| Introduction | Computational Framework | <b>Data</b><br>00 | Experiments | Conclusions |
|--------------|-------------------------|-------------------|-------------|-------------|
|              |                         |                   |             |             |

Modeling Vocal Interaction for Text-Independent Classification of Conversation Type

Kornel Laskowski<sup>1,3</sup>, Mari Ostendorf<sup>1,2</sup> & Tanja Schultz<sup>1,3</sup>

<sup>1</sup>interACT, Universität Karlsruhe
 <sup>2</sup>Dept. Electrical Engineering, University of Washington
 <sup>3</sup>interACT, Carnegie Mellon University

September 2, 2007

| Introduction<br>●○○○ | Computational Framework | Data<br>oo | Experiments | Conclusions |
|----------------------|-------------------------|------------|-------------|-------------|
|                      |                         |            |             |             |
| What Is <b>\</b>     | Vocal Interaction?      |            |             |             |

- the patterns of vocal activity for all participants to a conversation
  - no words —> a text-independent representation of multi-party conversation
- as used in psycholinguistics (Dabbs & Ruback, 1987)
- in telecommunications: "on-off patterns" (Brady, 1967)
- studied since the 1930s



(人間) (人) (人) (人) (人)

| Introduction<br>●○○○ | Computational Framework | Data<br>oo | Experiments | Conclusions |
|----------------------|-------------------------|------------|-------------|-------------|
|                      |                         |            |             |             |
| What Is              | Vocal Interaction?      |            |             |             |

- the patterns of vocal activity for all participants to a conversation
  - no words —> a text-independent representation of multi-party conversation
- as used in psycholinguistics (Dabbs & Ruback, 1987)
- in telecommunications: "on-off patterns" (Brady, 1967)
- studied since the 1930s



<回> < 回> < 回> < 回>

| Introduction<br>●○○○ | Computational Framework | Data<br>oo | Experiments | Conclusions |
|----------------------|-------------------------|------------|-------------|-------------|
|                      |                         |            |             |             |
| What Is              | Vocal Interaction?      |            |             |             |

- the patterns of vocal activity for all participants to a conversation
  - no words → a text-independent representation of multi-party conversation
- as used in psycholinguistics (Dabbs & Ruback, 1987)
- in telecommunications: "on-off patterns" (Brady, 1967)
- studied since the 1930s



| Introduction | Computational Framework | Data<br>oo | Experiments | Conclusions |
|--------------|-------------------------|------------|-------------|-------------|
|              |                         |            |             |             |
| What Is      | Vocal Interaction?      |            |             |             |

- the patterns of vocal activity for all participants to a conversation
  - no words —> a text-independent representation of multi-party conversation
- as used in psycholinguistics (Dabbs & Ruback, 1987)
- in telecommunications: "on-off patterns" (Brady, 1967)
- studied since the 1930s



| Introduction | Computational Framework | Data<br>oo | Experiments | Conclusions |
|--------------|-------------------------|------------|-------------|-------------|
|              |                         |            |             |             |
| What Is      | Vocal Interaction?      |            |             |             |

- the patterns of vocal activity for all participants to a conversation
  - no words —> a text-independent representation of multi-party conversation
- as used in psycholinguistics (Dabbs & Ruback, 1987)
- in telecommunications: "on-off patterns" (Brady, 1967)
- studied since the 1930s



| Introduction | Computational Framework | Data | Experiments | Conclusions |
|--------------|-------------------------|------|-------------|-------------|
| 0000         |                         |      |             |             |
|              |                         |      |             |             |

#### a basic competence in conversation understanding

- type is most often taken for granted
  - ie. "My project is about cocktail parties. Why would I ever need to know that a cocktail party is not a business meeting?"
- searching & indexing in heterogenous multi-party conversation recordings (or portions)
- text-independence: pre-ASR availability of type hypothesis/prior

伺下 イヨト イヨ

| Introduction | Computational Framework | Data | Experiments | Conclusions |
|--------------|-------------------------|------|-------------|-------------|
| 0000         |                         |      |             |             |
|              |                         |      |             |             |

- a basic competence in conversation understanding
- type is most often taken for granted
  - ie. "My project is about cocktail parties. Why would I ever need to know that a cocktail party is not a business meeting?"
- searching & indexing in heterogenous multi-party conversation recordings (or portions)
- text-independence: pre-ASR availability of type hypothesis/prior

伺下 イヨト イヨ

| Introduction<br>○●○○ | Computational Framework | <b>Data</b><br>00 | Experiments | Conclusions |
|----------------------|-------------------------|-------------------|-------------|-------------|
|                      |                         |                   |             |             |

- a basic competence in conversation understanding
- type is most often taken for granted
  - ie. "My project is about cocktail parties. Why would I ever need to know that a cocktail party is not a business meeting?"
- searching & indexing in heterogenous multi-party conversation recordings (or portions)
- text-independence: pre-ASR availability of type hypothesis/prior
  - may contribute to optimal selection of ASR components.
     type classification possible where no ASR or upstream processing possible.

| Introduction<br>○●○○ | Computational Framework | <b>Data</b><br>00 | Experiments | Conclusions |
|----------------------|-------------------------|-------------------|-------------|-------------|
|                      |                         |                   |             |             |

- a basic competence in conversation understanding
- type is most often taken for granted
  - ie. "My project is about cocktail parties. Why would I ever need to know that a cocktail party is not a business meeting?"
- searching & indexing in heterogenous multi-party conversation recordings (or portions)
- **text-independence**: pre-ASR availability of type hypothesis/prior
  - may contribute to optimal selection of ASR components
     type classification possible where no ASR or upstream processing possible

| Introduction<br>○●○○ | Computational Framework | <b>Data</b><br>00 | Experiments | Conclusions |
|----------------------|-------------------------|-------------------|-------------|-------------|
|                      |                         |                   |             |             |

- a basic competence in conversation understanding
- type is most often taken for granted
  - ie. "My project is about cocktail parties. Why would I ever need to know that a cocktail party is not a business meeting?"
- searching & indexing in heterogenous multi-party conversation recordings (or portions)
- text-independence: pre-ASR availability of type hypothesis/prior
  - may contribute to optimal selection of ASR components
  - type classification possible where no ASR or upstream processing possible

A (B) + A (B) + A (B) +

| Introduction<br>○●○○ | Computational Framework | <b>Data</b><br>00 | Experiments | Conclusions |
|----------------------|-------------------------|-------------------|-------------|-------------|
|                      |                         |                   |             |             |

- a basic competence in conversation understanding
- type is most often taken for granted
  - ie. "My project is about cocktail parties. Why would I ever need to know that a cocktail party is not a business meeting?"
- searching & indexing in heterogenous multi-party conversation recordings (or portions)
- text-independence: pre-ASR availability of type hypothesis/prior
  - may contribute to optimal selection of ASR components
  - type classification possible where no ASR or upstream processing possible

A (B) + A (B) + A (B) +

| Introduction<br>○●○○ | Computational Framework | <b>Data</b><br>00 | Experiments | Conclusions |
|----------------------|-------------------------|-------------------|-------------|-------------|
|                      |                         |                   |             |             |

- a basic competence in conversation understanding
- type is most often taken for granted
  - ie. "My project is about cocktail parties. Why would I ever need to know that a cocktail party is not a business meeting?"
- searching & indexing in heterogenous multi-party conversation recordings (or portions)
- text-independence: pre-ASR availability of type hypothesis/prior
  - may contribute to optimal selection of ASR components
  - type classification possible where no ASR or upstream processing possible

| Introduction<br>○○●○ | Computational Framework | Data<br>oo | Experiments | Conclusions |
|----------------------|-------------------------|------------|-------------|-------------|
|                      |                         |            |             |             |
| Defining (           | Conversation Type       |            |             |             |

• Sacks (1974) viewed conversation as one of several normative **speech-exchange system** types



< ∃⇒

| Introduction | Computational Framework | Data<br>oo | Experiments | Conclusions |
|--------------|-------------------------|------------|-------------|-------------|
|              |                         |            |             |             |
| Defining     | Conversation Type       |            |             |             |

- Sacks (1974) viewed conversation as one of several normative **speech-exchange system** types
- others include: lectures



| Introduction | Computational Framework | Data<br>oo | Experiments | Conclusions |
|--------------|-------------------------|------------|-------------|-------------|
|              |                         |            |             |             |
| Defining     | Conversation Type       |            |             |             |

- Sacks (1974) viewed conversation as one of several normative **speech-exchange system** types
- others include: lectures, rituals



| Introduction | Computational Framework | Data<br>oo | Experiments | Conclusions |
|--------------|-------------------------|------------|-------------|-------------|
|              |                         |            |             |             |
| Defining     | Conversation Type       |            |             |             |

- Sacks (1974) viewed conversation as one of several normative **speech-exchange system** types
- others include: lectures, rituals, debates, etc.



