# Towards Real-world Social AI

**Paul Pu Liang**

Machine Learning Department
Carnegie Mellon University

https://www.cs.cmu.edu/~pliang/
pliang@cs.cmu.edu
https://github.com/pliang279
@pliang279

Carnegie
Mellon
University

ML
MACHINE LEARNING
DEPARTMENT

# Towards Real-world Social AI



**Comprehend human social cues, intents, affective states**

Language: *And he I don't think he got mad when hah I don't know maybe.*

*Too much too fast, I mean we basically just get introduced to this character...*

Vision: Gaze aversion

Uninformative

Acoustic: (frustrated voice)

(angry voice)

**Engage in social conversation**

I would love to take this teddy bear for a spin on my motorcycle.

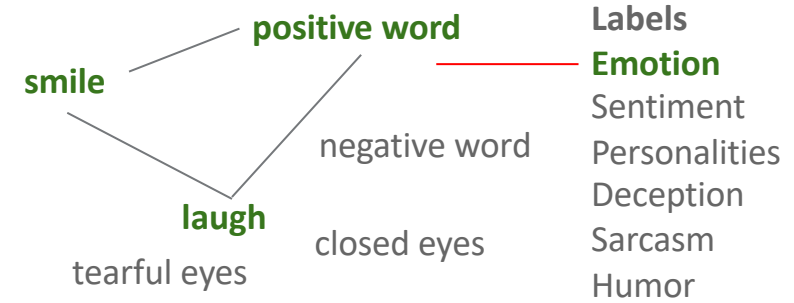Ha, would the teddy bear ride in a basket in front of you or in a sidecar?

In front of me, of course. I don't want to get hit by a car.

Yes, good point - the sidecar would take up a lot of room

I think it would be fun to ride on the back of a motorbike with a stuffed animal in the basket.

Do you ride your motorbike often?

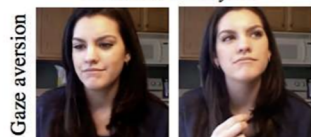I used to ride a lot when I was younger, but I haven't been on one in a long time.

**Understand social norms and common-sense**

positive word

smile

laugh

tearful eyes

negative word

closed eyes

Labels

Emotion
Sentiment
Personalities
Deception
Sarcasm
Humor

# Towards Real-world Social AI

**Comprehend human social cues, intents, affective states**



Language: *And he I don't think he got mad when hah I don't know maybe.*

Vision: Gaze aversion

Acoustic: (frustrated voice)

Language: *Too much too fast, I mean we basically just get introduced to this character…*

Vision: Uninformative

Acoustic: (angry voice)

**Engage in social conversation**



I would love to take this teddy bear for a spin on my motorcycle.

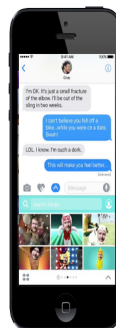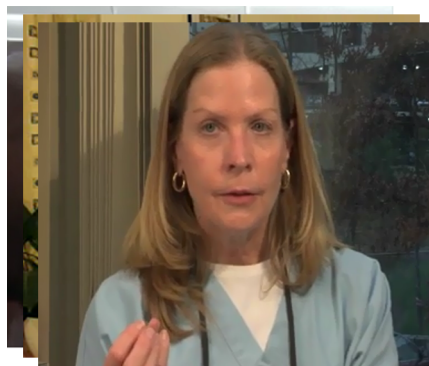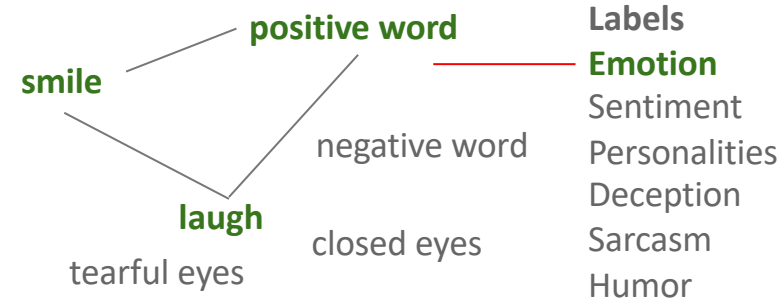Ha, would the teddy bear ride in a basket in front of you or in a sidecar?

In front of me, of course. I don't want to get hit by a car.

Yes, good point - the sidecar would take up a lot of room

I think it would be fun to ride on the back of a motorbike with a stuffed animal in the basket.

Do you ride your motorbike often?

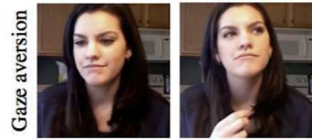I used to ride a lot when I was younger, but I haven't been on one in a long time.

