# Real-Time Estimation of Overt Attention from Dynamic Features of the Face using Deep Learning
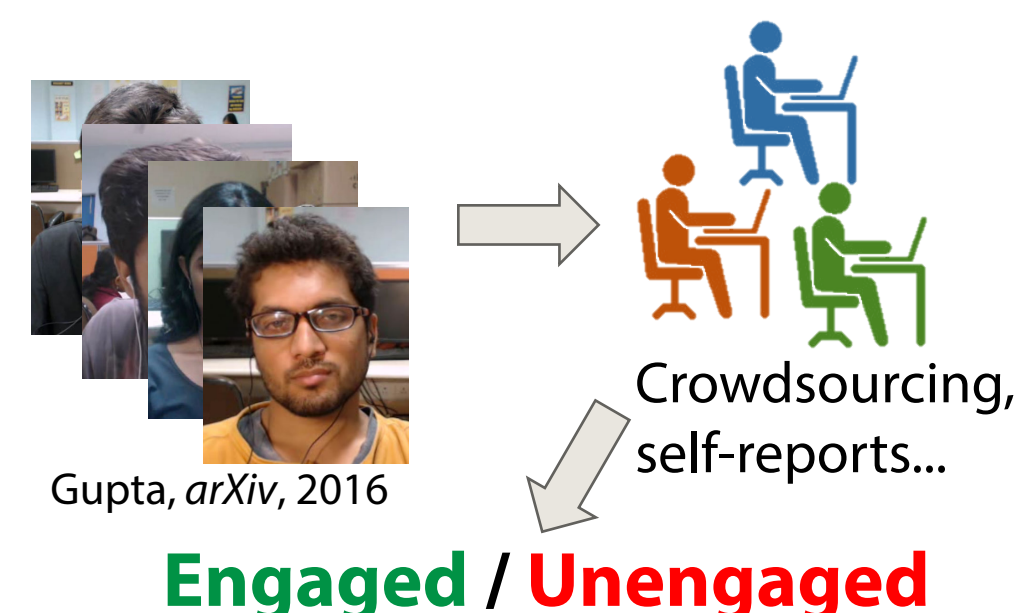
Aimar Silvan, Lucas C. Parra, Jens Madsen

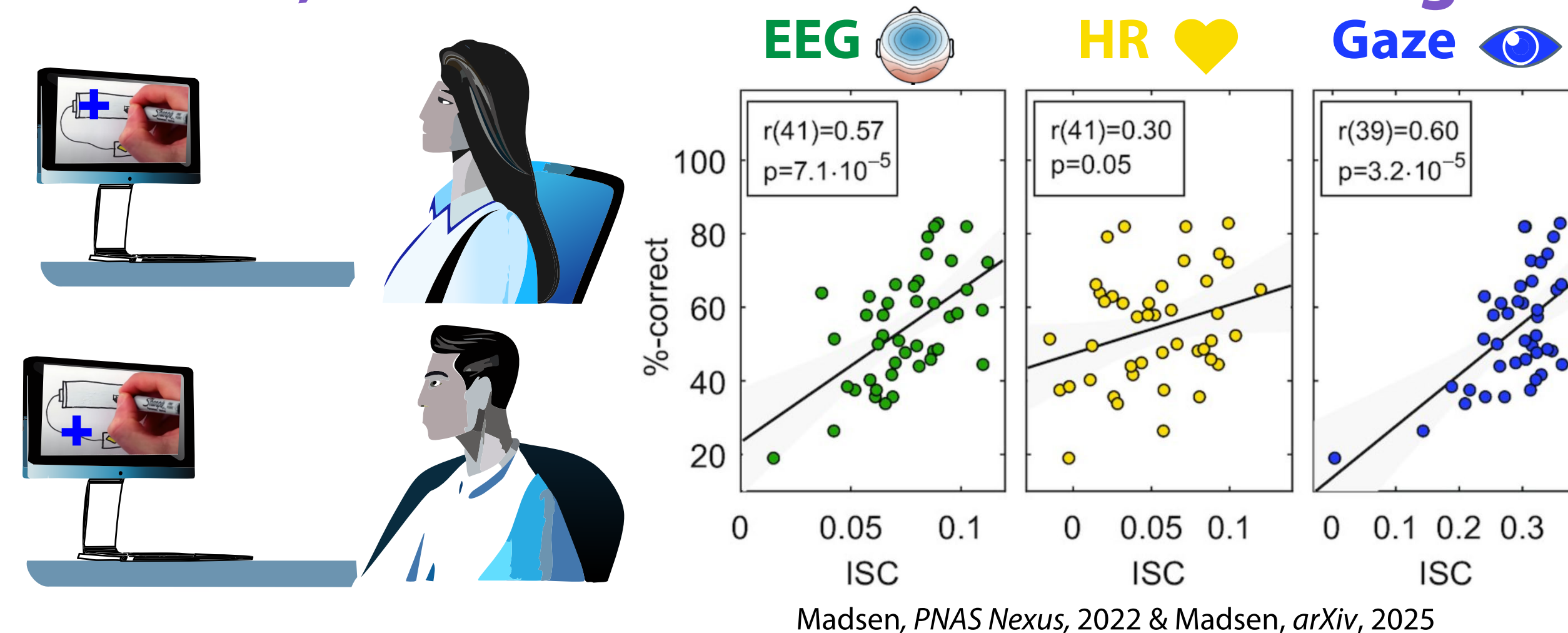- Department of Biomedical Engineering, City College of New York -

