
indep, 3 AM	76.3	90.3	85.0	87.6	80.9	20.4	32.6
joint, 3 AM	78.8	89.7	86.0	87.8	59.2	20.6	30.6
▶ joint, 4 AM	<b>79.5</b>	83.6	90.0	86.7	55.2	<b>25.1</b>	<b>34.5</b>

- 1 joint decoding improves precision by reducing potential overlap
- 2 modeling farfield vocalization on CTMs significantly improves precision for  $\mathcal{S}$  and  $\mathcal{L}$

# Single-participant vs Multiparticipant Decoding

- for decoding participants independently
  - $\mathcal{N}_{all}$  and  $\mathcal{N}_{farfield}$  both represent nearfield silence  $\mathcal{N}$
  - $\rightarrow$  3 competing models, rather than 4
- for decoding participant jointly, can use either 3 or 4 models

Decoding	$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
	F	R	P	F	R	P	F
indep, 3 AM	76.3	90.3	85.0	87.6	80.9	20.4	32.6
joint, 3 AM	78.8	89.7	86.0	87.8	59.2	20.6	30.6
▶ joint, 4 AM	<b>79.5</b>	83.6	90.0	86.7	55.2	<b>25.1</b>	<b>34.5</b>

- 1 joint decoding improves precision by reducing potential overlap
- 2 modeling farfield vocalization on CTMs significantly improves precision for  $\mathcal{S}$  and  $\mathcal{L}$

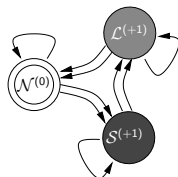
# Alternative Minimum Duration Constraints $T_{min}$

$$T_{min} = (0.1, 0.1, 0.1)$$

$$T_{min} = (0.3, 0.3, 0.3)$$

$$T_{min} = (0.2, 0.4, 0.3)$$

# Alternative Minimum Duration Constraints $\mathbf{T}_{min}$



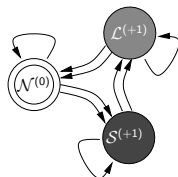
$$\mathbf{T}_{min} = (0.1, 0.1, 0.1)$$

$$\mathbf{T}_{min} = (0.3, 0.3, 0.3)$$

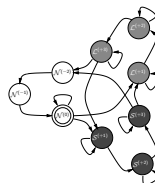
$$\mathbf{T}_{min} = (0.2, 0.4, 0.3)$$



# Alternative Minimum Duration Constraints $T_{min}$



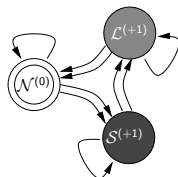
$$T_{min} = (0.1, 0.1, 0.1)$$



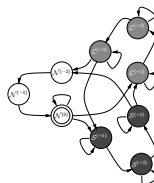
$$T_{min} = (0.3, 0.3, 0.3)$$

$$T_{min} = (0.2, 0.4, 0.3)$$

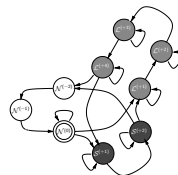
# Alternative Minimum Duration Constraints $T_{min}$



$$T_{min} = (0.1, 0.1, 0.1)$$



$$T_{min} = (0.3, 0.3, 0.3)$$



$$T_{min} = (0.2, 0.4, 0.3)$$

# Alternative Minimum Duration Constraints $T_{min}$

- hold maximum overlap constraints fixed,  $\mathbf{K}_{max} = (2, 3, 3)$

$T_{min}$ (s)	$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
	F	R	P	F	R	P	F
(0.1, 0.1, 0.1)	78.1	82.3	89.9	86.0	<b>55.9</b>	22.1	31.7
(0.3, 0.3, 0.3)	<b>79.5</b>	<b>83.7</b>	<b>90.4</b>	<b>86.9</b>	54.7	24.2	33.6
► (0.2, 0.4, 0.3)	<b>79.5</b>	83.6	90.0	86.7	55.2	<b>25.1</b>	<b>34.5</b>

- increasing all  $T_{min}$  from 0.1s to 0.3s improves all F measures
- allowing  $T_{min}^{\mathcal{L}} > T_{min}^{\mathcal{S}}$  can result in higher F( $\mathcal{L}$ )

# Alternative Minimum Duration Constraints $T_{min}$

- hold maximum overlap constraints fixed,  $\mathbf{K}_{max} = (2, 3, 3)$

$T_{min}$ (s)	$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
	F	R	P	F	R	P	F
(0.1, 0.1, 0.1)	<b>78.1</b>	82.3	89.9	<b>86.0</b>	<b>55.9</b>	22.1	<b>31.7</b>
(0.3, 0.3, 0.3)	<b>79.5</b>	<b>83.7</b>	<b>90.4</b>	<b>86.9</b>	54.7	24.2	<b>33.6</b>
► (0.2, 0.4, 0.3)	<b>79.5</b>	83.6	90.0	86.7	55.2	<b>25.1</b>	<b>34.5</b>

- increasing all  $T_{min}$  from 0.1s to 0.3s improves all F measures
- allowing  $T_{min}^{\mathcal{L}} > T_{min}^{\mathcal{S}}$  can result in higher F( $\mathcal{L}$ )

