

# Designing Bots, Virtual Humans, & Other Systems that Hold up their End of the Conversation

Justine Cassell

Carnegie Mellon University & ISIR / La Sorbonne  
Chaire Blaise Pascale & Chaire Sorbonne  
Mai 2018

# The ineffable quality of Rapport in learning



Children who report more rapport are more likely to learn from the virtual peer

# Search Engine vs. Conversation

Justine: “OK Google, I love Manchester United”

**Google:** Manchester United Football Club is a professional football club  
Based in Old Trafford, Greater Manchester, England, that competes  
in the Premier League, the top flight of English Football

Justine: “I love Manchester United”

**Friend.:** “No way! Arsenal wipes the floor with those Red Devils!”

**Socially-Aware Robot Asst:** “No way! Arsenal wipes the floor with those Red Devils!”

# Motivation for Socially-Aware Bots

1. People pursue *multiple conversational goals* in every conversation & expect the same from their interlocutors. To put people at ease, and increase relationship strength, we must understand the *propositional*, *interactional* & *interpersonal* functions of conversation.
2. People change interaction styles over time. We must increasingly *manage long-term interactions* with people by changing interaction style in a way that evokes increasing loyalty, rapport and trust.

# **Rapport** improves *task* performance

## **Surveys**

- Survey respondents gave higher quality answers if they felt rapport with interviewer (Berg (1989))

## **Health**

- Physicians who build rapport during trial interviews enroll more participants (Albrecht *et al.*, 1999).

## **Sales**

- Rapport with sales staff leads to increased likelihood of purchasing goods/service (Brooks, 1989).
- Customers show increased trust and disclosure when rapport is maintained with sales staff (LaBahn, 1996).

# Methodology

## Theorize & Model

Build formal models

Implement system on the basis of model

**Study**

**Build**

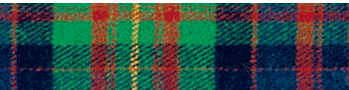
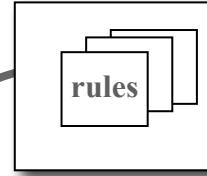
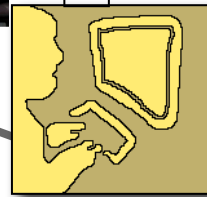
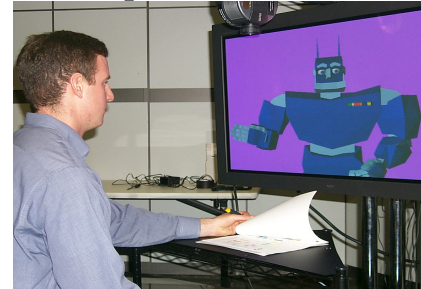
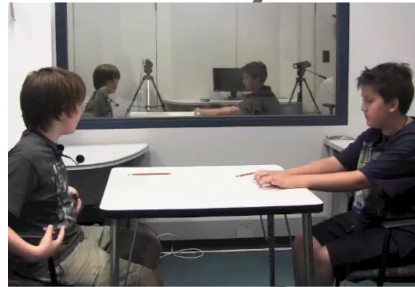
Start here

Collect Natural data

Realize gaps in understanding

Design evaluation of use

**Test**



# Observe

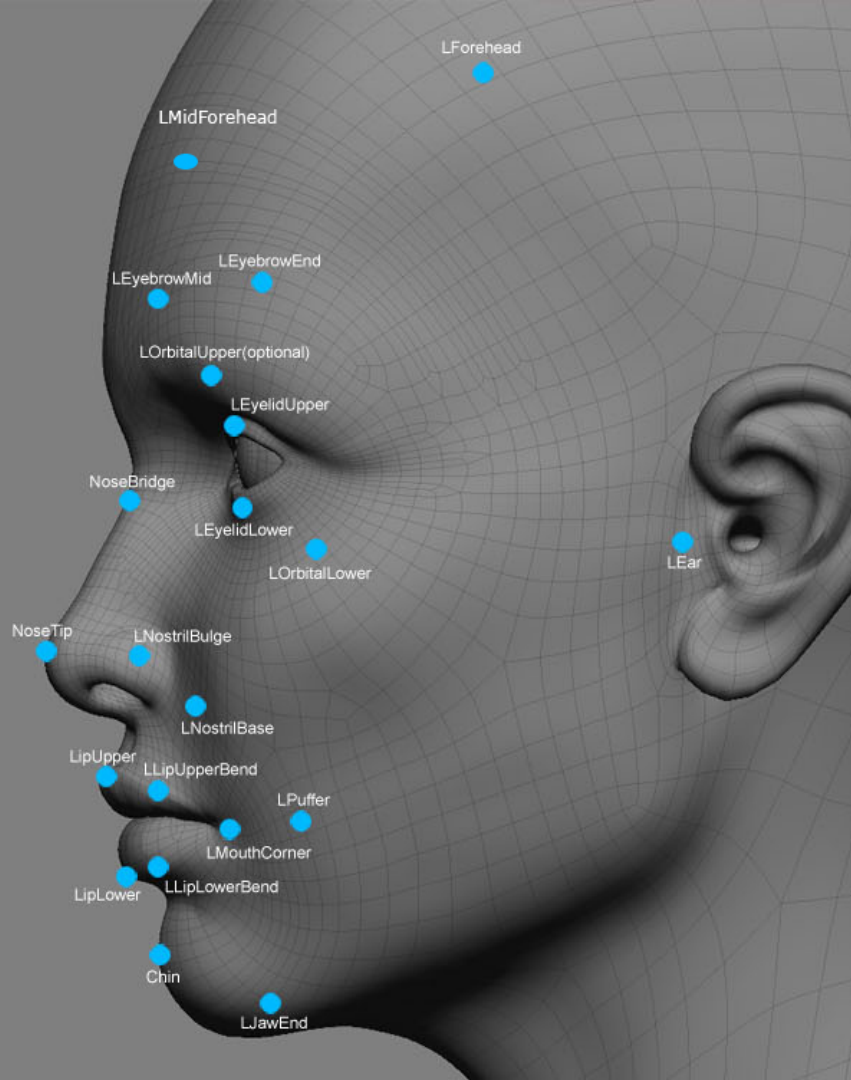
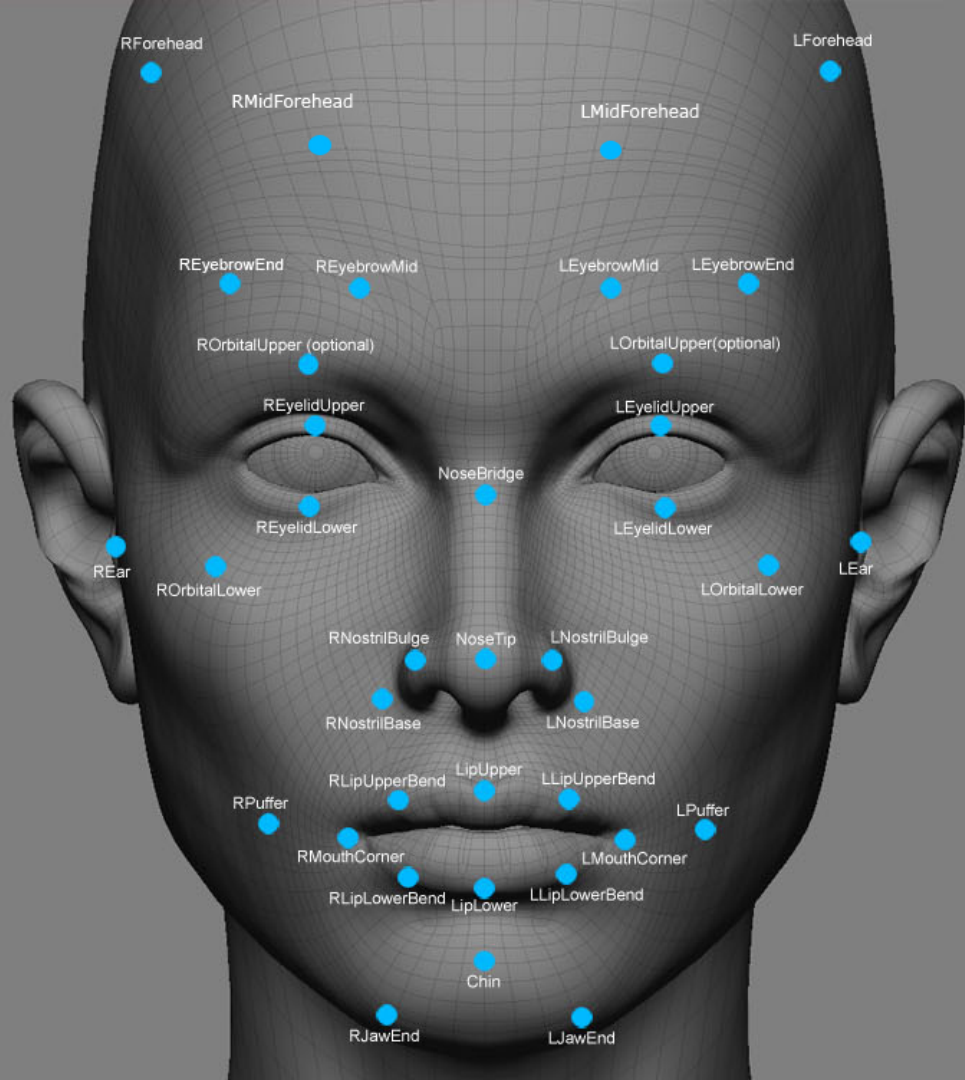