## There is a need for measuring attention objectively

- **Challenge**: Measuring task engagement is crucial across neuroscience, education, and psychology.

- **Current AI Limitations:** Existing methods for attention estimation often rely on self-reports, subjective ratings, and distributing face videos, raising privacy concerns.
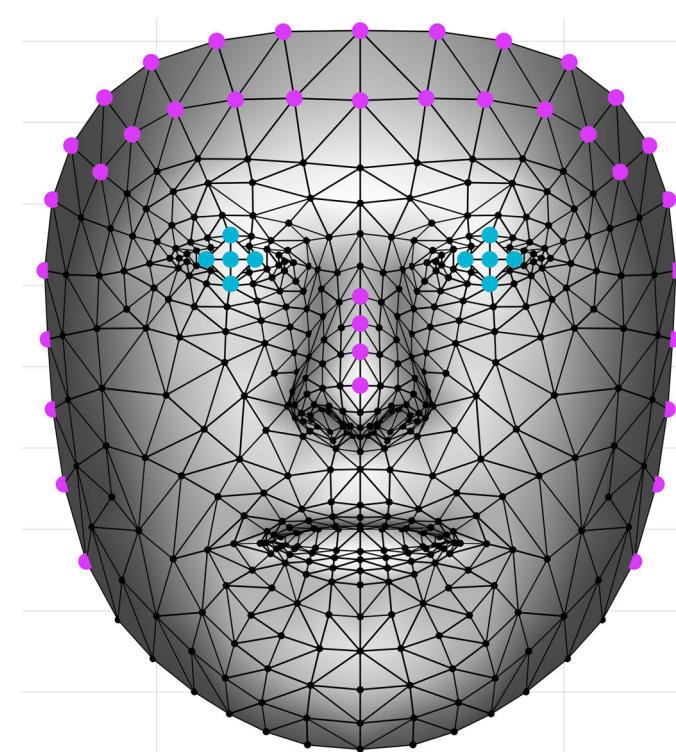

Gupta, *arXiv*, 2016

Crowdsourcing, self-reports...

**Engaged** / **Unengaged**

## Brains, hearts, and eyes synchronize when attending to videos, but are hard to measure on a large scale

**EEG**    **HR**    **Gaze**



r(41)=0.57  p=7.1·10⁻⁵
r(41)=0.30  p=0.05
r(39)=0.60  p=3.2·10⁻⁵

ISC    ISC    ISC

Madsen, *PNAS Nexus*, 2022 & Madsen, *arXiv*, 2025

- **Attentively watching videos synchronizes EEG, HR, gaze, and pupil**, and this Inter-Subject Correlation (**ISC**) strongly **correlates with their performance** when tested on the contents.
- But measuring these synchronized signals typically requires **complex sensors** that **are inaccessible or impractical for large-scale experiments.**