| Introduction | Computational Framework | Data<br>oo | Experiments | Conclusions |
|--------------|-------------------------|------------|-------------|-------------|
|              |                         |            |             |             |
| Defining     | Conversation Type       |            |             |             |

- Sacks (1974) viewed conversation as one of several normative speech-exchange system types
- others include: lectures, rituals, debates, etc.
- here, type of conversation = subtype of work-related conversation (meeting)
- implicitly assume that specific activities and specific participant groups and/or roles give rise to vocal interactions which are subtype-specific

| Introduction | Computational Framework | Data<br>oo | Experiments | Conclusions |
|--------------|-------------------------|------------|-------------|-------------|
|              |                         |            |             |             |
| Defining     | Conversation Type       |            |             |             |

- Sacks (1974) viewed conversation as one of several normative speech-exchange system types
- others include: lectures, rituals, debates, etc.
- here, type of conversation = subtype of work-related conversation (meeting)
- implicitly assume that specific activities and specific participant groups and/or roles give rise to vocal interactions which are subtype-specific

| Introduction<br>○○○● | Computational Framework | Data<br>00 | Experiments | Conclusions |
|----------------------|-------------------------|------------|-------------|-------------|
|                      |                         |            |             |             |
| Related V            | Vork                    |            |             |             |

#### • none on conversation type classification

- various, on evolving conversation state
  - (Banerjee & Rudnicky, 2004)
  - (McCowan et al, 2005)
  - (Zancanaro et al, 2006)
- several related text-independent tasks
  - participant dominance detection (Rienks et al, 2005), 4-party
  - interaction group recognition (Brdiczka et al, 2005), 4-party
  - conversational pair detection (Basu, 2002), 2-party
- modeling vocal interaction for vocal activity detection
  - meetings (Laskowski & Schultz, 2006)
  - ambulatory data (Wyatt et al, 2007)

<回> < 回> < 回> < 回>

| Introduction<br>○○○● | Computational Framework | Data<br>oo | Experiments | Conclusions |
|----------------------|-------------------------|------------|-------------|-------------|
|                      |                         |            |             |             |
| Related W            | /ork                    |            |             |             |

- none on conversation type classification
- various, on evolving conversation state
  - (Banerjee & Rudnicky, 2004)
  - (McCowan et al, 2005)
  - (Zancanaro et al, 2006)
- several related text-independent tasks
  - participant dominance detection (Rienks et al, 2005), 4-party
  - interaction group recognition (Brdiczka et al, 2005), 4-party
  - conversational pair detection (Basu, 2002), 2-party
- modeling vocal interaction for vocal activity detection
  - meetings (Laskowski & Schultz, 2006)
  - ambulatory data (Wyatt et al, 2007)

| Introduction<br>○○○● | Computational Framework | Data<br>oo | Experiments | Conclusions |
|----------------------|-------------------------|------------|-------------|-------------|
|                      |                         |            |             |             |
| Related Wo           | ork                     |            |             |             |

- none on conversation type classification
- various, on evolving conversation state
  - (Banerjee & Rudnicky, 2004)
  - (McCowan et al, 2005)
  - (Zancanaro et al, 2006)
- several related text-independent tasks
  - participant dominance detection (Rienks et al, 2005), 4-party
  - interaction group recognition (Brdiczka et al, 2005), 4-party
  - conversational pair detection (Basu, 2002), 2-party
- modeling vocal interaction for vocal activity detection
  - meetings (Laskowski & Schultz, 2006)
  - ambulatory data (Wyatt et al, 2007)

白 ト イヨ ト イヨト

| Introduction<br>○○○● | Computational Framework | Data<br>00 | Experiments | Conclusions |
|----------------------|-------------------------|------------|-------------|-------------|
|                      |                         |            |             |             |
| Related Wo           | rk                      |            |             |             |

- none on conversation type classification
- various, on evolving conversation state
  - (Banerjee & Rudnicky, 2004)
  - (McCowan et al, 2005)
  - (Zancanaro et al, 2006)
- several related text-independent tasks
  - participant dominance detection (Rienks et al, 2005), 4-party
  - interaction group recognition (Brdiczka et al, 2005), 4-party
  - conversational pair detection (Basu, 2002), 2-party
- modeling vocal interaction for vocal activity detection
  - meetings (Laskowski & Schultz, 2006)
  - ambulatory data (Wyatt et al, 2007)

A B K A B K

| Introduction | Computational Framework | Data | Experiments | Conclusions |
|--------------|-------------------------|------|-------------|-------------|
| 0000         | ●○○○○○○                 | oo   | 000         |             |
| Observable   | S                       |      |             |             |

 the vocal interaction record of a conversation C, of type T (of N<sub>T</sub> possible conversation types)



- at time t, each of K participants is in one of 2 discrete states, vocalizing (V) or not vocalizing (N)
- therefore, at time *t*, the state **q**<sub>t</sub> of *C*, as a whole, has one of  $2^{K}$  discrete values

- 4 回 ト 4 ヨ ト 4 ヨ ト

3

| Introduction | Computational Framework | Data | Experiments | Conclusions |
|--------------|-------------------------|------|-------------|-------------|
| 0000         | ●○○○○○○                 | oo   | 000         |             |
| Observable   | S                       |      |             |             |

 the vocal interaction record of a conversation C, of type T (of N<sub>T</sub> possible conversation types)



- at time t, each of K participants is in one of 2 discrete states, vocalizing (V) or not vocalizing (N)
- therefore, at time t, the state qt of C, as a whole, has one of 2<sup>K</sup> discrete values

・ 同 ト ・ ヨ ト ・ ヨ ト

-

| Introduction | <b>Computational Framework</b> | Data<br>oo | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Observabl    | es                             |            |             |             |

 the vocal interaction record of a conversation C, of type T (of N<sub>T</sub> possible conversation types)



- at time t, each of K participants is in one of 2 discrete states, vocalizing (V) or not vocalizing (N)
- therefore, at time t, the state  $\mathbf{q}_t$  of  $\mathcal{C}$ , as a whole, has one of  $2^K$  discrete values

向下 イヨト イヨト

3

| Introduction | <b>Computational Framework</b> | Data<br>00 | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Modeling     | Groups                         |            |             |             |



participants are drawn from a known population *P* of size ||*P*||
the number of distinct groups of size ||*G*|| ≤ ||*P*|| is



A (B) + A (B) + A (B) +

| Introduction | <b>Computational Framework</b> | Data<br>oo | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Modeling     | Groups                         |            |             |             |



participants are drawn from a known population *P* of size ||*P*||
the number of distinct groups of size ||*G*|| ≤ ||*P*|| is



A (B) + A (B) + A (B) +

| Introduction | Computational Framework | Data | Experiments | Conclusions |
|--------------|-------------------------|------|-------------|-------------|
| 0000         | ○●○○○○○○                | oo   | 000         |             |
| Modeling     | Groups                  |      |             |             |



participants are drawn from a known population *P* of size ||*P*||
the number of distinct groups of size ||*G*|| ≤ ||*P*|| is

$$N_{\mathcal{G}} = \frac{\|\mathcal{P}\|!}{(\|\mathcal{P}\| - \|\mathcal{G}\|)!}$$

| Introduction | <b>Computational Framework</b><br>○●○○○○○○ | Data<br>00 | Experiments<br>000 | Conclusions |
|--------------|--|------------|--------------------|-------------|
| Modeling     | Groups                                     |            |                    |             |



- participants are drawn from a known population  ${\mathcal P}$  of size  $\|{\mathcal P}\|$
- the number of distinct groups of size  $\|\mathcal{G}\| \leq \|\mathcal{P}\|$  is

$$N_{\mathcal{G}} = \frac{\|\mathcal{P}\|!}{(\|\mathcal{P}\| - \|\mathcal{G}\|)!}$$

伺下 イヨト イヨト

| Introduction | <b>Computational Framework</b><br>○●○○○○○○ | Data<br>00 | Experiments<br>000 | Conclusions |
|--------------|--|------------|--------------------|-------------|
| Modeling     | Groups                                     |            |                    |             |



- participants are drawn from a known population  $\mathcal P$  of size  $\|\mathcal P\|$
- $\bullet$  the number of distinct groups of size  $\|\mathcal{G}\| \leq \|\mathcal{P}\|$  is

$$N_{\mathcal{G}} = \frac{\|\mathcal{P}\|!}{(\|\mathcal{P}\| - \|\mathcal{G}\|)!}$$

| Introduction | <b>Computational Framework</b> | Data<br>00 | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Conversa     | tion Type Classifica           | tion       |             |             |

