THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY

DEPARTMENT OF
COMPUTER SCIENCE & ENGINEERING

PhD Qualifying Exam

# A Survey on Intelligent User Interfaces for the Learning of Verbal Communication Skills

Presenter: Xingbo Wang
Supervisor: Prof. Huamin Qu

2020 June 9

# Outline

# Motivation

## Background

**Verbal communication skills** :Proper usage of words and sounds to deliver message

What are other words for verbal communication?

speech, speaking, talking, articulation, dialogue, talk, conversation … …

Introduction | Automatic Assessment | User Interfaces | Conclusion

# Motivation

**Background**




Public speaking


Everyday conversation


Job interview

What are other words for verbal communication?

speech, speaking, talking, articulation, dialogue, talk, conversation ... ...

## Verbal communication skills

Adaptive speech content

Engaging vocal delivery

THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY

# Motivation

## Background

What are other words for verbal communication?

speech, speaking, talking, articulation, dialogue, talk, conversation … …

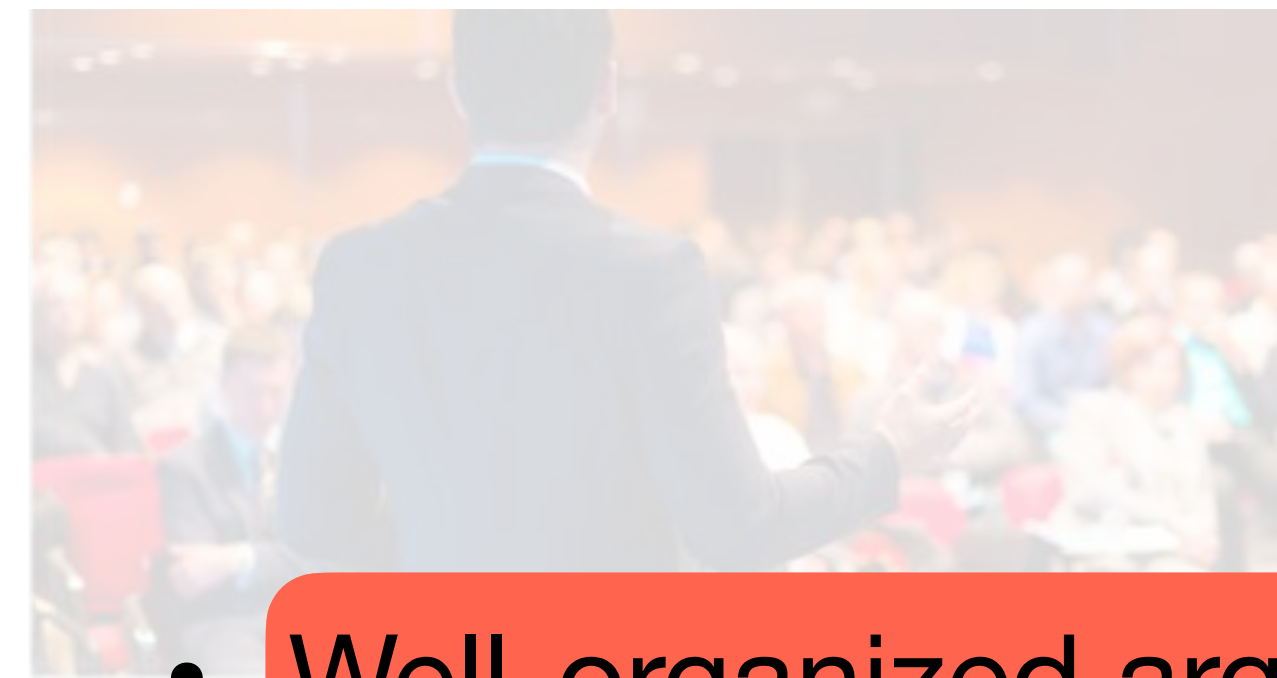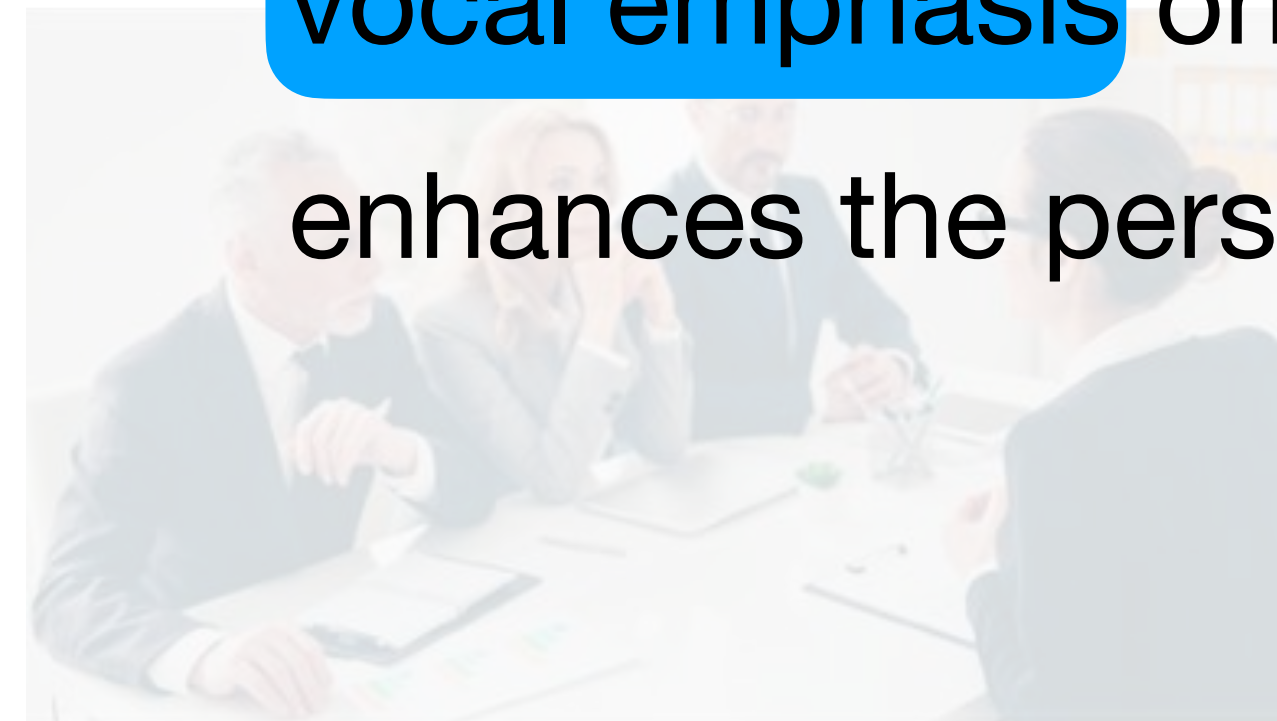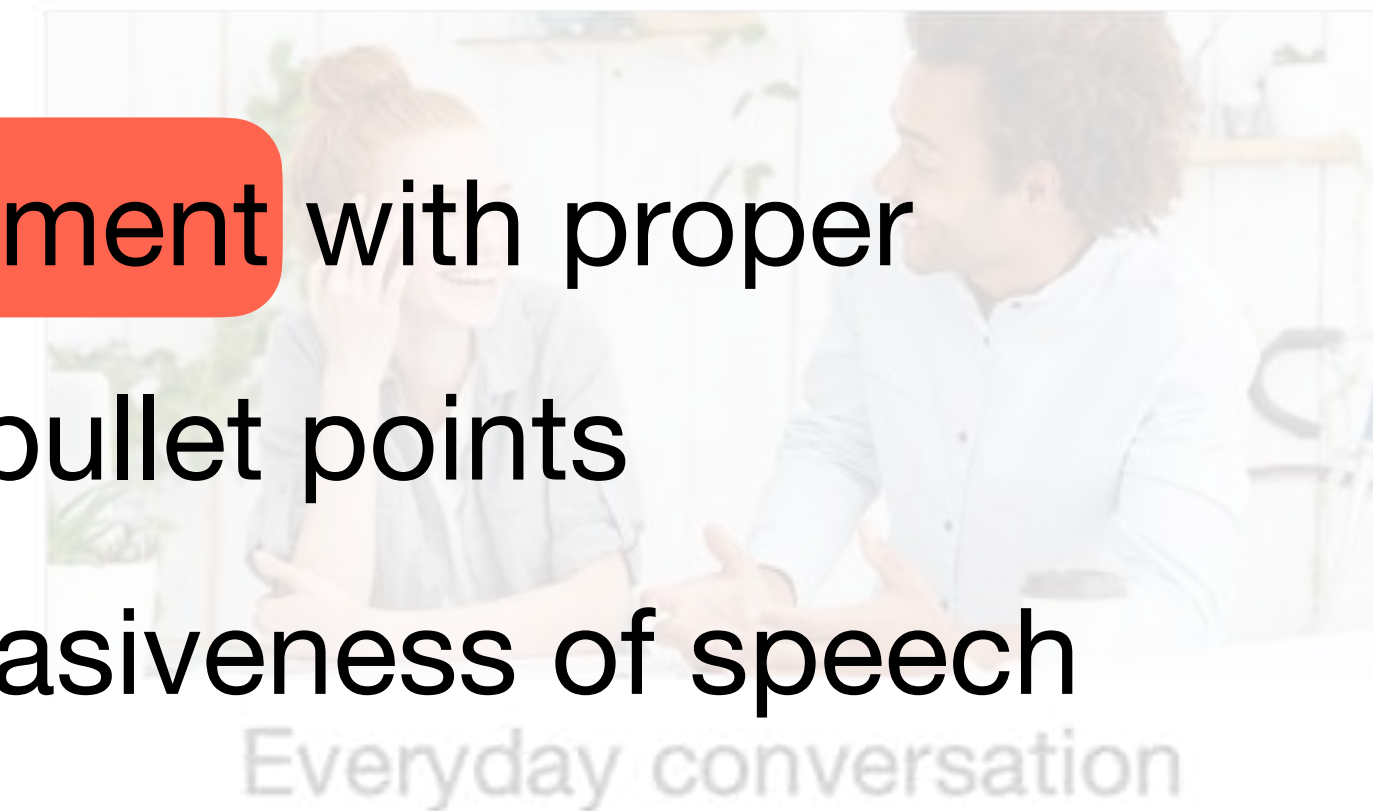**Verbal communication skills**