## Modern AI methods allow real-time, on-device and privacy-preserving face tracking



- Google MediaPipe enables **real-time** facial landmark and movement extraction from **standard webcams.**
- The model can **run on the user's device**, in-browser, on our experiment platform **Elicit.**
- **Privacy is preserved** by only trasmitting face landmarks and blendshapes, no webcam video.
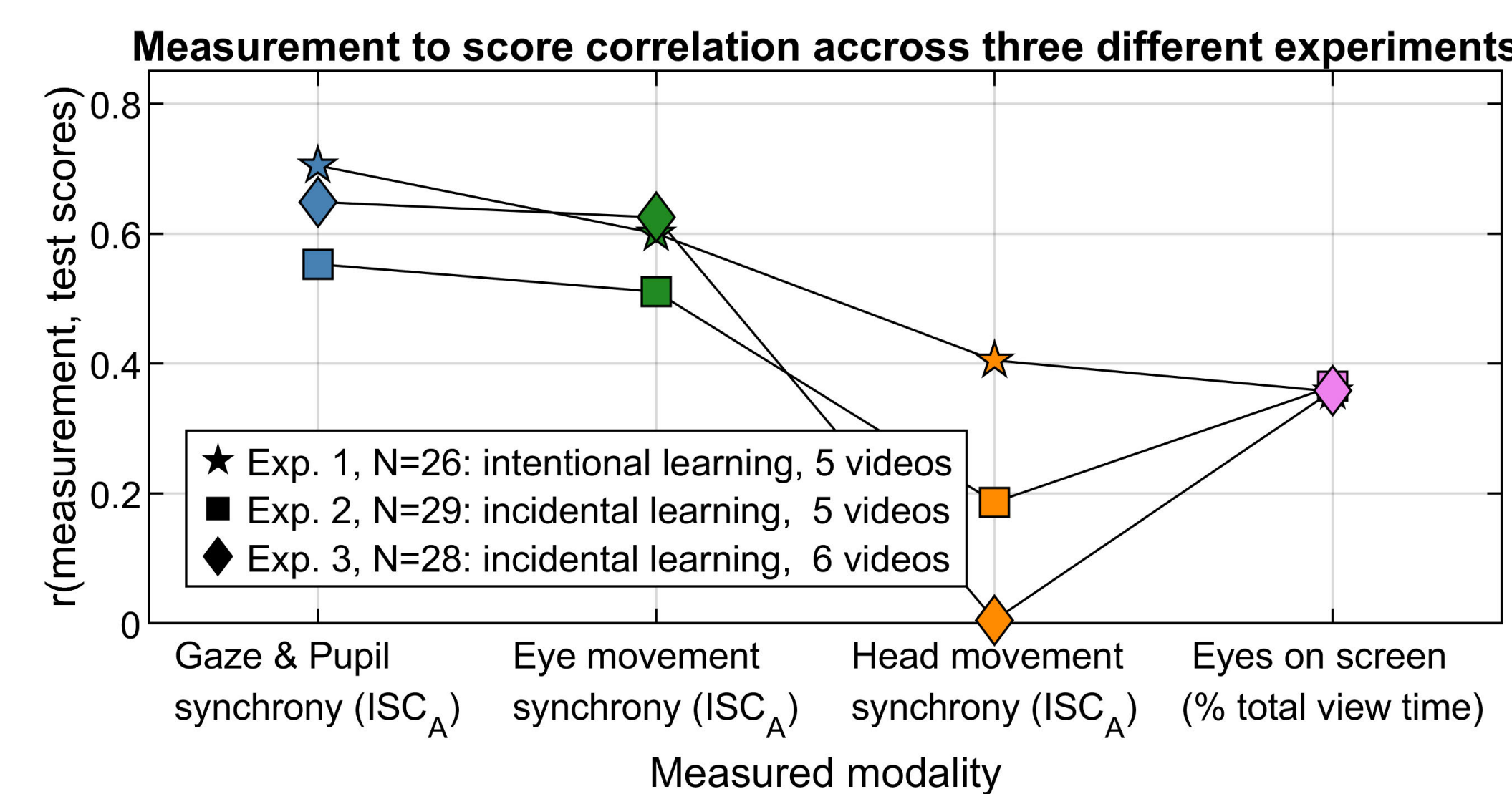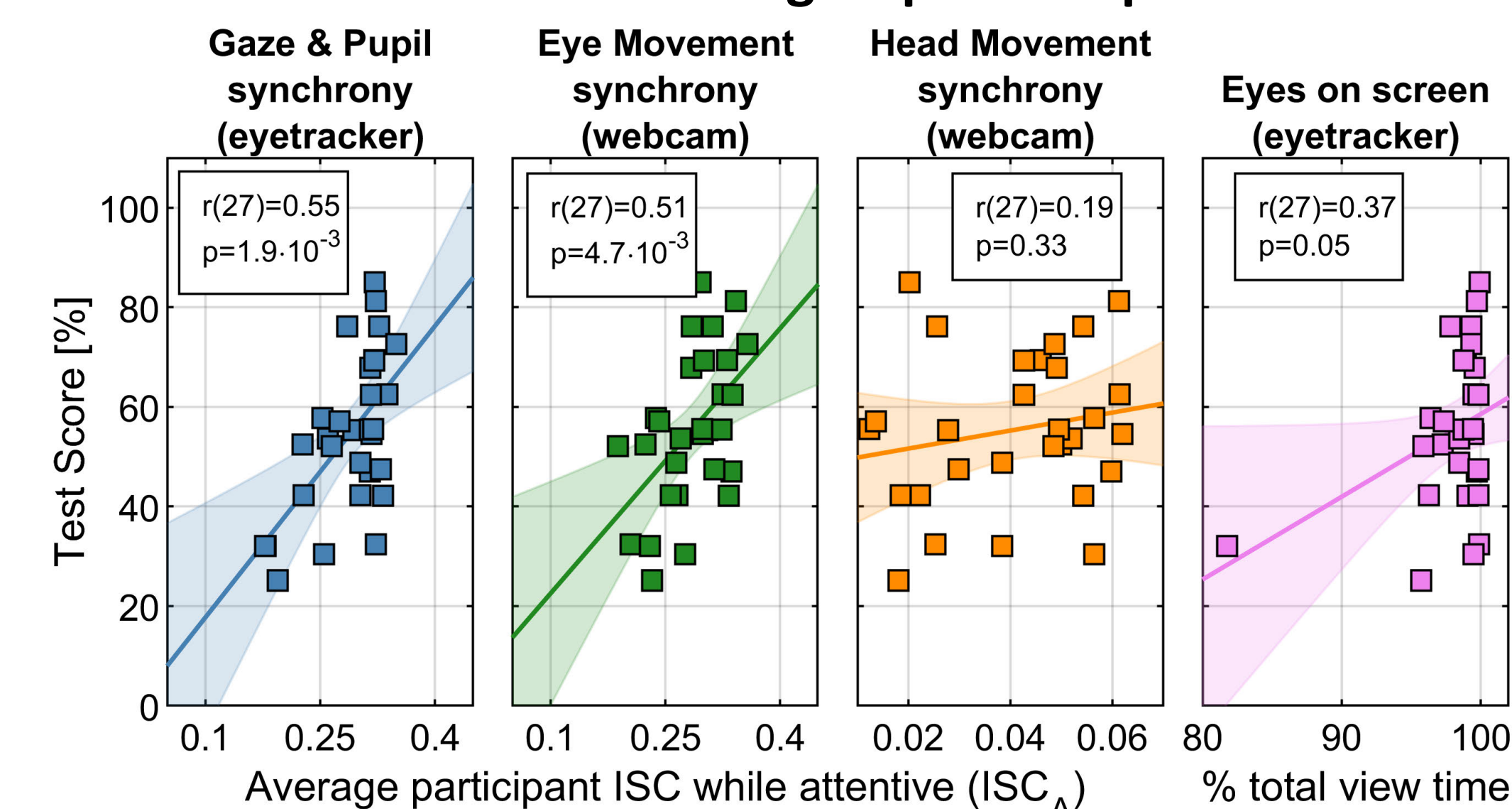
MediaPipe  Elicit  https://elicit-experiment.com/

