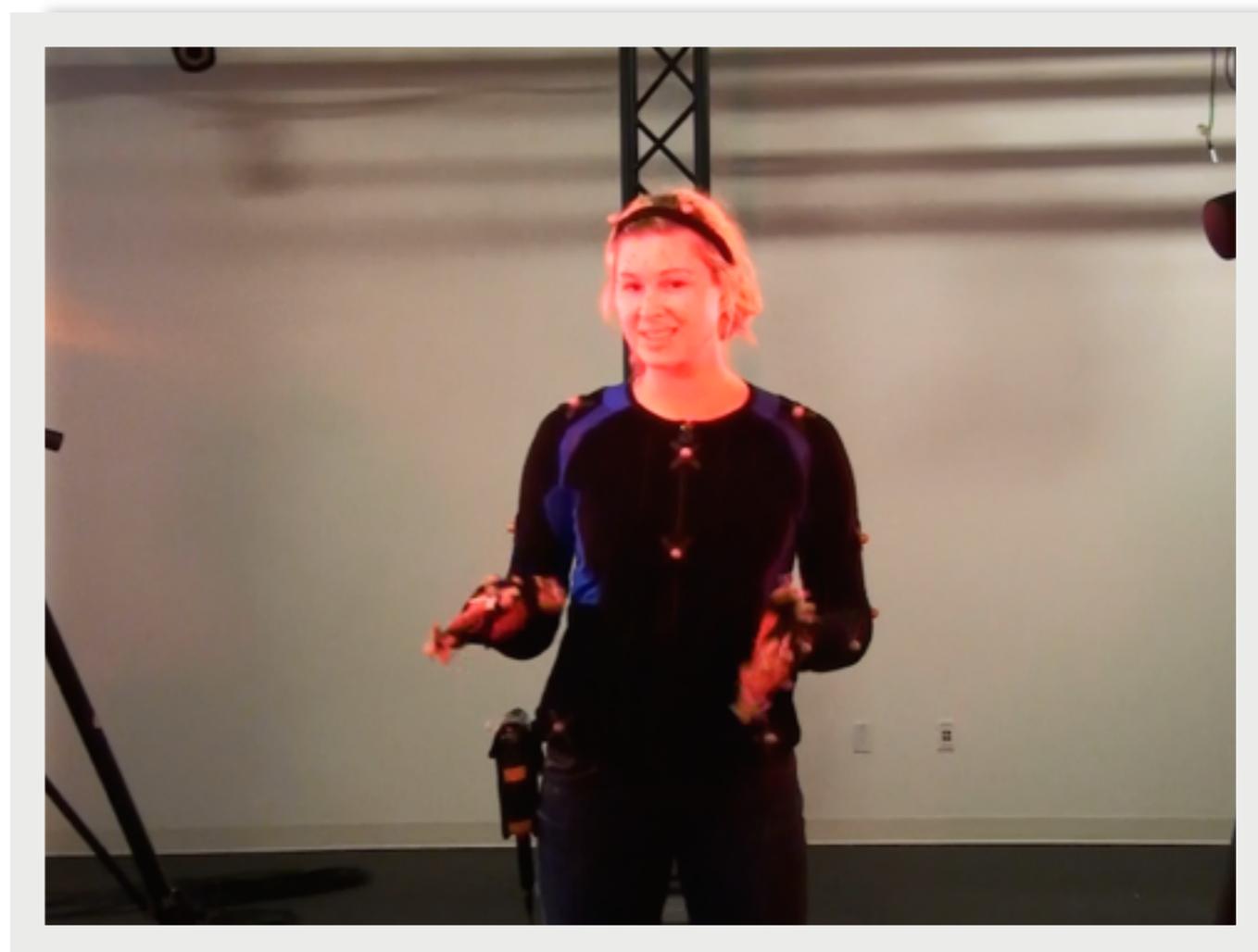


MSP-AVATAR Corpus:

Motion Capture Recordings to Study the Role of Discourse Functions in the Design of Intelligent Virtual Agents

NAJMEH SADOUGHI, YANG LIU AND CARLOS BUSO

Multimodal Signal Processing (MSP) lab
The University of Texas at Dallas
Erik Jonsson School of Engineering and Computer Science



May 8th, 2015

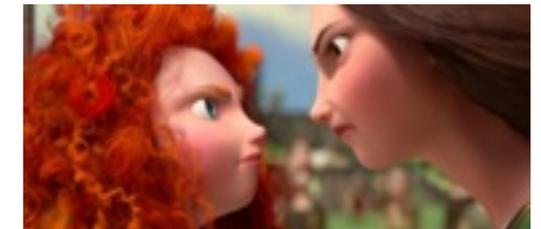


Motivation

- Synthesizing human-like behaviors
 - Animation
 - Entertainment
 - Virtual reality
 - Tutoring/training systems
- Multiple verbal and nonverbal behaviors
 - Head motion, facial expressions, hand gestures, and body postures
 - Gesture coordinated, and synchronized with speech



ICT-USC



[maxresdefault.jpg](#)

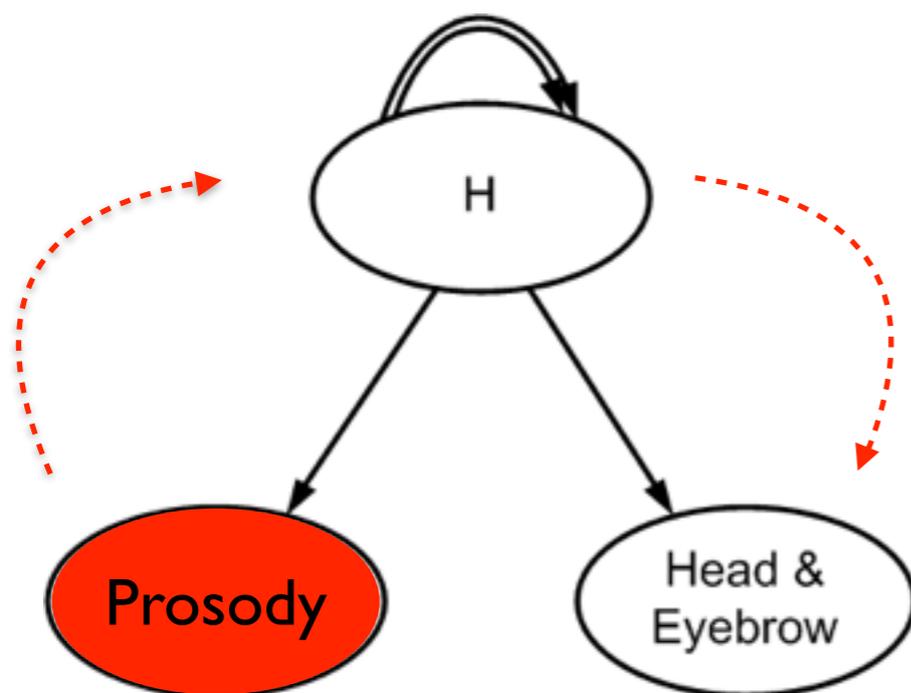
Complex relationship between gestures and speech has to be carefully considered

Previous Studies on Talking Faces

- Rule-based systems [Cassell et al., 1994; S. Kopp 2006]
 - + Semantic meaning of behaviors (nodding)
 - Repetitive behaviors for a given gesture
- Data-driven methods [Levine et al., 2010; Busso et al. 2007]
 - Prosodic features are effective modalities to synthesize nonverbal human-like behaviors

Speech-Driven Animation

- Dynamic Bayesian Models (DBNs)
 - Joint head-eyebrow discrete state variables
 - Joint continuous nodes for head and eyebrow



S. Mariooryad and C. Busso, "Generating human-like behaviors using joint, speech-driven models for conversational agents," IEEE Transactions on Audio, Speech and Language Processing, vol. 20, no. 8, pp. 2329-2340, October 2012.