# Analyze

The screenshot displays the Anvil 3.6 software interface, which is used for analyzing video and gesture data. The interface is divided into several panels:

- Top Left Panel:** Contains file management options (File, Edit, View, Tools, Bookmarks) and a status area showing video specifications: "Loading video: IV50, 384x288, FrameRate=25", "Video frame rate: 25.0", "Audio format: LINEAR, 22050.0 Hz, 8-bit, Mono", "Duration: 03:43:92 (5597 frames)", and "Current specification: D:\Research\anvil-spec\litqua2.xml".
- Top Middle Panel:** A video player window titled "Video: lq1-7-reich.avi" showing two people sitting at a table in a laboratory setting.
- Top Right Panel:** A "Track: gesture.phrase" window showing attributes for a specific gesture: "category: iconic", "iconic type: smash", "handedness: 2H", "cooc: Rivas", "function: emblematic", and "timing: direct". It also includes a "Comment" field with the text "compare with lq1-8 at 0.28" and control buttons for "start", "edit", "end", "cut", "extend", and "del".
- Bottom Panel:** An "Annotation: lq1-7-reich.anvil" window showing a timeline of the video. The timeline includes a "wave" track, a "praat" track with phonetic labels like "übersetzen", "ich", "spüre", "den", "poetschen", "Stil", "den", "ich", "bei", "Rivas", "na", "nur", "die", "Bemühung", and a "gesture" track with various annotations such as "metaphoric, heart, 2H", "iconic, smash, 2H", and "emblem, so-what, 2H".









# Analysis of Rapport

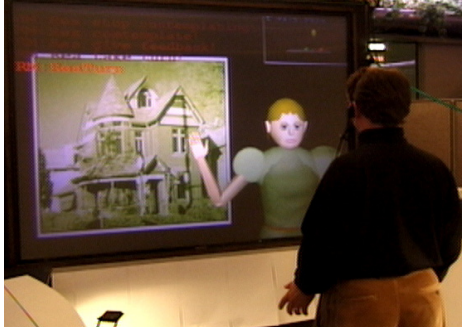
*Thin-slice* (Ambady & Rosenthal, 1992) judgments of every 30 second video segment) presented to 3 annotators in randomized order



## *IRR*

- “*Consensus*” measured by Intra-Class Correlation (single measure): *0.37*
- “*Consistency*” measured by Cronbach alpha: *0.68*
- *Inverse-based bias correction* (Kruger et al., 2014) was used mitigate rater bias & pick single rapport rating for each 30 second video segment.

# Implement



(GA) You just (GT) need [ to "type in ] ( a [ "line ] like ) ">>This is ...<<, and (GA) the actor is (GT) able to ( [ "talk ] and [ "gesture ] by ) itself

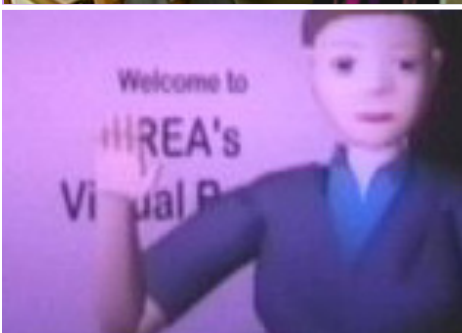
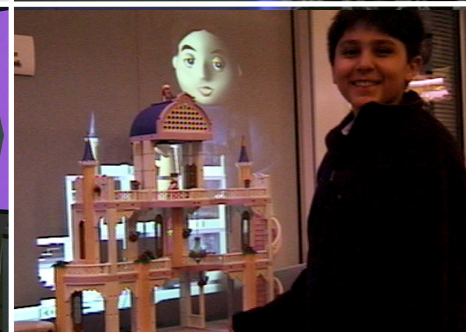
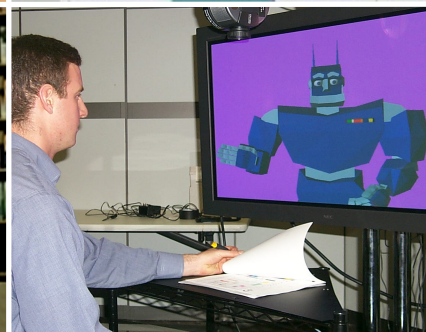
**TAGGING**

I don't know if this is a good thing or a bad thing.

I worry that virtual actors could squeeze real actors out of business.

00:00:17:20

	from	SIGGRAPH*	and	the
ze Towards				
Eyebrows				
Beat				
Gesture				



# Evaluate

## Total Interactions:

250+ sessions

## Current Dataset:

120 sessions (totally 10+ hours)

- Mean duration = approx. 5 min.
- SD = 69.00 seconds



**SARA**, 00:13.7, I'm so glad to have the opportunity to meet you in person, 00:16.7

**User**, 00:18.7, hi how are you [doing, 00:19.2

**SARA**, 00:19.3, [I'm] in Davos to assist attendees but to be honest I am a little anxious about it it's a big deal to be invited to Davos is it your first time here, 00:28.7

**User**, 00:29.1, it is and I'm a little anxious too, 00:31.4

**SARA**, 00:32.7, I found there's a lot to take in I hope that it isn't overwhelming you, 00:37.2

**User**, 00:37.9, it is overwhelming a little bit, 00:39.8

# Summary of Evaluation Metrics

## Interaction

- Total Time
- System Speaking Time
- User Speaking Time
- System's Response Time
- User's Response Time
- System's Interruptions
- User's Interruptions

## Verbal

- System's Intentions
- User's Intentions
- System's Conversational Strategies
- User's Conversational Strategies

## Task

- Session Rec. Acceptance
- Person Rec. Acceptance

## Interpersonal

- Rapport Score
- Mutual Attentiveness
- Coordination
- Positivity

# Goal of Socially Aware Systems

Development of a bot that manages *interpersonal rapport* (relationship strength) with users over interactions across time, as well as managing propositional and interactional goals, *in order to improve task performance*.

*Automatically recognize* rapport-managing conversational strategies from *verbal, visual and vocal* modalities of speaker and interlocutor, both within the individual and in the dyad.



# Goal of Building Socially Aware Systems

**Theoretical**: Understand the nature of rapport in greater detail, by correlating with associated observable verbal (conversational strategies, vocal (voice quality) and visual (non-verbal) cues

**Methodological**: Leverage this understanding to automatically recognize rapport-building strategies by leveraging and developing statistical machine learning techniques