### The Data:

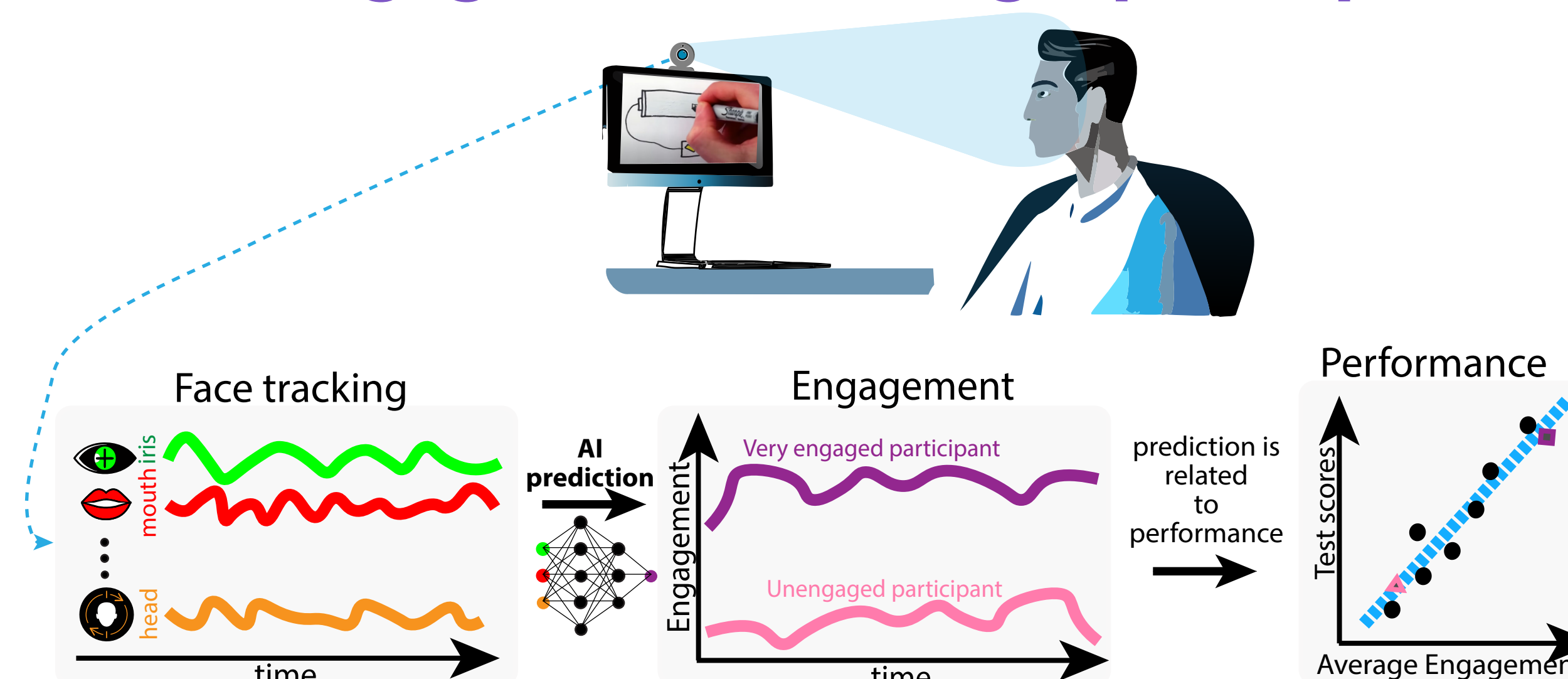| Dataset | Subjects | | Stimuli | | Webcam Data |
|---|---|---|---|---|---|
| | N | N | Set | | Duration (hours) |
| Experiment 1 | 26 (10M 16F) | 5 | A | | 9.18 |
| Experiment 2 | 29 (10M 19F) | 5 | A | | 11.14 |
| Experiment 3 | 28 (8M 20F) | 6 | B | | 15.42 |