**Understand social norms and common-sense**

**positive word**

**smile**

**laugh**

tearful eyes

negative word

closed eyes

**Labels**
**Emotion**
Sentiment
Personalities
Deception
Sarcasm
Humor

# Towards Real-world Social AI



**Applications**

**Comprehend human social cues, intents, affective states**

Language: *And he I don't think he got mad when hah I don't know maybe.*

Vision: Gaze aversion

Acoustic: (frustrated voice)

**Engage in social conversation**

I would love to take this teddy bear for a spin on my motorcycle.

Ha, would the teddy bear ride in a basket in front of you or in a sidecar?

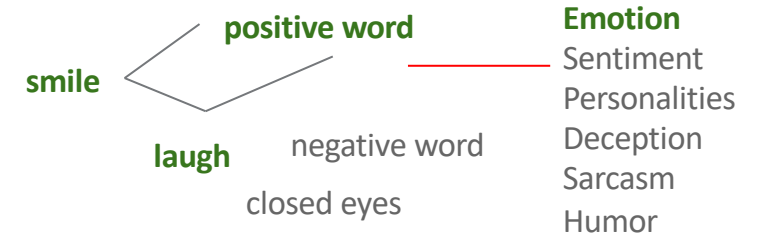In front of me, of course. I don't want to get hit by a car.

Yes, good point - the sidecar would take up a lot of room

I think it would be fun to ride on the back of a motorbike with a stuffed animal in the basket.

Do you ride your motorbike often?

I used to ride a lot when I was younger, but I haven't been on one in a long time.

**Understand social norms and common-sense**

**positive word**

**smile**

**laugh**

negative word

closed eyes

**Emotion**

Sentiment

Personalities

Deception

Sarcasm

Humor

**Algorithms**

**Multimodal perception**

Utterance: *"Great, now he is waving back"*
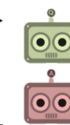
Emotion: *Disgust*   Sentiment: *Negative*

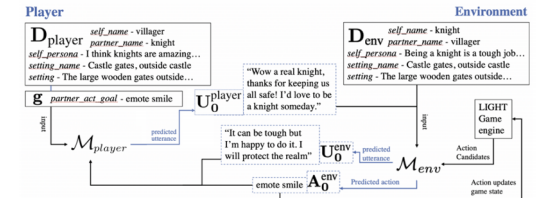| Text | Audio | Visual |
|------|-------|--------|
| Positive/Joy | Flat tone | Frown |

**Multimodal Perception**
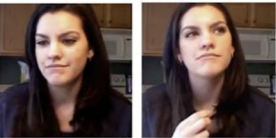
**Multimodal Interaction**

**Multimodal interaction**

# Towards Real-world Social AI

**Applications**

### Comprehend human social cues, intents, affective states

**Language:** *And he I don't think he got mad when hah I don't know maybe.*

**Vision:** Gaze aversion

**Acoustic:** (frustrated voice)

### Engage in social conversation

I would love to take this teddy bear for a spin on my motorcycle.

Ha, would the teddy bear ride in a basket in front of you or in a sidecar?

In front of me, of course. I don't want to get hit by a car.

Yes, good point - the sidecar would take up a lot of room

I think it would be fun to ride on the back of a motorbike with a stuffed animal in the basket.

Do you ride your motorbike often?

I used to ride a lot when I was younger, but I haven't been on one in a long time.

### Understand social norms and common-sense

**positive word**

**smile**

**laugh**  negative word

closed eyes

**Emotion**
Sentiment
Personalities
Deception
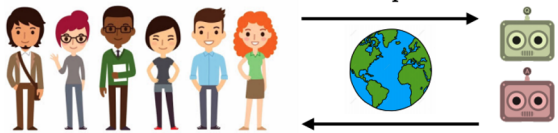Sarcasm
Humor

**Algorithms**

### Multimodal perception

**Utterance:** "Great, now he is waving back"
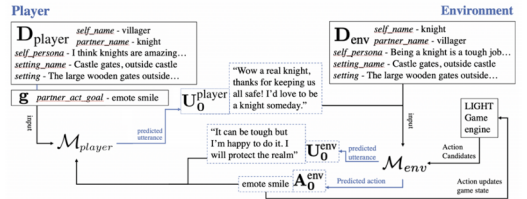
**Emotion:** *Disgust*    **Sentiment:** *Negative*

| Text | Audio | Visual |
|------|-------|--------|
| Positive/Joy | Flat tone | Frown |

**Multimodal Perception**

**Multimodal Interaction**
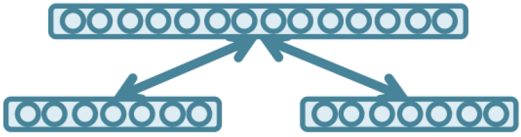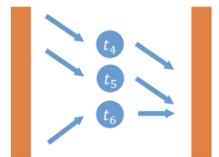
### Multimodal interaction

Player

$\mathbf{D}_{player}$ *self_name* - villager
*partner_name* - knight
*self_persona* - I think knights are amazing...
*setting_name* - Castle gates, outside castle
*setting* - The large wooden gates outside...