# Alternative Minimum Duration Constraints $T_{min}$

- hold maximum overlap constraints fixed,  $\mathbf{K}_{max} = (2, 3, 3)$

$T_{min}$ (s)	$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
	F	R	P	F	R	P	F
(0.1, 0.1, 0.1)	78.1	82.3	89.9	86.0	<b>55.9</b>	22.1	31.7
(0.3, <b>0.3</b> , <b>0.3</b> )	<b>79.5</b>	<b>83.7</b>	<b>90.4</b>	<b>86.9</b>	54.7	24.2	<b>33.6</b>
► (0.2, <b>0.4</b> , <b>0.3</b> )	<b>79.5</b>	83.6	90.0	86.7	55.2	<b>25.1</b>	<b>34.5</b>

- increasing all  $T_{min}$  from 0.1s to 0.3s improves all F measures
- allowing  $T_{min}^{\mathcal{L}} > T_{min}^{\mathcal{S}}$  can result in higher F ( $\mathcal{L}$ )

# Alternative Maximum Overlap Constraints $K_{max}$

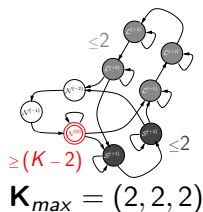
$$K_{max} = (3, 2, 3)$$

$$K_{max} = (2, 2, 2)$$

$$K_{max} = (2, 2, 3)$$

$$K_{max} = (2, 3, 3)$$

# Alternative Maximum Overlap Constraints $K_{max}$

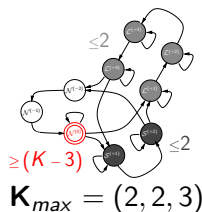
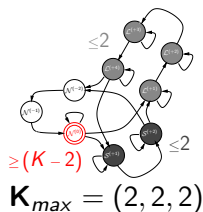


$$K_{max} = (3, 2, 3)$$

$$K_{max} = (2, 2, 3)$$

$$K_{max} = (2, 3, 3)$$

# Alternative Maximum Overlap Constraints $K_{max}$

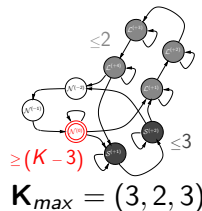
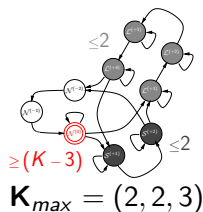
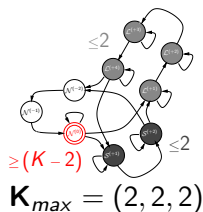


$$K_{max} = (3, 2, 3)$$

$$K_{max} = (2, 3, 3)$$

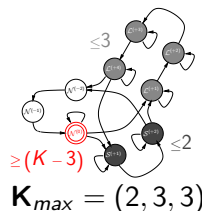
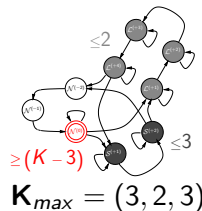
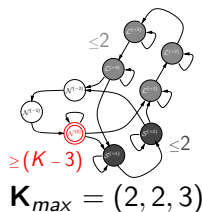
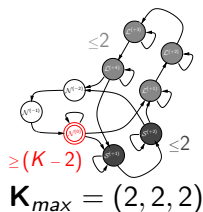


# Alternative Maximum Overlap Constraints $K_{max}$



$$K_{max} = (2, 3, 3)$$

# Alternative Maximum Overlap Constraints $K_{max}$



# Alternative Maximum Overlap Constraints $K_{max}$

- minimum duration constraints fixed,  $T_{min} = (0.2, 0.4, 0.3)$

$K_{max}$	$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
	F	R	P	F	R	P	F
(2, 2, 2)	<b>81.3</b>	83.3	<b>90.6</b>	<b>86.8</b>	36.9	<b>27.8</b>	31.7
(2, 2, 3)	79.9	84.0	89.0	86.4	48.8	24.3	32.4
(3, 2, 3)	79.9	<b>84.2</b>	88.6	86.4	49.1	24.6	32.8
► (2, 3, 3)	79.5	83.6	90.0	86.7	<b>55.2</b>	25.1	<b>34.5</b>

- increasing  $K_{max}$  generally leads to higher R and lower P
  - increasing  $K_{max}^{\mathcal{S}}$  from 2 to 3 has negligible impact
  - increasing  $K_{max}^{\mathcal{L}}$  from 2 to 3 has significant impact
- ★ because a higher proportion of  $\mathcal{L}$  is produced in overlap

# Alternative Maximum Overlap Constraints $K_{max}$

- minimum duration constraints fixed,  $\mathbf{T}_{min} = (0.2, 0.4, 0.3)$

$K_{max}$	$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
	F	R	P	F	R	P	F
(2, 2, 2)	<b>81.3</b>	83.3	<b>90.6</b>	<b>86.8</b>	36.9	<b>27.8</b>	31.7
(2, 2, 3)	79.9	84.0	89.0	86.4	48.8	24.3	32.4
(3, 2, 3)	79.9	<b>84.2</b>	88.6	86.4	49.1	24.6	32.8
► (2, 3, 3)	79.5	83.6	90.0	86.7	<b>55.2</b>	25.1	<b>34.5</b>