## Highlights

- **Engagement tracking from a standard webcam is feasible.**
- **Attention can be measured privately and remotely.**
- **The prediction generalizes to unseen participants and stimuli.**

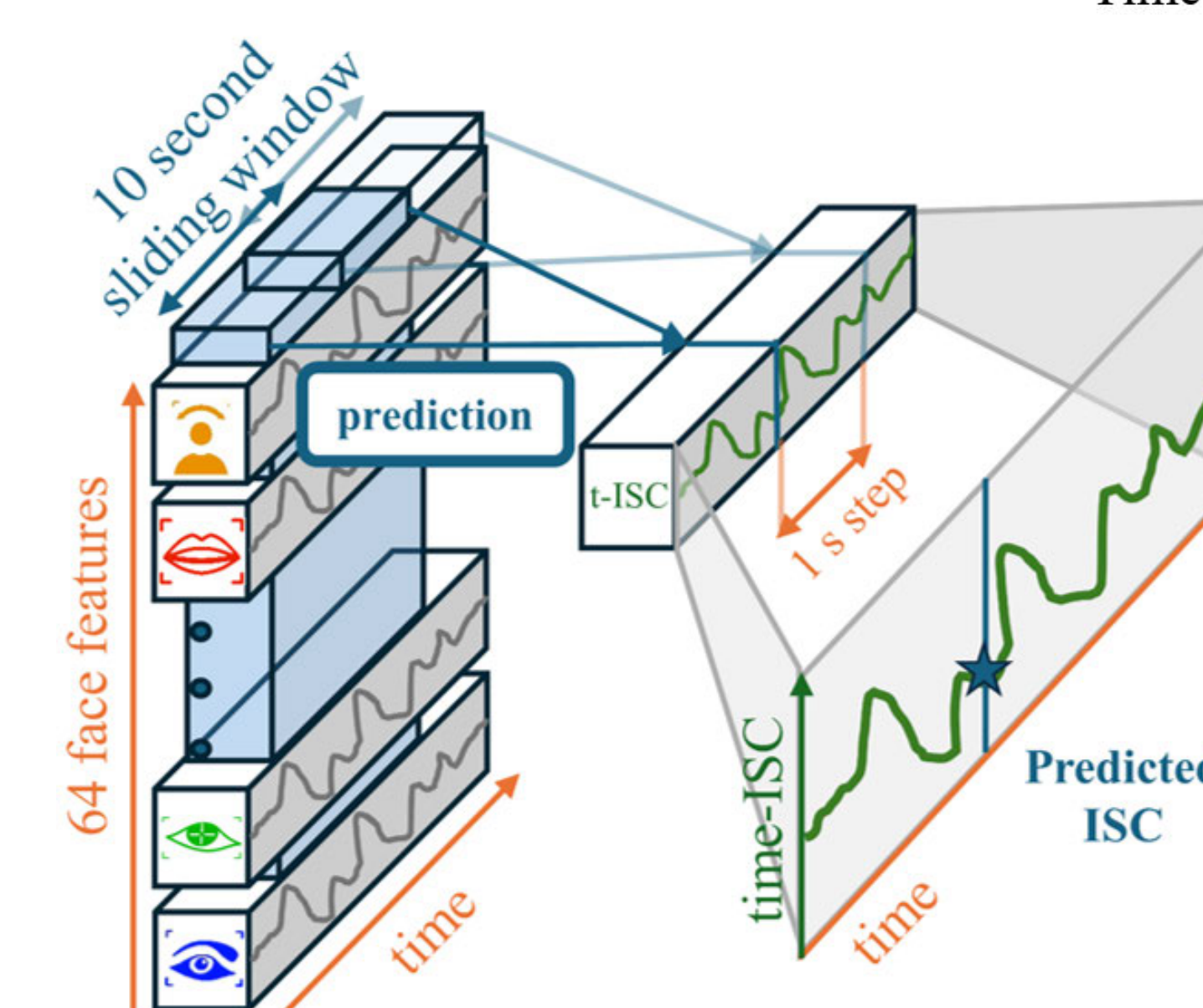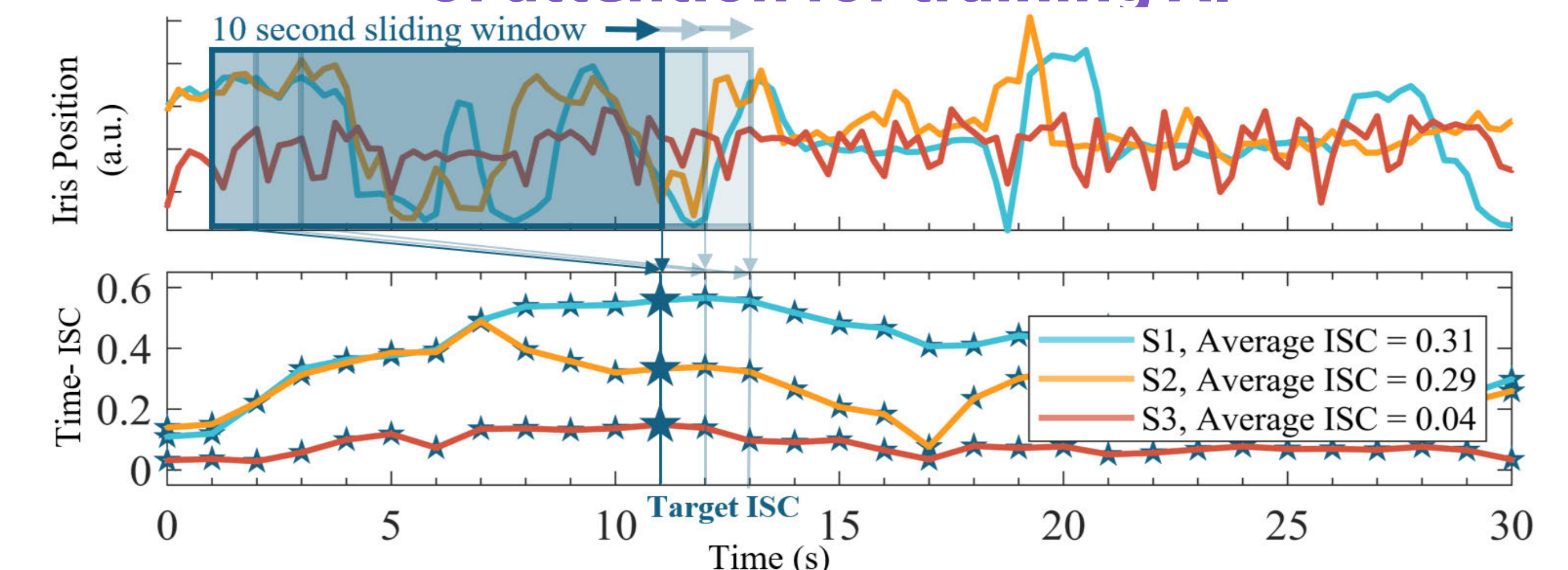## Engagement tracking from the face is comparable to eyetracking

- This is **consistent across participants, experimental conditions, and tested videos.**
- But we **still need a reference group to compute ISC.**



Gaze & Pupil synchrony (eyetracker)  r(27)=0.55 p=1.9·10⁻³
Eye Movement synchrony (webcam)  r(27)=0.51 p=4.7·10⁻³
Head Movement synchrony (webcam)  r(27)=0.19 p=0.33
Eyes on screen (eyetracker)  r(27)=0.37 p=0.05

Test Score [%]
Average participant ISC while attentive (ISC_A)
% total view time


Measurement to score correlation accross three different experiments

- ★ Exp. 1, N=26: intentional learning, 5 videos
- ■ Exp. 2, N=29: incidental learning, 5 videos
- ◆ Exp. 3, N=28: incidental learning, 6 videos

r(measurement, test scores)
Gaze & Pupil synchrony (ISC_A) | Eye movement synchrony (ISC_A) | Head movement synchrony (ISC_A) | Eyes on screen (% total view time)
Measured modality