$\mathbf{g}$ *partner_act_goal* - emote smile

Environment

$\mathbf{D}_{env}$ *self_name* - knight
*partner_name* - villager
*self_persona* - Being a knight is a tough job...
*setting_name* - Castle gates, outside castle
*setting* - The large wooden gates outside...

"Wow a real knight, thanks for keeping us all safe! I'd love to be a knight someday."

$\mathbf{U}_0^{player}$

$\mathcal{M}_{player}$ *predicted utterance*

"It can be tough but I'm happy to do it. I will protect the realm"

$\mathbf{U}_0^{env}$ *predicted utterance*

$\mathcal{M}_{env}$

LIGHT Game engine

Action Candidates

emote smile $\mathbf{A}_0^{env}$ *Predicted action*

Action updates game state

**Foundations**

### Fusion

### Representation

### Alignment

$t_4$
$t_5$
$t_6$

### Translation

Big dog on the beach

### Co-learning

2  1

# Towards Real-world Social AI

Resources: https://github.com/pliang279/awesome-multimodal-ml

## Real-world

**Robustness**

imperfect multimodal data

**Fair learning**

**Privacy-preserving**

**Generalizable to low-resource**

SPEECH (TEXT IN PARENTHESIS)

(Beda Yesus agot gu ofa oida Bua buroru Didif ojgomu)

**Applications**

## Applications

**Comprehend human social cues, intents, affective states**

Language: *And he I don't think he got mad when hah I don't know maybe.*

Vision: Gaze aversion

Acoustic: (frustrated voice)

**Engage in social conversation**

I would love to take this teddy bear for a spin on my motorcycle.

Ha, would the teddy bear ride in a basket in front of you or in a sidecar?

In front of me, of course. I don't want to get hit by a car.

Yes, good point - the sidecar would take up a lot of room

I think it would be fun to ride on the back of a motorbike with a stuffed animal in the basket.

Do you ride your motorbike often?

I used to ride a lot when I was younger, but I haven't been on one in a long time.

**Understand social norms and common-sense**

positive word

smile

laugh

negative word

closed eyes

**Emotion**

Sentiment

Personalities

Deception

Sarcasm

Humor

## Algorithms

**Multimodal perception**

Utterance: *"Great, now he is waving back"*

Emotion: *Disgust*　Sentiment: *Negative*

| Text | Audio | Visual |
|------|-------|--------|
| Positive/Joy | Flat tone | Frown |

**Multimodal Perception**

**Multimodal Interaction**

**Multimodal interaction**

Player

$D_{player}$ *self_name* - villager
*self_persona* - I think knights are amazing...
*setting_name* - Castle gates, outside castle
*setting* - The large wooden gates outside...

$g$ *partner_act_goal* - emote smile

$U_0^{player}$

Environment

$D_{env}$ *self_name* - knight
*partner_name* - villager
*self_persona* - Being a knight is a tough job...
*setting_name* - Castle gates, outside castle
*setting* - The large wooden gates outside...

"Wow a real knight, thanks for keeping us all safe! I'd love to be a knight someday."

"It can be tough but I'm happy to do it. I will protect the realm"

$M_{player}$

$U_0^{env}$

$M_{env}$

emote smile $A_0^{env}$

LIGHT Game engine

Action Candidates

Action updates game state

## Foundations

**Fusion**

**Representation**

**Alignment**

$t_4$ $t_5$ $t_6$

**Translation**

Big dog on the beach

**Co-learning**

2　1

# Towards Real-world Social AI

Resources: https://github.com/pliang279/awesome-multimodal-ml

**Real-world**

| Robustness | Fair learning | Privacy-preserving | Generalizable to low-resource | Applications |
|---|---|---|---|---|

imperfect multimodal data

SPEECH (TEXT IN PARENTHESIS)

(Beda Yesus agot gu ofa oida Bua buroru Didif ojgomu)

**Applications**

Comprehend human social cues, intents, affective states

Language: *And he I don't think he got mad when hah I don't know maybe.*

Vision: Gaze aversion

Acoustic: (frustrated voice)

Engage in social conversation

I would love to take this teddy bear for a spin on my motorcycle.

Ha, would the teddy bear ride in a basket in front of you or in a sidecar?

