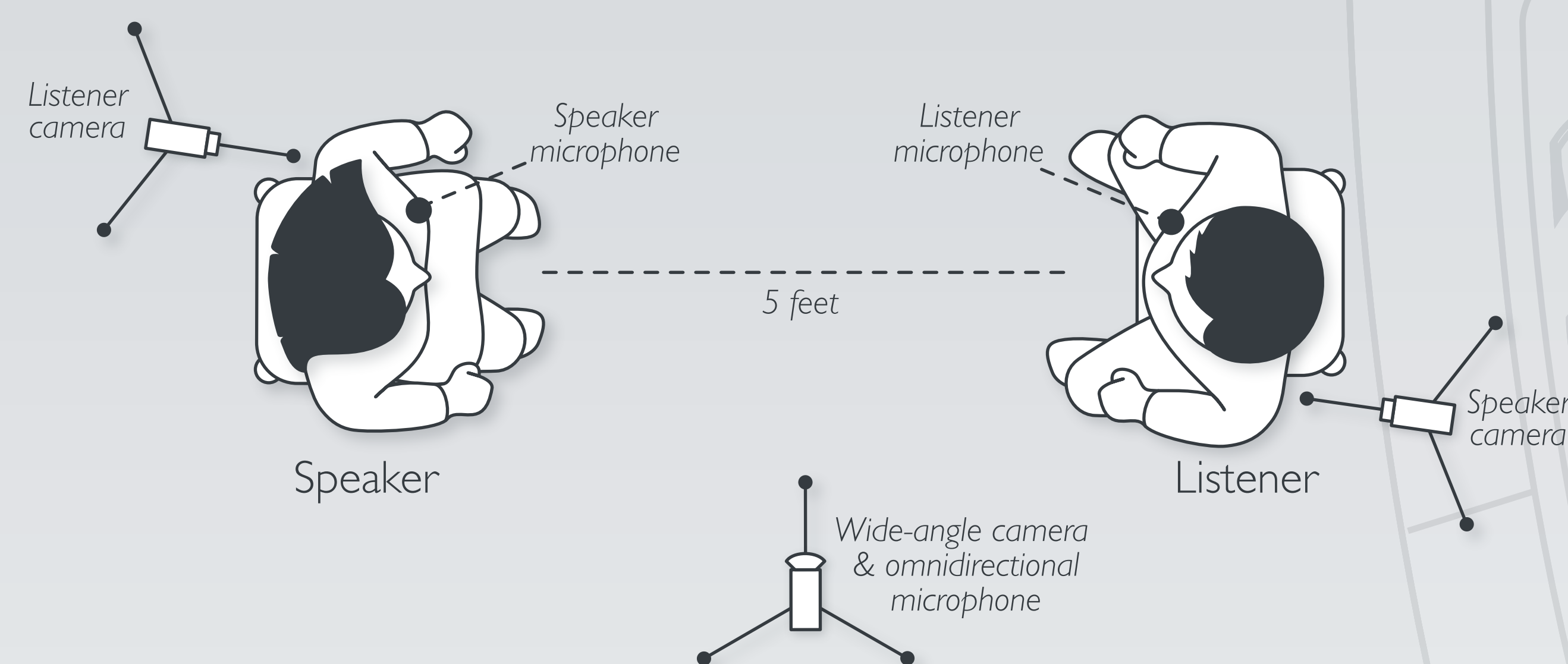# MODELING SOCIAL CUES:
# EFFECTIVE FEATURES FOR PREDICTING LISTENER NODS

FAISAL KHAN, BILGE MUTLU, XIAOJIN ZHU | DEPARTMENT OF COMPUTER SCIENCES | UNIVERSITY OF WISCONSIN–MADISON

## RESEARCH GOALS

Gaining a computational understanding of human social behavior

Building socially interactive systems such as agents and robots

Current study explores:

*Using of a small set of real-time features to predict listener nods*

## DATA COLLECTION SETUP



Listener camera

Speaker microphone

Listener microphone

5 feet

Speaker

Listener

Speaker camera

Wide-angle camera & omnidirectional microphone



Data collection with 24 dyads

Equal number of MM, FM, MF, and FF gender combinations

Perform a "storytelling task"

Participants were paid $10

## SPONSORS

## *RAW* FEATURES

A "raw" set of features extracted automatically from multimodal data

Speech segmentation (speech/pause)

Speaker classification (speaker/listener)

Pitch values and slopes (rising/falling intonation)

Speaker head movements including *nodding*

## *DERIVED* FEATURES

Temporal dependencies between raw of features and listener nods captured by a derived set of features based on multiple windows of averages of raw features and differences across window averages