Adaptive speech content

Engaging vocal delivery

- Well-organized argument with proper vocal emphasis on bullet points enhances the persuasiveness of speech
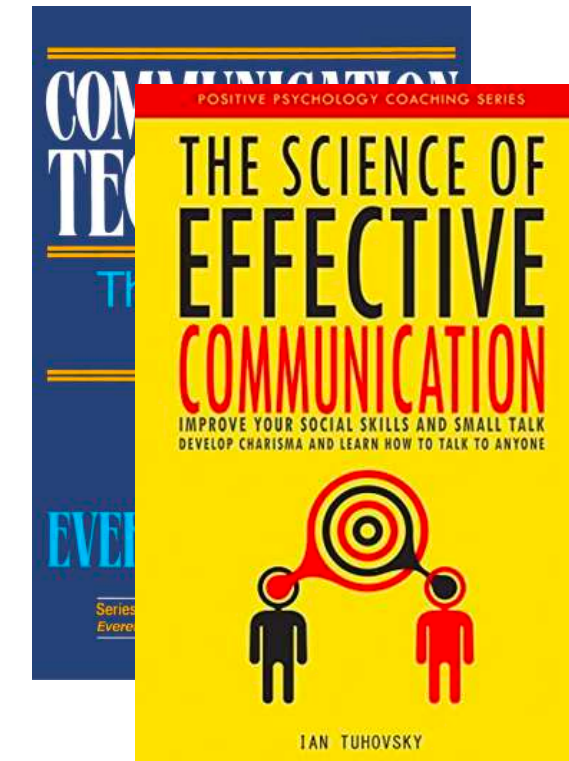
Everyday conversation

