

A Prototype-Based Generalized Zero-Shot Learning Framework for Hand Gesture Recognition

Jinting Wu

Institute of Automation, Chinese Academy of Sciences

University of Chinese Academy of Sciences

- **Introduction**
- **Methods**
- **Results**
- **Conclusion**

□ Motivation

- Most existing works can only recognize a limited number of categories that have been seen during training.
- Generalized Zero-Shot Learning (GZSL) provides a solution for tackling the above challenges. However, GZSL approaches for dynamic hand gesture recognition are less explored.

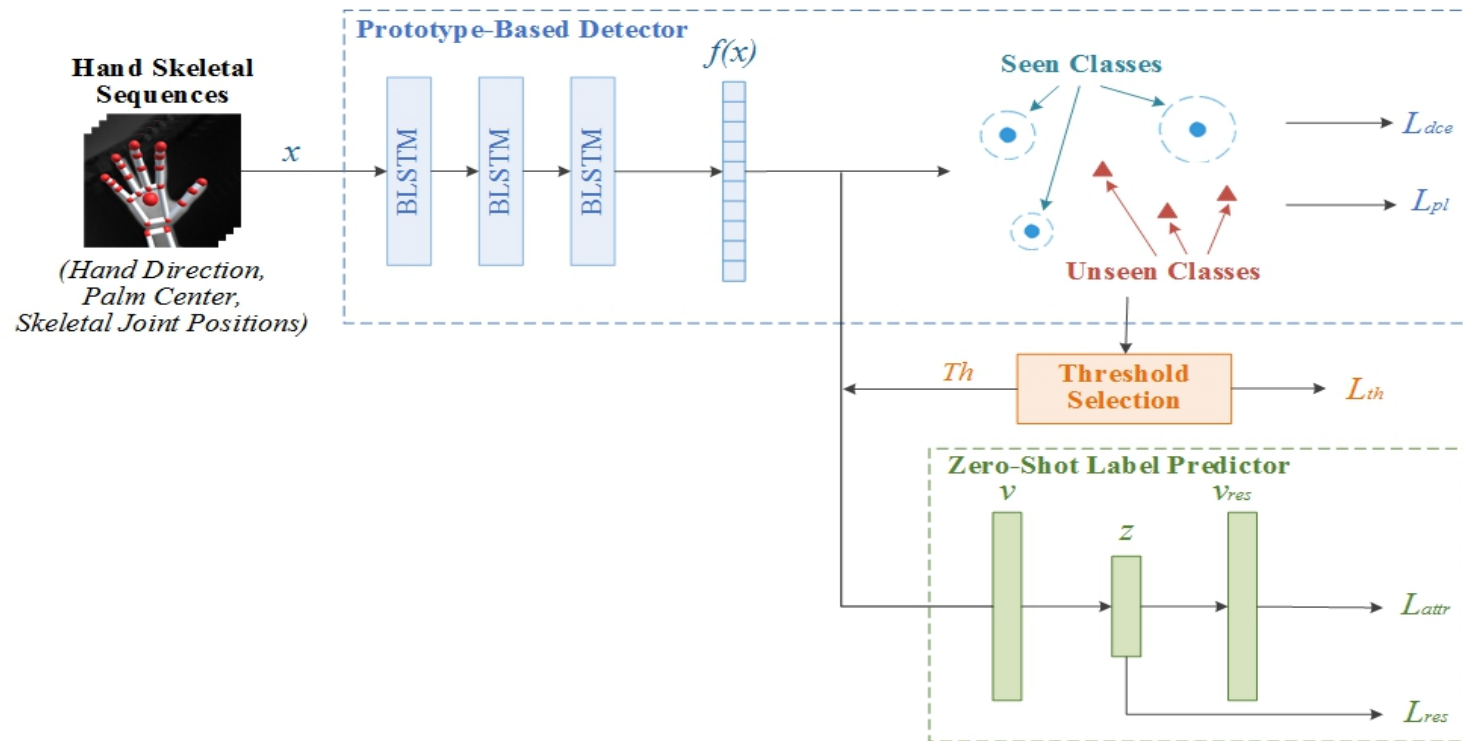
□ Contributions

- We propose an end-to-end prototype-based GZSL framework for hand gesture recognition which consists of two branches.
- We establish a hand gesture dataset that specifically targets this GZSL task.