Previous Studies on Talking Faces

- Rule-based systems [Cassell et al., 1994; S. Kopp 2006]
 - + Semantic meaning of behaviors (nodding)
 - Repetitive behaviors for a given gesture
- Data-driven methods [Levine et al., 2010; Busso et al. 2007]
 - Prosodic features are effective modalities to synthesize nonverbal human-like behaviors
 - + Capturing more variability, and synchrony
 - Behaviors may not convey the semantics



Our Vision

Rule-based systems



Data-driven systems

- Considering the underlying discourse function to bridge the gap between data driven and rule-based systems.
 - Synthesizing behaviors that are timely aligned with speech
 - Synthesizing behaviors that convey the right meaning



Data driven approach
requires a corpus rich in
discourse functions

MSP-AVATAR Corpus

- Multimodal database comprising:

- Motion capture data
- Video camera
- Speech recordings



- Four dyadic interaction between actors

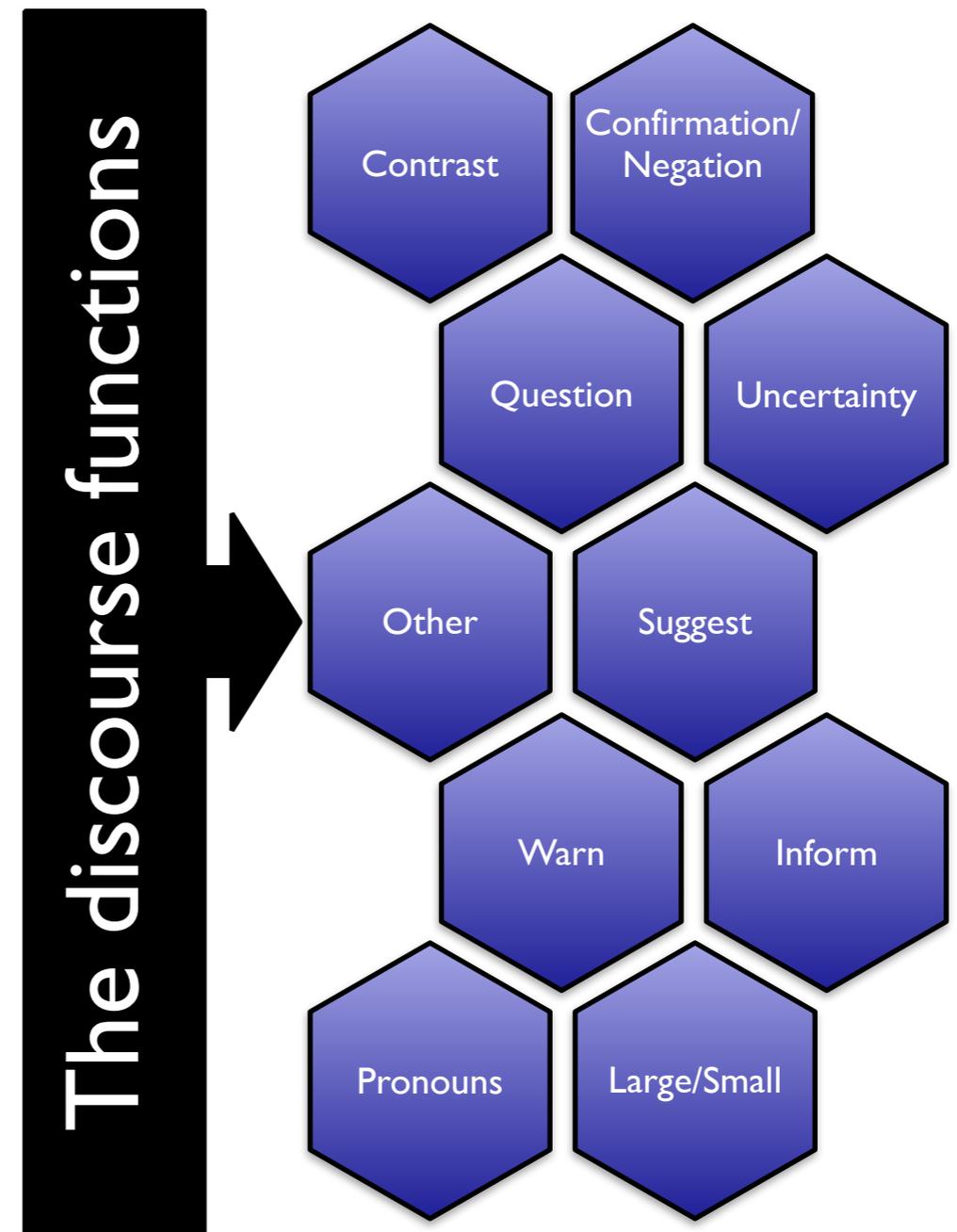
- We motion captured one of the actors

- Database rich in terms of discourse functions

- There are no corpora available to explore the role of discourse functions

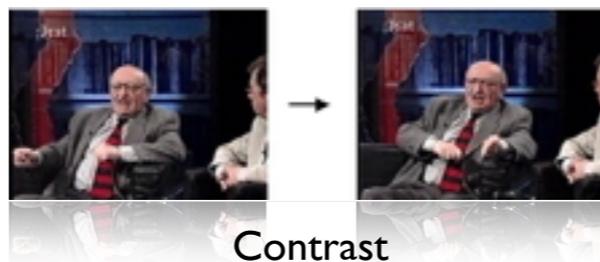
Selection of Discourse Functions

- We look for discourse functions that elicit specific gestural behaviors
- Selection guided by previous studies
 - Poggi et al [2005]
 - Marsella et al. [2013]
- 2-5 scenarios per discourse function