Job interview

# Motivation

## Verbal communication skills learning

- **Self-learning**: guidelines from books

  No feedback

- **Professional training**: feedback from coaches
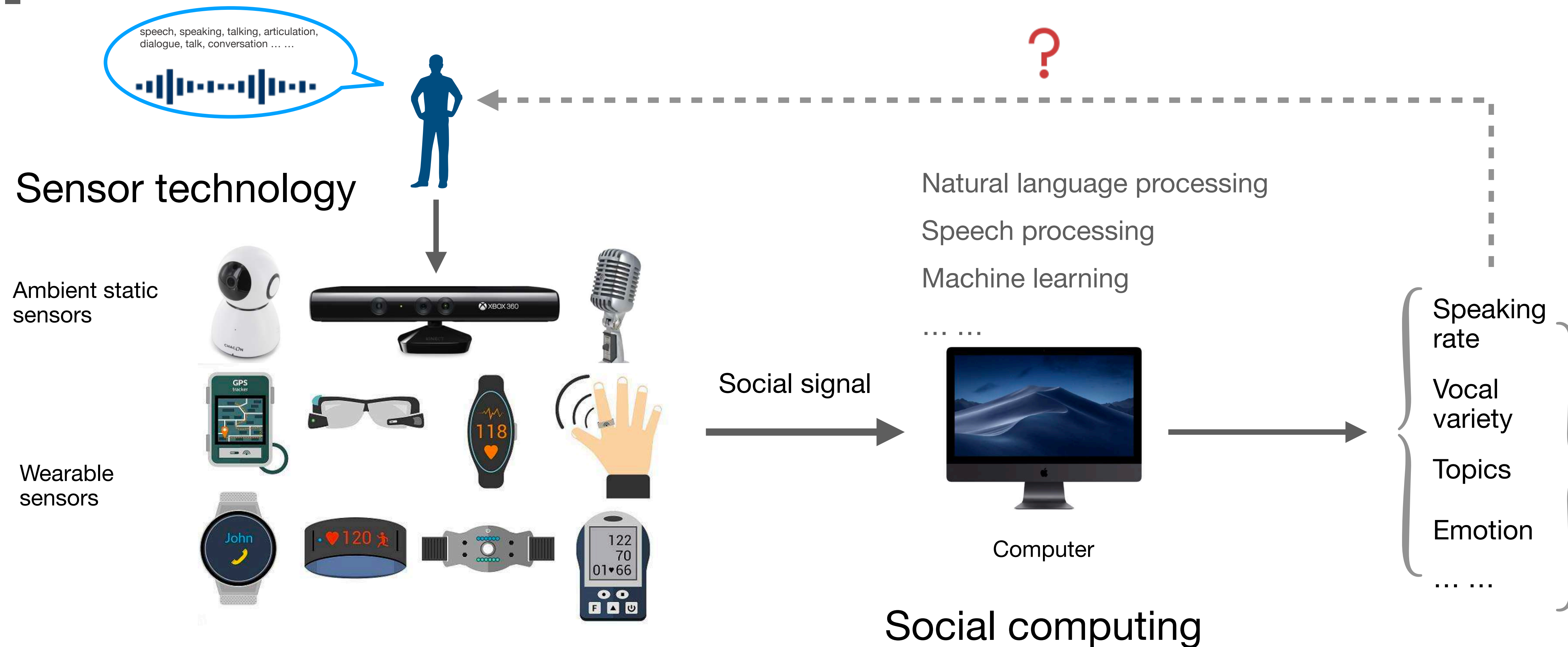
  - Qualitative
  - Inflexible

There is a lack of

- Quantitative & automated feedback

- Tool support for effective learning

Introduction | Automatic Assessment | User Interfaces | Conclusion

# Motivation

THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY

## Quantitative automatic feedback



speech, speaking, talking, articulation, dialogue, talk, conversation … …

Sensor technology

Ambient static sensors

Wearable sensors

Natural language processing

Speech processing

Machine learning

… …

Social signal

Computer

Social computing

Speaking rate

Vocal variety

Topics

Emotion

… …

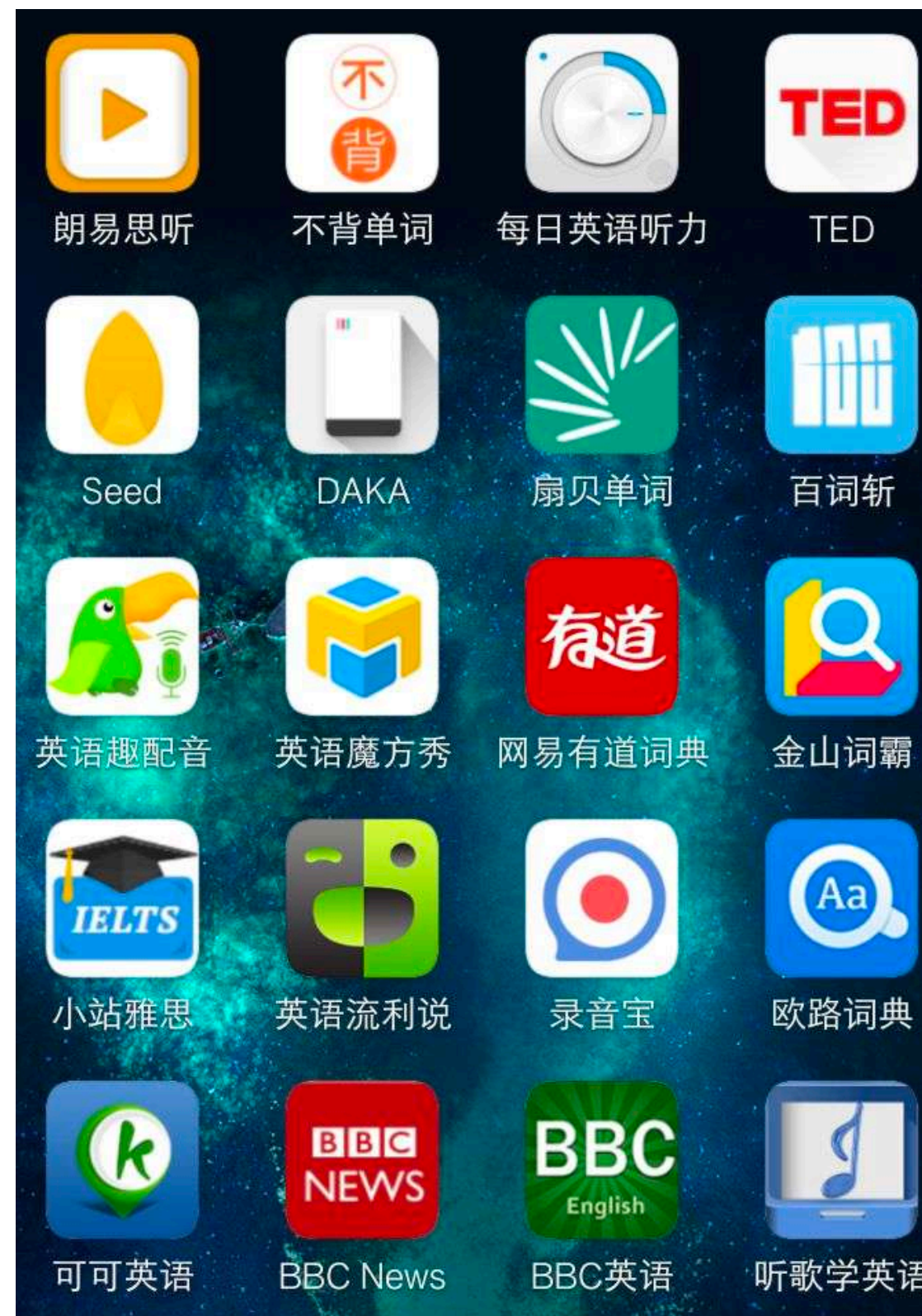Introduction | Automatic Assessment | User Interfaces | Conclusion

# Motivation

## User interfaces for learning



Knowledge learning platforms



Language learning softwares

Not for "soft skills"

- They have more clear criteria of what are "correct"

- Their feedback is mostly offline

- Their feedback is mostly in visual forms

# Motivation

## Challenges

### Intelligence

Quantitative automatic feedback

- Derive quantitative descriptors about speech behavior

- Assess multimodal speech behavior

### Learning

User interfaces for effective learning

- Offer valuable guidance on users' behavior

  - When ? (e.g., realtime or post hoc feedback)

  - How ? (e.g., through visuals or other forms)