- $\bullet$  participant identities, and therefore  $\mathcal{G},$  are hidden variables
- given a set of features F extracted from C,

$$T^{*} = \arg \max_{T} P(T | \mathbf{F})$$
  
=  $\arg \max_{T} \sum_{\mathcal{G}} P(\mathcal{G}, \mathcal{T}, \mathbf{F})$   
=  $\arg \max_{T} \sum_{\mathcal{G}} P(\mathcal{T}) \times \underbrace{P(\mathcal{G} | \mathcal{T})}_{Membership} \times \underbrace{P(\mathbf{F} | \mathcal{G}, \mathcal{T})}_{Behavior}$   
Model Model

• hypothesis testing: cycle through  $N_T$  types and  $N_G$  groups

・ロン ・回と ・ヨン ・ヨン

| Introduction | Computational Framework | Data<br>oo | Experiments | Conclusions |
|--------------|-------------------------|------------|-------------|-------------|
|              |                         |            |             |             |
| Conversa     | tion Type Classifica    | tion       |             |             |

- $\bullet$  participant identities, and therefore  $\mathcal{G},$  are hidden variables
- $\bullet$  given a set of features  ${\bm F}$  extracted from  ${\cal C},$

$$\begin{aligned} \mathcal{T}^* &= \arg \max_{\mathcal{T}} P(\mathcal{T} | \mathbf{F}) \\ &= \arg \max_{\mathcal{T}} \sum_{\mathcal{G}} P(\mathcal{G}, \mathcal{T}, \mathbf{F}) \\ &= \arg \max_{\mathcal{T}} \sum_{\mathcal{G}} P(\mathcal{T}) \times \underbrace{P(\mathcal{G} | \mathcal{T})}_{\text{Membership}} \times \underbrace{P(\mathbf{F} | \mathcal{G}, \mathcal{T})}_{\text{Behavior}} \\ & \text{Behavior} \\ & \text{Model} \end{aligned}$$

• hypothesis testing: cycle through  $N_T$  types and  $N_G$  groups

▲ □ → ▲ □ → ▲ □ →

| Introduction | Computational Framework | Data<br>oo | Experiments | Conclusions |
|--------------|-------------------------|------------|-------------|-------------|
|              |                         |            |             |             |
| Conversa     | tion Type Classifica    | tion       |             |             |

- $\bullet$  participant identities, and therefore  $\mathcal{G},$  are hidden variables
- $\bullet$  given a set of features  ${\bm F}$  extracted from  ${\cal C},$

$$\begin{aligned} \mathcal{T}^* &= \arg \max_{\mathcal{T}} P(\mathcal{T} | \mathbf{F}) \\ &= \arg \max_{\mathcal{T}} \sum_{\mathcal{G}} P(\mathcal{G}, \mathcal{T}, \mathbf{F}) \\ &= \arg \max_{\mathcal{T}} \sum_{\mathcal{G}} P(\mathcal{T}) \times \underbrace{P(\mathcal{G} | \mathcal{T})}_{\text{Membership}} \times \underbrace{P(\mathbf{F} | \mathcal{G}, \mathcal{T})}_{\text{Behavior}} \\ & \text{Behavior} \end{aligned}$$

 $\bullet$  hypothesis testing: cycle through  $N_{\mathcal{T}}$  types and  $N_{\mathcal{G}}$  groups

白 と く ヨ と く ヨ と

э

| Introduction | Computational Framework | Data<br>00 | Experiments | Conclusions |
|--------------|-------------------------|------------|-------------|-------------|
|              |                         |            |             |             |
| Features     |                         |            |             |             |

 probability, when no-one else is vocalizing, that k initiates vocalization (VI) and that k continues vocalization (VC)

 probability, when j is vocalizing, that k initiates vocalization overlap (OI) and that k continues vocalization overlap (OC)

| Introduction | Computational Framework | Data<br>00 | Experiments | Conclusions |
|--------------|-------------------------|------------|-------------|-------------|
|              |                         |            |             |             |
| Features     |                         |            |             |             |

 probability, when no-one else is vocalizing, that k initiates vocalization (VI) and that k continues vocalization (VC)



 probability, when j is vocalizing, that k initiates vocalization overlap (OI) and that k continues vocalization overlap (OC)
| Introduction | <b>Computational Framework</b> | Data<br>oo | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Features     |                                |            |             |             |



 probability, when j is vocalizing, that k initiates vocalization overlap (OI) and that k continues vocalization overlap (OC)

(4) (3) (4) (4) (4)

| Introduction | <b>Computational Framework</b> | Data<br>oo | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Features     |                                |            |             |             |

$$f_{k}^{VI} = P(\mathbf{q}_{t+1}[k] = \mathcal{V} | \mathbf{q}_{t}[i] = \mathcal{N} \quad \forall i)$$
  
$$f_{k}^{VC} = P(\mathbf{q}_{t+1}[k] = \mathcal{V} | \mathbf{q}_{t}[k] = \mathcal{V}, \mathbf{q}_{t}[i] = \mathcal{N} \quad \forall i \neq k)$$

 probability, when j is vocalizing, that k initiates vocalization overlap (OI) and that k continues vocalization overlap (OC)

向下 イヨト イヨト

| Introduction | <b>Computational Framework</b> | Data<br>oo | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Features     |                                |            |             |             |

$$\begin{aligned} f_k^{VI} &= P\left(\mathbf{q}_{t+1}\left[k\right] = \mathcal{V} \,|\, \mathbf{q}_t\left[i\right] = \mathcal{N} \quad \forall i \right) \\ f_k^{VC} &= P\left(\mathbf{q}_{t+1}\left[k\right] = \mathcal{V} \,|\, \mathbf{q}_t\left[k\right] = \mathcal{V}, \, \mathbf{q}_t\left[i\right] = \mathcal{N} \quad \forall i \neq k \right) \end{aligned}$$

 probability, when j is vocalizing, that k initiates vocalization overlap (OI) and that k continues vocalization overlap (OC)

• E • • E •

| Introduction | <b>Computational Framework</b> | Data<br>oo | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Features     |                                |            |             |             |

$$\begin{aligned} f_k^{VI} &= P\left(\mathbf{q}_{t+1}\left[k\right] = \mathcal{V} \mid \mathbf{q}_t\left[i\right] = \mathcal{N} \quad \forall i \right) \\ f_k^{VC} &= P\left(\mathbf{q}_{t+1}\left[k\right] = \mathcal{V} \mid \mathbf{q}_t\left[k\right] = \mathcal{V}, \, \mathbf{q}_t\left[i\right] = \mathcal{N} \quad \forall i \neq k \right) \end{aligned}$$

 probability, when j is vocalizing, that k initiates vocalization overlap (OI) and that k continues vocalization overlap (OC)

OI 
$$t + 1$$
  
 $j \times k$ 

| Introduction | <b>Computational Framework</b> | Data<br>oo | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Features     |                                |            |             |             |

$$\begin{aligned} f_k^{VI} &= P\left(\mathbf{q}_{t+1}\left[k\right] = \mathcal{V} \,|\, \mathbf{q}_t\left[i\right] = \mathcal{N} \quad \forall i \right) \\ f_k^{VC} &= P\left(\mathbf{q}_{t+1}\left[k\right] = \mathcal{V} \,|\, \mathbf{q}_t\left[k\right] = \mathcal{V}, \, \mathbf{q}_t\left[i\right] = \mathcal{N} \quad \forall i \neq k \right) \end{aligned}$$

 probability, when j is vocalizing, that k initiates vocalization overlap (OI) and that k continues vocalization overlap (OC)



A B K A B K

| Introduction | <b>Computational Framework</b> | Data<br>oo | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Features     |                                |            |             |             |

$$\begin{aligned} f_k^{VI} &= P\left(\mathbf{q}_{t+1}\left[k\right] = \mathcal{V} \,|\, \mathbf{q}_t\left[i\right] = \mathcal{N} \quad \forall i \right) \\ f_k^{VC} &= P\left(\mathbf{q}_{t+1}\left[k\right] = \mathcal{V} \,|\, \mathbf{q}_t\left[k\right] = \mathcal{V}, \,\mathbf{q}_t\left[i\right] = \mathcal{N} \quad \forall i \neq k \right) \end{aligned}$$

 probability, when j is vocalizing, that k initiates vocalization overlap (OI) and that k continues vocalization overlap (OC)

$$f_{k,j}^{OI} = P(\mathbf{q}_{t+1}[k] = \mathcal{V} | \mathbf{q}_t[j] = \mathcal{V}, \mathbf{q}_t[i] = \mathcal{N} \quad \forall i \neq j)$$
  

$$f_{k,j}^{OC} = P(\mathbf{q}_{t+1}[k] = \mathcal{V} | \mathbf{q}_t[k] = \mathbf{q}_t[j] = \mathcal{V}, \mathbf{q}_t[i] = \mathcal{N}$$
  

$$\forall i \neq j, i \neq k)$$

• E • • E •

| Introduction | <b>Computational Framework</b> | Data<br>oo | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Feature E    | Estimation                     |            |             |             |

$$\mathbf{F} = \bigcup_{k=1}^{K} \left\{ f_k^{VI}, f_k^{VC}, \bigcup_{j \neq k}^{K} \left\{ f_{k,j}^{OI}, f_{k,j}^{OC} \right\} \right\}$$

discretize the vocal interaction record using 200 ms frames

estimate features using maximum likelihood (ML)

probabilities with unseen conditioning contexts are set to 0.5

・ロン ・回と ・ヨン ・ヨン

| Introduction | <b>Computational Framework</b> | Data<br>00 | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Feature E    | Estimation                     |            |             |             |

$$\mathbf{F} = \bigcup_{k=1}^{K} \left\{ f_k^{VI}, f_k^{VC}, \bigcup_{j \neq k}^{K} \left\{ f_{k,j}^{OI}, f_{k,j}^{OC} \right\} \right\}$$