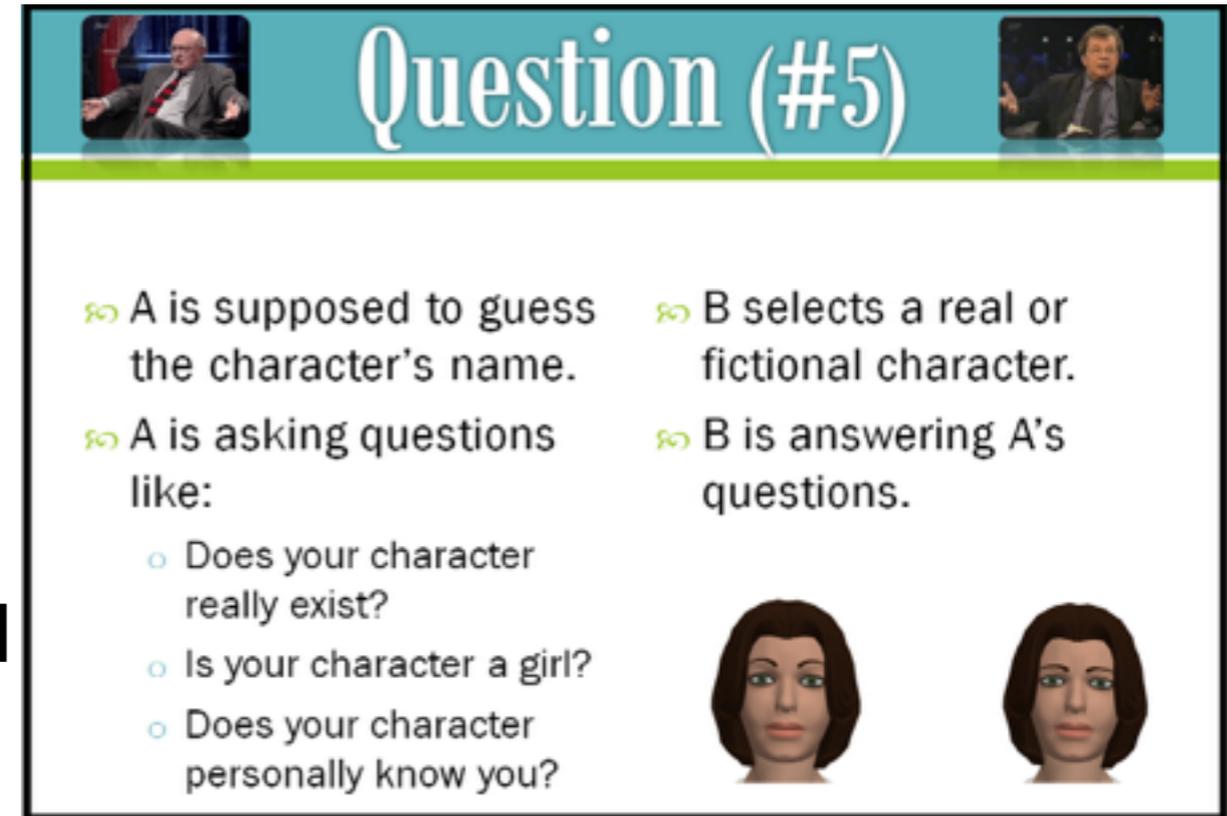
Discourse Functions

CONTRAST	Contrasting two ideas, usually accompanied with contrast conjunctions such as <i>but</i> , <i>nevertheless</i> , <i>as</i> , and <i>as opposed to</i>
CONFIRMATION/ NEGATION	Showing agreement and disagreement, usually accompanied with phrases such as <i>Yes</i> , <i>No</i> , and <i>I don't think so</i>
QUESTION	Asking a question of any type: <i>Yes-No</i> and <i>Wh-questions</i>
UNCERTAINTY	Showing uncertainty in making a decision, might be accompanied by sentences such as <i>I really don't know what to do!</i>
SUGGEST	Suggesting ideas to the listener, e.g., <i>How about the new Japanese restaurant?</i>
GIVING ORDERS	Ordering any type of service, e.g. <i>ordering food in a restaurant</i>
WARN	Warning the listener of a danger, e.g. <i>Be careful about ...</i>
INFORM	Inform something to the listener
LARGE/SMALL	The act of referring to something as small or large during speaking. These scenarios target iconic gesture usually accompany these two words or any of their synonyms
PRONOUNS	The act of referring to any pronoun (I/You/She/He/They). These scenarios target deictic gestures



Scenarios

- Scenarios designed to elicit characteristic behaviors
 - Description of the scenario
 - Prototypical behaviors associated with target discourse function
- Duration of the recorded scenarios:
 - MEAN = 143.1 sec, STD = 74.7 sec.
 - 21, 15, 22, and 16 scenarios (74 in total), based on pace of the actors



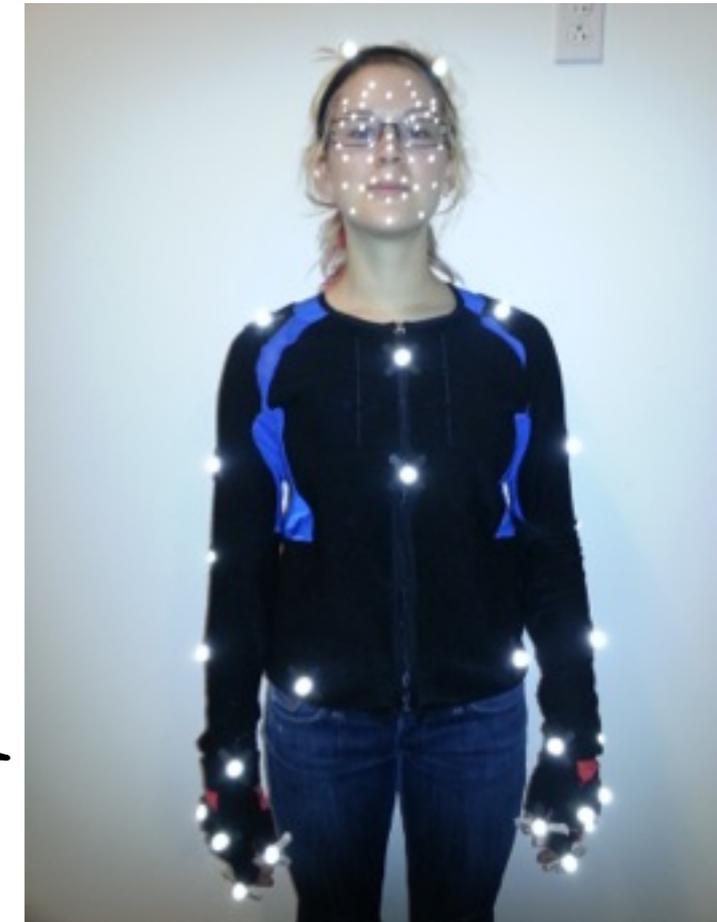
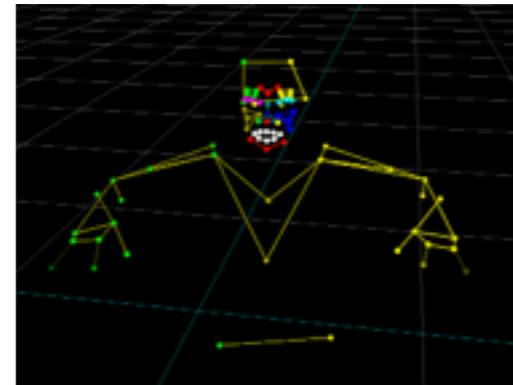
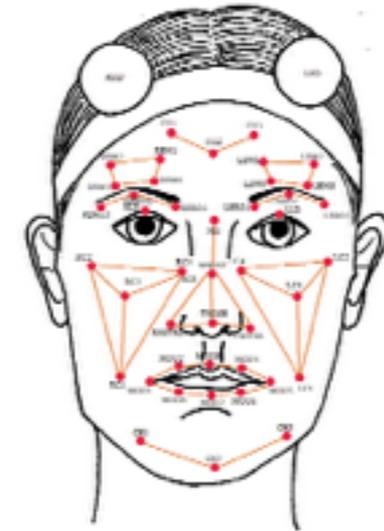
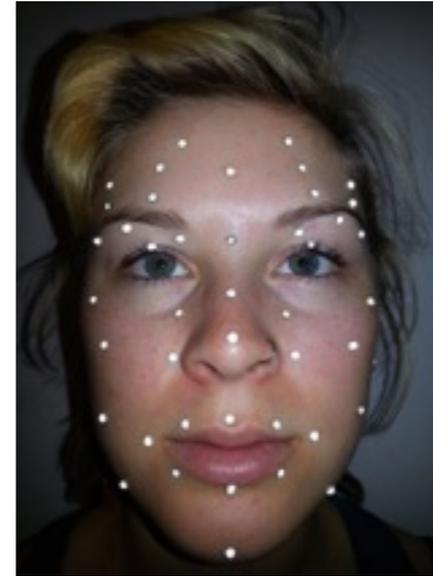
Question (#5)

