**ETH**

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

MagicLeap

Distributed Computing

Prof. R. Wattenhofer

# Multi-modal Contrastive Learning for Emotion Recognition

HMDs today have a lot of sensors mounted on them. There is a huge diversity among the sensor types such as cameras, depth sensors, microphones, IMUs, etc. Data coming from these sensors is rich with information. This data can be used to infer different aspects of the person's physiological and physical condition.

The goal of this thesis is to create a multi-modal pipeline that can extract facial and body expressions using machine learning approaches. The idea is to use images from world, eyes, and face cameras, data from IMUs, and data from other sensors mounted on the HMD. Data coming from the sensors is usually processed in isolation, and used to infer things like, head pose tracking, eye tracking, hand tracking but they are rarely used multi-modally to provide a full description, and emotional and expression state of the person wearing the HMD.

The pipeline should be able to do voice stress analysis in combination with semantic speech analysis, eye rapid movement analysis, eye expressions analysis and extract different emotional states of the person that can be used to be displayed on an avatar and/or used in different high risk environments.

Some of the goals include:

- Develop a self-supervised approach to cluster emotions for different modalities.

- Establish a baseline for all of those expressions and emotions and figure out the minimal subset of information needed to describe a person.

- Put all that together and try to improve the results by combining the data to achieve a multimodal approach.

In this project, you will have the opportunity to collaborate with Magic Leap.

**Requirements:**
Knowledge in Deep Learning, or solid background in Machine Learning.
Implementation experience with TensorFlow or PyTorch is an advantage.

**Interested? Please contact us for more details!**

**Contact**

- Ard Kastrati: kard@ethz.ch, ETZ G61.3

- Dushan Vasilevski: dvasilevski@magicleap.com, Magic Leap