  - …

Introduction → Automatic Assessment → User Interfaces → Conclusion

# Outline

Introduction

## Automatic Assessment

Competence rubrics

Computational features

Performance assessment

User Interfaces

Conclusion

# Automatic Assessment

## Competence rubrics

Researchers have conducted a number of studies to identify core communication competency and its rubrics for the practice of communication skills

(Quianthy, 1990; Lucas, 2007; Morreale et al., 2007; Rhodes, 2010; Thomson & Rucker, 2002)

### Core aspects of communication proficiency

**Adaptive speech content**

Topic selection, support material usage, idea organization, word choices

**Engaging vocal delivery**

Vocal variety, articulation, non-verbal behavior

Introduction | Automatic Assessment | User Interfaces | Conclusion

# Automatic Assessment

## Competence rubrics

Researchers have conducted a number of studies to identify core communication competency and its rubrics for the practice of communication skills

(Quianthy, 1990; Lucas, 2007; Morreale et al., 2007; Rhodes, 2010; Thomson & Rucker, 2002)

Adaptive speech content

Engaging vocal delivery

| Performance standard The student … | Assessment Criteria | | | | |
|---|---|---|---|---|---|
| | Advanced 4 | Proficient 3 | Basic 2 | Minimal 1 | Deficient 0 |
| Demonstrates a careful choice of words | Language is exceptionally clear, imaginative and vivid; completely free from bias, grammar errors and inappropriate usage | Language appropriate to the goals of the presentation; no conspicuous errors in grammar; no evidence of bias | Language selection adequate; some errors in grammar; language at times misused (e.g., jargon, slang, awkward structure) | Grammar and syntax need to be improved as can level of language sophistication; occasionally biased | Many errors in grammar and syntax; extensive use of jargon, slang, sexist/racist terms or mispronunciations |
| Effectively uses vocal expression and paralanguage to engage the audience | Excellent use of vocal variation, intensity and pacing; vocal expression natural and enthusiastic; avoids fillers | Good vocal variation and pace; vocal expression suited to assignment; few if any fillers | Demonstrates some vocal variation; enunciates clearly and speaks audibly; generally avoids fillers (e.g., um, uh, like) | Sometimes uses a voice too soft or articulation too indistinct for listeners to comfortably hear; often uses fillers | Speaks inaudibly; enunciates poorly; speaks in monotone; poor pacing; distracts listeners with fillers |

Public Speaking Competency Rubric (PSCR) (Schreiber et al., 2012)

Introduction → Automatic Assessment → User Interfaces → Conclusion

# Automatic Assessment

## Computational features

Vocal delivery          Speech content

| | Prosodic Features | Language Features | Other Features |
|---|---|---|---|
| low-level Features | tempo, loudness, pitch | Syntactic features (e.g., Part-of-Speech features), semantic meaning | - |
| High-level Features | speaking rate, liveliness, fluency | word choices (e.g., PMI, LIWC), topics, content features, presentation state, content structure | emotion, stage atmosphere (e.g., laughter, applause) |
| Feature Encodings | freq., max., min., avg., std, multi., quot., add, minus | | |

Computational features summarized from previous work

*Tempo, loudness, pitch* -> auditory perception of a speech

**Articulation**
- *Speaking rate: syllables/ words/sentences per minute*
- *fluency: smoothness of speech -> filled pauses, filler words (e.g., "em","hmm")*

**Vocal variety**
- *liveliness: expressiveness of voice -> intonation -> variation of pitch and volume*

| Introduction | Automatic Assessment | User Interfaces | Conclusion |

13

# Automatic Assessment

## Computational features

Vocal delivery          Speech content

| | Prosodic Features | Language Features | Other Features |
|---|---|---|---|
| low-level Features | tempo, loudness, pitch | Syntactic features (e.g., Part-of-Speech features), semantic meaning | - |
| High-level Features | speaking rate, liveliness, fluency | word choices (e.g., PMI, LIWC), topics, content features, presentation state, content structure | emotion, stage atmosphere (e.g., laughter, applause) |
| Feature Encodings | freq., max., min., avg., std, multi., quot., add, minus | | |

Computational features summarized from previous work

**Word usage**

=> commonness: *PMI*

=> psychology: *LIWC, sentiment/subjectivity lexicons*

**Topics** => LDA

**Content organization** => BoW, Word2Vec

**Adaptation**

=> *Presentation state: presentation & QA*

Introduction | Automatic Assessment | User Interfaces | Conclusion

# Automatic Assessment

## Computational features

Vocal delivery          Speech content

| | Prosodic Features | Language Features | Other Features |
|---|---|---|---|
| low-level Features | tempo, loudness, pitch | Syntactic features (e.g., Part-of-Speech features), semantic meaning | - |
| High-level Features | speaking rate, liveliness, fluency | word choices (e.g., PMI, LIWC), topics, content features, presentation state, content structure | emotion, stage atmosphere (e.g., laughter, applause) |
| Feature Encodings | freq., max., min., avg., std, multi., quot., add, minus | | |

Computational features summarized from previous work

**Speaker engagement**

=> emotion

**Audience engagement**

=> stage atmosphere (e.g., laughter, applause, booing)

# Automatic Assessment

## Performance assessment

### Rule-based methods

Based on the statistical properties of features to set the range of "good"/"bad" performance
(e.g., std, freq., mean)

### Vocal delivery

Pitch variety, speech speed => mean+std of words/sentences

### Speech content

Content coverage => spotted keywords in speech and their weights (tf.idf) / text in slides

It is simple and useful for basic features. However,

- it is intricate to decide thresholds for complex features (e.g., emotion)
- it fails to adapt to different speakers and different speaking scenarios

Introduction | Automatic Assessment | User Interfaces | Conclusion

# Automatic Assessment

## Performance assessment

Machine learning (supervised learning)



Performance is much inferior to human evaluation

Rubrics → Human raters

speech, speaking, talking, articulation, …

$f$

Computational behavior descriptors

Machine learning models

Verbal communication

Human evaluation examples

Quality control: inter-agreement

Supervised learning models

**Q1:** Models <=?=> judgements

- SVM/SVR, L1/L2 Regularized Logistic Regression, Lasso, tree-based models (e.g., RF)
- BN, HMM

**Q2:** Features <=?=> judgements

- Correlation Coefficients

Introduction | Automatic Assessment | User Interfaces | Conclusion

# Automatic Assessment

## Performance assessment

Deep learning

CNN (Krizhevsky et al., 2012), LSTM (Hochreiter and Schmidhuber, 1997), Transformer (Vaswani et al., 2017) and their variations achieve impressing results on complex analytical tasks of human communication understanding