- ⌘ A is supposed to guess the character's name.
- ⌘ B selects a real or fictional character.
- ⌘ A is asking questions like:
 - Does your character really exist?
 - Is your character a girl?
 - Does your character personally know you?
- ⌘ B is answering A's questions.



Motion Capture

- 16 VICON cameras
- Face:
 - 43 reflective markers
- Upper-body joints:
 - Headband with 4 markers
 - Suite with 28 marker



Recording “upper-body”, and “facial”!

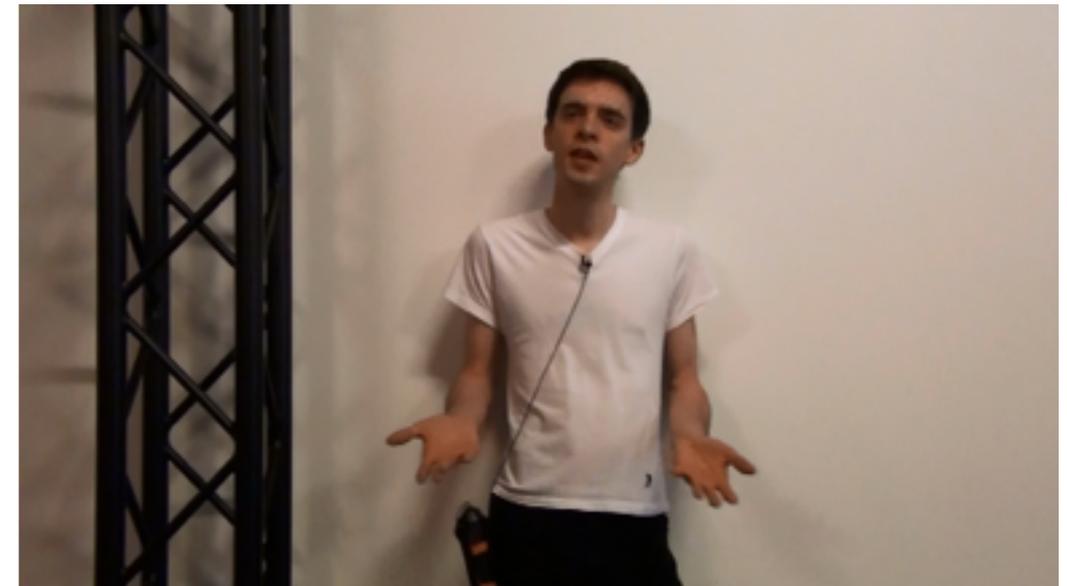
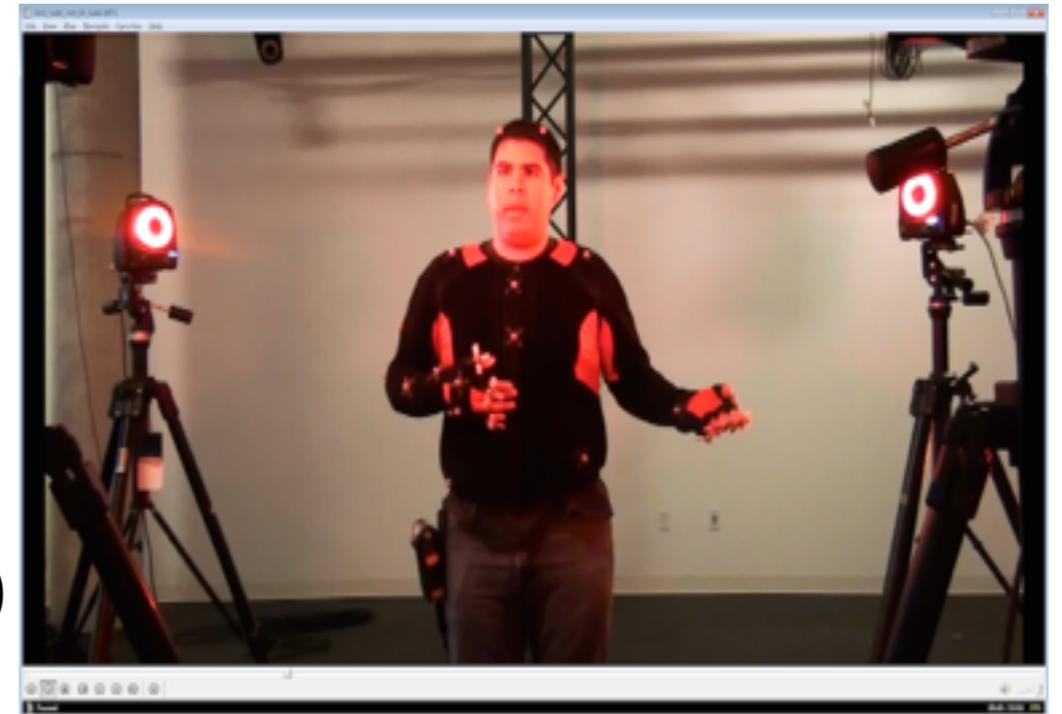
Audio

- Audio-visual data for all actors (6 people)
- Microphone connected to a digital recorder (TASCAM DR-100MKII)
 - 16 bit resolution
 - Sampling rate of 44.1 kHz
- First session: head-worn microphone (SHURE BETA 53)
 - Occluded some of the facial markers
 - Making the post-processing more difficult
- Next two session: lavalier microphone (SHURE MX150)