• discretize the vocal interaction record using 200 ms frames



• estimate features using maximum likelihood (ML)

• probabilities with unseen conditioning contexts are set to 0.5

(日) (同) (E) (E) (E)

| Introduction | <b>Computational Framework</b> | Data<br>oo | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Feature E    | Estimation                     |            |             |             |

$$\mathbf{F} = \bigcup_{k=1}^{K} \left\{ f_k^{VI}, f_k^{VC}, \bigcup_{j \neq k}^{K} \left\{ f_{k,j}^{OI}, f_{k,j}^{OC} \right\} \right\}$$

• discretize the vocal interaction record using 200 ms frames





$$\mathbf{F} = \bigcup_{k=1}^{K} \left\{ f_k^{VI}, f_k^{VC}, \bigcup_{j \neq k}^{K} \left\{ f_{k,j}^{OI}, f_{k,j}^{OC} \right\} \right\}$$

• discretize the vocal interaction record using 200 ms frames

estimate features using maximum likelihood (ML)
probabilities with unseen conditioning contexts are set to 0.5

| Introduction | <b>Computational Framework</b> | Data<br>00 | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Feature E    | Estimation                     |            |             |             |

$$\mathbf{F} = \bigcup_{k=1}^{K} \left\{ f_k^{VI}, f_k^{VC}, \bigcup_{j \neq k}^{K} \left\{ f_{k,j}^{OI}, f_{k,j}^{OC} \right\} \right\}$$

• discretize the vocal interaction record using 200 ms frames

estimate features using maximum likelihood (ML)

probabilities with unseen conditioning contexts are set to 0.5

イロト イポト イヨト イヨト

| Introduction | <b>Computational Framework</b> | Data<br>oo | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Feature E    | Estimation                     |            |             |             |

$$\mathbf{F} = \bigcup_{k=1}^{K} \left\{ f_k^{VI}, f_k^{VC}, \bigcup_{j \neq k}^{K} \left\{ f_{k,j}^{OI}, f_{k,j}^{OC} \right\} \right\}$$

• discretize the vocal interaction record using 200 ms frames



イロト イポト イヨト イヨト

| Introduction | <b>Computational Framework</b> | Data<br>oo | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Feature E    | Estimation                     |            |             |             |

$$\mathbf{F} = \bigcup_{k=1}^{K} \left\{ f_k^{VI}, f_k^{VC}, \bigcup_{j \neq k}^{K} \left\{ f_{k,j}^{OI}, f_{k,j}^{OC} \right\} \right\}$$

• discretize the vocal interaction record using 200 ms frames



• estimate features using maximum likelihood (ML)

probabilities with unseen conditioning contexts are set to 0.5

イロト イポト イヨト イヨト

| Introduction | <b>Computational Framework</b> | Data<br>oo | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Feature E    | Estimation                     |            |             |             |

$$\mathbf{F} = \bigcup_{k=1}^{K} \left\{ f_k^{VI}, f_k^{VC}, \bigcup_{j \neq k}^{K} \left\{ f_{k,j}^{OI}, f_{k,j}^{OC} \right\} \right\}$$

• discretize the vocal interaction record using 200 ms frames



estimate features using maximum likelihood (ML)

• probabilities with unseen conditioning contexts are set to 0.5

・ 同下 ・ ヨト ・ ヨト

| Introduction | <b>Computational Framework</b> | Data<br>oo | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Feature E    | Estimation                     |            |             |             |

$$\mathbf{F} = \bigcup_{k=1}^{K} \left\{ f_k^{VI}, f_k^{VC}, \bigcup_{j \neq k}^{K} \left\{ f_{k,j}^{OI}, f_{k,j}^{OC} \right\} \right\}$$

• discretize the vocal interaction record using 200 ms frames



• estimate features using maximum likelihood (ML)

• probabilities with unseen conditioning contexts are set to 0.5

向下 イヨト イヨト

| Introduction | <b>Computational Framework</b> | Data<br>oo | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Feature E    | Estimation                     |            |             |             |

$$\mathbf{F} = \bigcup_{k=1}^{K} \left\{ f_k^{VI}, f_k^{VC}, \bigcup_{j \neq k}^{K} \left\{ f_{k,j}^{OI}, f_{k,j}^{OC} \right\} \right\}$$

• discretize the vocal interaction record using 200 ms frames



- estimate features using maximum likelihood (ML)
- probabilities with unseen conditioning contexts are set to 0.5

| Introduction | Computational Framework | Data | Experiments | Conclusions |
|--------------|-------------------------|------|-------------|-------------|
|              | 0000000                 |      |             |             |
|              |                         |      |             |             |

- use a variant of a model from stochastic dynamics, the **Ising model** (Glauber, 1963)
- assume the conditional probability of vocal activity state transition, for each k

$$P(\mathbf{q}_{t+1}[k] = \mathcal{V} | \mathbf{q}_t = \mathbf{S}_i) = y_k(\mathbf{S}_i)$$

where

$$y_k(\mathbf{x}) = \frac{1}{1 + e^{-\beta \left(\sum_{j=1}^K w_{k,j} \times_j + b_k\right)}}$$

- not coincidentally, this is a one-layer neural network
- obviates the need for designing a back-off/smoothing strategy in ML estimation of features

| Introduction | Computational Framework | Data | Experiments | Conclusions |
|--------------|-------------------------|------|-------------|-------------|
|              | 0000000                 |      |             |             |
|              |                         |      |             |             |

- use a variant of a model from stochastic dynamics, the **Ising model** (Glauber, 1963)
- assume the conditional probability of vocal activity state transition, for each k

$$P(\mathbf{q}_{t+1}[k] = \mathcal{V} | \mathbf{q}_t = \mathbf{S}_i) = y_k(\mathbf{S}_i)$$

where

$$y_k(\mathbf{x}) = \frac{1}{1 + e^{-\beta \left(\sum_{j=1}^{K} w_{k,j} x_j + b_k\right)}}$$

- not coincidentally, this is a one-layer neural network
- obviates the need for designing a back-off/smoothing strategy in ML estimation of features

| Introduction | Computational Framework | Data | Experiments | Conclusions |
|--------------|-------------------------|------|-------------|-------------|
|              | 0000000                 |      |             |             |
|              |                         |      |             |             |

- use a variant of a model from stochastic dynamics, the Ising model (Glauber, 1963)
- assume the conditional probability of vocal activity state transition, for each k

$$P(\mathbf{q}_{t+1}[k] = \mathcal{V} | \mathbf{q}_t = \mathbf{S}_i) = y_k(\mathbf{S}_i)$$

where

$$y_k(\mathbf{x}) = \frac{1}{1 + e^{-\beta \left(\sum_{j=1}^{K} w_{k,j} \times_j + b_k\right)}}$$

• not coincidentally, this is a one-layer neural network

• obviates the need for designing a back-off/smoothing strategy in ML estimation of features

| Introduction | Computational Framework | Data | Experiments | Conclusions |
|--------------|-------------------------|------|-------------|-------------|
|              | 0000000                 |      |             |             |
|              |                         |      |             |             |

- use a variant of a model from stochastic dynamics, the Ising model (Glauber, 1963)
- assume the conditional probability of vocal activity state transition, for each k

$$P(\mathbf{q}_{t+1}[k] = \mathcal{V} | \mathbf{q}_t = \mathbf{S}_i) = y_k(\mathbf{S}_i)$$

where

$$y_k(\mathbf{x}) = \frac{1}{1 + e^{-\beta \left(\sum_{j=1}^{\kappa} w_{k,j} \times j + b_k\right)}}$$

- not coincidentally, this is a one-layer neural network
- obviates the need for designing a back-off/smoothing strategy in ML estimation of features

| Introduction | Computational Framework | Data | Experiments | Conclusions |
|--------------|-------------------------|------|-------------|-------------|
|              | 0000000                 |      |             |             |
|              |                         |      |             |             |