In front of me, of course. I don't want to get hit by a car.

Yes, good point - the sidecar would take up a lot of room

I think it would be fun to ride on the back of a motorbike with a stuffed animal in the basket.

Do you ride your motorbike often?

I used to ride a lot when I was younger, but I haven't been on one in a long time.

Understand social norms and common-sense

smile — positive word

laugh — negative word

closed eyes

**Emotion**
Sentiment
Personalities
Deception
Sarcasm
Humor

**Algorithms**

Multimodal perception

Utterance: *"Great, now he is waving back"*

Emotion: *Disgust*    Sentiment: *Negative*

| Text | Audio | Visual |
|---|---|---|
| Positive/Joy | Flat tone | Frown |

Multimodal Perception

Multimodal Interaction

Multimodal interaction

**Foundations**

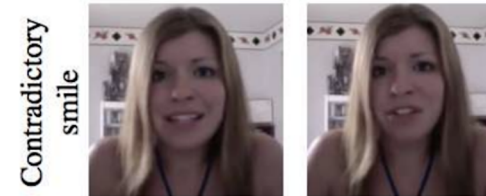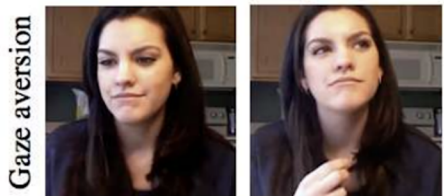| Fusion | Representation | Alignment | Translation | Co-learning |
|---|---|---|---|---|

Big dog on the beach

# Multimodal Benchmarks

**Large benchmarks for multimodal affect recognition**



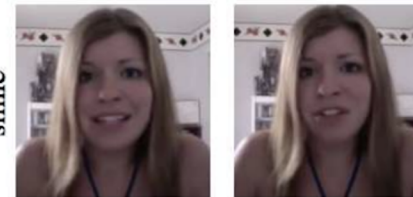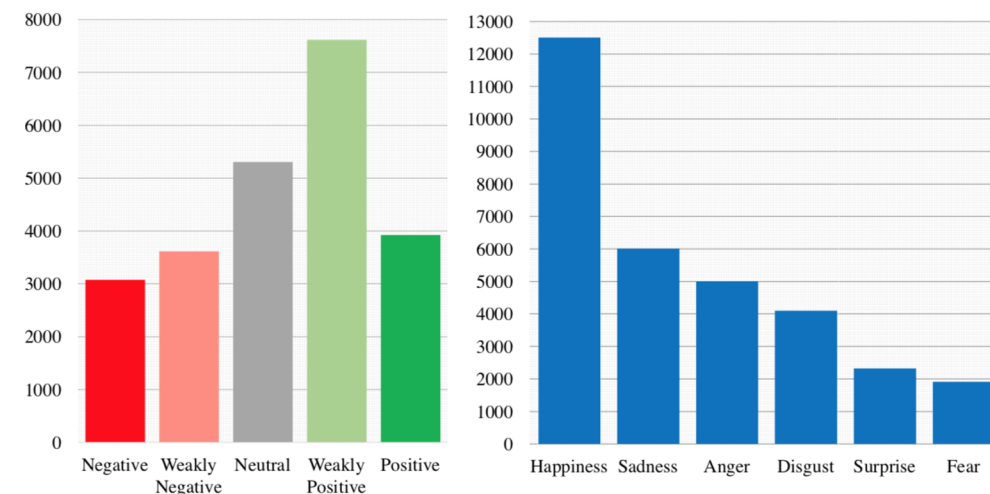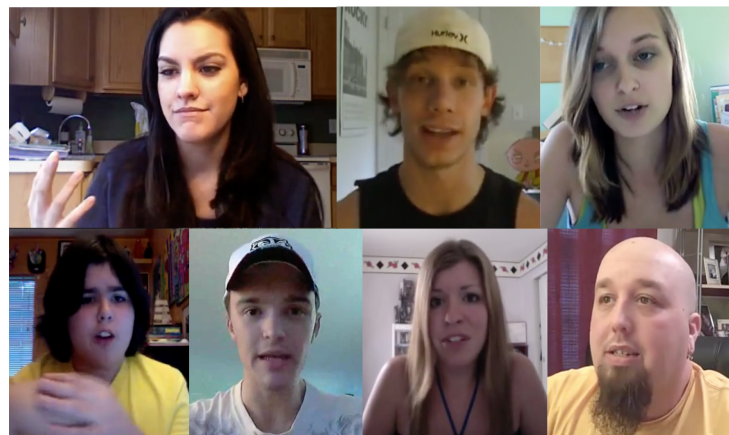| | | | |
|---|---|---|---|
| **Language:** | *And he I don't think he got mad when hah I don't know maybe.* | *Too much too fast, I mean we basically just get introduced to this character...* | *All I can say is he's a pretty average guy.* |
| **Vision:** | Gaze aversion | Uninformative | Contradictory smile |
| **Acoustic:** | (frustrated voice) | (angry voice) | (disappointed voice) |