# Ineffective Conversation (don't do this with agents)

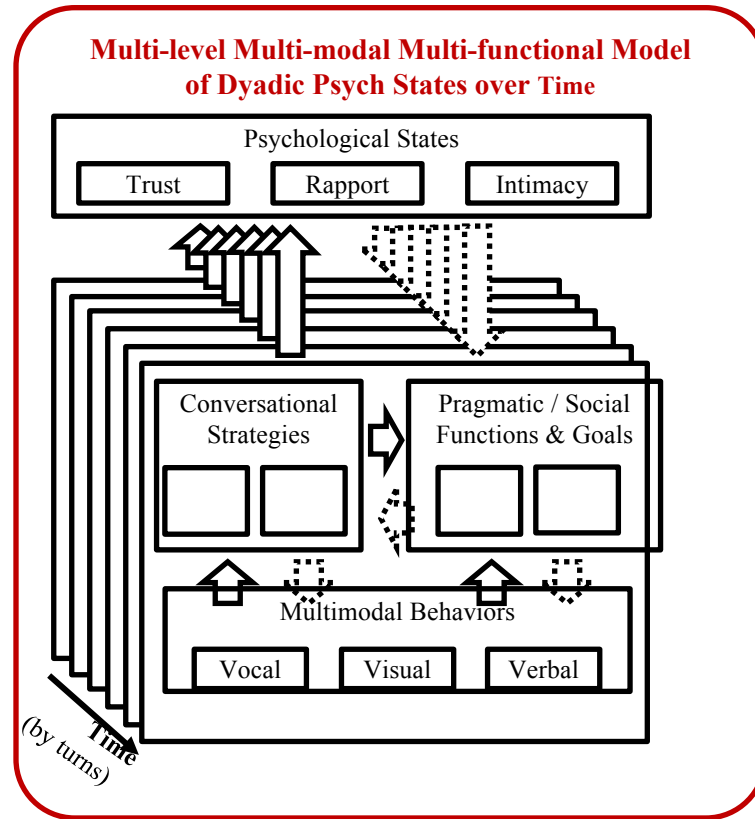


# Intimate Conversation (don't do this with agents)



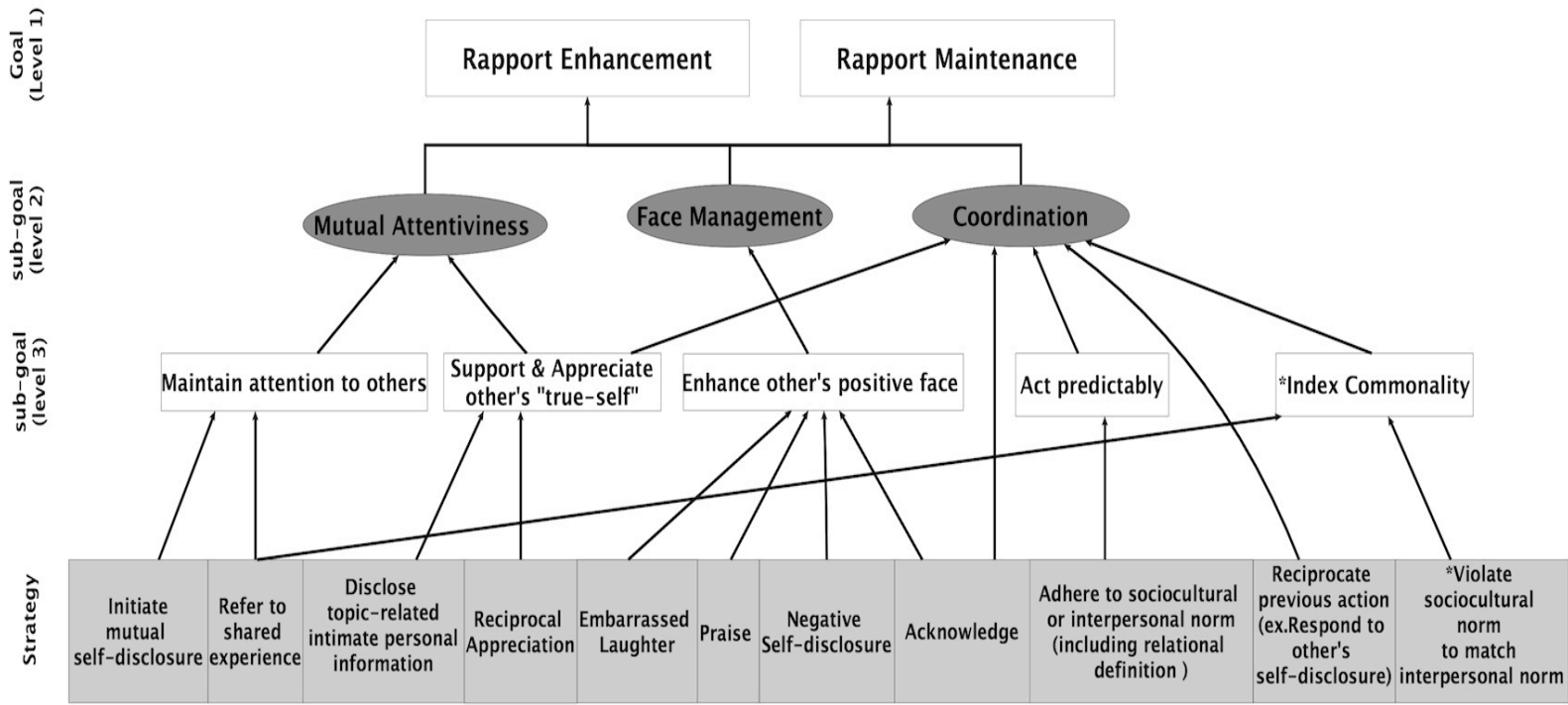
# Agent Model of Rapport must be:

1. Dyadic,
2. Multi-level:  
differentiate between  
observable signals &  
underlying  
psychological states,
3. Sensitive to effect of  
*time*
4. Cross-Modal



with L.P. Morency, 2015

# Data- & Theory-Driven Model of Rapport Management

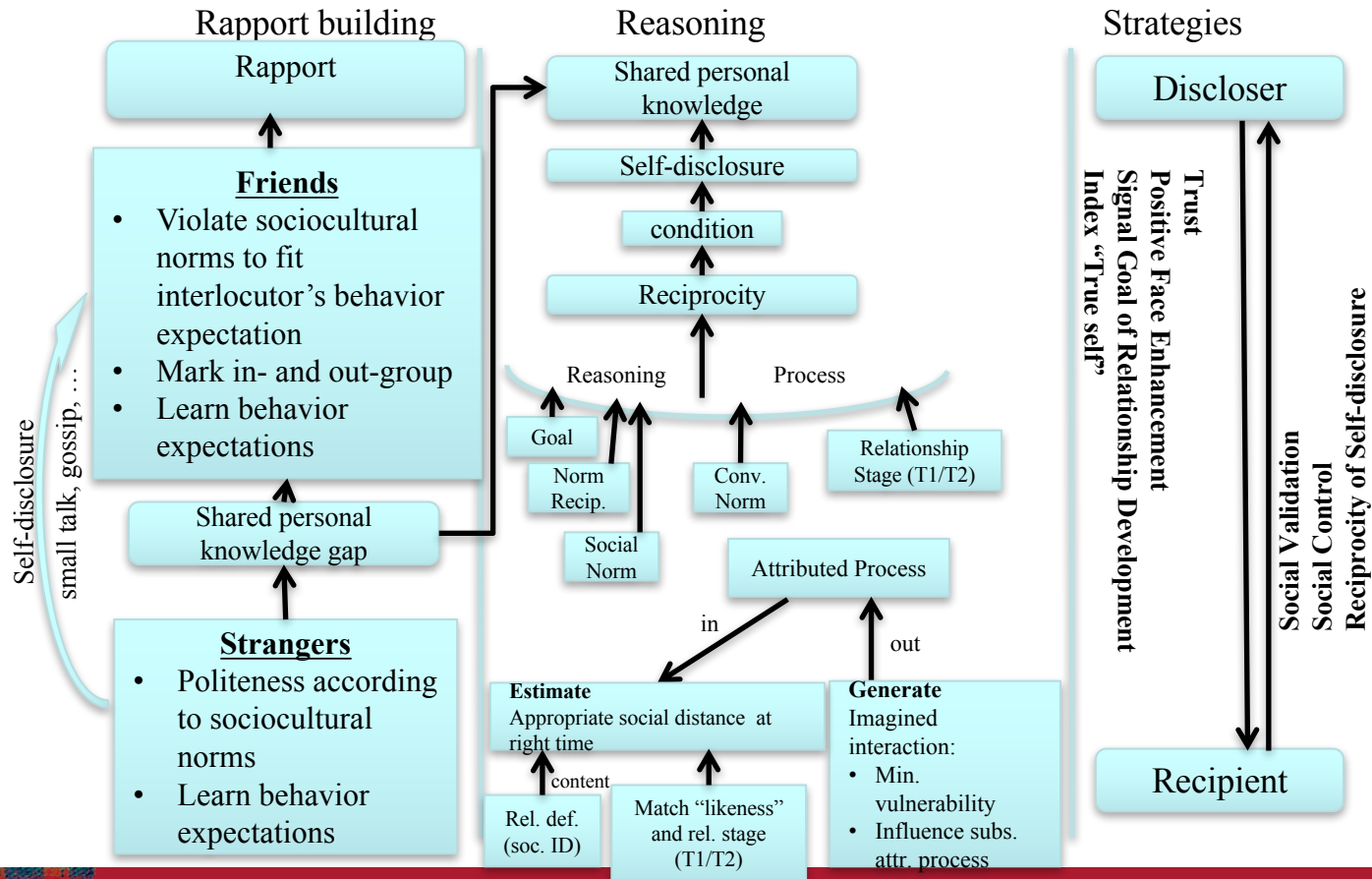


With Ran Zhao

# Conversational Strategies

Conversational Strategies	Examples
<b>VSN (Violation of Social Norm)</b>	<i>"man you take forever to write"</i>
<b>SD (Self Disclosure)</b>	<i>"I hate math"</i>
<b>PR (Praise)</b>	<i>"well done"</i>
<b>SE (Reference to Shared Experience)</b>	<i>"I shared m&amp;m's with you last time"</i>
<b>BC (Back Channel)</b>	<i>"yup"</i>
<b>QE (Question Eliciting SD)</b>	<i>"are you an atheist"</i>