- use a variant of a model from stochastic dynamics, the Ising model (Glauber, 1963)
- assume the conditional probability of vocal activity state transition, for each k

$$P(\mathbf{q}_{t+1}[k] = \mathcal{V} | \mathbf{q}_t = \mathbf{S}_i) = y_k(\mathbf{S}_i)$$

where

$$y_k(\mathbf{x}) = \frac{1}{1 + e^{-\beta \left(\sum_{j=1}^{K} w_{k,j} x_j + b_k\right)}}$$

- not coincidentally, this is a one-layer neural network
- obviates the need for designing a back-off/smoothing strategy in ML estimation of features

| Introduction | <b>Computational Framework</b> | Data<br>oo | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Behavior M   | lodel                          |            |             |             |

• for each conversation type  $\mathcal{T}$  and each group  $\mathcal{G}$ , require the likelihood of **F** (as estimated from the observed vocal interaction record)

$$P(\mathbf{F} | \mathcal{G}, \mathcal{T}) = \prod_{k=1}^{K} P\left(f_{k}^{VI} | \theta_{T,\mathcal{G}[k]}^{VI}\right) P\left(f_{k}^{VC} | \theta_{T,\mathcal{G}[k]}^{VC}\right) \\ \times \prod_{j \neq k}^{K} P\left(f_{k,j}^{OI} | \theta_{T,\mathcal{G}[k],\mathcal{G}[j]}^{OI}\right) P\left(f_{k,j}^{OC} | \theta_{T,\mathcal{G}[k],\mathcal{G}[j]}^{OC}\right)$$

• each  $\theta$  represents a single one-dimensional Gaussian mean  $\mu$  and variance  $\Sigma$  pair

コント イヨン イヨン

| Introduction | <b>Computational Framework</b> | Data<br>oo | Experiments | Conclusions |
|--------------|--------------------------------|------------|-------------|-------------|
|              |                                |            |             |             |
| Behavior M   | lodel                          |            |             |             |

• for each conversation type  $\mathcal{T}$  and each group  $\mathcal{G}$ , require the likelihood of **F** (as estimated from the observed vocal interaction record)

$$P(\mathbf{F} | \mathcal{G}, \mathcal{T}) = \prod_{k=1}^{K} P\left(f_{k}^{VI} | \theta_{T,\mathcal{G}[k]}^{VI}\right) P\left(f_{k}^{VC} | \theta_{T,\mathcal{G}[k]}^{VC}\right) \\ \times \prod_{j \neq k}^{K} P\left(f_{k,j}^{OI} | \theta_{T,\mathcal{G}[k],\mathcal{G}[j]}^{OI}\right) P\left(f_{k,j}^{OC} | \theta_{T,\mathcal{G}[k],\mathcal{G}[j]}^{OC}\right)$$

• each  $\theta$  represents a single one-dimensional Gaussian mean  $\mu$  and variance  $\Sigma$  pair

4 B K 4 B K

| Introduction | Computational Framework | Data<br>00 | Experiments | Conclusions |
|--------------|-------------------------|------------|-------------|-------------|
|              |                         |            |             |             |
| Members      | hip Model               |            |             |             |

• for each conversation type  ${\cal T},$  require the probability of group  ${\cal G}$  (as hypothesized)

$$P(\mathcal{G} \mid \mathcal{T}) = \frac{1}{Z_{\mathcal{G}}} \prod_{k=1}^{K} P(\mathcal{G}[k] \mid \mathcal{T})$$

•  $Z_{\mathcal{G}}$  is a normalization constant,  $\sum_{N_{\mathcal{G}}} P\left(\mathcal{G} | \mathcal{T} 
ight) = 1$ 

白 と く ヨ と く ヨ と

| Introduction | Computational Framework | Data<br>00 | Experiments | Conclusions |
|--------------|-------------------------|------------|-------------|-------------|
|              |                         |            |             |             |
| Members      | hip Model               |            |             |             |

for each conversation type *T*, require the probability of group *G* (as hypothesized)

$$P(\mathcal{G} \mid \mathcal{T}) = \frac{1}{Z_{\mathcal{G}}} \prod_{k=1}^{K} P(\mathcal{G}[k] \mid \mathcal{T})$$

•  $Z_{\mathcal{G}}$  is a normalization constant,  $\sum_{N_{\mathcal{G}}} P\left(\mathcal{G} | \mathcal{T} 
ight) = 1$ 

. . . . . . . .

| Introduction | Computational Framework | Data<br>●○ | Experiments | Conclusions |
|--------------|-------------------------|------------|-------------|-------------|
|              |                         |            |             |             |

## The ICSI Meeting Corpus

#### (Janin et al, 2003), (Shriberg et al, 2004)

- naturally occurring project-oriented conversations
- for our purposes, 4 types of longitudinal collections:

- "other" contains types of which there are ≤3 meetings
- rarely, meetings contain additional, uninstrumented participants (whose contributions we ignore)

- 4 同 5 - 4 日 5 - 4 日

| Introduction | Computational Framework | Data<br>●○ | Experiments | Conclusions |
|--------------|-------------------------|------------|-------------|-------------|
|              |                         |            |             |             |
| The ICSI     | Meeting Corpus          |            |             |             |

- naturally occurring project-oriented conversations
- for our purposes, 4 types of longitudinal collections:

- "other" contains types of which there are  $\leq$ 3 meetings
- rarely, meetings contain additional, uninstrumented participants (whose contributions we ignore)

(4) (3) (4) (4) (4)

| Introduction            | Computational Framework | Data<br>●○ | Experiments | Conclusions |  |  |
|-------------------------|-------------------------|------------|-------------|-------------|--|--|
|                         |                         |            |             |             |  |  |
| The ICSI Meeting Corpus |                         |            |             |             |  |  |

- naturally occurring project-oriented conversations
- for our purposes, 4 types of longitudinal collections:

| type  | # of     | # of possible $\#$ of partic |     | partici | pants |
|-------|----------|------------------------------|-----|---------|-------|
| type  | meetings | participants                 | mod | min     | max   |
| Bed   | 15       | 13                           | 6   | 4       | 7     |
| Bmr   | 29       | 15                           | 7   | 3       | 9     |
| Bro   | 23       | 10                           | 6   | 4       | 8     |
| other | 8        | 27                           | 6   | 5       | 8     |

● "other" contains types of which there are ≤3 meetings

 rarely, meetings contain additional, uninstrumented participants (whose contributions we ignore)

(3)

| Introduction            | Computational Framework | Data<br>●○ | Experiments | Conclusions |  |  |
|-------------------------|-------------------------|------------|-------------|-------------|--|--|
|                         |                         |            |             |             |  |  |
| The ICSI Meeting Corpus |                         |            |             |             |  |  |

- naturally occurring project-oriented conversations
- for our purposes, 4 types of longitudinal collections:

| type  | # of     | # of possible $\#$ of partic |     | partici | pants |
|-------|----------|------------------------------|-----|---------|-------|
| type  | meetings | participants                 | mod | min     | max   |
| Bed   | 15       | 13                           | 6   | 4       | 7     |
| Bmr   | 29       | 15                           | 7   | 3       | 9     |
| Bro   | 23       | 10                           | 6   | 4       | 8     |
| other | 8        | 27                           | 6   | 5       | 8     |

- "other" contains types of which there are ≤3 meetings
- rarely, meetings contain additional, uninstrumented participants (whose contributions we ignore)

| Introduction            | Computational Framework | Data<br>●○ | Experiments | Conclusions |  |
|-------------------------|-------------------------|------------|-------------|-------------|--|
|                         |                         |            |             |             |  |
| The ICSI Meeting Corpus |                         |            |             |             |  |

- naturally occurring project-oriented conversations
- for our purposes, 4 types of longitudinal collections:

| type  | # of     | # of possible $\#$ of partic |     | partici | pants |
|-------|----------|------------------------------|-----|---------|-------|
| type  | meetings | participants                 | mod | min     | max   |
| Bed   | 15       | 13                           | 6   | 4       | 7     |
| Bmr   | 29       | 15                           | 7   | 3       | 9     |
| Bro   | 23       | 10                           | 6   | 4       | 8     |
| other | 8        | 27                           | 6   | 5       | 8     |

- "other" contains types of which there are  $\leq$ 3 meetings
- rarely, meetings contain additional, uninstrumented participants (whose contributions we ignore)

| Introduction | Computational Framework      | Data<br>○● | Experiments | Conclusions |
|--------------|------------------------------|------------|-------------|-------------|
|              |                              |            |             |             |
| Differences  | Between Meeting <sup>-</sup> | Fypes      |             |             |

• 36-minute excerpts (from 1000 sec to 2000 sec)



| Introduction | Computational Framework | Data<br>oo | Experiments<br>●○○ | Conclusions |
|--------------|-------------------------|------------|--------------------|-------------|
|              |                         |            |                    |             |
| Baseline     |                         |            |                    |             |