- increasing  $K_{max}$  generally leads to higher R and lower P
  - increasing  $K_{max}^{\mathcal{S}}$  from 2 to 3 has negligible impact
  - increasing  $K_{max}^{\mathcal{L}}$  from 2 to 3 has significant impact
- ★ because a higher proportion of  $\mathcal{L}$  is produced in overlap

# Alternative Maximum Overlap Constraints $K_{max}$

- minimum duration constraints fixed,  $\mathbf{T}_{min} = (0.2, 0.4, 0.3)$

$K_{max}$	$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
	F	R	P	F	R	P	F
(2, 2, 2)	<b>81.3</b>	83.3	<b>90.6</b>	<b>86.8</b>	36.9	<b>27.8</b>	31.7
(2, 2, 3)	79.9	84.0	89.0	86.4	48.8	24.3	32.4
(3, 2, 3)	79.9	<b>84.2</b>	88.6	86.4	49.1	24.6	32.8
► (2, 3, 3)	79.5	83.6	90.0	86.7	<b>55.2</b>	25.1	<b>34.5</b>

- increasing  $K_{max}$  generally leads to higher R and lower P
  - increasing  $K_{max}^{\mathcal{S}}$  from 2 to 3 has negligible impact
  - increasing  $K_{max}^{\mathcal{L}}$  from 2 to 3 has significant impact
- ★ because a higher proportion of  $\mathcal{L}$  is produced in overlap

# Alternative Maximum Overlap Constraints $K_{max}$

- minimum duration constraints fixed,  $\mathbf{T}_{min} = (0.2, 0.4, 0.3)$

$K_{max}$	$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
	F	R	P	F	R	P	F
(2, 2, 2)	<b>81.3</b>	83.3	<b>90.6</b>	<b>86.8</b>	36.9	<b>27.8</b>	31.7
(2, <b>2</b> , 3)	<b>79.9</b>	84.0	89.0	<b>86.4</b>	<b>48.8</b>	24.3	<b>32.4</b>
(3, 2, 3)	79.9	<b>84.2</b>	88.6	86.4	49.1	24.6	32.8
► (2, <b>3</b> , 3)	<b>79.5</b>	83.6	90.0	86.7	<b>55.2</b>	25.1	<b>34.5</b>

- increasing  $K_{max}$  generally leads to higher R and lower P
- increasing  $K_{max}^{\mathcal{S}}$  from 2 to 3 has negligible impact
- increasing  $K_{max}^{\mathcal{L}}$  from 2 to 3 has significant impact

★ because a higher proportion of  $\mathcal{L}$  is produced in overlap

# Alternative Maximum Overlap Constraints $K_{max}$

- minimum duration constraints fixed,  $\mathbf{T}_{min} = (0.2, 0.4, 0.3)$

$K_{max}$	$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
	F	R	P	F	R	P	F
(2, 2, 2)	<b>81.3</b>	83.3	<b>90.6</b>	<b>86.8</b>	36.9	<b>27.8</b>	31.7
(2, 2, 3)	79.9	84.0	89.0	86.4	48.8	24.3	32.4
(3, 2, 3)	79.9	<b>84.2</b>	88.6	86.4	49.1	24.6	32.8
► (2, 3, 3)	79.5	83.6	90.0	86.7	<b>55.2</b>	25.1	<b>34.5</b>

- increasing  $K_{max}$  generally leads to higher R and lower P
  - increasing  $K_{max}^{\mathcal{S}}$  from 2 to 3 has negligible impact
  - increasing  $K_{max}^{\mathcal{L}}$  from 2 to 3 has significant impact
- ★ because a higher proportion of  $\mathcal{L}$  is produced in overlap

# Generalization to Other Meetings

Test data $p_V(\mathcal{L})$			$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
			F	R	P	F	R	P	F
Bmr	train	10.91	<b>80.1</b>	83.4	89.8	86.5	53.0	19.4	28.4
	test	14.94	79.5	83.6	90.0	<b>86.7</b>	55.2	<b>25.1</b>	<b>34.5</b>
Bro	(all)	5.94	78.1	81.1	<b>90.6</b>	85.6	57.8	11.4	19.0
Bed	(all)	7.53	75.1	<b>84.6</b>	85.7	85.2	<b>58.7</b>	10.0	17.0

- ①  $F(\mathcal{V})$ : Bmr(train) > Bmr(test) > Bro > Bed
  - Bmr(train) and Bmr(test) have lots of participants in common
  - Bmr(train) and Bmr(test) are very similar
- ②  $F(\mathcal{L})$  on Bmr(test) higher than on Bmr(train)



# Generalization to Other Meetings

Test data $p_V(\mathcal{L})$			$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
			F	R	P	F	R	P	F
Bmr	train	10.91	<b>80.1</b>	83.4	89.8	86.5	53.0	19.4	28.4
	test	14.94	<b>79.5</b>	83.6	90.0	<b>86.7</b>	55.2	<b>25.1</b>	<b>34.5</b>
Bro	(all)	5.94	<b>78.1</b>	81.1	<b>90.6</b>	85.6	57.8	11.4	19.0
Bed	(all)	7.53	<b>75.1</b>	<b>84.6</b>	85.7	85.2	<b>58.7</b>	10.0	17.0

❶  $F(\mathcal{V})$ :  $\text{Bmr}(\text{train}) > \text{Bmr}(\text{test}) > \text{Bro} > \text{Bed}$

