# Zero-shot Emotion Classification via Reinforced Self-training

@進捗報告
M2 LIU YI(21860638)
情報科学域
指導教員 下川原英理

Yamaguchi Lab.,
**TOKYO METROPOLITAN UNIVERSITY**

# Zero-shot Learning

- **Zero-shot learning (ZSL)**

ZSL is a challenging task as no labeled data is available
for unseen classes during training.

- **When we could use ZSL?**

If we need   **A more generalizable AI**   that can even recognize non-observed classes
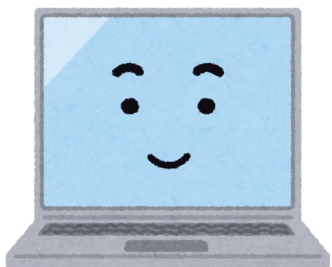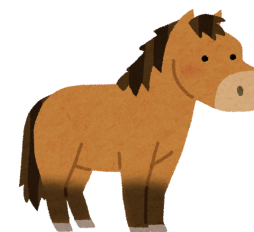
If we are   **Lack of labled training data**   labeling is a pain
or we even don't have the data at all
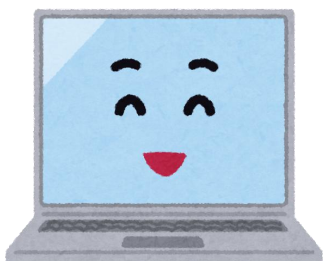
# Zero-shot Learning

# Zero-shot Learning



**Feature space
extracted from Image**

**Semantic space
transformed from auxiliary information**

**So we do not need Zebra's image
we only need Semantic space about zebra**

Yamaguchi Lab
**TOKYO METROPOLITAN UNIVERSITY**

# Zero-shot Learning

**We need some form of auxiliary information**

and this type of information can be of several types:

1) Attributes

2) Textual description

3) Class-class similarity

Can be converted to **semantic space**

Yamaguchi Lab
TOKYO METROPOLITAN UNIVERSITY

# Zero-shot Learning

## Issues in ZSL

- **How to accurately define the description of the Zero-shot class**

how about learning it from the (unseen)test dataset?

- **Domain shift problem**



The same 'hasTail' attribute different visual appearance

(a) visual space     (b) attribute space

- **Hubness problem**

In high-dimensional space, some points will be the nearest neighbors of most points

- **Semantic gap**



- **Semantic loss**

Yamaguchi Lab
TOKYO METROPOLITAN UNIVERSITY

# Zero-shot Learning for emotion classification

**Complex, compounded** emotional expressions are common!
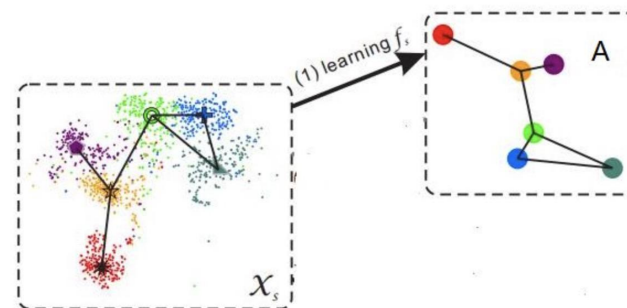
*Ex. **happily surprised** and **angrily surprised**

**Because of**

- Individual differences
- environmental influences
- diversity of emotional expressions

**it is difficult!** →

- enumerate all emotional biometric data
- collect enough samples for each category

For the relatively rare samples of emotional expressions
**It's just like unseen class!** (So can we use Class-class similarity or sth..?)