- use inclusive-OR of "talk-spurt" (Shriberg et al, 2001) and "laugh-bout" (Laskowski & Burger, 2007) segmentations
- compute a single feature f<sup>I</sup><sub>k</sub>: vocalizing time proportion
   employed for assessing speaker diarization performance (Jin et al., 2004). (Mirghafori & Wooters, 2006)
   e continue (Hamed of speaker diarization distribution across speaker)
- leave-one-out classification

- cluster participants for training the behavior model
- accuracy: 65.7% (random guessing: 43%)

( ) < </p>

| Introduction | Computational Framework | Data<br>00 | Experiments<br>●○○ | Conclusions |
|--------------|-------------------------|------------|--------------------|-------------|
|              |                         |            |                    |             |
| Baseline     |                         |            |                    |             |

- use inclusive-OR of "talk-spurt" (Shriberg et al, 2001) and "laugh-bout" (Laskowski & Burger, 2007) segmentations
- compute a single feature  $f_k^T$ : vocalizing time proportion
  - employed for assessing speaker diarization performance (Jin et al, 2004), (Mirghafori & Wooters, 2006)
  - captures "flatness" of speaking-time distribution across speakers
- leave-one-out classification

- cluster participants for training the behavior model
- accuracy: 65.7% (random guessing: 43%)

| Introduction | Computational Framework | Data<br>00 | Experiments<br>●○○ | Conclusions |
|--------------|-------------------------|------------|--------------------|-------------|
|              |                         |            |                    |             |
| Baseline     |                         |            |                    |             |

- use inclusive-OR of "talk-spurt" (Shriberg et al, 2001) and "laugh-bout" (Laskowski & Burger, 2007) segmentations
- compute a single feature  $f_k^T$ : vocalizing time proportion
  - employed for assessing speaker diarization performance (Jin et al, 2004), (Mirghafori & Wooters, 2006)
  - captures "flatness" of speaking-time distribution across speakers
- leave-one-out classification
  - train on 65 meetings, test on 1 meeting, rotate
  - too little data for a true, unseen evaluation set
- cluster participants for training the behavior model
- accuracy: 65.7% (random guessing: 43%)

| Introduction | Computational Framework | Data<br>00 | Experiments<br>●○○ | Conclusions |
|--------------|-------------------------|------------|--------------------|-------------|
|              |                         |            |                    |             |
| Baseline     |                         |            |                    |             |

- use inclusive-OR of "talk-spurt" (Shriberg et al, 2001) and "laugh-bout" (Laskowski & Burger, 2007) segmentations
- compute a single feature  $f_k^T$ : vocalizing time proportion
  - employed for assessing speaker diarization performance (Jin et al, 2004), (Mirghafori & Wooters, 2006)
  - captures "flatness" of speaking-time distribution across speakers
- leave-one-out classification
  - train on 65 meetings, test on 1 meeting, rotate
  - too little data for a true, unseen evaluation set
- cluster participants for training the behavior model
- accuracy: 65.7% (random guessing: 43%)

( ) < </p>

| Introduction | Computational Framework | Data<br>00 | Experiments<br>●○○ | Conclusions |
|--------------|-------------------------|------------|--------------------|-------------|
|              |                         |            |                    |             |
| Baseline     |                         |            |                    |             |

- use inclusive-OR of "talk-spurt" (Shriberg et al, 2001) and "laugh-bout" (Laskowski & Burger, 2007) segmentations
- compute a single feature  $f_k^T$ : vocalizing time proportion
  - employed for assessing speaker diarization performance (Jin et al, 2004), (Mirghafori & Wooters, 2006)
  - captures "flatness" of speaking-time distribution across speakers
- leave-one-out classification
  - train on 65 meetings, test on 1 meeting, rotate
  - too little data for a true, unseen evaluation set
- cluster participants for training the behavior model
  - o side-effect: renders impact of membership model negligible
- accuracy: 65.7% (random guessing: 43%)

向下 イヨト イヨ
| Introduction | Computational Framework | Data<br>00 | Experiments<br>●○○ | Conclusions |
|--------------|-------------------------|------------|--------------------|-------------|
|              |                         |            |                    |             |
| Baseline     |                         |            |                    |             |

- use inclusive-OR of "talk-spurt" (Shriberg et al, 2001) and "laugh-bout" (Laskowski & Burger, 2007) segmentations
- compute a single feature  $f_k^T$ : vocalizing time proportion
  - employed for assessing speaker diarization performance (Jin et al, 2004), (Mirghafori & Wooters, 2006)
  - captures "flatness" of speaking-time distribution across speakers
- leave-one-out classification
  - train on 65 meetings, test on 1 meeting, rotate
  - too little data for a true, unseen evaluation set
- cluster participants for training the behavior model

side-effect: renders impact of membership model negligible.

accuracy: 65.7% (random guessing: 43%)

| Introduction | Computational Framework | Data<br>00 | Experiments<br>●○○ | Conclusions |
|--------------|-------------------------|------------|--------------------|-------------|
|              |                         |            |                    |             |
| Baseline     |                         |            |                    |             |

- use inclusive-OR of "talk-spurt" (Shriberg et al, 2001) and "laugh-bout" (Laskowski & Burger, 2007) segmentations
- compute a single feature  $f_k^T$ : vocalizing time proportion
  - employed for assessing speaker diarization performance (Jin et al, 2004), (Mirghafori & Wooters, 2006)
  - captures "flatness" of speaking-time distribution across speakers
- leave-one-out classification
  - train on 65 meetings, test on 1 meeting, rotate
  - too little data for a true, unseen evaluation set
- cluster participants for training the behavior model
  - side-effect: renders impact of membership model negligible
- accuracy: 65.7% (random guessing: 43%)

| Introduction | Computational Framework | Data<br>00 | Experiments<br>●○○ | Conclusions |
|--------------|-------------------------|------------|--------------------|-------------|
|              |                         |            |                    |             |
| Baseline     |                         |            |                    |             |

- use inclusive-OR of "talk-spurt" (Shriberg et al, 2001) and "laugh-bout" (Laskowski & Burger, 2007) segmentations
- compute a single feature  $f_k^T$ : vocalizing time proportion
  - employed for assessing speaker diarization performance (Jin et al, 2004), (Mirghafori & Wooters, 2006)
  - captures "flatness" of speaking-time distribution across speakers
- leave-one-out classification
  - train on 65 meetings, test on 1 meeting, rotate
  - too little data for a true, unseen evaluation set
- cluster participants for training the behavior model
  - side-effect: renders impact of membership model negligible

accuracy: 65.7% (random guessing: 43%)

(4月) (4日) (4日)

| Introduction | Computational Framework | Data<br>00 | Experiments<br>●○○ | Conclusions |
|--------------|-------------------------|------------|--------------------|-------------|
|              |                         |            |                    |             |
| Baseline     |                         |            |                    |             |

- use inclusive-OR of "talk-spurt" (Shriberg et al, 2001) and "laugh-bout" (Laskowski & Burger, 2007) segmentations
- compute a single feature  $f_k^T$ : vocalizing time proportion
  - employed for assessing speaker diarization performance (Jin et al, 2004), (Mirghafori & Wooters, 2006)
  - captures "flatness" of speaking-time distribution across speakers
- leave-one-out classification
  - train on 65 meetings, test on 1 meeting, rotate
  - too little data for a true, unseen evaluation set
- cluster participants for training the behavior model
  - side-effect: renders impact of membership model negligible
- accuracy: 65.7% (random guessing: 43%)

向下 イヨト イヨト

| Introduction | Computational Framework | <b>Data</b><br>00 | Experiments | Conclusions |
|--------------|-------------------------|-------------------|-------------|-------------|
|              |                         |                   |             |             |

| Feature(s)                       | ML Esti              | ML Estimation |               | NN Estimation |  |
|----------------------------------|----------------------|---------------|---------------|---------------|--|
| reactive(3)                      | w/o $f_k^{\ \prime}$ | w/ $f_k^I$    | w/o $f_k^{I}$ | w/ $f_k^I$    |  |
| baseline                         | —                    | 65.7          | —             | 65.7          |  |
| $f_k^{VI}$                       | 59.7                 | 67.2          | 56.7          | 65.7          |  |
| $f_k^{VC}$                       | 62.7                 | 77.6          | 56.7          | 71.6          |  |
| $\langle f_{k,i}^{OI} \rangle_i$ | 35.8                 | 52.2          | 64.2          | 67.2          |  |
| $\langle f_{k,i}^{OC} \rangle_i$ | 53.7                 | 67.2          | 64.2          | 80.6          |  |
| $f_{k,i}^{Oi}$                   | 41.8                 | 46.3          | 67.2          | 64.2          |  |
| $f_{k,j}^{OC}$                   | 61.2                 | 68.7          | 73.1          | 79.1          |  |