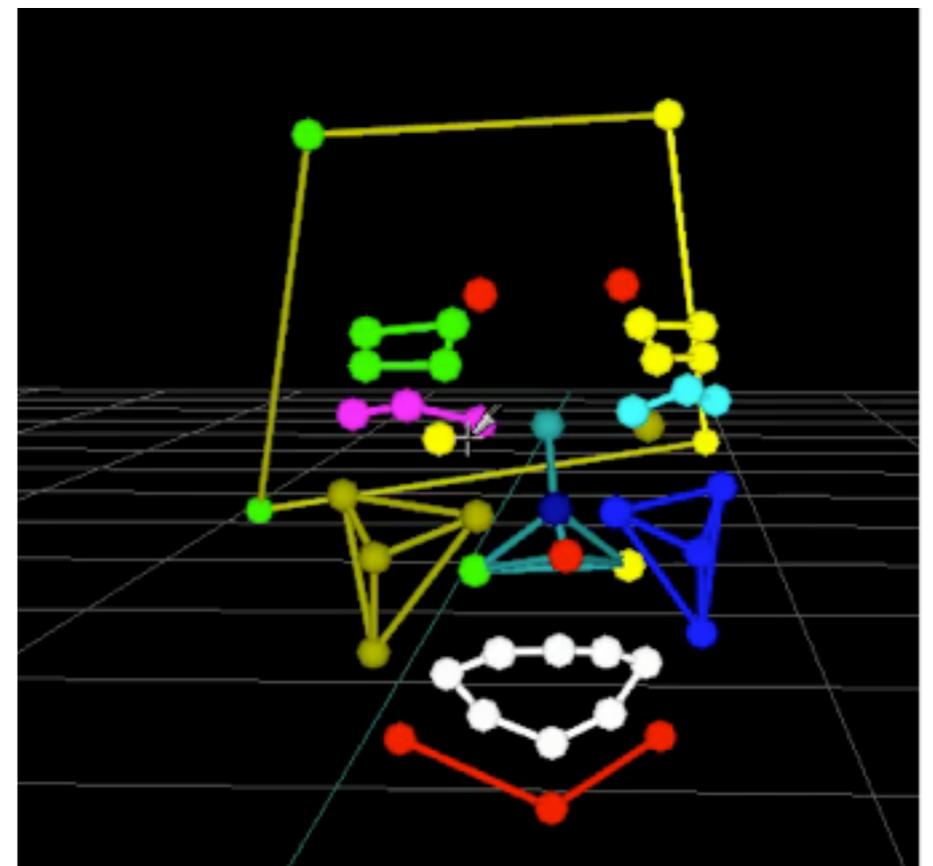
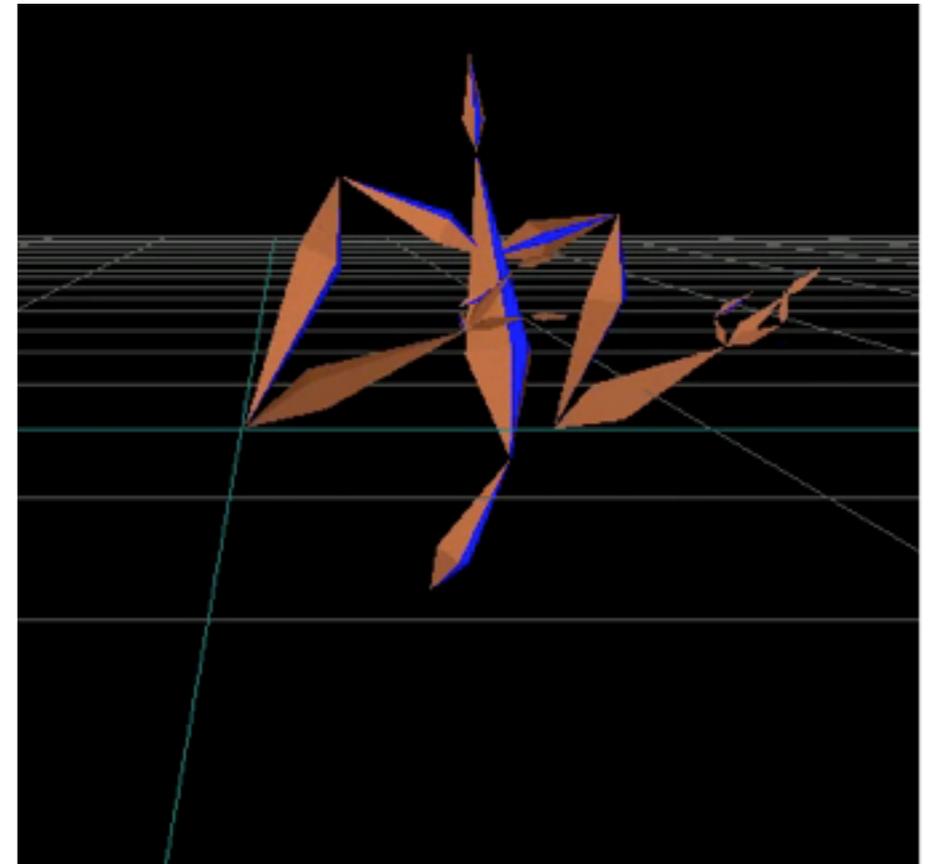
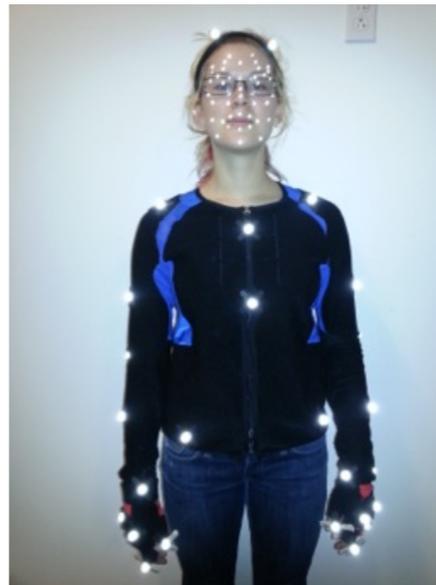
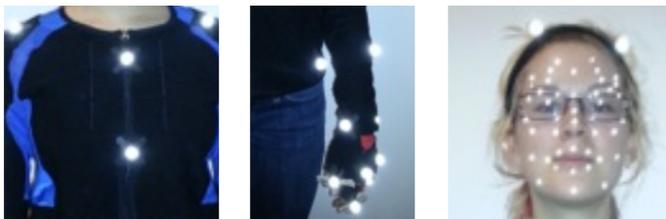
Video

- Frontal view of both actors
- Two Sony Handycams HDR-XR100
 - 1920x1080 resolution in Full HD
- We use these videos to annotate the behaviors
- They can be used in extracting visual features from both actors
- We use a clapboard with two reflective markers to synchronize audio, video and motion capture data