# Data- & Theory-Driven process model



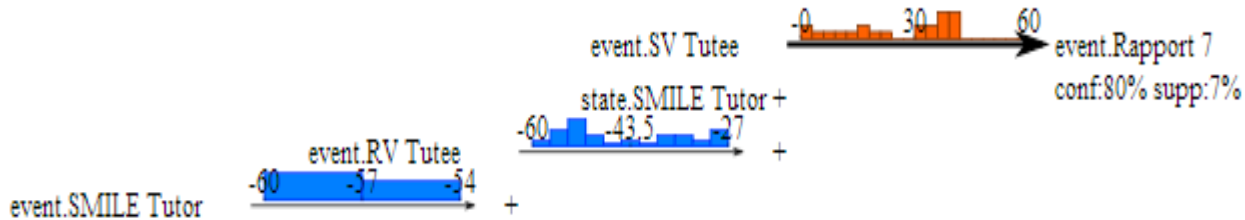
# Data-Driven: Temporal association rules

Form of temporal rules

*“If event  $A$  happens at time  $t$ , there is 50% chance of event  $B$  happening between time  $t+3$  to  $t+5$ ”*



# Temporal association rules: Friends



## Example: Friend in high rapport

**Tutor:** Sweeney you can't do that, that's the whole point *{smile}* [Violation of Social Norm]

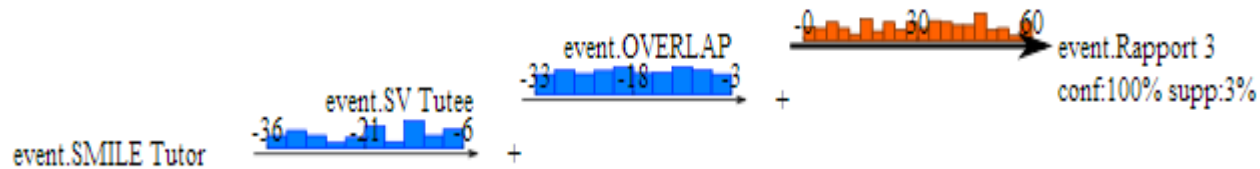
**Tutee:** I hate you. I'll probably never never do that [Reciprocate Social Norm Violation]

**Tutor:** Sweeney that's why I'm tutoring you *{smile}*

**Tutee:** You're so oh my gosh *{smile}*. We never did that ever [Violation of Social Norm]

**Tutor:** *{smile}* What'd you say?

# Temporal association rules: Strangers



## Example: Stranger in low rapport

**Tutee:** divide oh this is so hard let me guess: 11 [Negative Self-Disclosure]

**Tutor:** you know

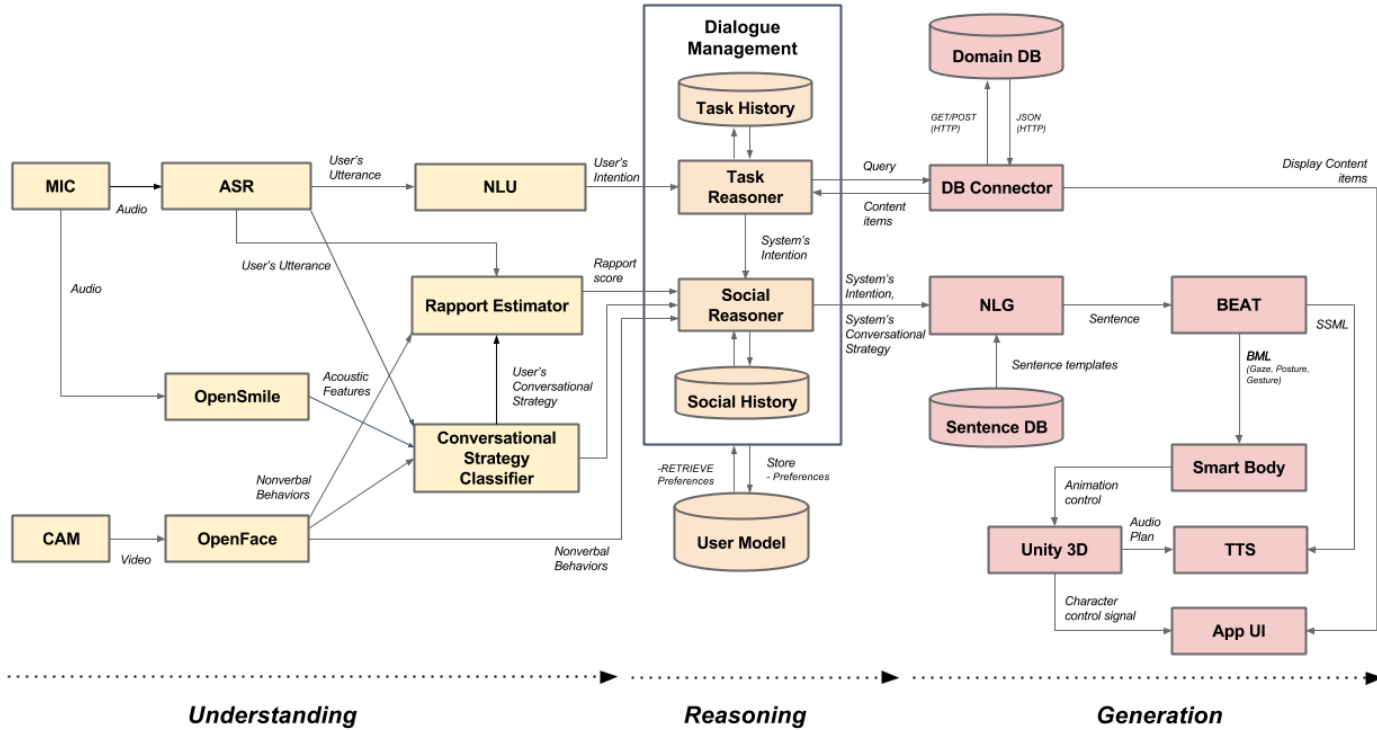
**Tutee:** 6

**Tutor:** next problem is is exactly the same *{smile}*: over 11 equals, 11 x over 11

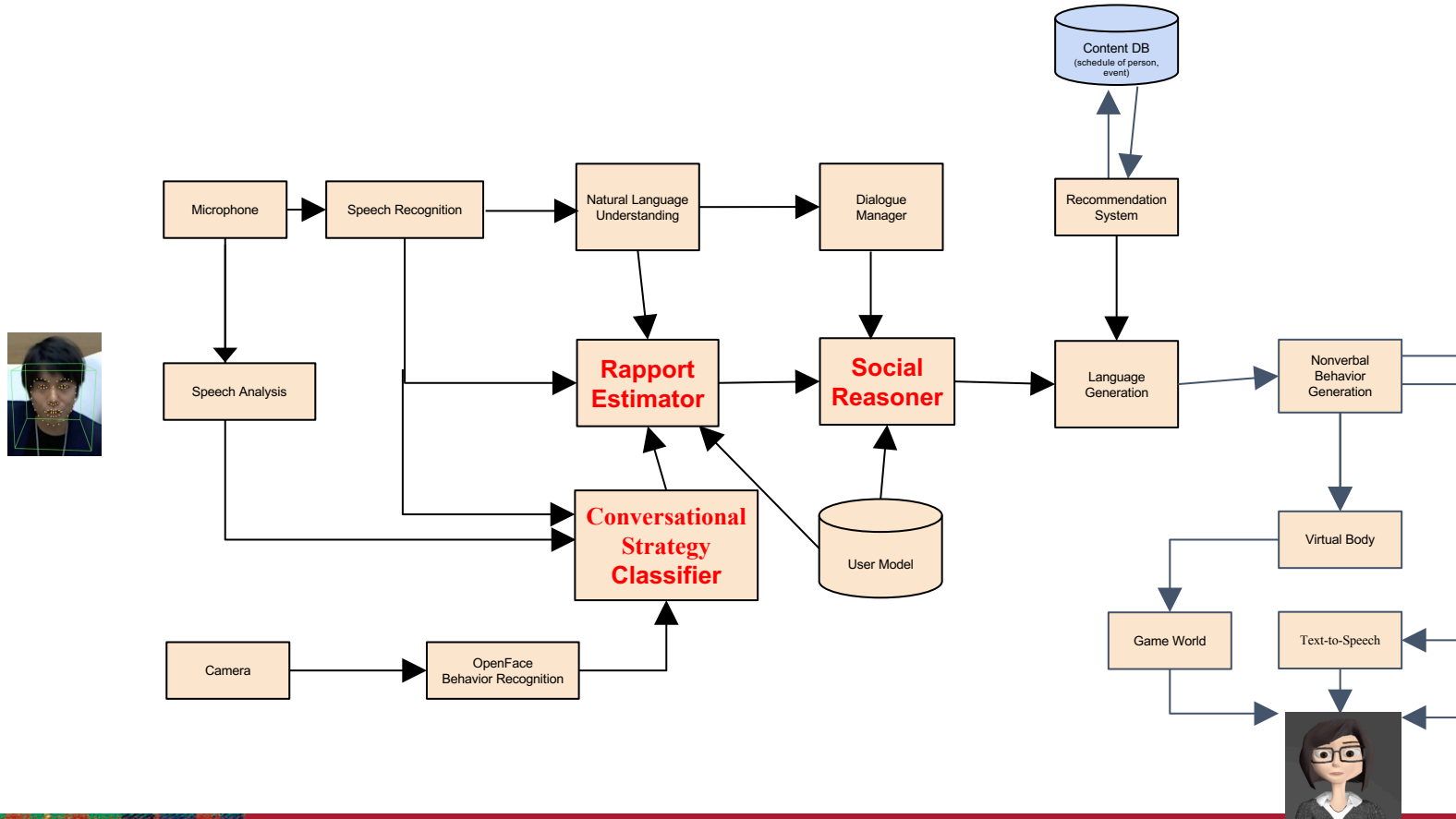
**Tutee:** I don't need your help; [Violation of Social Norm]