4

э

< ≣ >

| Introduction | Computational Framework | <b>Data</b><br>00 | Experiments | Conclusions |
|--------------|-------------------------|-------------------|-------------|-------------|
|              |                         |                   |             |             |

| Feature(s)                               | ML Estimation |            | NN Estimation |            |
|--|---------------|------------|---------------|------------|
| r cature(3)                              | w/o $f_k^{I}$ | w/ $f_k^I$ | w/o $f_k^I$   | w/ $f_k^I$ |
| baseline                                 | —             | 65.7       |               | 65.7       |
| $f_k^{VI}$                               | 59.7          | 67.2       | 56.7          | 65.7       |
| $f_k^{VC}$                               | 62.7          | 77.6       | 56.7          | 71.6       |
| $\langle \hat{f}_{k,i}^{OI} \rangle_{i}$ | 35.8          | 52.2       | 64.2          | 67.2       |
| $\langle f_{k,i}^{OC} \rangle_i$         | 53.7          | 67.2       | 64.2          | 80.6       |
| $f_{k,i}^{Oi}$                           | 41.8          | 46.3       | 67.2          | 64.2       |
| $f_{k,j}^{OC}$                           | 61.2          | 68.7       | 73.1          | 79.1       |

• the baseline feature  $f_k^T$  outperforms most other features

| Introduction | Computational Framework | <b>Data</b><br>00 | Experiments | Conclusions |
|--------------|-------------------------|-------------------|-------------|-------------|
|              |                         |                   |             |             |

| Feature(s)                         | ML Esti              | ML Estimation |               | NN Estimation |  |
|------------------------------------|----------------------|---------------|---------------|---------------|--|
| reactive(3)                        | w/o $f_k^{\ \prime}$ | w/ $f_k^I$    | w/o $f_k^{I}$ | w/ $f_k^I$    |  |
| baseline                           | —                    | 65.7          | —             | 65.7          |  |
| $f_k^{VI}$                         | 59.7                 | 67.2          | 56.7          | 65.7          |  |
| $f_k^{VC}$                         | 62.7                 | 77.6          | 56.7          | 71.6          |  |
| $\langle f_{k,j}^{OI} \rangle_{j}$ | 35.8                 | 52.2          | 64.2          | 67.2          |  |
| $\langle f_{k,i}^{OC} \rangle_i$   | 53.7                 | 67.2          | 64.2          | 80.6          |  |
| $f_{k,j}^{OI}$                     | 41.8                 | 46.3          | 67.2          | 64.2          |  |
| $f_{k,j}^{OC}$                     | 61.2                 | 68.7          | 73.1          | 79.1          |  |

• by themselves, specific participant-pair features outperform each participant's average participant-pair features

| Introduction | Computational Framework | <b>Data</b><br>00 | Experiments | Conclusions |
|--------------|-------------------------|-------------------|-------------|-------------|
|              |                         |                   |             |             |

| Feature(s)                         | ML Estimation        |            | NN Estimation |            |
|------------------------------------|----------------------|------------|---------------|------------|
| r cature(3)                        | w/o $f_k^{\ \prime}$ | w/ $f_k^I$ | w/o $f_k^{I}$ | w/ $f_k^I$ |
| baseline                           | —                    | 65.7       | —             | 65.7       |
| $f_k^{VI}$                         | 59.7                 | 67.2       | 56.7          | 65.7       |
| $f_k^{VC}$                         | 62.7                 | 77.6       | 56.7          | 71.6       |
| $\langle f_{k,j}^{OI} \rangle_{i}$ | 35.8                 | 52.2       | 64.2          | 67.2       |
| $\langle f_{k,i}^{OC} \rangle_i$   | 53.7                 | 67.2       | 64.2          | 80.6       |
| $f_{k,i}^{OI}$                     | 41.8                 | 46.3       | 67.2          | 64.2       |
| $f_{k,j}^{OC}$                     | 61.2                 | 68.7       | 73.1          | 79.1       |

• most features, when combined with  $f_k^T$ , lead to improved performance

(B) (B)

| Introduction | Computational Framework | <b>Data</b><br>00 | Experiments | Conclusions |
|--------------|-------------------------|-------------------|-------------|-------------|
|              |                         |                   |             |             |

| Feature(s)                               | ML Estimation        |            | NN Estimation |            |
|--|----------------------|------------|---------------|------------|
| reactive(3)                              | w/o $f_k^{\ \prime}$ | w/ $f_k^I$ | w/o $f_k^{I}$ | w/ $f_k^I$ |
| baseline                                 | —                    | 65.7       | —             | 65.7       |
| $f_k^{VI}$                               | 59.7                 | 67.2       | 56.7          | 65.7       |
| $f_k^{VC}$                               | 62.7                 | 77.6       | 56.7          | 71.6       |
| $\langle \hat{f}_{k,j}^{OI} \rangle_{j}$ | 35.8                 | 52.2       | 64.2          | 67.2       |
| $\langle f_{k,i}^{OC} \rangle_i$         | 53.7                 | 67.2       | 64.2          | 80.6       |
| $f_{k,i}^{Oi}$                           | 41.8                 | 46.3       | 67.2          | 64.2       |
| $f_{k,j}^{OC}$                           | 61.2                 | 68.7       | 73.1          | 79.1       |

- all NN-estimated features together yield 82.1%
- an optimal NN-estimated feature subset (forward selection) yields 83.6%

| Introduction | Computational Framework | <b>Data</b><br>00 | Experiments | Conclusions |
|--------------|-------------------------|-------------------|-------------|-------------|
|              |                         |                   |             |             |

| Feature(s)                               | ML Estimation        |            | NN Estimation |            |
|--|----------------------|------------|---------------|------------|
| reactive(3)                              | w/o $f_k^{\ \prime}$ | w/ $f_k^I$ | w/o $f_k^{I}$ | w/ $f_k^I$ |
| baseline                                 | —                    | 65.7       | —             | 65.7       |
| $f_k^{VI}$                               | 59.7                 | 67.2       | 56.7          | 65.7       |
| $f_k^{VC}$                               | 62.7                 | 77.6       | 56.7          | 71.6       |
| $\langle \hat{f}_{k,j}^{OI} \rangle_{j}$ | 35.8                 | 52.2       | 64.2          | 67.2       |
| $\langle f_{k,i}^{OC} \rangle_i$         | 53.7                 | 67.2       | 64.2          | 80.6       |
| $f_{k,i}^{Oi}$                           | 41.8                 | 46.3       | 67.2          | 64.2       |
| $f_{k,j}^{OC}$                           | 61.2                 | 68.7       | 73.1          | 79.1       |

- all NN-estimated features together yield 82.1%
- an optimal NN-estimated feature subset (forward selection) yields **83.6%**

| Introduction | Computational Framework | Data<br>00 | Experiments<br>○○● | Conclusions |
|--------------|-------------------------|------------|--------------------|-------------|
|              |                         |            |                    |             |
| Discussion   |                         |            |                    |             |

| Estimated | Actual Type |                |     |  |
|-----------|-------------|----------------|-----|--|
| LStimated | Bed         | $\mathtt{Bmr}$ | Bro |  |
| Bed       | 11          | 1              | 3   |  |
| Bmr       | 2           | 26             | 1   |  |
| Bro       | 3           | 1              | 19  |  |

- Bmr (discussions among peers) is the most distinct type
- Bed and Bro (both more structured meetings) are more similar to each other than to Bmr
- miss-classification pattern reflects intuition

回 と く ヨ と く ヨ と

| Introduction | Computational Framework | Data<br>00 | Experiments<br>○○● | Conclusions |
|--------------|-------------------------|------------|--------------------|-------------|
|              |                         |            |                    |             |
| Discussion   |                         |            |                    |             |

| Estimated      | Actual Type |                |     |  |
|----------------|-------------|----------------|-----|--|
| LStimated      | Bed         | $\mathtt{Bmr}$ | Bro |  |
| Bed            | 11          | 1              | 3   |  |
| $\mathtt{Bmr}$ | 2           | 26             | 1   |  |
| Bro            | 3           | 1              | 19  |  |

- Bmr (discussions among peers) is the most distinct type
- Bed and Bro (both more structured meetings) are more similar to each other than to Bmr
- miss-classification pattern reflects intuition

| Introduction | Computational Framework | Data<br>00 | Experiments<br>○○● | Conclusions |
|--------------|-------------------------|------------|--------------------|-------------|
|              |                         |            |                    |             |
| Discussion   |                         |            |                    |             |

| Estimated      | Actual Type |                |     |  |
|----------------|-------------|----------------|-----|--|
| LStimated      | Bed         | $\mathtt{Bmr}$ | Bro |  |
| Bed            | 11          | 1              | 3   |  |
| $\mathtt{Bmr}$ | 2           | 26             | 1   |  |
| Bro            | 3           | 1              | 19  |  |