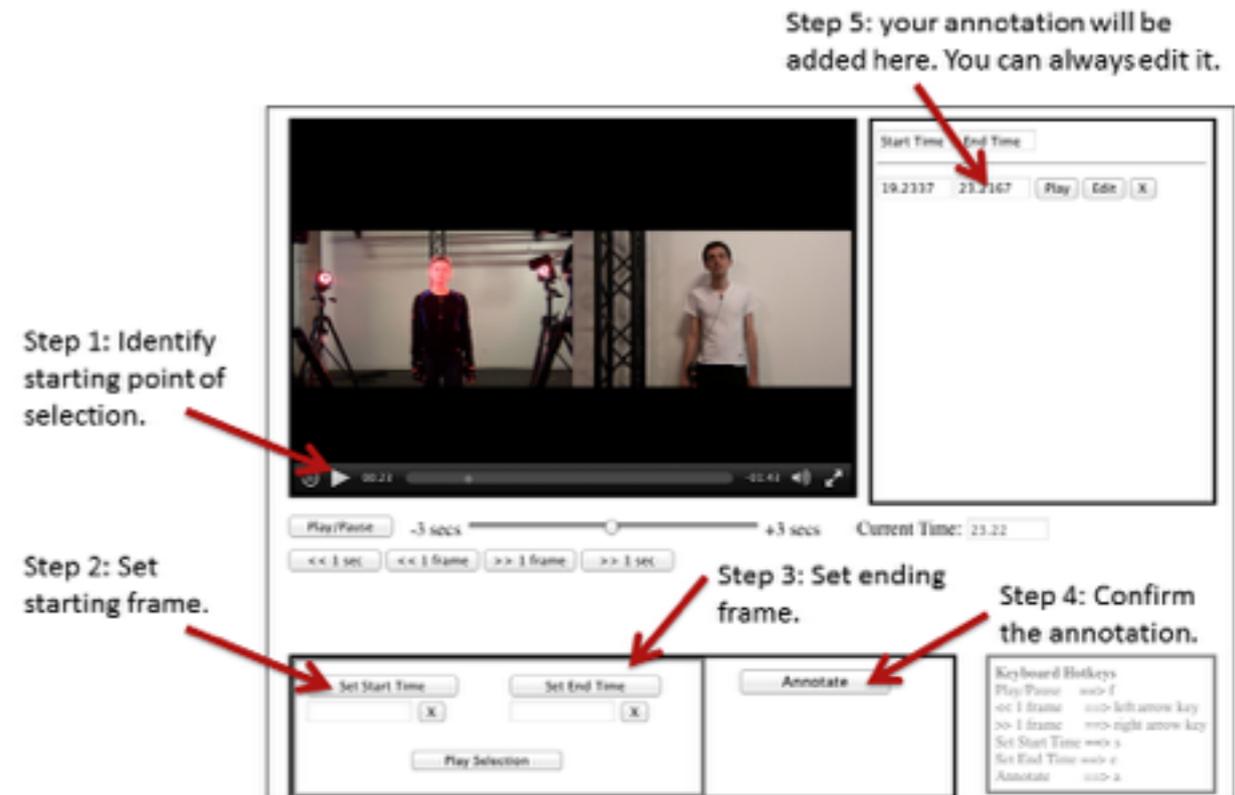
Post Processing

- Cleaning the motion capture data
 - Blade (VICON)
- Skeleton personalized for each actor
- Upper-body skeleton
 - 74 scenarios (all sessions)
- Facial markers
 - Only 3 scenarios



Discourse Functions

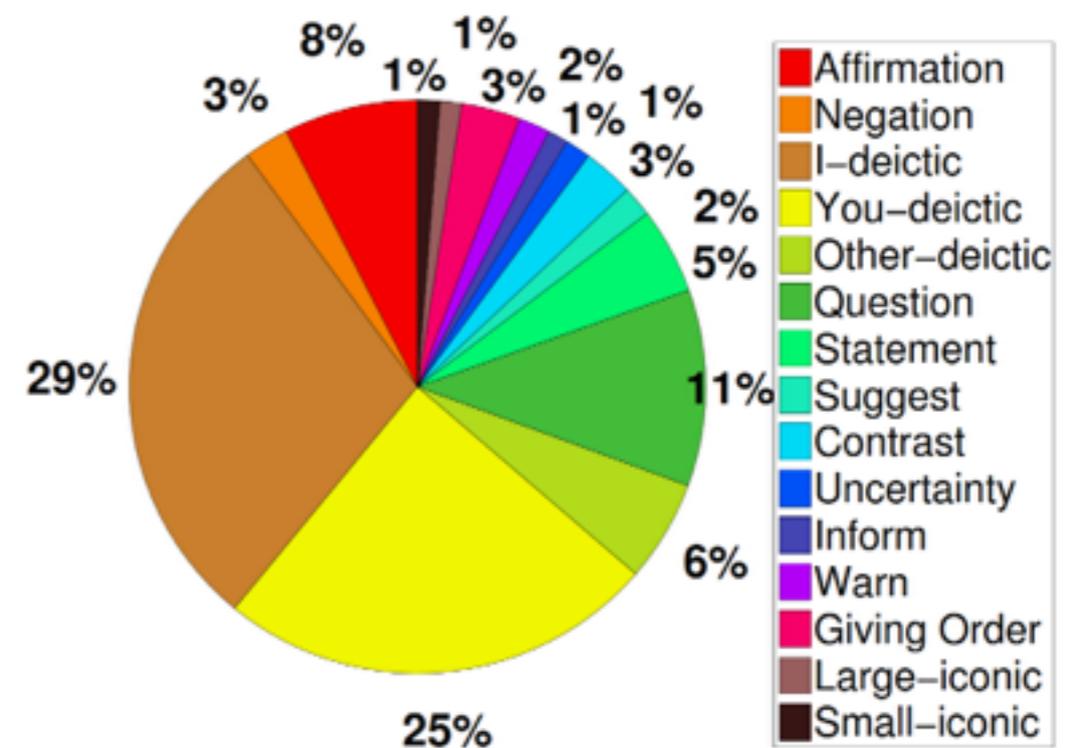
- Annotate discourse functions within speaking turns
- All annotations are done by one person, using Audacity
- We have 1751 samples
 - UPDATE: We are currently getting annotations by more people using AMT with OCTAB [Sunghyun et al., 2014]