# Self-training and Reinforcement Learning Model

# Experiment Design

**APP:** VRChat VR CHAT
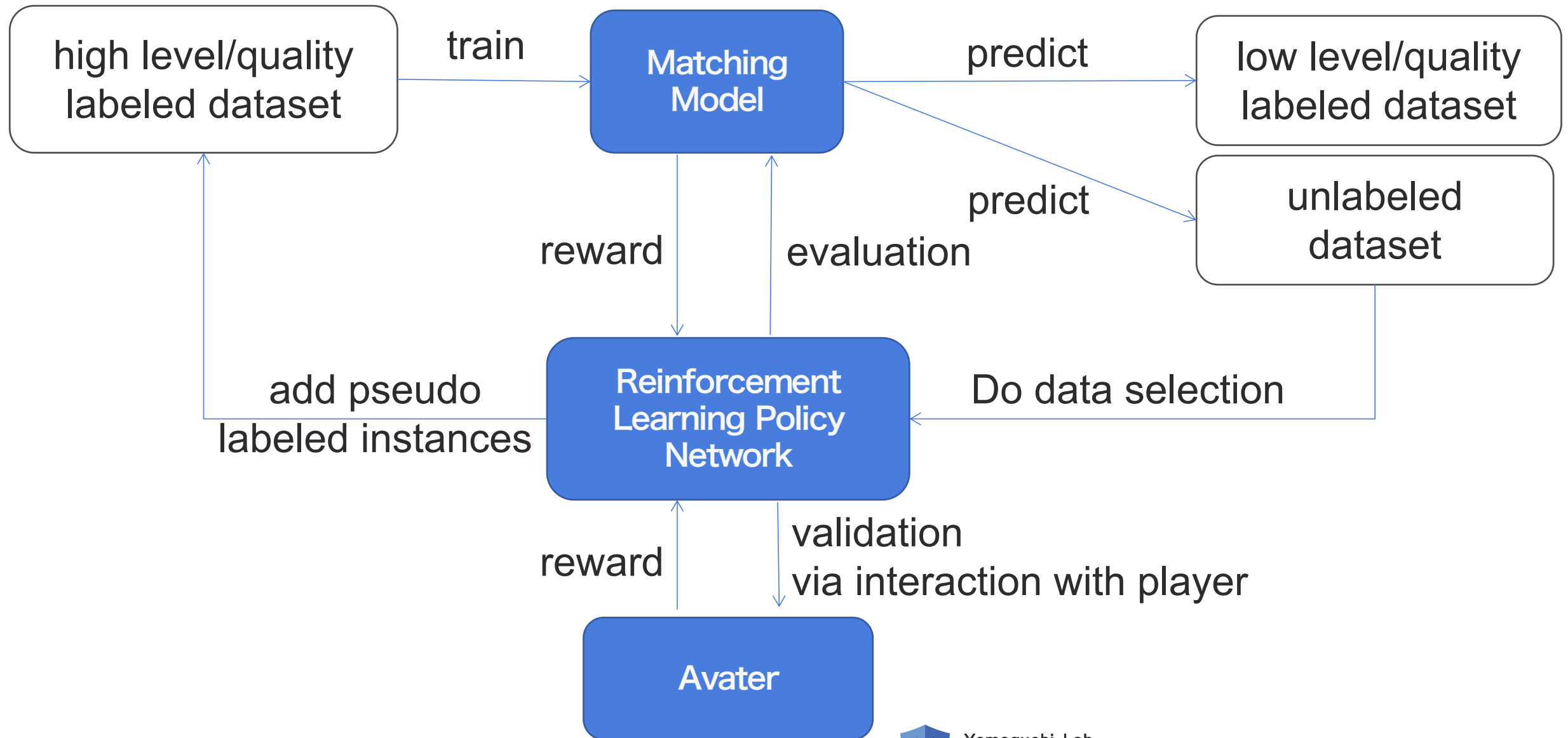
**Participants:** more than 20 groups of collaborators (3 people a group)

**Used Raw Data:** ECG   Audio(wav)   Eye tracking   Head position&rotation

**Details:**

One participant wear sensors and VR headset as a main talker
other two participants just wear VR headset and talk
Each conversation lasted approximately **three minutes**