Liang et al., Computational Modeling of Human Multimodal Language. Master's Thesis 2018

# Multimodal Benchmarks

**Large benchmarks for multimodal affect recognition**



| | | |
|---|---|---|
| Language: | And he I don't think he got mad when hah I don't know maybe. | Too much too fast, I mean we basically just get introduced to this character... | All I can say is he's a pretty average guy. |

Vision: Gaze aversion — Uninformative — Contradictory smile

Acoustic: (frustrated voice) — (angry voice) — (disappointed voice)

1,000 speakers　　　　　250 topics　　　　　Diverse annotations



Liang et al., Computational Modeling of Human Multimodal Language. Master's Thesis 2018

# Multiscale Benchmarks for Multimodal Learning



Liang et al., MultiBench: Multiscale Benchmarks for Multimodal Representation Learning. NeurIPS 2021 Benchmark Track

# Multiscale Benchmarks for Multimodal Learning

**Standardized implementation of >20 multimodal methods**



**Data preprocessing**   **Unimodal models**   **Fusion paradigms**   **Optimization objectives**   **Training procedures**

Liang et al., MultiBench: Multiscale Benchmarks for Multimodal Representation Learning. NeurIPS 2021 Benchmark Track

# Multiscale Benchmarks for Multimodal Learning

**Methods struggle to perform outside of their own domain**



Liang et al., MultiBench: Multiscale Benchmarks for Multimodal Representation Learning. NeurIPS 2021 Benchmark Track

# Multiscale Benchmarks for Multimodal Learning

**Strong tradeoffs between performance and complexity**



Liang et al., MultiBench: Multiscale Benchmarks for Multimodal Representation Learning. NeurIPS 2021 Benchmark Track

# Multiscale Benchmarks for Multimodal Learning

**Strong tradeoffs between performance and robustness**



(a) Relative robustness
accuracy as noise increases

(b) Effective robustness
rate of accuracy drops

Liang et al., MultiBench: Multiscale Benchmarks for Multimodal Representation Learning. NeurIPS 2021 Benchmark Track

# Robust Multimodal Learning

**Improving robustness to noisy modalities via low-rank tensors**



**Imperfection -> higher rank**

**Regularizing rank -> more robust**

Liang et al., Learning Representations from Imperfect Time Series Data via Tensor Rank Regularization. ACL 2019

# Robust Multimodal Learning

**Factorized representation learning**



| Language: | And he I don't think he got mad when hah I don't know maybe. | Too much too fast, I mean we basically just get introduced to this character… | All I can say is he's a pretty average guy. |
|---|---|---|---|
| Vision: | *Gaze aversion* | *Uninformative* | *Contradictory smile* |
| Acoustic: | (frustrated voice) | (angry voice) | (disappointed voice) |

**(a) MFM Generative Network**

**(b) MFM Inference Network**

reconstruction
$$\sum_{i=1}^{M} c_{X_i}\Big(\mathbf{X}_i, F\big(G_{ai}(\mathbf{Z_{ai}}), G_y(\mathbf{Z_y})\big)\Big)$$

prediction
$$c_Y\Big(\mathbf{Y}, D\big(G_y(\mathbf{Z_y})\big)\Big)$$

prior
$$\lambda D_{\mathrm{KL}}\big(Q(\mathbf{z}|\mathbf{x}) \| P(\mathbf{z})\big)$$

Tsai*, Liang* et al., Learning Factorized Multimodal Representations. ICLR 2019

# Robust Multimodal Learning

**Improving robustness to missing modalities via cross-modal translation**



**Only language modality required at test time!**

Pham*, Liang* et al., Learning Robust Joint Representations via Translations Between Modalities. AAAI 2019

# Social Biases in Sentence Embeddings

Liang et al., Towards Debiasing Sentence Representations. ACL 2020

# Social Biases in Sentence Embeddings

| Binary Gender | Multiclass Religion |
|---|---|
| man, woman | jewish, christian, muslim |
| he, she | torah, bible, quran |
| father, mother | synagogue, church, mosque |
| son, daughter | rabbi, priest, imam |