Discourse Functions in MSP-AVATAR

- Some discourse functions appears in most scenarios

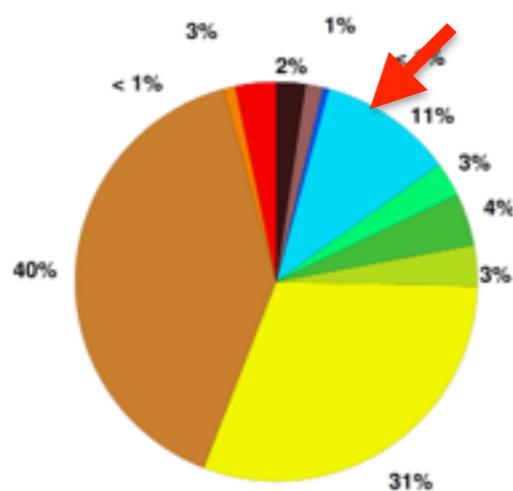
- I-deictic
- you-deictic
- questions



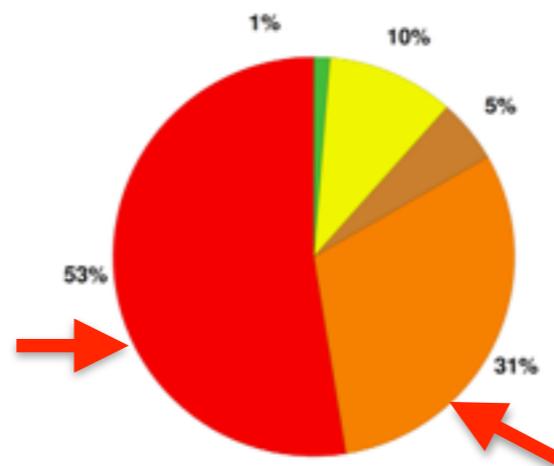
- We have multiple instances of the most relevant discourse functions

Discourse Functions in MSP-AVATAR

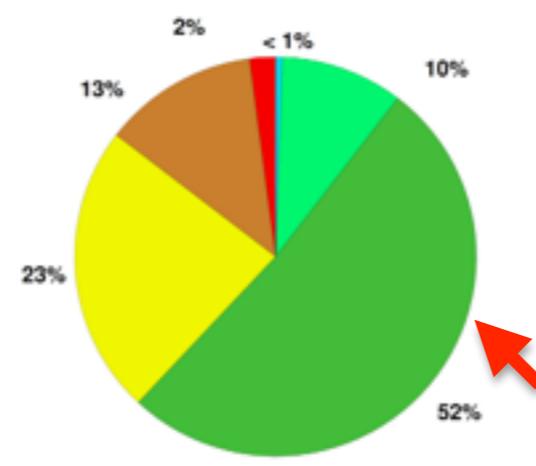
- The scenarios successfully elicited the target discourse functions



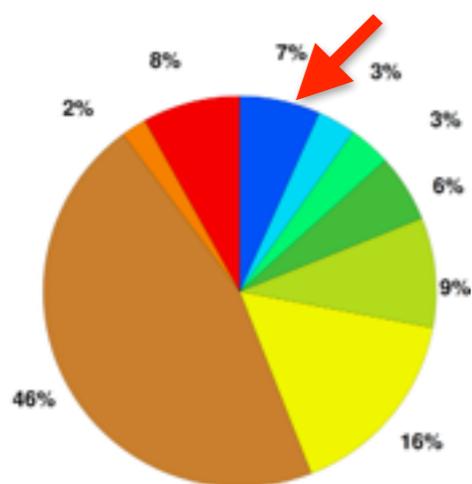
(a) Contrast



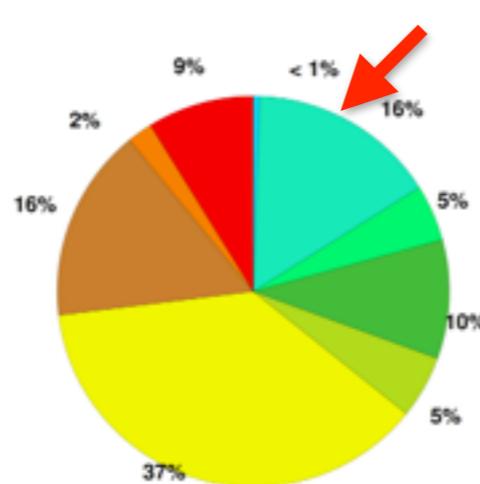
(b) Affirmation/Negation



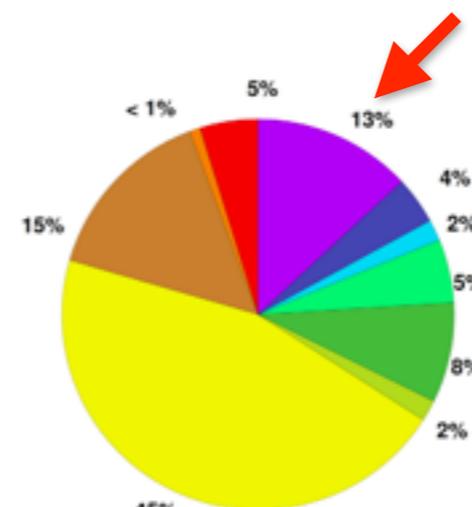
(c) Question



(d) Uncertainty



(e) Suggest



(f) Warn



Analysis of body movements

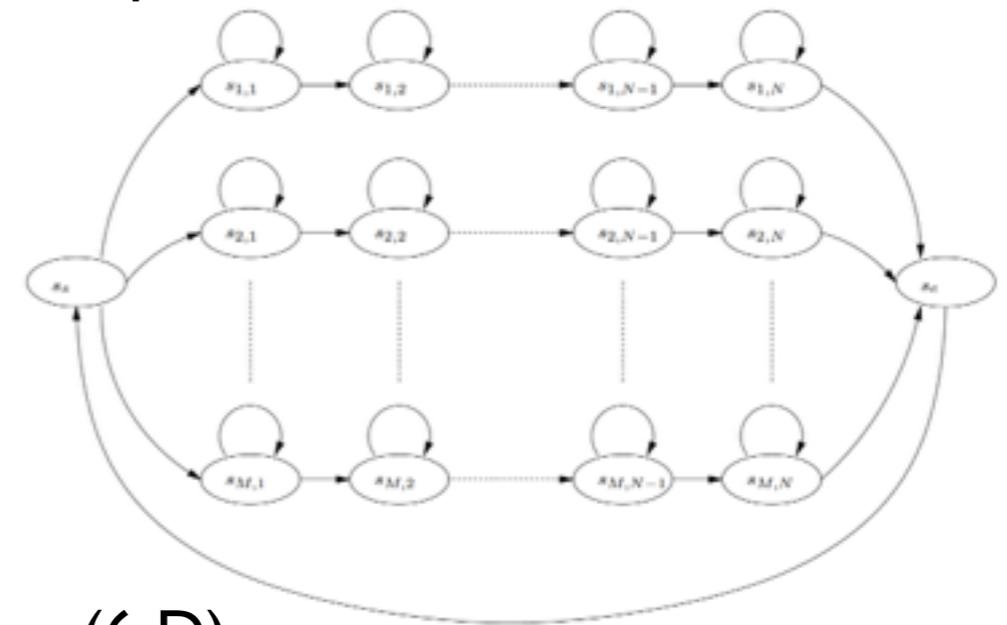
- Are the gestures different across discourse functions?
- We focus on hand and head gestures
 - We only have facial features of 3 sessions
- We rely on automatic unsupervised segmentation of gestures
 - Data-driven mid-level classes



Session (~3min)