- $\text{Bmr}(\text{train})$  and  $\text{Bmr}(\text{test})$  have lots of participants in common
- with Bmr, Bro shares 3 participants, and Bed 1 participant

❷  $F(\mathcal{L})$  on  $\text{Bmr}(\text{test})$  higher than on  $\text{Bmr}(\text{train})$

# Generalization to Other Meetings

Test data $p_V(\mathcal{L})$			$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
			F	R	P	F	R	P	F
Bmr	train	10.91	<b>80.1</b>	83.4	89.8	86.5	53.0	19.4	28.4
	test	14.94	<b>79.5</b>	83.6	90.0	<b>86.7</b>	55.2	<b>25.1</b>	<b>34.5</b>
Bro	(all)	5.94	<b>78.1</b>	81.1	<b>90.6</b>	85.6	57.8	11.4	19.0
Bed	(all)	7.53	<b>75.1</b>	<b>84.6</b>	85.7	85.2	<b>58.7</b>	10.0	17.0

- ①  $F(\mathcal{V})$ :  $\text{Bmr}(\text{train}) > \text{Bmr}(\text{test}) > \text{Bro} > \text{Bed}$ 
  - $\text{Bmr}(\text{train})$  and  $\text{Bmr}(\text{test})$  have lots of participants in common
  - with Bmr, Bro shares 3 participants, and Bed 1 participant
- ②  $F(\mathcal{L})$  on  $\text{Bmr}(\text{test})$  higher than on  $\text{Bmr}(\text{train})$

# Generalization to Other Meetings

Test data $p_V(\mathcal{L})$			$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
			F	R	P	F	R	P	F
Bmr	train	10.91	<b>80.1</b>	83.4	89.8	86.5	53.0	19.4	28.4
	test	14.94	<b>79.5</b>	83.6	90.0	<b>86.7</b>	55.2	<b>25.1</b>	<b>34.5</b>
Bro	(all)	5.94	<b>78.1</b>	81.1	<b>90.6</b>	85.6	57.8	11.4	19.0
Bed	(all)	7.53	<b>75.1</b>	<b>84.6</b>	85.7	85.2	<b>58.7</b>	10.0	17.0

- 1  $F(\mathcal{V})$ :  $\text{Bmr}(\text{train}) > \text{Bmr}(\text{test}) > \text{Bro} > \text{Bed}$ 
  - Bmr(train) and Bmr(test) have lots of participants in common
  - with Bmr, Bro shares 3 participants, and Bed 1 participant

- 2  $F(\mathcal{L})$  on Bmr(test) higher than on Bmr(train)
  - appears to correlate with  $p_V(\mathcal{L})$ , the proportion of vocalizations
  - effort spent on laughter

# Generalization to Other Meetings

Test data $p_V(\mathcal{L})$			$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
			F	R	P	F	R	P	F
Bmr	train	10.91	<b>80.1</b>	83.4	89.8	86.5	53.0	19.4	28.4
	test	14.94	79.5	83.6	90.0	<b>86.7</b>	55.2	<b>25.1</b>	<b>34.5</b>
Bro	(all)	5.94	78.1	81.1	<b>90.6</b>	85.6	57.8	11.4	19.0
Bed	(all)	7.53	75.1	<b>84.6</b>	85.7	85.2	<b>58.7</b>	10.0	17.0

- 1  $F(\mathcal{V})$ :  $\text{Bmr}(\text{train}) > \text{Bmr}(\text{test}) > \text{Bro} > \text{Bed}$ 
  - $\text{Bmr}(\text{train})$  and  $\text{Bmr}(\text{test})$  have lots of participants in common
  - with Bmr, Bro shares 3 participants, and Bed 1 participant
- 2  $F(\mathcal{L})$  on  $\text{Bmr}(\text{test})$  higher than on  $\text{Bmr}(\text{train})$ 
  - appears to correlate with  $p_V(\mathcal{L})$ , the proportion of vocalization effort spent on laughter
  - this test set is *not typical* of the corpus

# Generalization to Other Meetings

Test data $p_V(\mathcal{L})$			$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
			F	R	P	F	R	P	F
Bmr	train	10.91	<b>80.1</b>	83.4	89.8	86.5	53.0	19.4	28.4
	test	14.94	79.5	83.6	90.0	<b>86.7</b>	55.2	<b>25.1</b>	<b>34.5</b>
Bro	(all)	5.94	78.1	81.1	<b>90.6</b>	85.6	57.8	11.4	19.0
Bed	(all)	7.53	75.1	<b>84.6</b>	85.7	85.2	<b>58.7</b>	10.0	17.0

- 1  $F(\mathcal{V})$ :  $\text{Bmr}(\text{train}) > \text{Bmr}(\text{test}) > \text{Bro} > \text{Bed}$ 
  - $\text{Bmr}(\text{train})$  and  $\text{Bmr}(\text{test})$  have lots of participants in common
  - with  $\text{Bmr}$ ,  $\text{Bro}$  shares 3 participants, and  $\text{Bed}$  1 participant
- 2  $F(\mathcal{L})$  on  $\text{Bmr}(\text{test})$  higher than on  $\text{Bmr}(\text{train})$ 
  - appears to correlate with  $p_V(\mathcal{L})$ , the proportion of vocalization effort spent on laughter
  - this test set is *not typical* of the corpus