- CNN (Hershey et al., 2017) => Audio event detection and classification

- BERT (Devlin et al., 2018) => Various NLP tasks

- MFN (Zadeh et al. 2018) => Multimodal feature fusion

It is difficult for human to understand and interpret the model results

# Outline

Introduction

Automatic Assessment

**Intelligent User Interface**

Taxonomy

Prior feedback

Live feedback

Posterior feedback

Conclusion

# Intelligent User Interface

## Taxonomy

# Intelligent User Interface

## Taxonomy

**1** Considering learning scenarios,



Speaking anxiety

Content organization

Vocal delivery

Stage management

Group interactions & dynamics

Design focus

Introduction | Automatic Assessment | User Interfaces | Conclusion

# Intelligent User Interface

## Taxonomy

**2** Considering the learning process,

Where am I going?
(learning goals)

**Perception**

Learning cycle

(Dewey, 1993), (Lewin, 1946),
(Kolb, 1975), (Mumford, 1995)

Where do I go next?
(self-adjustment)

**Action**

**Reflection**

How am I going?
(self-awareness)

Introduction | Automatic Assessment | User Interfaces | Conclusion

22

# Intelligent User Interface

## Taxonomy

**2** Considering the learning process,



Exploration of knowledge base

Reflection of vocal/verbal behavior

Putting knowledge into action

# Intelligent User Interface

## Taxonomy

Feedback has been considered as an <u>effective intervention</u> in skills learning and a <u>key consideration</u> of learning interfaces

**3** Considering the channels of feedback,



Most widely used for feedback on speech content and vocal delivery

Reduce cognitive load

Introduction  Automatic Assessment  User Interfaces  Conclusion

24

# Intelligent User Interface

## Taxonomy

**4** Considering the feedback strategy (in the learning cycle),

# Intelligent User Interface

## Taxonomy

**4** Considering the feedback strategy (in the learning cycle),



Prior feedback
- content–based cues
- delivery–based cues

Feedback strategy

Live feedback
- implicit feedback
- explicit feedback
  - simple verification
  - elaborated feedback
  - termination

Posterior feedback
- summary feedback
- focused feedback

**Evidences & consequences**

Excellent learning examples
↓
Good speech delivery needs varied intonation

**Feedback complexity**

Avoid information overload

**Focus of self-reflection**

Exploration from coarse to fine

Introduction | Automatic Assessment | User Interfaces | Conclusion

# Intelligent User Interface

**Prior feedback:** **Delivery-based cues** ≫ **Content-based cues**

Exploring narration strategies (pitch, pause, volume)



SpeechLens (Yuan et al., 2019)

Sentence-level

A: context+focus design for visualization of prosodic features

B: Structural query

Word-level

C: Word clouds for phrase intonation

Introduction → Automatic Assessment → User Interfaces → Conclusion

27

# Intelligent User Interface

**Prior feedback:** **Delivery-based cues** ≫ Content-based cues

Exploring emotion coherence in presentation



A: video barcode charts & coherence curves → overview

B: an augmented Sankey diagram → channel connection

C: sentence clustering view → emotion combinations & its distribution

EmoCo (Zeng et al., 2019)

Introduction | Automatic Assessment | User Interfaces | Conclusion

# Intelligent User Interface

**Prior feedback:** **Delivery-based cues** ≫ Content-based cues

From feature exploration ➡ Speech styles generation for voice-over



It can automatically modify the pitch, volume and duration curves to generate desired emphasis and flow.

It is difficult for novice speakers to identify the words for resynthesis

NarrationCoach (Rubin et al., 2015)

Introduction | Automatic Assessment | User Interfaces | Conclusion

# Intelligent User Interface

**Prior feedback:** **Delivery-based cues** » Content-based cues

Data-driven recommendation of voice modulation techniques



Overview of VoiceCoach (Wang et al., 2020)

Introduction | Automatic Assessment | User Interfaces | Conclusion

# Intelligent User Interface

## Prior feedback: Delivery-based cues ≫ Content-based cues

Data-driven recommendation of voice modulation techniques



N-gram based hierarchical summary

VoiceCoach (Wang et al., 2020)

# Intelligent User Interface

## Prior feedback: Delivery-based cues ≫ Content-based cues

Data-driven recommendation of voice modulation techniques



N-gram based hierarchical summary

VoiceCoach (Wang et al., 2020)

## **Prior feedback:** **Delivery based cues** » Content-based cues

### Data-driven recommendation of modulation techniques

**Query**: Tact is … making **an enemy**



One-line mode

Multi-line mode

VoiceCoach        l., 2020)

### **Speech examples**

- modulation of interest is highlighted

- context for phrase of interest

- listening to original audio clips

Introduction | Automatic Assessment | User Interfaces | Conclusion

33

# Intelligent User Interface

**Prior feedback:** Delivery-based cues ≫ **Content-based cues**

Data-driven interaction for video navigation



LectureScape (Kim et al., 2014)

- Timeline (1)

- Search (2)

- Summarization (3)

Introduction → Automatic Assessment → User Interfaces → Conclusion

34

# Intelligent User Interface

**Prior feedback:** Delivery-based cues ≫ **Content-based cues**

Topic-based content summarization



MMToC (Biswas et al., 2015)

**Visual & spoken** word fusion

- Extend visual salient words with a group of spoken salient words based on the semantic similarity

Content segmentation

- Minimize the inner-group difference
- Maximize the inter-group difference

# Intelligent User Interface

**Prior feedback:** Delivery-based cues ≫ **Content-based cues**

Textbook-inspired chapter/section content organization



Video Digest (Pavel et al., 2014)

Chapter — Topically coherent sections

Section — A set of varying topics

Bayesian topic segmentation

Crowdsourcing summary and ranking

Introduction | Automatic Assessment | User Interfaces | Conclusion

# Intelligent User Interface

**Prior feedback:** Delivery-based cues ≫ **Content-based cues**

From content exploration ➜ Hierarchical content structure planning



HyperSlides (Edge et al., 2014)

| Action | Syntax | Explanation |
|---|---|---|
| Create scenes | [Scene 1 < image1.jpg]<br>[Scene 2 < image2.jpg] | Create scene slides with titles "Scene X" that have the background imageX.jpg. |
| Add details | [Scene 1 < image1.jpg]<br>[> Point A]<br>[>> Point A1]<br>[>> Point A2]<br>[> Point B] | Add Point A and Point B as details of Scene 1, with Point A1 and Point A2 sub-details of Point A. A third level of detail is possible using [>>>...], and so on. |
| Add hyperlinks | [> Point A >> http://url.tld]<br>[> Point B >> anyfile.ext] | Link from Point A to a URL.<br>Link from Point B to a file. |

Mark-up language to create hierarchically structured scenes

⬇

hyperlinked slides of a consistent and minimalist style

# Intelligent User Interface

**Prior feedback:** Delivery-based cues ≫ **Content-based cues**

From manual planning ➡ Automatic structure generation & path suggestions



NextSlidePlease (Spicer et al., 2012)