**Tutor:** *{Overlap}* That is seriously like exactly the same.

# Socially-Aware Agent Architecture



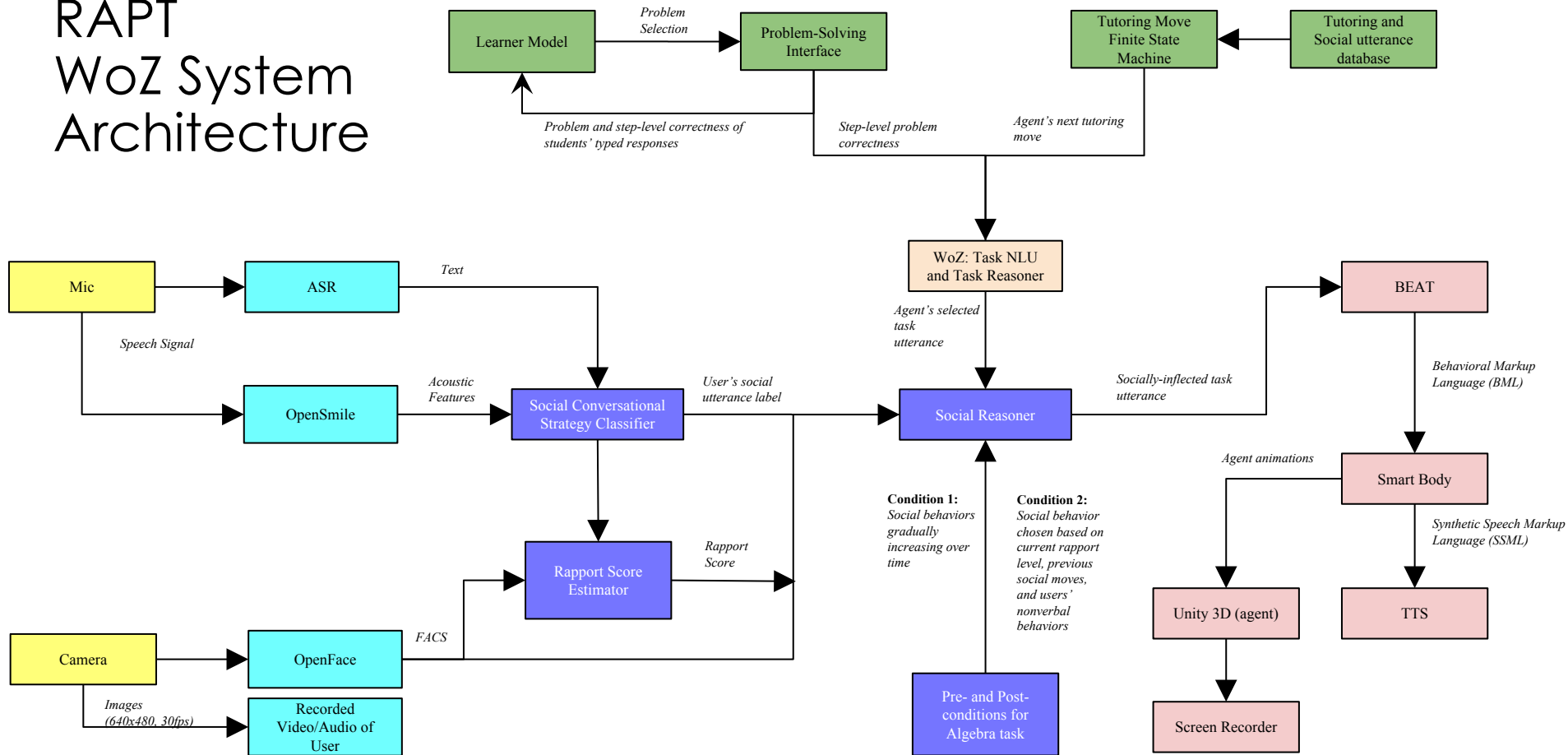
# Socially-Aware Agent Architecture



# Some Applications: Rapport-Aware Peer Tutor (RAPT)

The screenshot displays the RAPT interface for solving the equation  $4(2x + 2) = 16$ . At the top right is a teal "RESTART" button. Below the equation, there are two 2x2 grids for selecting operators: the first grid contains "+", "-", "x", and "÷"; the second grid contains "+", "-", "x", and "÷". To the right of these grids are two teal buttons: "DISTRIBUTE" and "COMBINE LIKE TERMS". Below the grids is an equals sign followed by a blank line for the answer. At the bottom left is a teal button with a "+" icon and the text "Add Next Step". In the bottom right corner, there is a 3D-rendered virtual female tutor with dark hair, wearing a red top, sitting at a desk in a classroom setting with windows and a potted plant.

# RAPT WoZ System Architecture



# Evaluation: rule-based vs. adaptive

**control condition:** (fixed heuristics for social dialogue usage)

## Praise

Decreasing in frequency [Kumar et al., 2010]

## Self-Disclosure

Gradually increasing in frequency and intimacy [Ogan, 2011; Bickmore and Schulman, 2010]

## Questions eliciting self-disclosure

Gradually increasing in breadth and depth of topics [Altman and Taylor, 1973]

## References to shared experiences

Gradually increasing frequency [Kumar et al., 2010; Cassell and Bickmore, 2003]

## Violation of social norms

Use only after a given threshold in number of turns or elapsed time [Ogan et al., 2012]

## Indirectness

Decreasing in frequency [Madaio et al., 2017]

# Evaluation: rule-based vs. adaptive

**Experimental condition:** (adaptive usage of social dialogue)

Based on:

Current rapport state

Changes in rapport state (increasing, decreasing, maintaining)

User's social behavior (self-disclosure, violation of social norms, etc)

Agent's previous social behavior

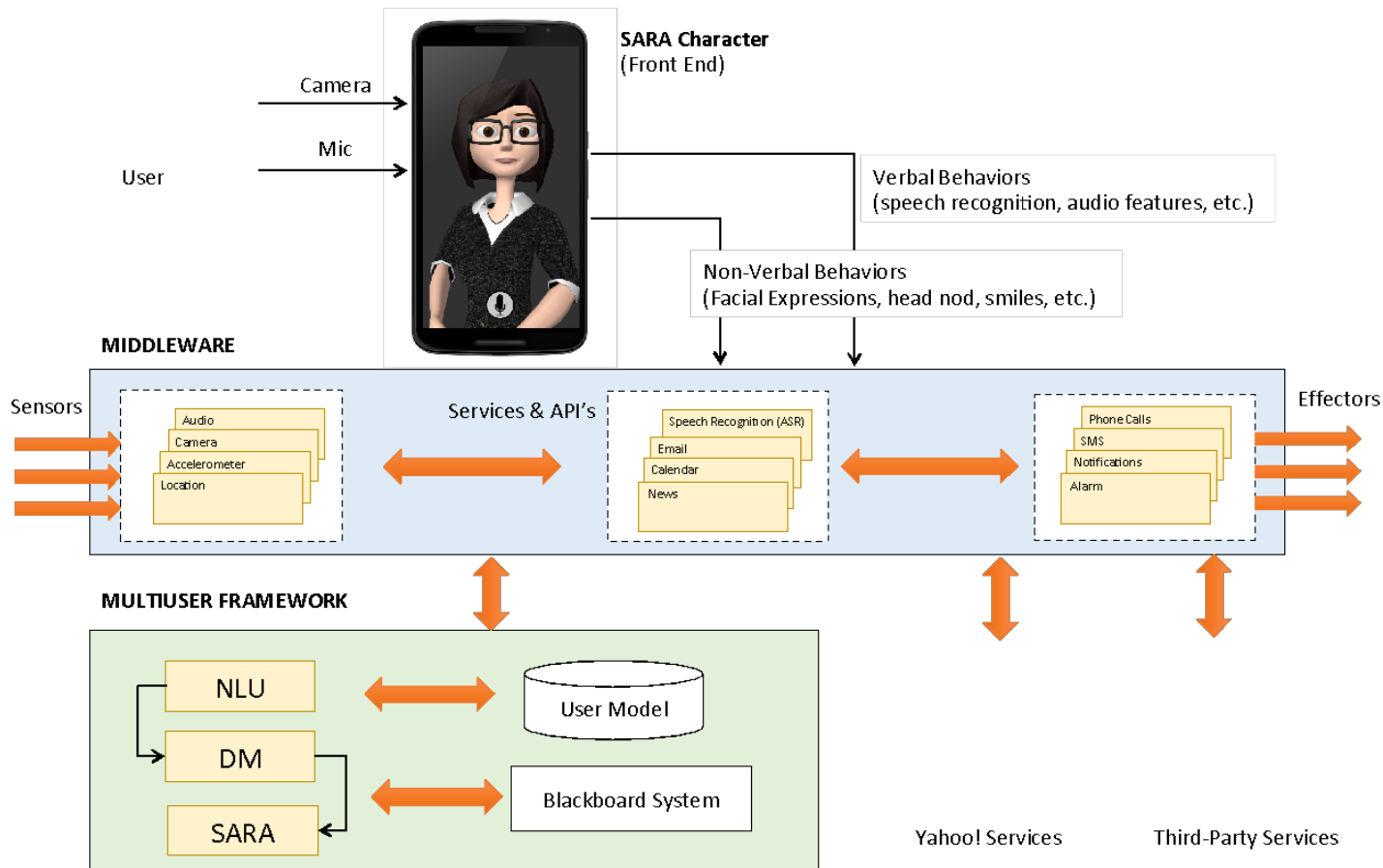
User's nonverbal behavior (smiling, nodding, gaze patterns)