- Bmr (discussions among peers) is the most distinct type
- Bed and Bro (both more structured meetings) are more similar to each other than to Bmr
- miss-classification pattern reflects intuition

| Introduction | Computational Framework | Data<br>00 | Experiments<br>○○● | Conclusions |
|--------------|-------------------------|------------|--------------------|-------------|
|              |                         |            |                    |             |
| Discussion   |                         |            |                    |             |

| Estimated | Ac  | Actual Type    |     |  |
|-----------|-----|----------------|-----|--|
| LStimated | Bed | $\mathtt{Bmr}$ | Bro |  |
| Bed       | 11  | 1              | 3   |  |
| Bmr       | 2   | 26             | 1   |  |
| Bro       | 3   | 1              | 19  |  |

- Bmr (discussions among peers) is the most distinct type
- Bed and Bro (both more structured meetings) are more similar to each other than to Bmr
- miss-classification pattern reflects intuition

| Introduction | Computational Framework | Data<br>00 | Experiments | Conclusions<br>●○○ |
|--------------|-------------------------|------------|-------------|--------------------|
|              |                         |            |             |                    |
| Conclusio    | ons                     |            |             |                    |

#### • classification paradigm with several novel elements:

- exclusively text-independent features, from vocal interaction patterns
- Participant groups, allowing for modeling multi-participant behaviors
- $\bigcirc$  Ising model assumption of C transition probabilities
- meeting sub-type classification accuracy: 83%
- relative error reduction of 52% over the baseline
- (specific) multi-participant interaction features play a large role in this improvement

(4回) (4回) (4回)

| Introduction | Computational Framework | Data<br>00 | Experiments | Conclusions<br>●○○ |
|--------------|-------------------------|------------|-------------|--------------------|
|              |                         |            |             |                    |
| Conclusio    | ons                     |            |             |                    |

- classification paradigm with several novel elements:
  - exclusively text-independent features, from vocal interaction patterns
  - participant groups, allowing for modeling multi-participant behaviors
  - Ising model assumption of C transition probabilities
- meeting sub-type classification accuracy: 83%
- relative error reduction of 52% over the baseline
- (specific) multi-participant interaction features play a large role in this improvement

(不同) とうり くうり

| Introduction | Computational Framework | Data<br>00 | Experiments | Conclusions<br>●○○ |
|--------------|-------------------------|------------|-------------|--------------------|
|              |                         |            |             |                    |
| Conclusio    | ons                     |            |             |                    |

- classification paradigm with several novel elements:
  - exclusively text-independent features, from vocal interaction patterns
  - Participant groups, allowing for modeling multi-participant behaviors
  - $\bigcirc$  Ising model assumption of  $\mathcal C$  transition probabilities
- meeting sub-type classification accuracy: 83%
- relative error reduction of 52% over the baseline
- (specific) multi-participant interaction features play a large role in this improvement

- 4 回 2 - 4 三 2 - 4 三 2

| Introduction | Computational Framework | Data<br>00 | Experiments | Conclusions<br>●○○ |
|--------------|-------------------------|------------|-------------|--------------------|
|              |                         |            |             |                    |
| Conclusio    | ons                     |            |             |                    |

- classification paradigm with several novel elements:
  - exclusively text-independent features, from vocal interaction patterns
  - Participant groups, allowing for modeling multi-participant behaviors
  - 3 Ising model assumption of  $\mathcal C$  transition probabilities
- meeting sub-type classification accuracy: 83%
- relative error reduction of 52% over the baseline
- (specific) multi-participant interaction features play a large role in this improvement

| Introduction | Computational Framework | Data<br>00 | Experiments | Conclusions<br>●○○ |
|--------------|-------------------------|------------|-------------|--------------------|
|              |                         |            |             |                    |
| Conclusio    | ons                     |            |             |                    |

- classification paradigm with several novel elements:
  - exclusively text-independent features, from vocal interaction patterns
  - Participant groups, allowing for modeling multi-participant behaviors
  - 3 Ising model assumption of  $\mathcal C$  transition probabilities
- meeting sub-type classification accuracy: 83%
- relative error reduction of 52% over the baseline
- (specific) multi-participant interaction features play a large role in this improvement

(人間) とうり くうり

| Introduction | Computational Framework | Data<br>00 | Experiments | Conclusions<br>●○○ |
|--------------|-------------------------|------------|-------------|--------------------|
|              |                         |            |             |                    |
| Conclusio    | ons                     |            |             |                    |

- classification paradigm with several novel elements:
  - exclusively text-independent features, from vocal interaction patterns
  - Participant groups, allowing for modeling multi-participant behaviors
  - 3 Ising model assumption of  $\mathcal C$  transition probabilities
- meeting sub-type classification accuracy: 83%
- relative error reduction of 52% over the baseline
- (specific) multi-participant interaction features play a large role in this improvement

- 4 回 ト 4 ヨ ト 4 ヨ ト

| Introduction | Computational Framework | Data<br>00 | Experiments | Conclusions<br>●○○ |
|--------------|-------------------------|------------|-------------|--------------------|
|              |                         |            |             |                    |
| Conclusio    | ons                     |            |             |                    |

- classification paradigm with several novel elements:
  - exclusively text-independent features, from vocal interaction patterns
  - Participant groups, allowing for modeling multi-participant behaviors
  - 3 Ising model assumption of  $\mathcal C$  transition probabilities
- meeting sub-type classification accuracy: 83%
- relative error reduction of 52% over the baseline
- (specific) multi-participant interaction features play a large role in this improvement

・ 同下 ・ ヨト ・ ヨト

| Introduction | Computational Framework | Data<br>oo | Experiments | Conclusions<br>○●○ |
|--------------|-------------------------|------------|-------------|--------------------|
|              |                         |            |             |                    |
| Future W     | /ork                    |            |             |                    |

#### • use automatic, rather than manual, segmentation

- include verbal (words, DAs) features
- explore the dual problem of role/participant detection:

$$\begin{aligned} \mathcal{G}^* &= \arg \max_{\mathcal{G}} P(\mathcal{G} | \mathbf{F}) \\ &= \arg \max_{\mathcal{G}} \sum_{\mathcal{T}} P(\mathcal{G}, \mathcal{T}, \mathbf{F}) \\ &= \arg \max_{\mathcal{G}} \sum_{\mathcal{T}} P(\mathcal{T}) \times P(\mathcal{G} | \mathcal{T}) \times P(\mathbf{F} | \mathcal{G}, \mathcal{T}) \end{aligned}$$

回 と く ヨ と く ヨ と

| Introduction | Computational Framework | Data<br>00 | Experiments<br>000 | Conclusions<br>○●○ |
|--------------|-------------------------|------------|--------------------|--------------------|
| Future W     | ork                     |            |                    |                    |

- use automatic, rather than manual, segmentation
- include verbal (words, DAs) features
- explore the dual problem of role/participant detection:

$$\begin{aligned} \mathcal{J}^* &= \arg \max_{\mathcal{G}} P(\mathcal{G} | \mathbf{F}) \\ &= \arg \max_{\mathcal{G}} \sum_{\mathcal{T}} P(\mathcal{G}, \mathcal{T}, \mathbf{F}) \\ &= \arg \max_{\mathcal{G}} \sum_{\mathcal{T}} P(\mathcal{T}) \times P(\mathcal{G} | \mathcal{T}) \times P(\mathbf{F} | \mathcal{G}, \mathcal{T}) \end{aligned}$$

• • = • • = •

| Introduction | Computational Framework | Data<br>00 | Experiments | Conclusions<br>○●○ |
|--------------|-------------------------|------------|-------------|--------------------|
|              |                         |            |             |                    |
| Future W     | /ork                    |            |             |                    |

- use automatic, rather than manual, segmentation
- include verbal (words, DAs) features
- explore the dual problem of role/participant detection:

$$\begin{aligned} \mathcal{G}^* &= \arg \max_{\mathcal{G}} P(\mathcal{G} | \mathbf{F}) \\ &= \arg \max_{\mathcal{G}} \sum_{\mathcal{T}} P(\mathcal{G}, \mathcal{T}, \mathbf{F}) \\ &= \arg \max_{\mathcal{G}} \sum_{\mathcal{T}} P(\mathcal{T}) \times P(\mathcal{G} | \mathcal{T}) \times P(\mathbf{F} | \mathcal{G}, \mathcal{T}) \end{aligned}$$

| Introduction | Computational Framework | Data<br>oo | Experiments | Conclusions<br>○○● |
|--------------|-------------------------|------------|-------------|--------------------|
|              |                         |            |             |                    |
| Thanks!      |                         |            |             |                    |

We'd also like to thank:

- Liz Shriberg
  - lots of helpful discussion
  - access to the ICSI MRDA annotation
- CHIL project for funding

2

- E