Final feature vector for a given frame:

$$\mathbf{f}_i = \left[\begin{array}{cccccccc} \mathbf{r}_i & \mathbf{g}_i^1 & . & . & . & \mathbf{g}_i^7 & \mathbf{h}_i^1 & . & . & . & \mathbf{h}_i^7 \end{array}\right]'$$

Where

$$\mathbf{r}_i = \left[\begin{array}{ccccccccc} speech & speaker & head_x & head_y & nodding & pitch & s_1 & . & . & . & s_9 \end{array}\right]'$$

$$\mathbf{g}_i^m = \frac{1}{2^m}\sum_{k=0}^{2^m-1}\mathbf{r}_{i-k} \text{ and } \mathbf{h}_i^m = \mathbf{g}_i^m - \mathbf{g}_{i-2^m}^m$$
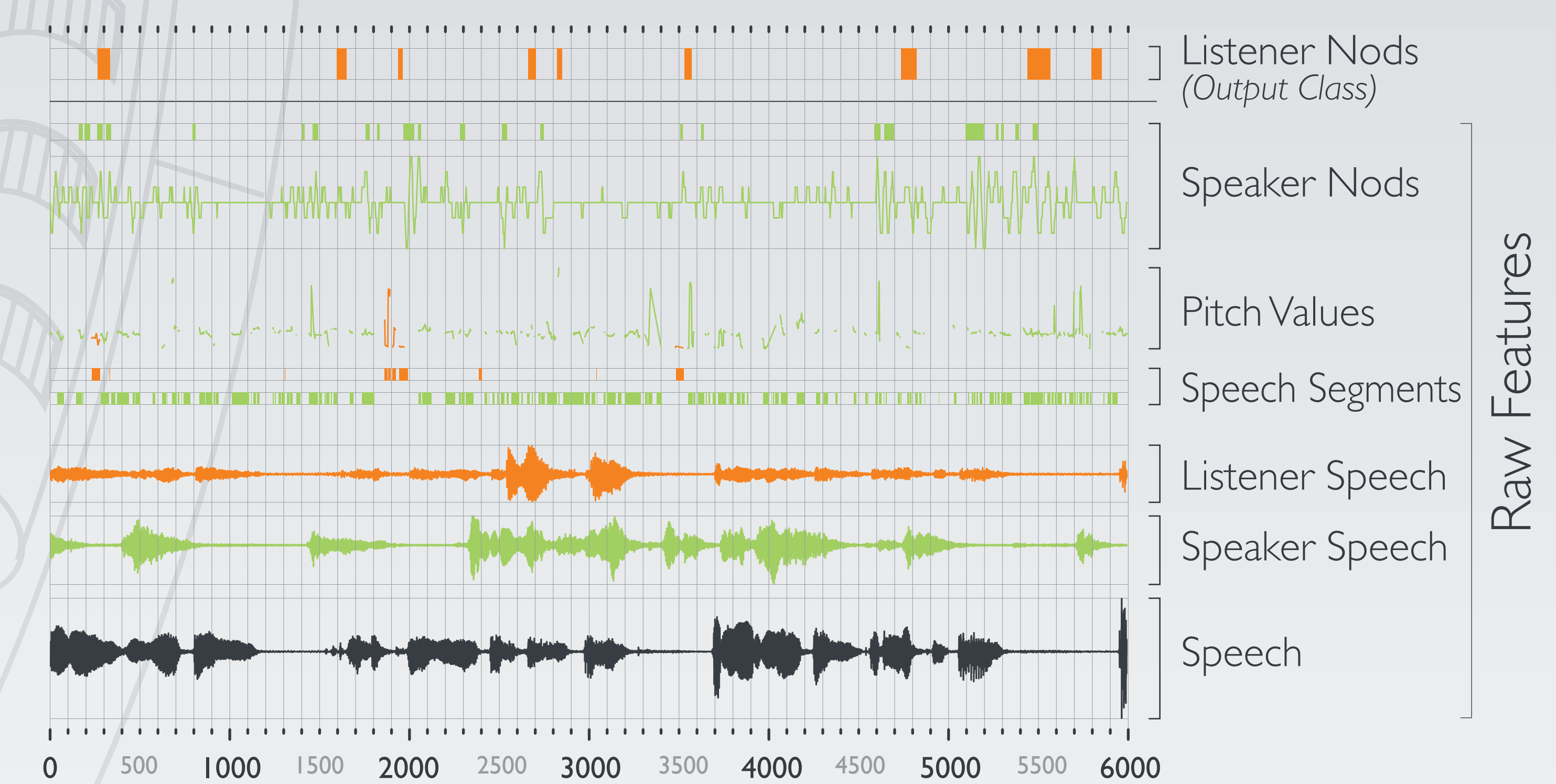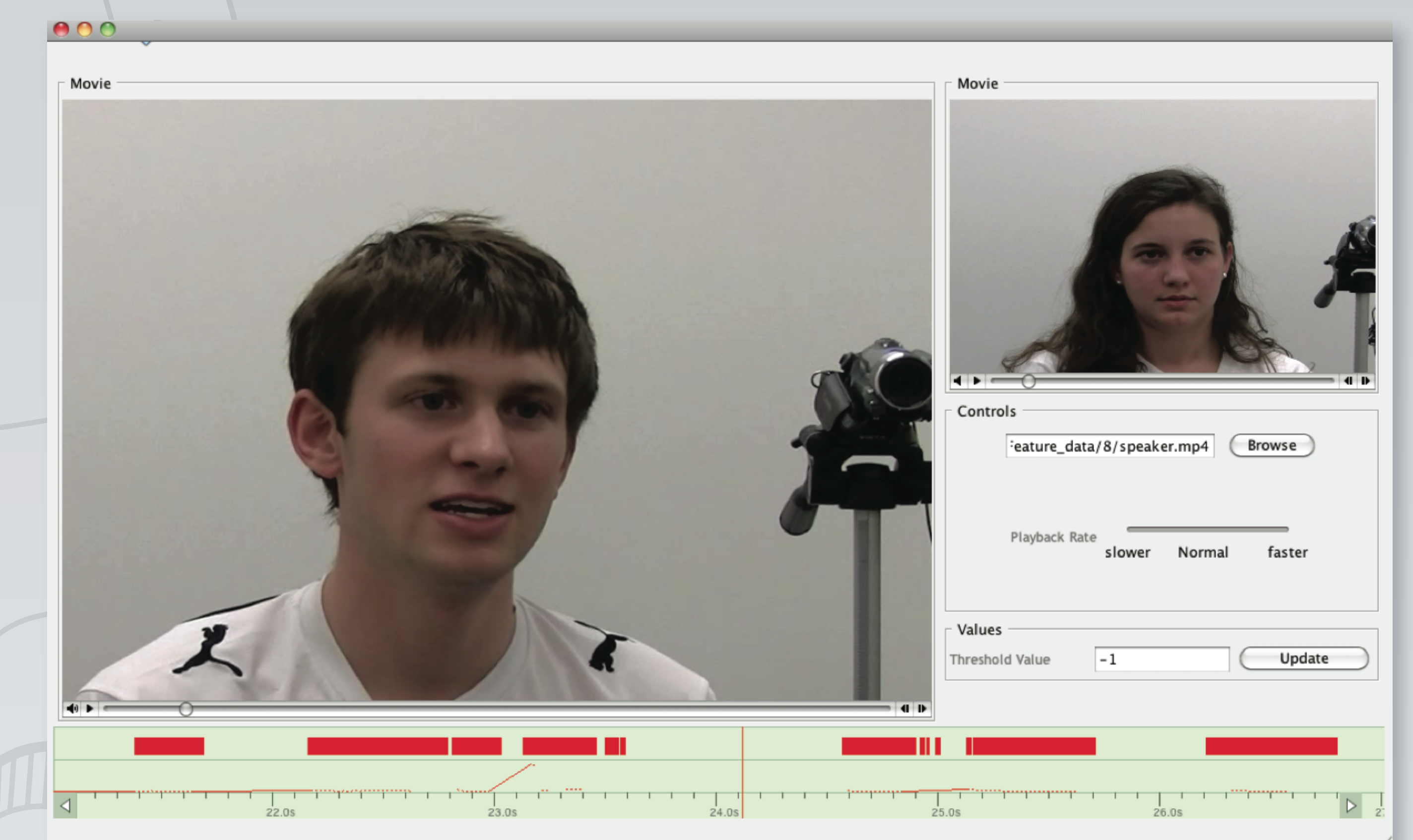
## CLASSIFICATION RESULTS

Predictions using a Support Vector Machine (SVN) classifier

Four-fold cross validation

Precision   =   **0.1083**

Recall   =   **0.3165**

F-measure =   **0.1605**

## DATA ANNOTATION & VISUALIZATION





Listener Nods *(Output Class)*

Speaker Nods

Pitch Values

Speech Segments

Listener Speech

Speaker Speech

Speech

Raw Features

0  500  1000  1500  2000  2500  3000  3500  4000  4500  5000  5500  6000

## NEXT STEPS

Improving the modeling of temporal dependencies using:

Encoding templates (Morency et al., 2010)

Sequential models (e.g., CRF, HMM)

Using model predictions to control a robot's nods

Conducting human-robot interaction studies to test effectiveness