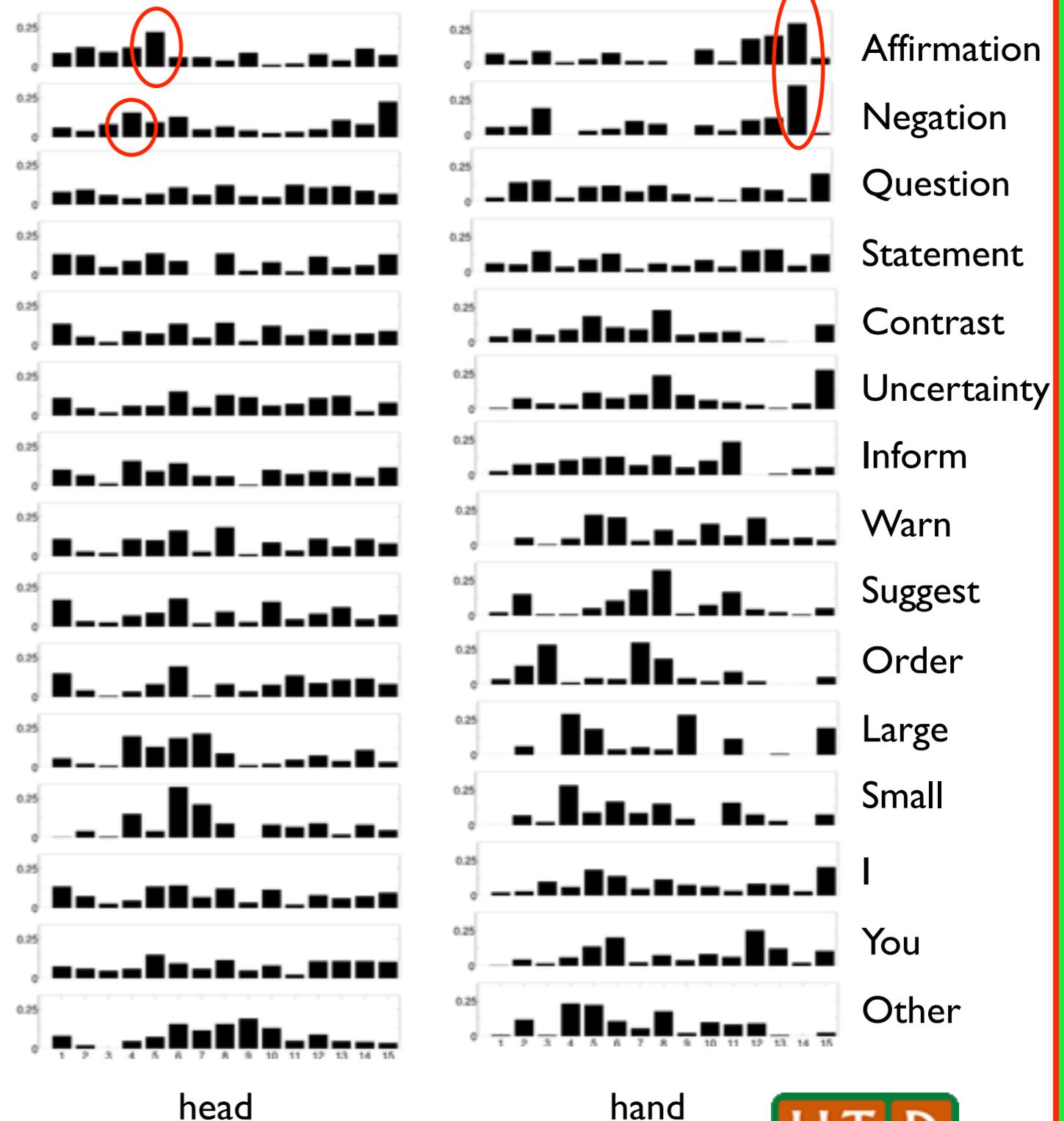
Parallel Hidden Markov Models

- Cluster the behaviors using PHMM
 - Individual left-to-right HMMs
 - States within the branches have self transitions
 - Model dynamic behaviors having different temporal durations
- #states per branch = 15
 - Minimum duration: 125 ms
- 15 branches for both head and hand
- Analysis on head, and arms movements
 - Head: three angular rotations + derivatives (6 D)
 - Arms: angular rotations + derivatives (20D)



Analysis of Head Motion and Hand Gestures

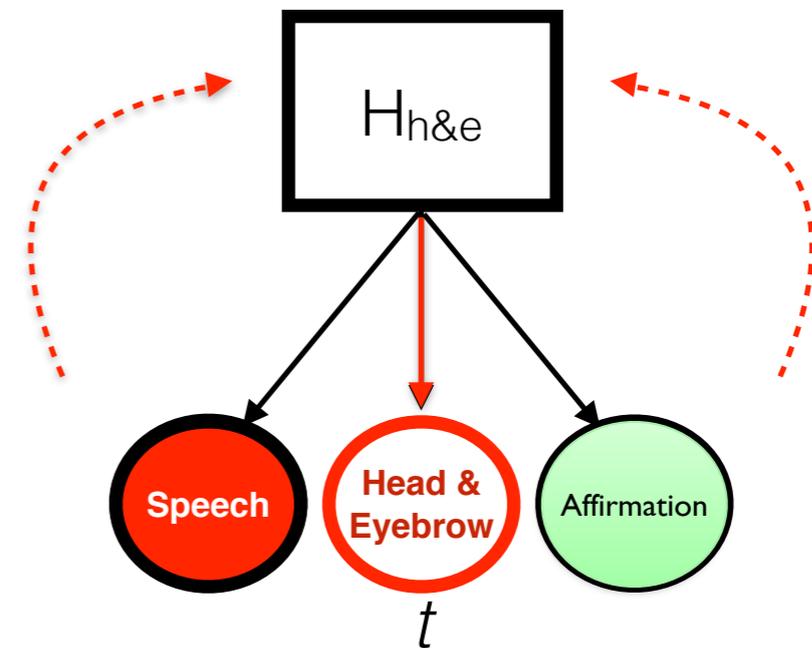
- Different gestures across discourse functions
- Hand: *affirmation* and *negation* have a distribution with a peak in cluster 14
- Hand gesture less active
- Head: peak in *affirmation* relate to head nod (cluster 5)
- Head: *negation* include head shakes (cluster 4)



Constraining DBNs with Discourse Functions



Constrained on Affirmation



Conclusions

- This paper introduces the MSP-AVATAR database
- This corpus comprises video and audio of 6 actors and motion capture recording of 4 actors in dyadic conversations
- Scenarios are designed such that they elicit the behaviors associated with discourse functions
- One drawback of this corpus is small number of subjects that were motion captured
- We expect to release the database after finishing the correcting process of motion capture recordings

<http://ecs.utdallas.edu/research/researchlabs/msp-lab/MSP-AVATAR.html>

Multimodal Signal Processing (MSP)

- Questions?



<http://msp.utdallas.edu/>