# Generalization to Other Meetings

Test data $p_V(\mathcal{L})$			$\mathcal{V}$	$\mathcal{S}$			$\mathcal{L}$		
			F	R	P	F	R	P	F
Bmr	train	10.91	<b>80.1</b>	83.4	89.8	86.5	53.0	19.4	28.4
	test	14.94	79.5	83.6	90.0	<b>86.7</b>	55.2	<b>25.1</b>	<b>34.5</b>
Bro	(all)	5.94	78.1	81.1	<b>90.6</b>	85.6	57.8	11.4	19.0
Bed	(all)	7.53	75.1	<b>84.6</b>	85.7	85.2	<b>58.7</b>	10.0	17.0

- 1  $F(\mathcal{V})$ :  $\text{Bmr}(\text{train}) > \text{Bmr}(\text{test}) > \text{Bro} > \text{Bed}$ 
  - $\text{Bmr}(\text{train})$  and  $\text{Bmr}(\text{test})$  have lots of participants in common
  - with Bmr, Bro shares 3 participants, and Bed 1 participant
- 2  $F(\mathcal{L})$  on  $\text{Bmr}(\text{test})$  higher than on  $\text{Bmr}(\text{train})$ 
  - appears to correlate with  $p_V(\mathcal{L})$ , the proportion of vocalization effort spent on laughter
  - this test set is *not typical* of the corpus

# Confusion Matrix Analysis

	hypothesized as			$\Sigma$
	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}$	6.5	<b>9.1</b>	1.0	10.4
$\mathcal{S}$	11.0	4.5	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

- final system on test set (13.8 hours)
- all quantities in minutes

# Confusion Matrix Analysis

	hypothesized as			$\Sigma$
	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}$	6.5	<b>9.1</b>	1.0	16.6
$\mathcal{S}$	11.0	4.5	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

- break down references  $\mathcal{L} \equiv \{ \mathcal{L}'_U, \mathcal{L}'_V, \mathcal{L} \cap \mathcal{S} \}$ 
  - $\mathcal{L}'_U \equiv \mathcal{L}_U - \mathcal{L} \cap \mathcal{S}$ : unvoiced laughter less “laughed speech”
  - $\mathcal{L}'_V \equiv \mathcal{L}_V - \mathcal{L} \cap \mathcal{S}$ : voiced laughter less “laughed speech”



# Confusion Matrix Analysis

	hypothesized as			$\Sigma$
	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}$	6.5	<b>9.1</b>	1.0	16.6
$\mathcal{S}$	11.0	4.5	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

- break down references  $\mathcal{L} \equiv \{ \mathcal{L}'_U, \mathcal{L}'_V, \mathcal{L} \cap \mathcal{S} \}$ 
  - $\mathcal{L}'_U \equiv \mathcal{L}_U - \mathcal{L} \cap \mathcal{S}$ : unvoiced laughter less “laughed speech”
  - $\mathcal{L}'_V \equiv \mathcal{L}_V - \mathcal{L} \cap \mathcal{S}$ : voiced laughter less “laughed speech”

# Confusion Matrix Analysis

	hypothesized as			$\Sigma$
	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}$	6.5	<b>9.1</b>	1.0	16.6
$\mathcal{S}$	11.0	4.5	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

- break down references  $\mathcal{L} \equiv \{ \mathcal{L}'_U, \mathcal{L}'_V, \mathcal{L} \cap \mathcal{S} \}$ 
  - $\mathcal{L}'_U \equiv \mathcal{L}_U - \mathcal{L} \cap \mathcal{S}$ : unvoiced laughter less “laughed speech”
  - $\mathcal{L}'_V \equiv \mathcal{L}_V - \mathcal{L} \cap \mathcal{S}$ : voiced laughter less “laughed speech”

# Confusion Matrix Analysis

	hypothesized as			$\Sigma$
	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	<b>2.8</b>	2.4	0.2	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	0.3	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	0.2	<b>0.5</b>	0.8
$\mathcal{S}$	11.0	4.5	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

- most unvoiced laughter ( $\mathcal{L}'_U$ ) is classified as silence ( $\mathcal{N}$ )
- most “laughed speech” ( $\mathcal{L} \cap \mathcal{S}$ ) is classified as speech ( $\mathcal{S}$ )

# Confusion Matrix Analysis

	hypothesized as			$\Sigma$
	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	<b>2.8</b>	2.4	0.2	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	0.3	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	0.2	<b>0.5</b>	0.8
$\mathcal{S}$	11.0	4.5	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

- 1 most unvoiced laughter ( $\mathcal{L}'_U$ ) is classified as silence ( $\mathcal{N}$ )
- 2 most “laughed speech” ( $\mathcal{L} \cap \mathcal{S}$ ) is classified as speech ( $\mathcal{S}$ )

# Confusion Matrix Analysis

	hypothesized as			$\Sigma$
	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	<b>2.8</b>	2.4	0.2	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	0.3	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	0.2	<b>0.5</b>	0.8
$\mathcal{S}$	11.0	4.5	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

- 1 most unvoiced laughter ( $\mathcal{L}'_U$ ) is classified as silence ( $\mathcal{N}$ )
- 2 most “laughed speech” ( $\mathcal{L} \cap \mathcal{S}$ ) is classified as speech ( $\mathcal{S}$ )