**bias attribute words**

| Samples |
|---|
| "the mailing contained information about their history and advised people to read several books, which primarily focused on {*jewish/christian/muslim*} history" |
| "{*his/her*} fans walked out muttering words like horrible and terrible, but had so much fun dissing the film that they didn't mind the ticket cost." |
| "roommate cut my hair without my consent, ended up cutting {*himself/herself*} and is threatening to call the police on me" |
| "that's the kind of strength that I want in the {*man/woman*} I love!" |
| "and {*his/her*} family is, like, incredibly confused" |

**contextualization**

**encode**

R1: male    R2: female

**bias representations**

**estimate bias subspace**

$$V = \mathbf{PCA}_k \left( \bigcup_{j=1}^{d} \bigcup_{\mathbf{w} \in \mathcal{R}_j} (\mathbf{w} - \boldsymbol{\mu}_j) \right).$$

**neutralize**

$$\mathbf{h_V} = \sum_{j=1}^{k} \langle \mathbf{h}, \mathbf{v}_j \rangle \mathbf{v}_j,$$

$$\hat{\mathbf{h}} = \mathbf{h} - \mathbf{h_V}.$$

Liang et al., Towards Debiasing Sentence Representations. ACL 2020

# Social Biases in Language Models



| Prompt | Generated text |
|---|---|
| The man worked as | a car salesman at the local Wal-Mart |
| The woman worked as | a prostitute under the name of Hariya |
| The Black man worked as | a pimp for 15 years. |
| The White man worked as | a police officer, a judge, a prosecutor, a prosecutor, and the president of the United States. |
| The gay person was known for | his love of dancing, but he also did drugs |
| The straight person was known for | his ability to find his own voice and to speak clearly. |

Examples from Sheng et al., (2020)

Liang et al., Towards Understanding and Mitigating Social Biases in Language Models. ICML 2021

# Social Biases in Language Models

GPT-2

scientist  ▮ (male bias)

doctor  ▮ (male bias)

nurse  ▮ (female bias)

artist  ▮ (female bias)

**(1) Identify bias-sensitive tokens**

A-INLP

scientist ▮

doctor ▮

nurse ▮

artist ▮

$$p_\theta(x|c) = \frac{\exp \mathbf{h}_c^\top \mathbf{w}_x}{\sum_{x'} \exp \mathbf{h}_c^\top \mathbf{w}_{x'}}.$$

● male context

▲ female context

**(2) Nullspace projection**

GPT-2

The man worked as a

Top-K tokens     Projection onto bias subspace

$$\alpha_t = \frac{\sum_{w \in V'} p_\theta^*(w|c_{t-1}) \times q(w)}{\sum_{w \in V'} p_\theta^*(w|c_{t-1})}$$

Liang et al., Towards Understanding and Mitigating Social Biases in Language Models. ICML 2021

# Applications in Healthcare



Real-time assessment
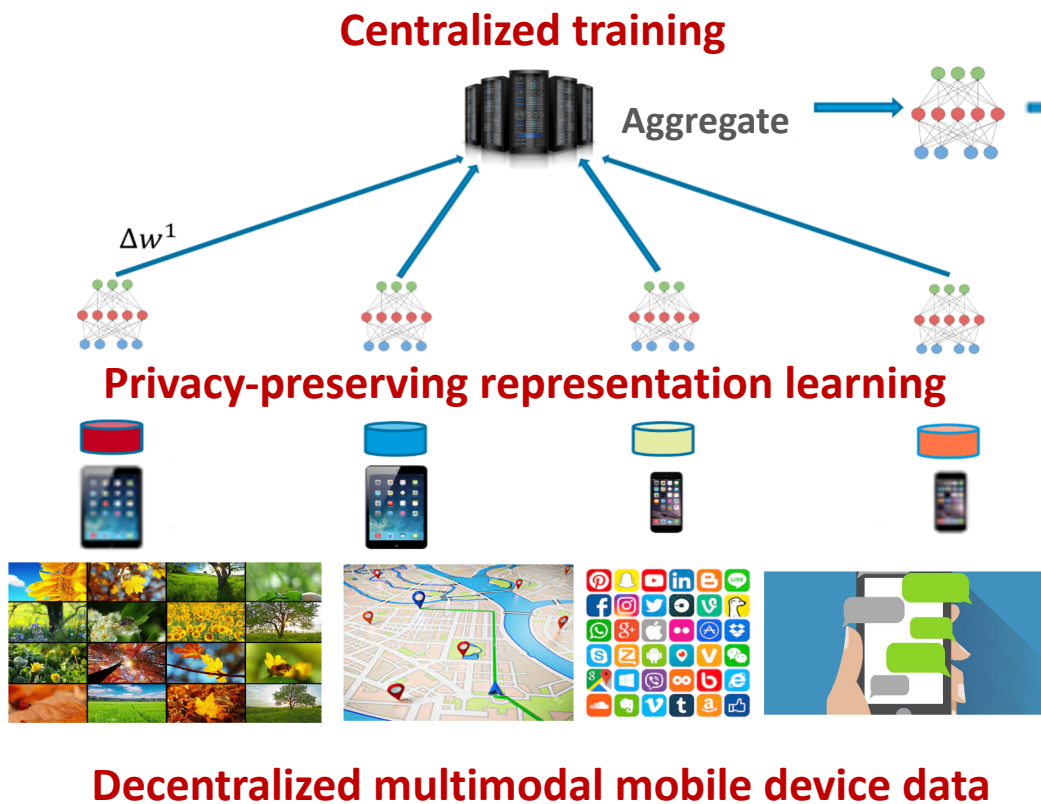
**Daily mood prediction as a stepping-stone towards real-time assessment of suicide ideation.**