## Using facial movements to predict engagement on a single participant



Face tracking | Engagement | Performance
iris | mouth | head | time
AI prediction
Very engaged participant
Unengaged participant
Engagement | time
prediction is related to performance
Test scores
Average Engagement

## Using time-resolved ISC as an objective index of attention for training AI



10 second sliding window
Iris Position (a.u.)
Time-ISC
Target ISC

- S1, Average ISC = 0.31
- S2, Average ISC = 0.29
- S3, Average ISC = 0.04

Time (s)



10 second sliding window
64 face features
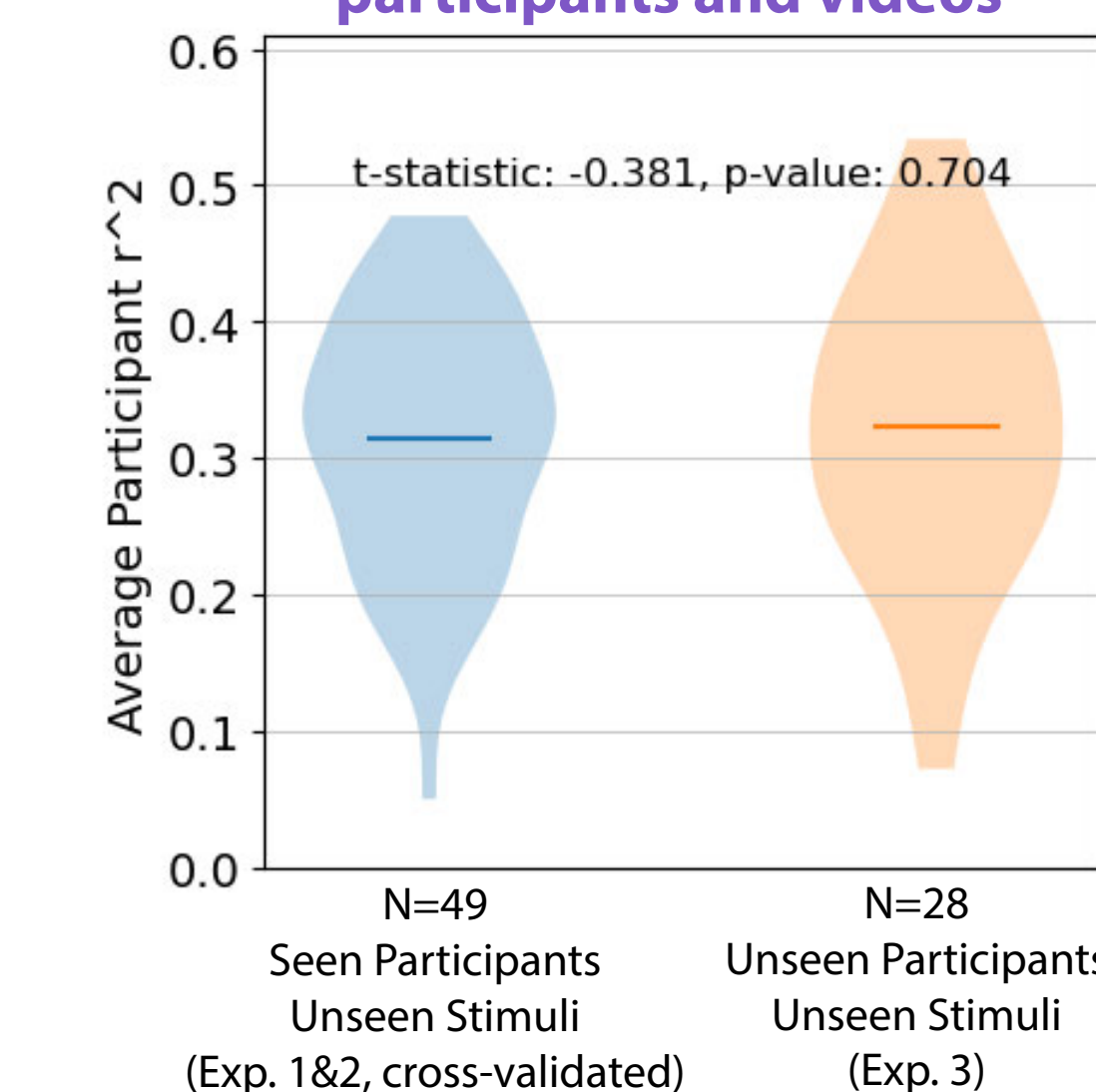prediction
t-ISC
1 s step
time-ISC
Predicted ISC