# Confusions Between $\mathcal{L}$ and $\mathcal{S}$

	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	$\Sigma$
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	2.8	<b>2.4</b>	<b>0.2</b>	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	<b>0.3</b>	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	<b>0.2</b>	<b>0.5</b>	0.8
$\mathcal{S}$	11.0	<b>4.5</b>	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

	$\mathcal{L}$	$\mathcal{S}$
$\mathcal{L}'$	8.9	0.5
$\mathcal{L} \cap \mathcal{S}$	0.2	0.5
$\mathcal{S}$	4.5	79.0

- looking at  $\mathcal{L}$  and  $\mathcal{S}$  only

# Confusions Between $\mathcal{L}$ and $\mathcal{S}$

	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	$\Sigma$
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	2.8	<b>2.4</b>	<b>0.2</b>	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	<b>0.3</b>	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	<b>0.2</b>	<b>0.5</b>	0.8
$\mathcal{S}$	11.0	<b>4.5</b>	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

Recall:

	$\mathcal{L}$	$\mathcal{S}$
$\mathcal{L}'$	<b>94.7</b>	5.3
$\mathcal{L} \cap \mathcal{S}$	28.6	<b>71.4</b>
$\mathcal{S}$	5.4	<b>93.6</b>

- 94% of speech is hypothesized as speech
- 95% of laughter (excluding "laughed speech") is hypothesized as laughter

# Confusions Between $\mathcal{L}$ and $\mathcal{S}$

	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	$\Sigma$
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	2.8	<b>2.4</b>	<b>0.2</b>	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	<b>0.3</b>	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	<b>0.2</b>	<b>0.5</b>	0.8
$\mathcal{S}$	11.0	<b>4.5</b>	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

Recall:

	$\mathcal{L}$	$\mathcal{S}$
$\mathcal{L}'$	<b>94.7</b>	5.3
$\mathcal{L} \cap \mathcal{S}$	28.6	<b>71.4</b>
$\mathcal{S}$	5.4	<b>93.6</b>

- 1 94% of speech is hypothesized as speech
- 2 95% of laughter (excluding "laughed speech") is hypothesized as laughter



# Confusions Between $\mathcal{L}$ and $\mathcal{S}$

	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	$\Sigma$
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	2.8	<b>2.4</b>	<b>0.2</b>	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	<b>0.3</b>	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	<b>0.2</b>	<b>0.5</b>	0.8
$\mathcal{S}$	11.0	<b>4.5</b>	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

Recall:

	$\mathcal{L}$	$\mathcal{S}$
$\mathcal{L}'$	<b>94.7</b>	5.3
$\mathcal{L} \cap \mathcal{S}$	28.6	<b>71.4</b>
$\mathcal{S}$	5.4	<b>93.6</b>

- 1 94% of speech is hypothesized as speech
- 2 95% of laughter (excluding "laughed speech") is hypothesized as laughter

# Confusions Between $\mathcal{L}$ and $\mathcal{S}$

	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	$\Sigma$
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	2.8	<b>2.4</b>	<b>0.2</b>	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	<b>0.3</b>	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	<b>0.2</b>	<b>0.5</b>	0.8
$\mathcal{S}$	11.0	<b>4.5</b>	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

Precision:

	$\mathcal{L}$	$\mathcal{S}$
$\mathcal{L}'$	<b>65.4</b>	0.6
$\mathcal{L} \cap \mathcal{S}$	1.5	0.6
$\mathcal{S}$	33.1	<b>98.8</b>

- 1 99% of hypothesized speech is speech
- 2 65% of hypothesized laughter is laughter

# Confusions Between $\mathcal{L}$ and $\mathcal{S}$

	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	$\Sigma$
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	2.8	<b>2.4</b>	<b>0.2</b>	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	<b>0.3</b>	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	<b>0.2</b>	<b>0.5</b>	0.8
$\mathcal{S}$	11.0	<b>4.5</b>	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

Precision:

	$\mathcal{L}$	$\mathcal{S}$
$\mathcal{L}'$	<b>65.4</b>	0.6
$\mathcal{L} \cap \mathcal{S}$	1.5	0.6
$\mathcal{S}$	33.1	<b>98.8</b>

- 1 99% of hypothesized speech is speech
- 2 65% of hypothesized laughter is laughter

# Confusions Between $\mathcal{L}$ and $\mathcal{S}$

	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	$\Sigma$
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	2.8	<b>2.4</b>	<b>0.2</b>	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	<b>0.3</b>	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	<b>0.2</b>	<b>0.5</b>	0.8
$\mathcal{S}$	11.0	<b>4.5</b>	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

Precision:

	$\mathcal{L}$	$\mathcal{S}$
$\mathcal{L}'$	<b>65.4</b>	0.6
$\mathcal{L} \cap \mathcal{S}$	1.5	0.6
$\mathcal{S}$	33.1	<b>98.8</b>

- 1 99% of hypothesized speech is speech
- 2 65% of hypothesized laughter is laughter

# Confusions Between $\mathcal{L}$ and $\mathcal{N}$

	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	$\Sigma$
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	2.8	<b>2.4</b>	0.2	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	0.3	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	<b>0.2</b>	0.5	0.8
$\mathcal{S}$	11.0	4.5	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