**Algorithm support**

- Semantic similarity between adjacent slides => presentation graph

- Time constraints, priority => path suggestions



Linear layout

Structure generation                    Post editing

# Intelligent User Interface

**Prior feedback:**

## Limitations

- Do not consider learning from <span style="color:red">BAD</span> examples

- Do not consider learning from <span style="color:red">multimodal</span> speech styles

**NEXT**: Live feedback (implicit feedback)

Introduction > Automatic Assessment > User Interfaces > Conclusion

# Intelligent User Interface

## Live feedback: **Implicit feedback** ≫ Explicit feedback

Simulate nonverbal behavior of virtual audience



- Posture (e.g., straight, relaxed, forward)
- Head orientation
- Gaze

Virtual audiences in Cicero (Batrinca et al., 2013)



MACH (Hoque et al., 2013)

**Acknowledgement**

- nodding, changing posture

- spoken acknowledgement "That's very interesting"

Introduction | Automatic Assessment | User Interfaces | Conclusion

40

# Intelligent User Interface

## Live feedback:  Implicit feedback ≫ Explicit feedback

Investigate the impact of nonverbal behavior of virtual audience



Two characters of virtual coaches (Gebhard et al., 2013)

|  | *Gestures* | *Facial expression* | *Pause* | *Gaze* | *Comments* |
|---|---|---|---|---|---|
| **Understanding** | Narrow | Positive | Shorter | Friendly | Many polite phrases |
| **Demanding** | Space-taking | Negative | Longer | Dominant | Few polite phrases |

**Within-subject study**

- Participants perceive the differences and they reported that **demanding** character induced higher level of stress

- Demanding condition: more breathing **pauses**, higher **movement** energy

# Intelligent User Interface

**Live feedback:** **Implicit feedback** ≫ Explicit feedback

## Limitations



- Most listening behaviors of virtual audience are controlled by finite state machines. There is a lack of more intelligent models to simulate affective states of listeners

**NEXT**: Explicit feedback (simple verification)

Introduction   Automatic Assessment   User Interfaces   Conclusion

42

# Intelligent User Interface

## Live feedback: Implicit feedback ≫ Explicit feedback - simple verification

Realtime behavioural checking on speech delivery with Google Glass



Logue (Hoque et al., 2013)

Feedback icon alternatives



Final design

Introduction ➤ Automatic Assessment ➤ User Interfaces ➤ Conclusion

43

# Intelligent User Interface

## **Live feedback:** Implicit feedback ≫ **Explicit feedback - simple verification**

Realtime behavioural checking on speech content with HMD



A wearable MC system (Okada et al., 2011)

**Manage stage**

- Realtime tracking of MC's speech
- Communication with operators
- Sensing atmosphere (e.g., buzzing, laugh)



Introduction | Automatic Assessment | User Interfaces | Conclusion

# Intelligent User Interface

## Live feedback: Implicit feedback ≫ Explicit feedback - elaborated feedback

Realtime instructions on speech delivery with visual feedback



Rhema (Okada et al., 2011)



VoiceAssist (Okada et al., 2011)

**User study**

- Verbal feedback is most favored

- Participants prefer sparse feedback to continuous feedback

Room acoustics: speech transmission index

Background noise: signal to noise ratio

Box: red/green indicates low/high audio quality

**Live feedback:** Implicit feedback ≫ **Explicit feedback - elaborated feedback**

From behavior awareness ➡ Dynamic time control during presentation



(a) Time target screen
(b) Time target setting
(c) Zone view (not late)
(d) Overview (not late)
(e) Planned-Overview (late)
(f) Adaptive-Overview (late)

Two types of timing support

**Less flexible**

- Planned rehearsal (e)

**More flexible**

- Adaptive guidance (f)

TalkZones (Saket et al., 2014)

Haptic feedback is enabled for redundant representation and reminder of lateness

Introduction | Automatic Assessment | User Interfaces | Conclusion

46

# Intelligent User Interface

## Live feedback: Implicit feedback ≫ Explicit feedback - termination

Interruption for improving specific skills



Presentation Trainer (Schneider et al., 2014)

**Action list** (e.g., voice modulation)

- Volume: loud, soft, normal
- Pause: long narration without pauses
- Filler sounds: "ehm", "hmm"

➡ **Corrective feedback**
(realtime visuals)

**Severe mistakes**

(vibration, pause sound, stops the program)

- Repetition of same mistakes
- Mistakes without being corrected for too long
- Predefined severe mistakes

➡ **Interruptive feedback**

Introduction | Automatic Assessment | User Interfaces | Conclusion

47

# Intelligent User Interface

THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY

## Live feedback

## Limitations

Implicit feedback

Explicit feedback

Simple verification

Elaborated feedback

Termination

- Most systems focus on providing timely suggestions about users' performance. They do not consider how to help them effectively and efficiently correct their mistakes

# Intelligent User Interface

**Posterior feedback:** **Summary feedback** ≫ Focus feedback

Summary of strengths & weaknesses



Aging and Engaging (Ali et al., 2018)



Automated Social Skills Trainer (Tanaka et al., 2015)
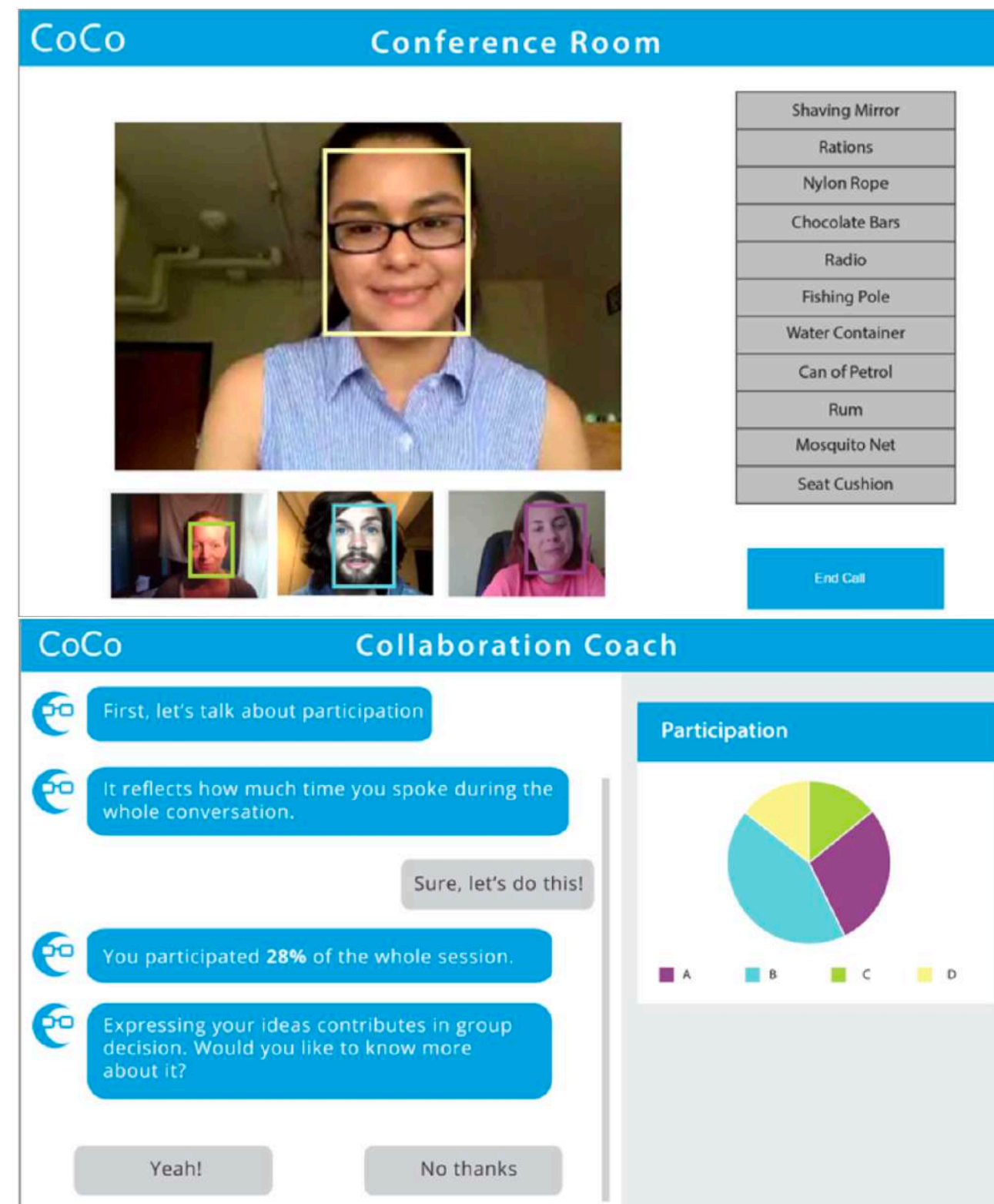
A: good points, B: bad points