### Target
**ISC** of eye movements measured in 10 second windows.
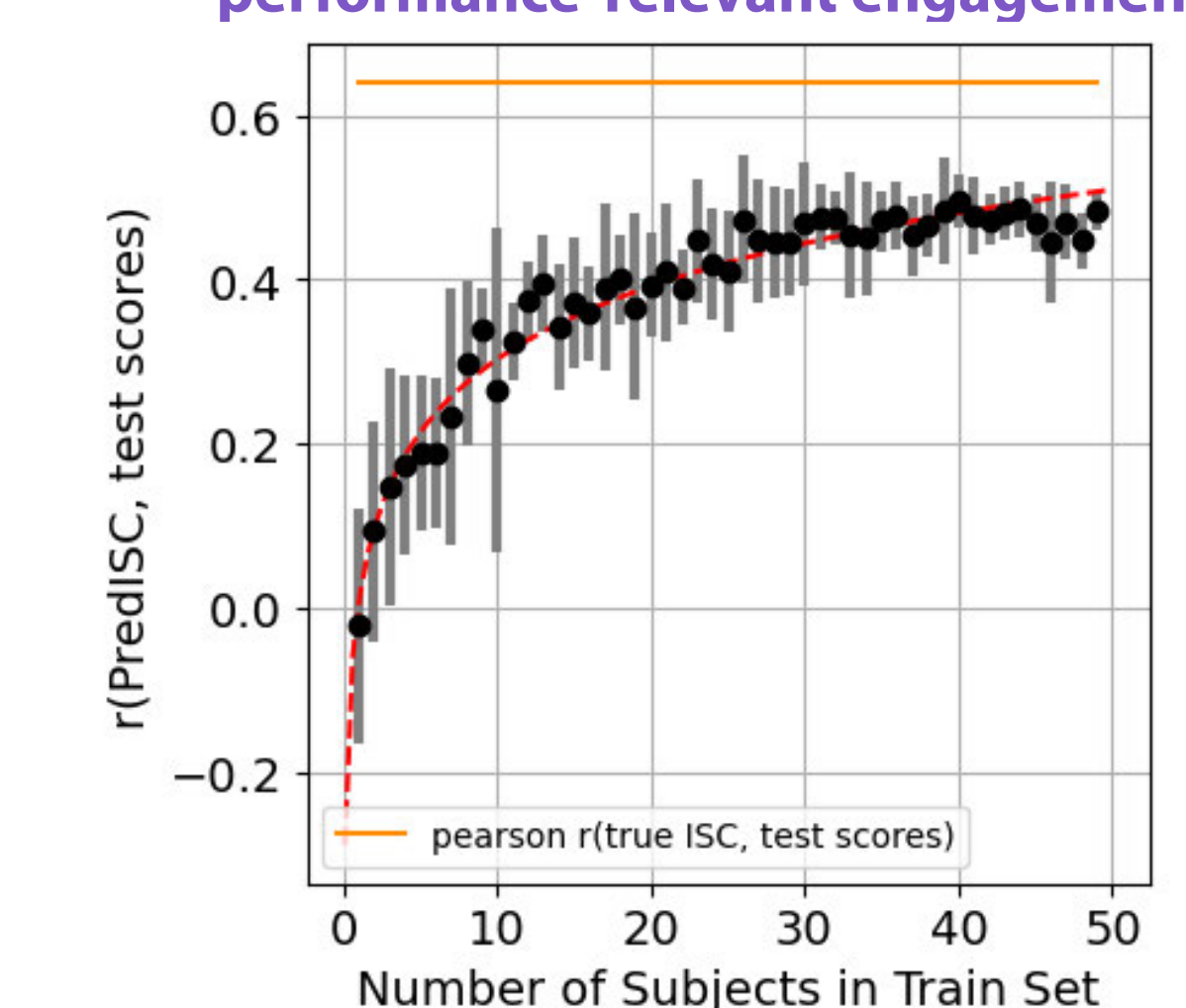
### Features
Mediapipe Blendshapes (facial movements) and Head Rotation features as predictors, over the preceeding 10 seconds.

## AI engagement predictions translate to unseen participants and videos, and correlate with scores

**Results generalize equally well to unseen participants and videos**



t-statistic: -0.381, p-value: 0.704

Average Participant r^2

N=49 Seen Participants Unseen Stimuli (Exp. 1&2, cross-validated)
N=28 Unseen Participants Unseen Stimuli (Exp. 3)

**How many participants are needed to capture performance-relevant engagement?**



r(PredISC, test scores)
pearson r(true ISC, test scores)
Number of Subjects in Train Set

### Check it out!
Article Preprint | Live Face Tracking demo | Time Resolved ISC

Code  github.com/asortubay/timeISC