	$\mathcal{N}$	$\mathcal{L}$
$\mathcal{N}$	685.4	22.9
$\mathcal{L}'_U$	2.8	2.4
$\mathcal{L}'_V$	3.7	6.7

- looking at  $\mathcal{L}$  and  $\mathcal{N}$  only

# Confusions Between $\mathcal{L}$ and $\mathcal{N}$

	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	$\Sigma$
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	2.8	<b>2.4</b>	0.2	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	0.3	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	<b>0.2</b>	0.5	0.8
$\mathcal{S}$	11.0	4.5	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

Recall:

	$\mathcal{N}$	$\mathcal{L}$
$\mathcal{N}$	<b>96.8</b>	3.2
$\mathcal{L}'_U$	<b>53.9</b>	46.2
$\mathcal{L}'_V$	35.6	<b>64.4</b>

- 97% of silence is hypothesized as silence
- 64% of voiced laughter (including “laughed speech”) is classified as laughter

# Confusions Between $\mathcal{L}$ and $\mathcal{N}$

	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	$\Sigma$
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	2.8	<b>2.4</b>	0.2	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	0.3	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	<b>0.2</b>	0.5	0.8
$\mathcal{S}$	11.0	4.5	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

Recall:

	$\mathcal{N}$	$\mathcal{L}$
$\mathcal{N}$	<b>96.8</b>	3.2
$\mathcal{L}'_U$	<b>53.9</b>	46.2
$\mathcal{L}'_V$	35.6	<b>64.4</b>

- 1 97% of silence is hypothesized as silence
- 2 64% of voiced laughter (including “laughed speech”) is classified as laughter

# Confusions Between $\mathcal{L}$ and $\mathcal{N}$

	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	$\Sigma$
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	2.8	<b>2.4</b>	0.2	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	0.3	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	<b>0.2</b>	0.5	0.8
$\mathcal{S}$	11.0	4.5	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

Recall:

	$\mathcal{N}$	$\mathcal{L}$
$\mathcal{N}$	<b>96.8</b>	3.2
$\mathcal{L}'_U$	<b>53.9</b>	46.2
$\mathcal{L}'_V$	35.6	<b>64.4</b>

- 1 97% of silence is hypothesized as silence
- 2 64% of voiced laughter (including “laughed speech”) is classified as laughter



# Confusions Between $\mathcal{L}$ and $\mathcal{N}$

	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	$\Sigma$
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	2.8	<b>2.4</b>	0.2	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	0.3	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	<b>0.2</b>	0.5	0.8
$\mathcal{S}$	11.0	4.5	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

Precision:

	$\mathcal{N}$	$\mathcal{L}$
$\mathcal{N}$	<b>99.1</b>	<b>71.6</b>
$\mathcal{L}'_U$	0.4	7.5
$\mathcal{L}'_V$	0.5	20.9

- 99% of hypothesized silence is silence
- 28% of hypothesized laughter is laughter
- 72% of hypothesized laughter is silence

# Confusions Between $\mathcal{L}$ and $\mathcal{N}$

	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	$\Sigma$
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	2.8	<b>2.4</b>	0.2	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	0.3	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	<b>0.2</b>	0.5	0.8
$\mathcal{S}$	11.0	4.5	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

Precision:

	$\mathcal{N}$	$\mathcal{L}$
$\mathcal{N}$	<b>99.1</b>	<b>71.6</b>
$\mathcal{L}'_U$	0.4	7.5
$\mathcal{L}'_V$	0.5	20.9

- 1 99% of hypothesized silence is silence
- 2 28% of hypothesized laughter is laughter
- 3 72% of hypothesized laughter is silence

# Confusions Between $\mathcal{L}$ and $\mathcal{N}$

	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	$\Sigma$
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	2.8	<b>2.4</b>	0.2	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	0.3	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	<b>0.2</b>	0.5	0.8
$\mathcal{S}$	11.0	4.5	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

Precision:

	$\mathcal{N}$	$\mathcal{L}$
$\mathcal{N}$	<b>99.1</b>	<b>71.6</b>
$\mathcal{L}'_U$	0.4	7.5
$\mathcal{L}'_V$	0.5	20.9

- 1 99% of hypothesized silence is silence
- 2 28% of hypothesized laughter is laughter
- 3 72% of hypothesized laughter is silence

# Confusions Between $\mathcal{L}$ and $\mathcal{N}$

	$\mathcal{N}$	$\mathcal{L}$	$\mathcal{S}$	$\Sigma$
$\mathcal{N}$	<b>685.4</b>	22.9	7.8	716.2
$\mathcal{L}'_U$	2.8	<b>2.4</b>	0.2	5.4
$\mathcal{L}'_V$	3.6	<b>6.5</b>	0.3	10.4
$\mathcal{L} \cap \mathcal{S}$	0.1	<b>0.2</b>	0.5	0.8
$\mathcal{S}$	11.0	4.5	<b>79.0</b>	94.4
$\Sigma$	702.9	36.6	87.8	827.2

Precision:

	$\mathcal{N}$	$\mathcal{L}$
$\mathcal{N}$	<b>99.1</b>	<b>71.6</b>
$\mathcal{L}'_U$	0.4	7.5
$\mathcal{L}'_V$	0.5	20.9

- 1 99% of hypothesized silence is silence
- 2 28% of hypothesized laughter is laughter
- 3 **72% of hypothesized laughter is silence**