C: overall score

D: pitch variation

E: comparison with model persons (pitch, power, energy, pause, WPM, 6 letters, fillers)

# Intelligent User Interface

## Posterior feedback: Summary feedback ≫ Focused feedback

Explaining affective behavioral performance on demand



CoCo (Samrose et al., 2017)



Performance graph

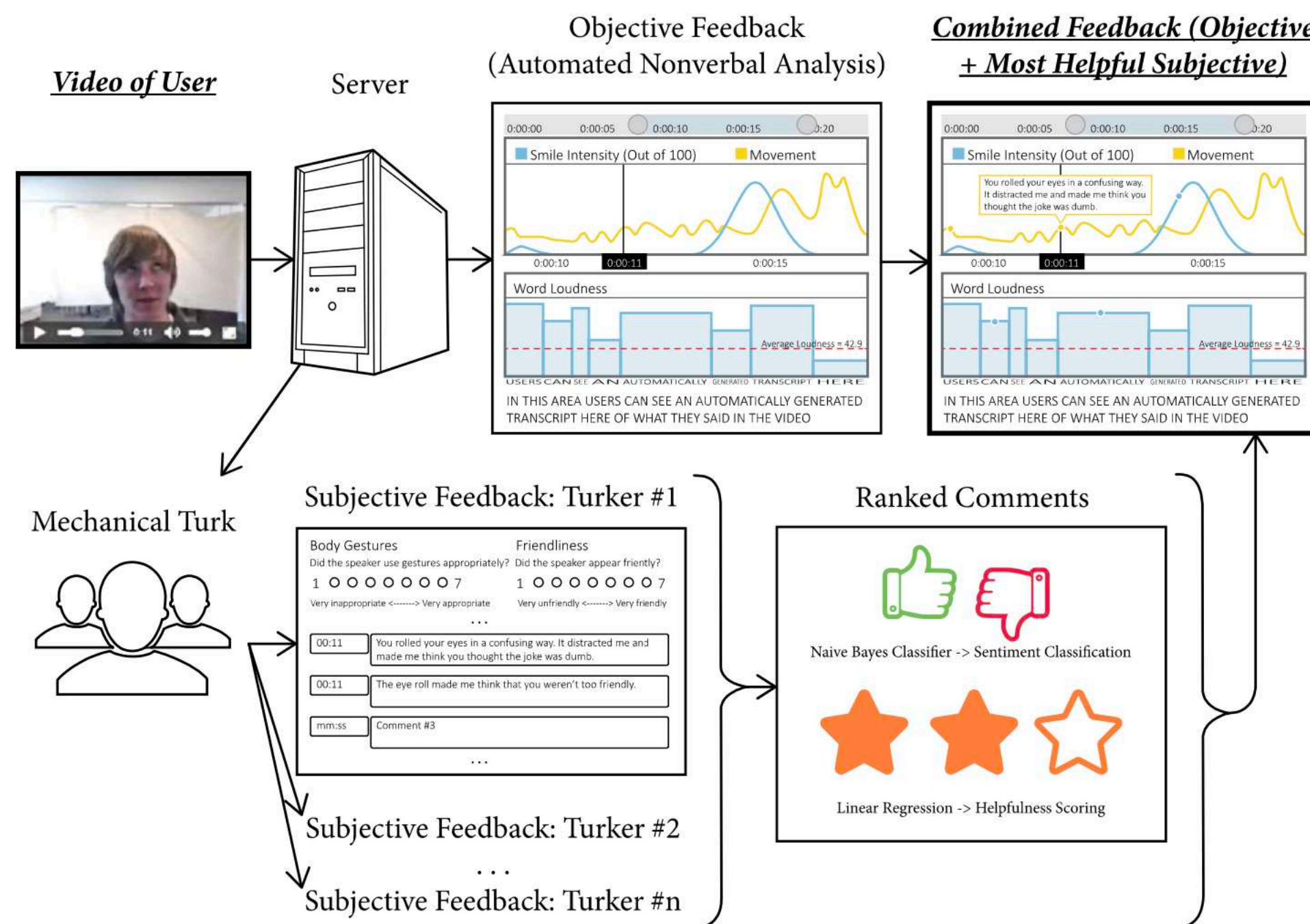Chat-based, post-conversation feedback

User study (with feedback)

- Balanced participation

- Skills awareness
  - How often they let others talk
  - Teammates's communicative skills

Introduction → Automatic Assessment → User Interfaces → Conclusion

50

# Intelligent User Interface

## Posterior feedback

Comprehensive feedback from the machine and crowdsourced workers



Overview of ROC Speak (Fung et al., 2015)

Motivation

**Machine**    Consistently & objectively sense subtle human behavior

**Human**    Interpreting contextual behavior

Gather human feedback

- score overall performance, voice modulation, friendliness, body gestures from 1 to 7