**Topics:**

Talk about specific topics that are likely to cause emotional fluctuations or produce opposing positions

# Experiment Design

I am not sure which questionnaire to use

# Data-processing and Matching Model



ECG → R-R Interval → Heart Rate Variability(HRV)

Audio(wav) → Time-domain Audio Features

Eye tracking Data → Pupil Diameter&EyeOpen

Pupil Position → Nearest Neighbor Index of pupil position

Head position&rotation

→ **Train class Matching Model**

# Reinforcement Learning for Self-training

## Time-series data selection

Biometric data of 3minutes Dialogue

Interval   Clip2                                    Clip Length

Time-series data

.......

Clip1                              ClipN

*Ex. If set interval 1s,
clip length 30s,
we can get about 270
clips

## Policy network
Train a policy network to decide whether select this clip or not

# Reinforcement Learning Network

## For Reinforcement Learning Network

**(1)State:**

preprocessed biological data and confidence point

**(2)Action:**

Two class, whether choose this sample

**(3)Reward:**

Based on the model's performance on validation set

**(4)Policy Network:**

Input is State

Output is Action's probability distribution

$$r_k = \frac{(F_k^s - \mu^s)}{\sigma^s} + \lambda \cdot \frac{(F_k^u - \mu^u)}{\sigma^u}$$

其中:

- $F^S$: 可以看见类型的序列
- $F^U$: 不可以看见类型的序列
- λ: 权重
- μ: 均值
- σ: 方差

**policy Network**: 使用多层感知机作为挑选策略网络，输入为state，输出为是否挑选当前实例的概率（action 的概率），计算公式如下，

$$z_t = \text{ReLU}\left(W_1^T c_{x,y^*} + W_2^T p_{x,y^*} + b_1\right)$$
$$P\left(a \mid s_t\right) = \text{softmax}\left(W_3^T z_t + b_2\right)$$

其中:

- $W_1, W_2, W_3, b_1, b_2$ 为多层感知机的参数
- P() 为 action 的概率

# Discussion

- How to collect **high quality data** to train a not bad Matching Model first

- What algorithm to use to train **Matching Model**

- What **Policy Network algorithm** to use

- How to **balance the reward and evaluation** of RL?

- How to design the **interaction between avater and player** in VR

Yamaguchi Lab
TOKYO METROPOLITAN UNIVERSITY

# Acknowledgment

**Thank you for listening**
**ご清聴ありがとうございます**

Yamaguchi Lab
TOKYO METROPOLITAN UNIVERSITY

# Related Work

There are extensive works proposed in **zero-shot image/text classification** task

## Related Work

Zero-shot Text Classification via Reinforced Self-training

A Generalized Zero-Shot Framework for Emotion Recognition from Body Gestures

# Why emotion and dialogue mood？

- VRにおける、複数人の対話の雰囲気や個人の感情を把握して適切な介入を行うことで、コミュニケーションを円滑させる対話支援アバターの開発を目指す。
- **The goal is to develop a dialogue support avatar in VR**

- 雰囲気工学では，多人数の会話場における雰囲気を分析することや，複数の会話エージェントや会話ロボットによる人工的な言語，非言語情報が作り出す会話場の雰囲気の分析を目指す

Yamaguchi Lab
TOKYO METROPOLITAN UNIVERSITY

# Shortcomings of the previous study

## Shortcomings of the previous study

- – Insufficient amount of experimental data
- – individual differences appeared
- – The means of feature extraction of the data needs to be improved
- – Difficulty in confirming whether self-report accurately describe their own emotions
- – Collaborators exposed to VR for the first time tend to show excitement

Yamaguchi Lab
**TOKYO METROPOLITAN UNIVERSITY**

# Question

如何保证收集到高质量的情感数据以训练出高质量的Matching Model

要用什么Policy Network算法
reward和evaluation的平衡怎么办

VR中avater和player的交互方式怎么设计
avater的evaluation/reward策略怎么设计