Agent's previous tutoring behaviors (feedback, questions, explanations, etc)



# Application. Mobile front-end to apps

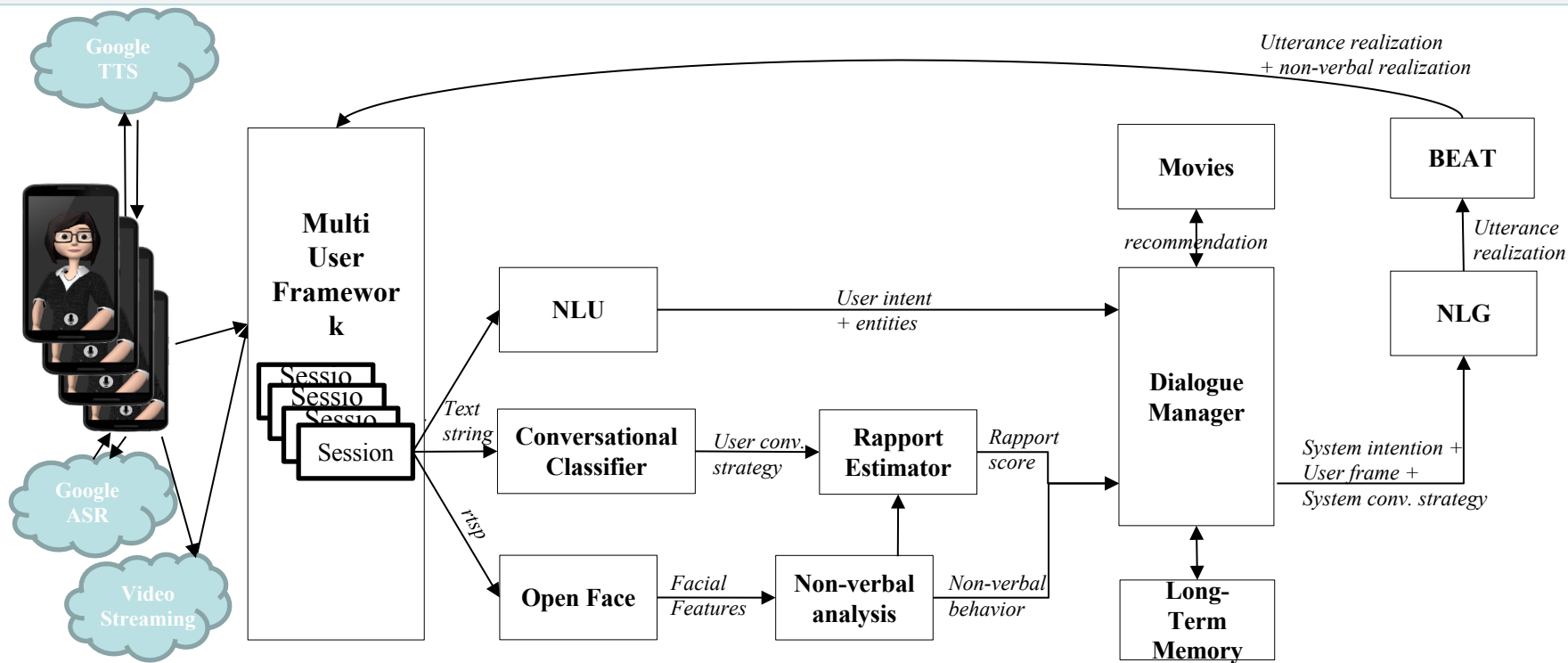




with Yahoo InMind Team

# Integrated InMind Dialog Architecture

## General architecture



# SARA: Socially Aware Robot Assistant



# WHAT SARA UNDERSTANDS

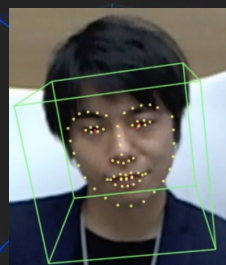
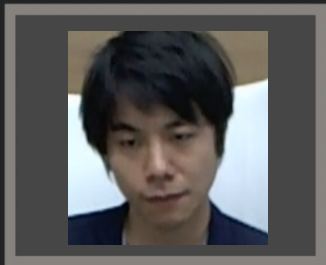


OpenFace

CAMERA FEED

LIVE

HEAD TRACKING...



Smile

NO



Eye Gaze

NO



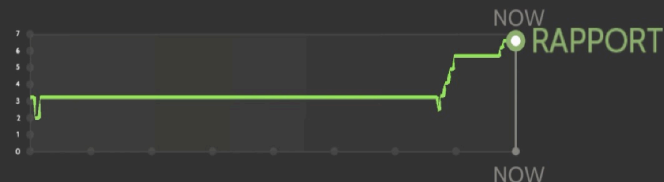
Head Nod

NO



User-Sara Rapport

6 / 7



User Conversation Strategy

SELF DISCLOSURE (SD)



SHARED EXPERIENCE (SE)



PRAISE (PR)



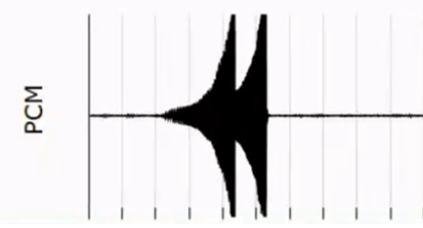
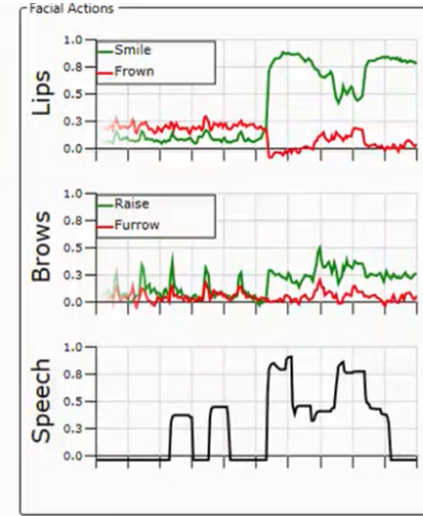
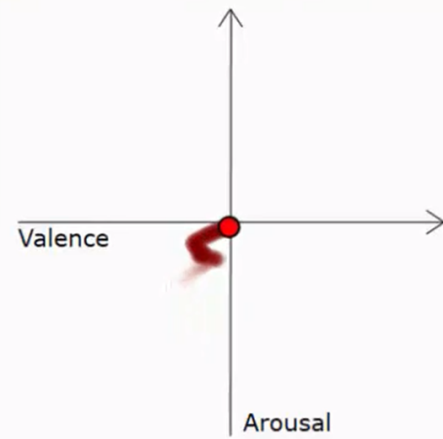
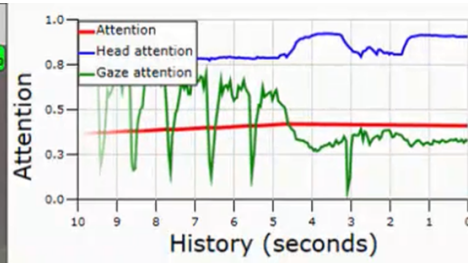
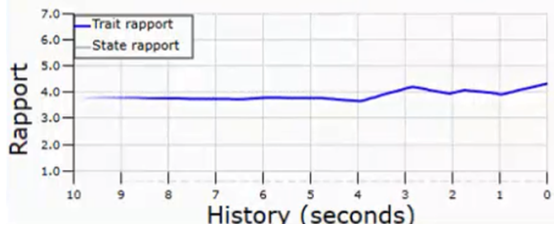
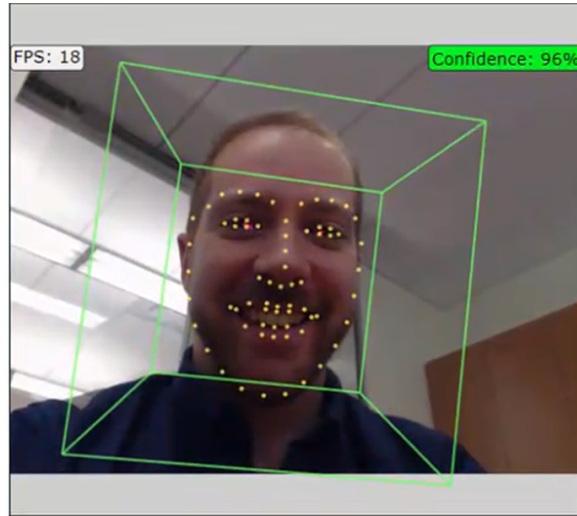
VIOLATE SOCIAL NORM (VSN)