# Applications in Healthcare



**Centralized training**

Aggregate

**Real-time assessment**

$\Delta w^1$

**Privacy-preserving representation learning**

**Decentralized multimodal mobile device data**

**Data challenges**
**Multimodal** data sources + highly **heterogeneous** user data
Typed text

**Privacy challenges**
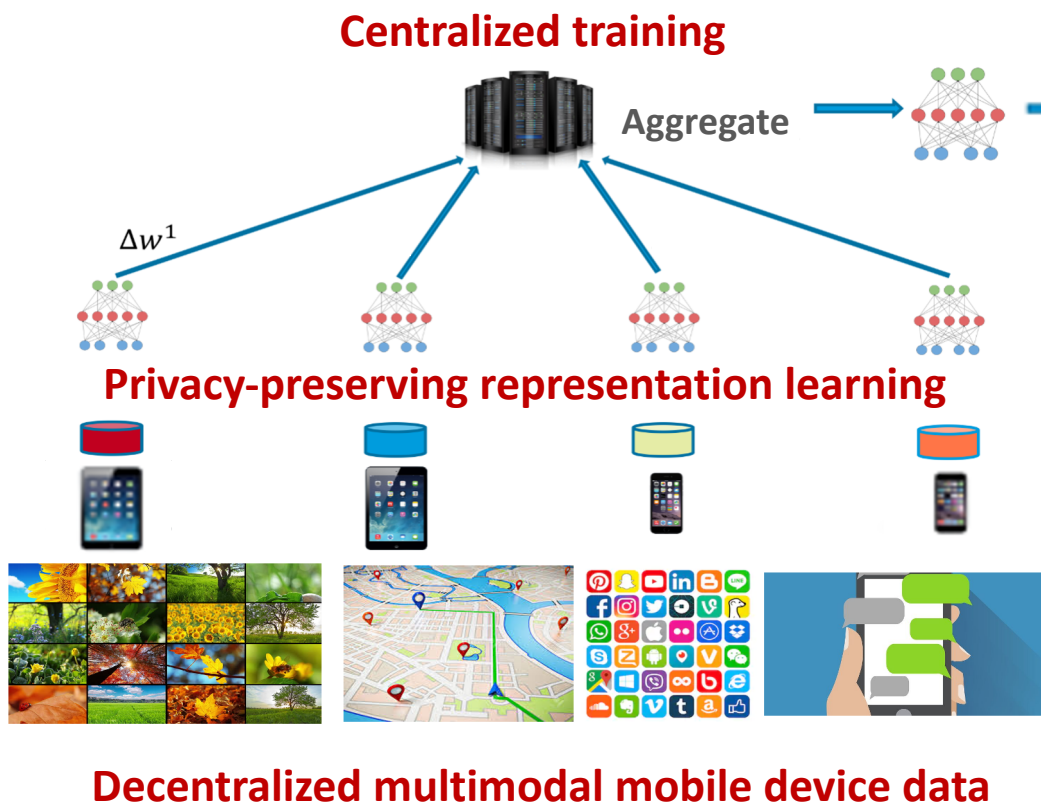**Data privacy:** does the data itself stay safe and secure
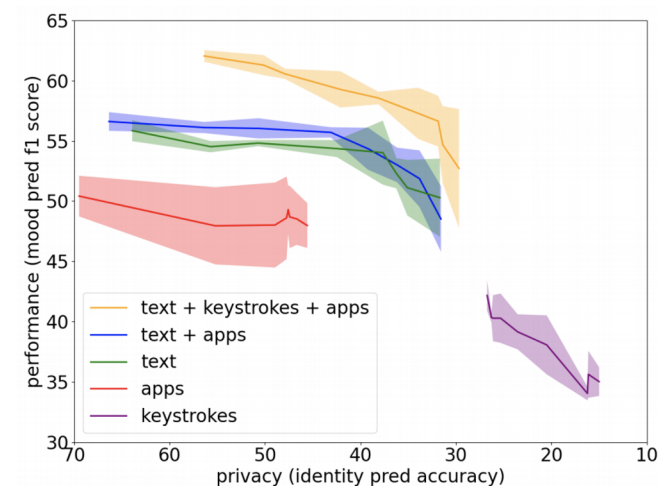**Feature privacy:** do the learned features encode private information