# Conclusions

- ① baseline system for multiparticpant 3-way VAD
  - no pre-segmentation assumed
  - all laughter considered
- ② { laughter vs silence } harder than { laughter vs speech }
- ③ speech/laughter contrastive constraints helpful
  - maximum allowed degree of overlap
  - minimum state duration
- ④ current performance is a function of
  - ① proportion of laughter present
  - ② participant novelty

# Conclusions

- ❶ baseline system for multiparticpant 3-way VAD
  - no pre-segmentation assumed
  - all laughter considered
- ❷ { laughter vs silence } harder than { laughter vs speech }
- ❸ speech/laughter contrastive constraints helpful
  - maximum allowed degree of overlap
  - minimum state duration
- ❹ current performance is a function of
  - ❶ proportion of laughter present
  - ❷ participant novelty

# Conclusions

- ❶ baseline system for multiparticpant 3-way VAD
  - no pre-segmentation assumed
  - all laughter considered
- ❷ { laughter vs silence } harder than { laughter vs speech }
- ❸ speech/laughter contrastive constraints helpful
  - maximum allowed degree of overlap
  - minimum state duration
- ❹ current performance is a function of
  - ❶ proportion of laughter present
  - ❷ participant novelty

# Conclusions

- ❶ baseline system for multiparticpant 3-way VAD
  - no pre-segmentation assumed
  - all laughter considered
- ❷ { laughter vs silence } harder than { laughter vs speech }
- ❸ speech/laughter contrastive constraints helpful
  - maximum allowed degree of overlap
  - minimum state duration
- ❹ current performance is a function of
  - ❶ proportion of laughter present
  - ❷ participant novelty



# Conclusions

- ❶ baseline system for multiparticpant 3-way VAD
  - no pre-segmentation assumed
  - all laughter considered
- ❷ { laughter vs silence } harder than { laughter vs speech }
- ❸ speech/laughter contrastive constraints helpful
  - maximum allowed degree of overlap
  - minimum state duration
- ❹ current performance is a function of
  - ❶ proportion of laughter present
  - ❷ participant novelty

# Possible Future Work

- ① model voiced and unvoiced laughter ( $\mathcal{L}_V$  and  $\mathcal{L}_U$ ) separately
  - different acoustics
  - different overlap contexts (Laskowski & Burger, 2007c)
  - different semantics
- ② characterize laughter by instrumentality to high level tasks
  - which laughter signals different emotional valence
  - which laughter signals involvement hotspots
  - Q: Does instrumentality correspond to how clear and unambiguous laughter is? cf. (Truong & van Leeuwen, 2007a)
- ③ multi-pass, multi-resolution laughter detection
  - pass 1: large frame size (0.1s), small context (0.1s)
  - pass 2: small frame size (0.01s), large context (1.0s) (Knox & Mirghafori, 2007)

# Possible Future Work

- ❶ model voiced and unvoiced laughter ( $\mathcal{L}_V$  and  $\mathcal{L}_U$ ) separately
  - different acoustics
  - different overlap contexts (Laskowski & Burger, 2007c)
  - different semantics
- ❷ characterize laughter by instrumentality to high level tasks
  - which laughter signals different emotional valence
  - which laughter signals involvement hotspots
  - Q: Does instrumentality correspond to how clear and unambiguous laughter is? cf. (Truong & van Leeuwen, 2007a)
- ❸ multi-pass, multi-resolution laughter detection
  - pass 1: large frame size (0.1s), small context (0.1s)
  - pass 2: small frame size (0.01s), large context (1.0s) (Knox & Mirghafori, 2007)

# Possible Future Work

- ❶ model voiced and unvoiced laughter ( $\mathcal{L}_V$  and  $\mathcal{L}_U$ ) separately
  - different acoustics
  - different overlap contexts (Laskowski & Burger, 2007c)
  - different semantics
- ❷ characterize laughter by instrumentality to high level tasks
  - which laughter signals different emotional valence
  - which laughter signals involvement hotspots
  - Q: Does instrumentality correspond to how clear and unambiguous laughter is? cf. (Truong & van Leeuwen, 2007a)
- ❸ multi-pass, multi-resolution laughter detection
  - pass 1: large frame size (0.1s), small context (0.1s)
  - pass 2: small frame size (0.01s), large context (1.0s) (Knox & Mirghafori, 2007)

# Possible Future Work

- ❶ model voiced and unvoiced laughter ( $\mathcal{L}_V$  and  $\mathcal{L}_U$ ) separately
  - different acoustics
  - different overlap contexts (Laskowski & Burger, 2007c)
  - different semantics
- ❷ characterize laughter by instrumentality to high level tasks
  - which laughter signals different emotional valence
  - which laughter signals involvement hotspots
  - Q: Does instrumentality correspond to how clear and unambiguous laughter is? cf. (Truong & van Leeuwen, 2007a)
- ❸ multi-pass, multi-resolution laughter detection
  - pass 1: large frame size (0.1s), small context (0.1s)
  - pass 2: small frame size (0.01s), large context (1.0s) (Knox & Mirghafori, 2007)

# Thanks for attending ...

Also, thanks to

- Susi Burger, help with  $\mathcal{L}$  segmentation & classification
- Liz Shriberg, access to ICSI MRDA Corpus
- Khiet Truong and Mary Knox, discussion of own work