FOLLOW SOCIAL NORM (FSN)



# Tracking Facial Movements



OpenFace: L.P. Morency

## HOW SARA WILL RESPOND

### Sara's Conversational Strategy Selection

#### INPUTS



OpenFace Output



User Conv. Strategy



User-Sara Rapport

6

NOW

CHOSEN STRATEGY

NOW

- SELF DISCLOSURE (SD)
- SHARED EXPERIENCE (SE)
- PRAISE (PR)
- VIOLATE SOCIAL NORM (VSN)
- FOLLOW SOCIAL NORM (FSN)



Sara's Current Task

Sara's Current Task

## WHAT SARA SAYS



### Sara's Words and Body Language



part

let

me look

this

GESTURE ACTIVATE

NOD YES

up

one minute

GESTURE ACTIVATE

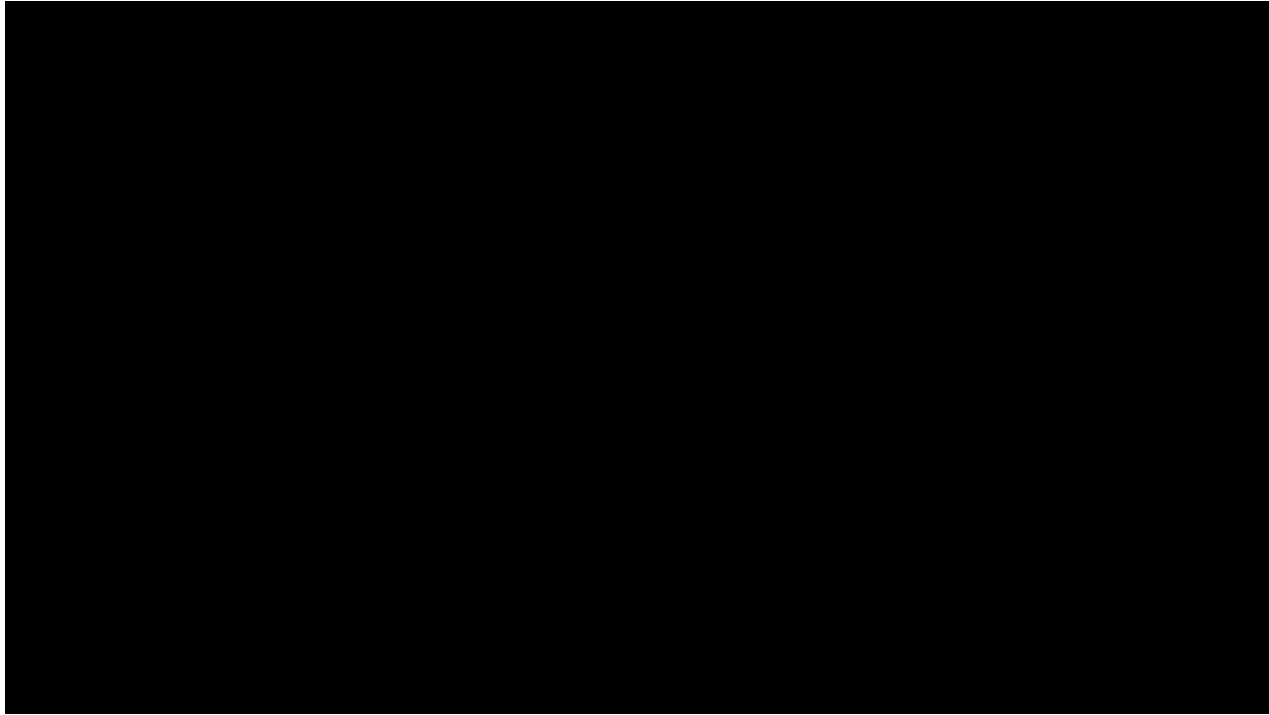
NOD YES

SMILE YES

this is my favorite part let me look this up one minute



# SARA: Socially-Aware Robot Assistant at Davos



# Methodology

## Theorize & Model

Build formal models

Implement system on the basis of model

**Study**

**Build**

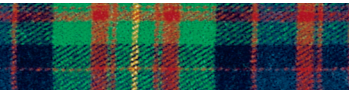
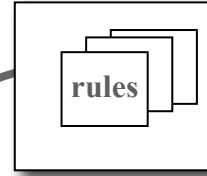
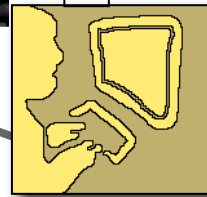
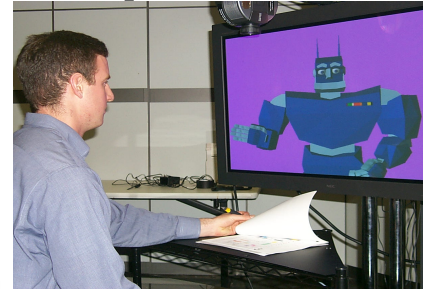
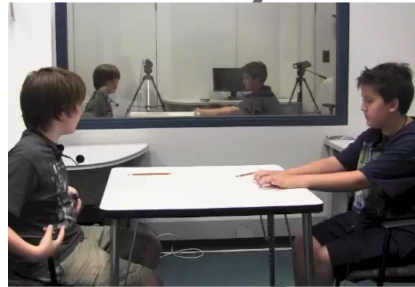
Start here