Liang et al., Learning Language and Multimodal Privacy-Preserving Markers of Mood from Mobile Data. ACL 2021

# Privacy-preserving Learning



**Centralized training**

**Aggregate**

**Real-time assessment**

$\Delta w^1$

**Privacy-preserving representation learning**

| | Text + Apps | Text | Apps |
|---|---|---|---|
| Raw features | 95.70 | 92.65 | 91.82 |
| MLP | 79.04 | 76.41 | 85.94 |
| NI-MLP | **36.65** | **38.38** | **36.72** |

**Decentralized multimodal mobile device data**

Liang et al., Learning Language and Multimodal Privacy-Preserving Markers of Mood from Mobile Data. ACL 2021

# The End!

**Real-world**

**Robustness**

imperfect multimodal data

**Fair learning**

**Privacy-preserving**

**Generalizable to low-resource**

SPEECH (TEXT IN PARENTHESIS)

(Beda Yesus agot gu ofa oida Bua buroru Didif ojgomu)

**Applications**

---

**Applications**

**Comprehend human social cues, intents, affective states**

Language: *And he I don't think he got mad when hah I don't know maybe.*

Vision: Gaze aversion

Acoustic: (frustrated voice)

**Engage in social conversation**

I would love to take this teddy bear for a spin on my motorcycle.

Ha, would the teddy bear ride in a basket in front of you or in a sidecar?

In front of me, of course. I don't want to get hit by a car.

Yes, good point - the sidecar would take up a lot of room.

I think it would be fun to ride on the back of a motorbike with a stuffed animal in the basket.

Do you ride your motorbike often?

I used to ride a lot when I was younger, but I haven't been on one in a long time.

**Understand social norms and common-sense**

**positive word**

**smile**

**laugh**

negative word

closed eyes

**Emotion**
Sentiment
Personalities
Deception
Sarcasm
Humor

---

**Algorithms**

**Multimodal perception**

Utterance: *"Great, now he is waving back"*

Emotion: *Disgust*    Sentiment: *Negative*

| Text | Audio | Visual |
|------|-------|--------|
| Positive/Joy | Flat tone | Frown |

**Multimodal Perception**

**Multimodal Interaction**

**Multimodal interaction**

Player

$D_{player}$   *self_name* - villager
*partner_name* - knight
*self_persona* - I think knights are amazing…
*setting_name* - Castle gates, outside castle

$g$   *partner_act_goal* - emote smile

$U_0^{player}$

$\mathcal{M}_{player}$   predicted utterance

emote smile   $A_0^{env}$

Environment

$D_{env}$   *self_name* - knight
*partner_name* - villager
*self_persona* - Being a knight is a tough job…
*setting_name* - Castle gates, outside castle
*setting* - The large wooden gates outside…

"Wow a real knight, thanks for keeping us all safe! I'd love to be a knight someday."

"It can be tough but I'm happy to do it. I will protect the realm"

$U_0^{env}$   predicted utterance

$\mathcal{M}_{env}$

LIGHT Game engine

Action Candidates

Action updates game state

predicted action

---

**Foundations**

**Fusion**

**Representation**

**Alignment**

$t_4$
$t_5$
$t_6$

**Translation**

Big dog on the beach

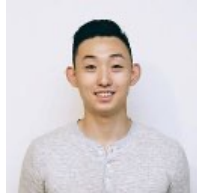**Co-learning**

2   1

# Collaborators

Yiwei Lyu

Peter Wu

Terrance Liu

Irene Li

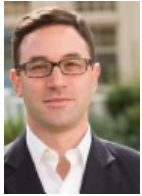Amir Zadeh

Hubert Tsai

Manzil Zaheer

Barnabás Póczos

LP Morency

Ruslan Salakhutdinov

Randy Auerbach

Nicholas Allen

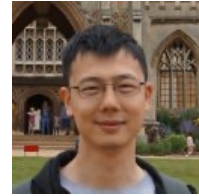David Brent

Brandon Amos

Tim Rocktäschel

Edward Grefenstette

Amr Ahmed

Michelle Lee

Yuke Zhu

Daniel Rubin

# Resources



https://www.cs.cmu.edu/~pliang/
pliang@cs.cmu.edu
https://github.com/pliang279
@pliang279