Automated ranking

- Label *helpfulness* & *sentiment*

- Train classifiers for prediction
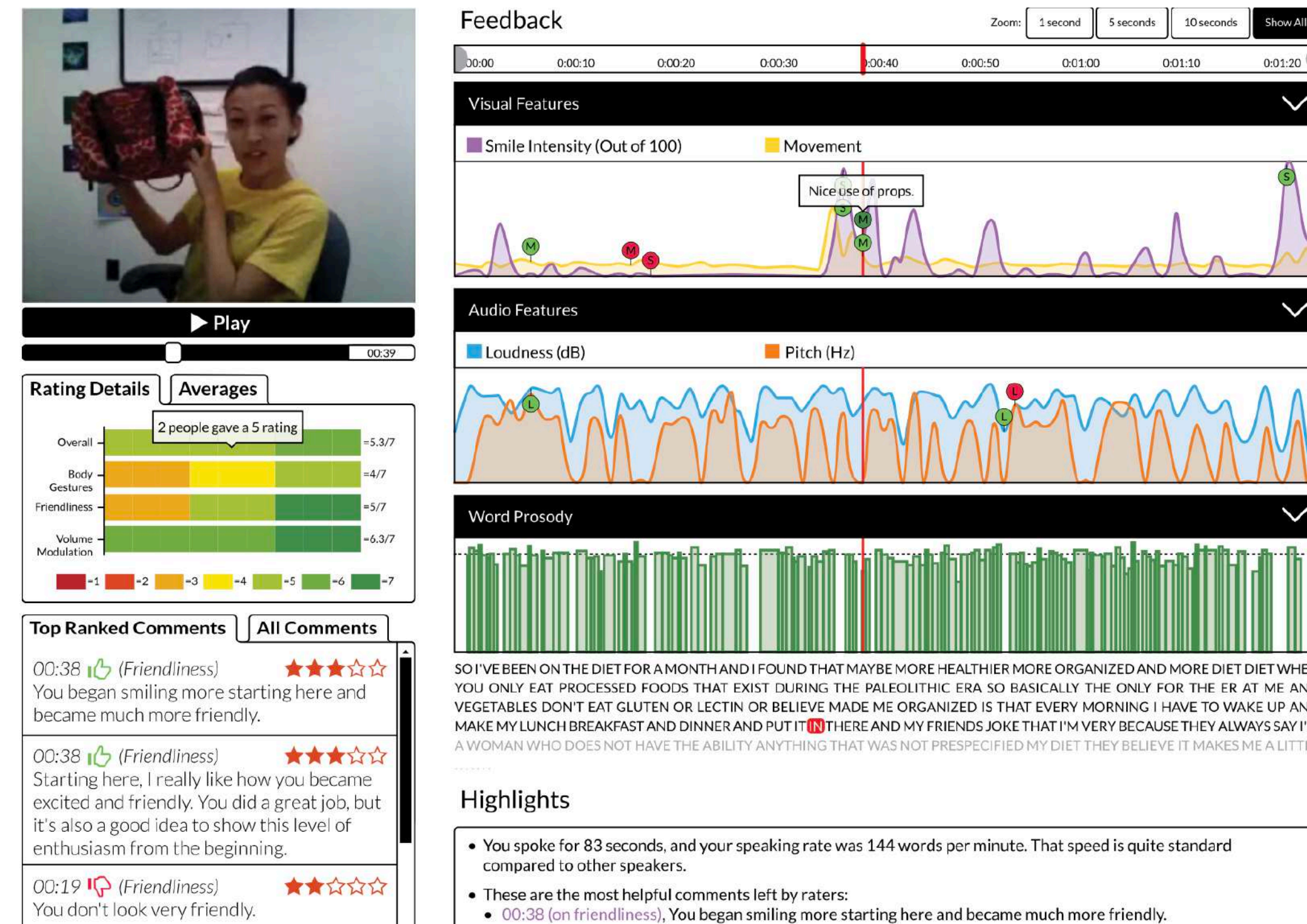
## Posterior feedback

Comprehensive feedback from the machine and crowdsourced workers



Overview of human feedback

Ranked comments

Quantitative visual graphs

**+**

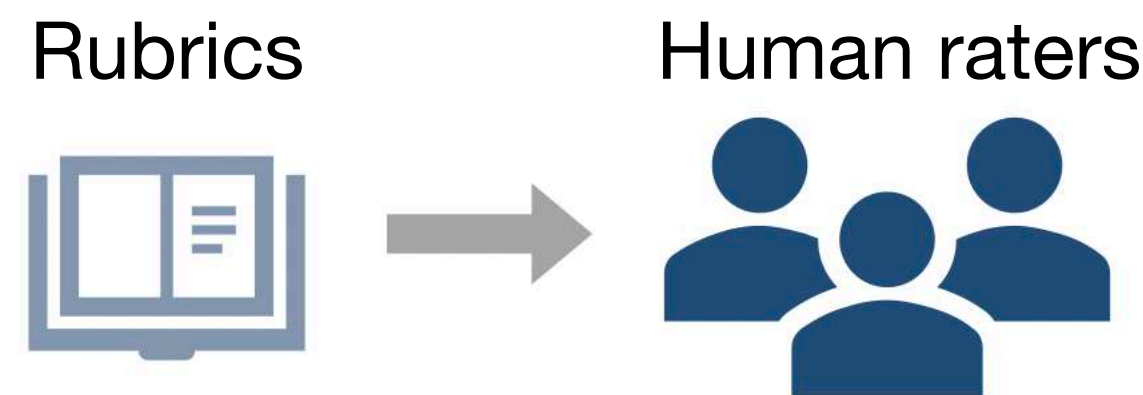Embedded human feedback
red: negative
green: positive

Most helpful comments for each category

ROC Speak (Fung et al., 2015)

Introduction → Automatic Assessment → User Interfaces → Conclusion

# Conclusion

## Summary & Future work **Machine Intelligence** ≫ Learning Interface

- Performance rubrics

- Computational features

- Machine learning models

Rubrics          Human raters



speech, speaking, talking, articulation, …

$f$

Computational behavior
descriptors

Machine learning models

Verbal
communication

- Developing more advanced and interpretable models for verbal communication assessment

- Investigating interactions among different modalities

Introduction | Automatic Assessment | User Interfaces | Conclusion

53

# Conclusion

**Summary & Future work** Machine Intelligence ≫ **Learning Interface**



**Prior feedback**
(Perception)

SpeechLens

EmoCo Video Digest

MMToC LectureScape

HyperSlides

Virtual Coach

NarrationCoach

VoiceCoach

Rhema

Logue

Cicero

VoiceAssist

MACH
TalkZones

PT

MC Wearable
System

ROC Speak

Aging & Engaging

Automated Social
Skills Trainer

CoCo

**Live feedback**
(action)

**Posterior feedback**
(reflection)

- Providing comprehensive feedback at all stages of learning cycle

- Engaging users in an iterative learning process

Introduction → Automatic Assessment → User Interfaces → Conclusion