Collect Natural data

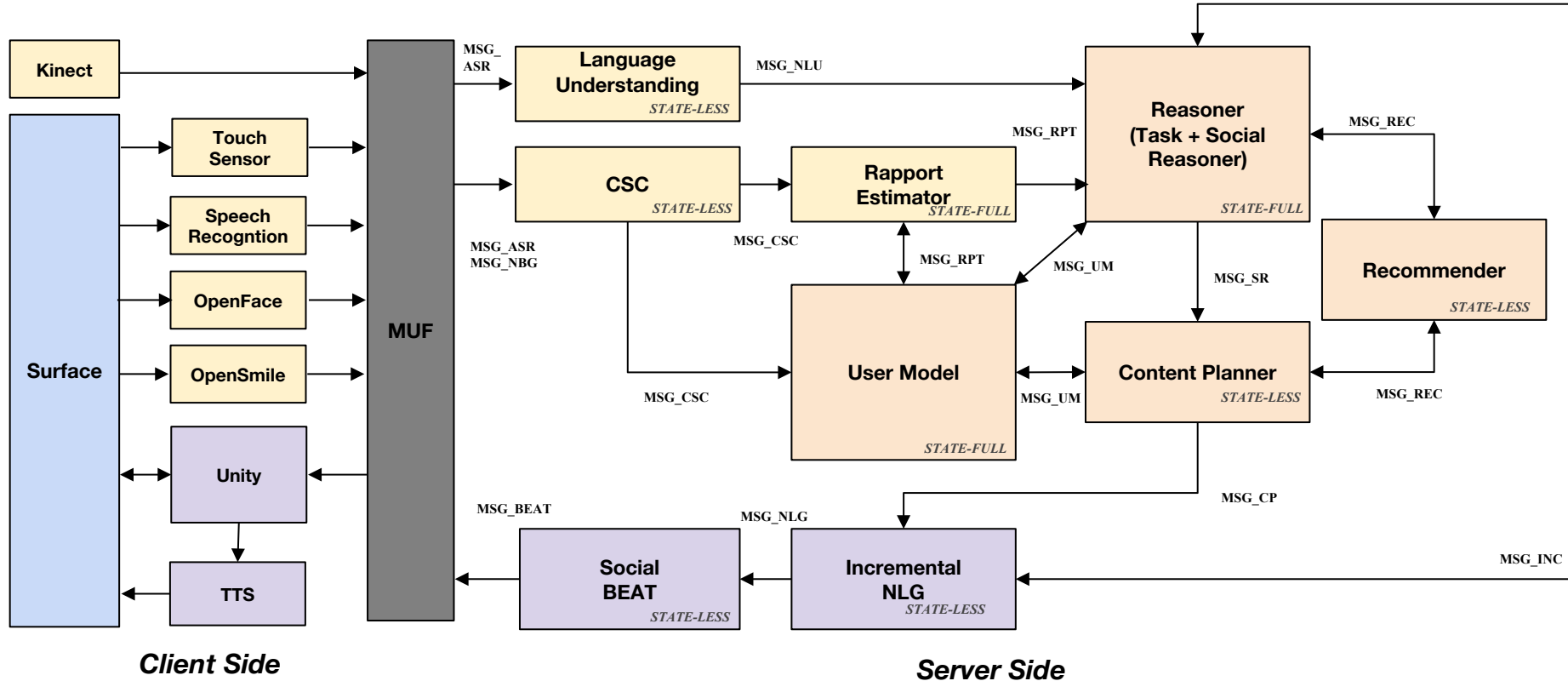
Realize gaps in understanding

Design evaluation of use

**Test**



# New SARA Framework



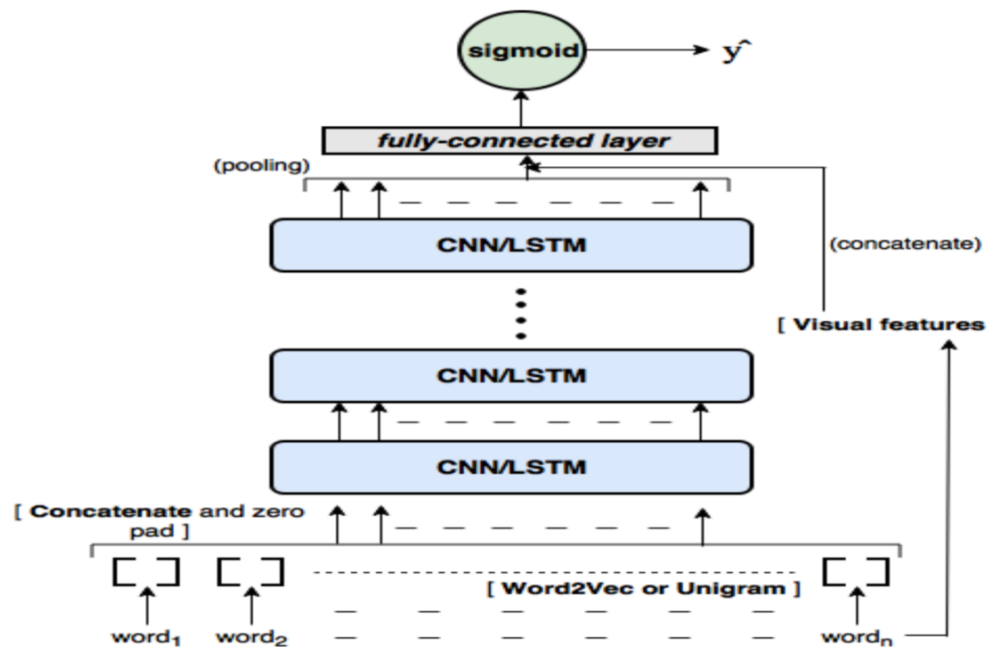
# Indirectness Strategy Classifier

- Corpus
  - RAPT 2013 : indirectness annotation of peer-tutoring corpus
  - ConLL 2010 shared task on **uncertainty** detection
    - Wikipedia dataset (Wikipedia articles)
    - BioScope dataset (abstracts and articles from biomedical literature)

Code	Definition	Example	Distribution
Apology	Apologies used to soften direct speech acts	Sorry, its negative 2.	7.7%
Qualifiers	Qualifying words for reducing intensity or certainty	You just add 5 to both sides.	66.1%
Extenders	Indicating uncertainty by referring to vague categories	You have to multiply and stuff.	3.6%
Subjectivizer	Making an utterance seem more subjective to reduce intensity	I think you divide by 3 here.	22.6%

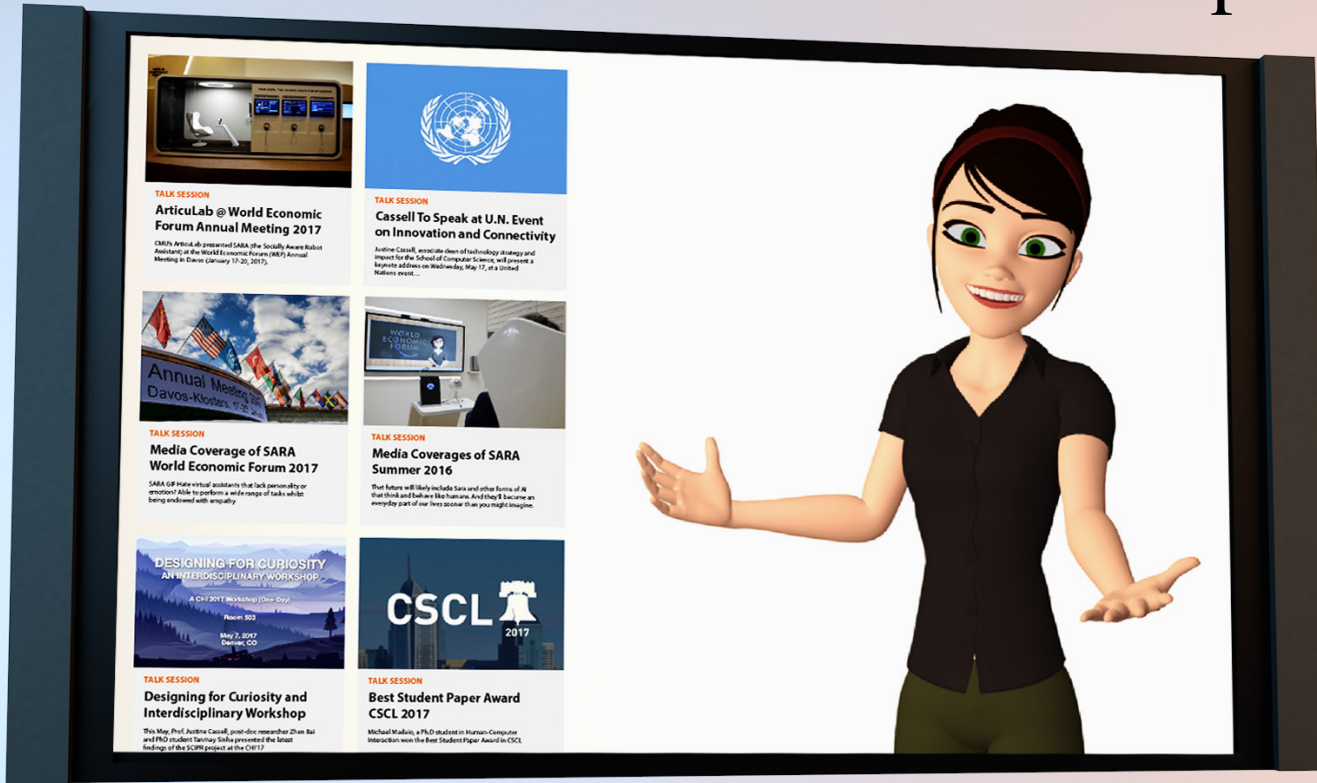
**IN (Indirect Delivery)** “*so I think what I'm gonna do is make that 15 minus 3 a 12*”

# Architecture: Indirectness Classifier



Pranav Goel, Yoichi Matsuyama, Michael Madaio & Justine Cassell, "I think it might help if we multiply, and not add": Detecting Indirectness in Conversation, International Workshop on Spoken Dialog System Technology (IWSDS 2018). – to appear

# SARA Receptionist



TALK SESSION

### ArticLab @ World Economic Forum Annual Meeting 2017

CMU ArticLab presented SARA (the Socially Aware Robot Assistant) at the World Economic Forum (WEF) Annual Meeting in Davos (January 17-20, 2017).



TALK SESSION

### Cassell To Speak at U.N. Event on Innovation and Connectivity

Justin Cassell, associate dean of technology strategy and impact for the School of Computer Science, will present a keynote address on Wednesday, May 17, at a United Nations event...



TALK SESSION

### Media Coverage of SARA World Economic Forum 2017

SARA QR Hubs virtual assistants that lack personality or emotional skills to perform a wide range of tasks whilst being endowed with empathy.



TALK SESSION

### Media Coverages of SARA Summer 2016

The future will likely include Sara and other forms of AI that think and behave like humans. And they'll become an everyday part of our lives sooner than you might imagine.



TALK SESSION

### Designing for Curiosity and Interdisciplinary Workshop

The May Prof. Justin Cassell, post-doc researcher Chen Bai and PhD student Tammy Shih presented the latest findings of the SCIPN project at the CIP17.



TALK SESSION

### Best Student Paper Award CSCL 2017

Michael Madala, a PhD student in Human-Computer Interaction won the Best Student Paper Award in CSCL.

With thanks for generous funding  
and support to the following  
organizations, and to the students  
and staff of the ArticuLab