- **Introduction**
- **Methods**
- **Results**
- **Conclusion**

□ Overview of the Proposed Framework

- Two branches
- Jointly training



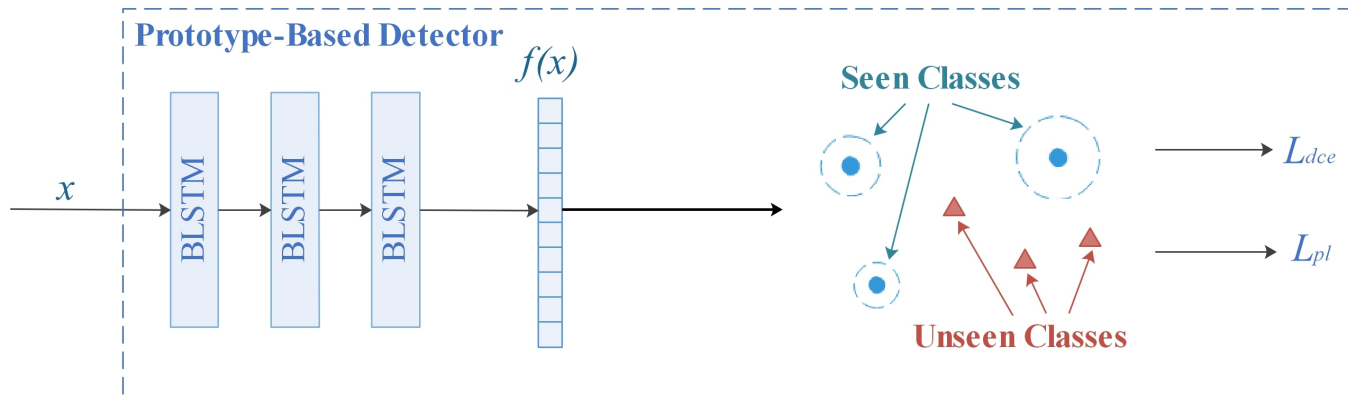
Methods

□ Prototype-Based Detector (PBD)

- Learning prototypes for each class
- Distance-based cross entropy loss and prototype loss

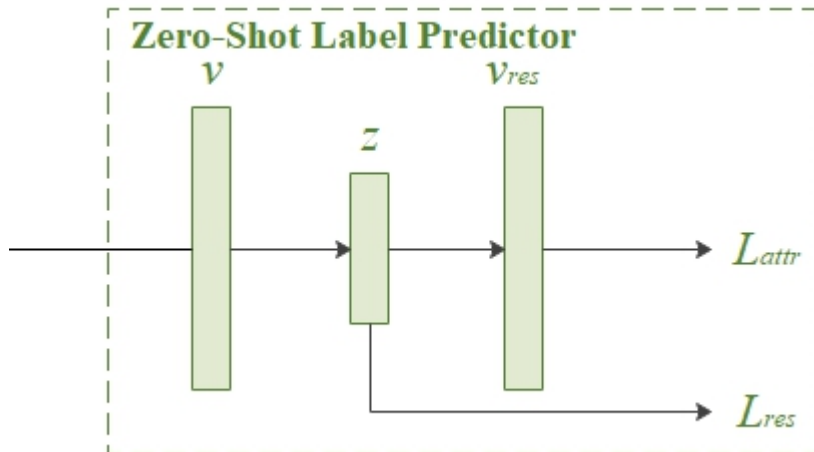
$$L_{dce}((x, y)|\theta, M) = -\log \sum_{j=1}^K \frac{e^{-\gamma \text{dis}(p_{pbd}(x), m_{yj})}}{\sum_{k=1}^C \sum_{l=1}^K e^{-\gamma \text{dis}(p_{pbd}(x), m_{kl})}}$$

$$L_{pl}((x, y)|\theta, M) = \|p_{pbd}(x) - m_{yj}\|_2^2$$



□ Zero-Shot Label Predictor

- Using a multi-layer Semantic Auto-Encoder (SAE) to predict the unseen gestures
- Attribute loss and reconstruction loss



$$L_{attr}((x, z_s) | \theta, \phi) = \|z - z_s\|_2^2$$

$$L_{res}((x, z_s) | \theta, \phi) = \|v - v_{res}\|_2^2$$

Methods:

□ End-to-End Learning Objective

- $L((x, y, z_s) | \theta, M, \phi) = L_{dce} + \lambda_1 L_{pl} + \lambda_2 L_{attr} + \lambda_3 L_{attr}$

□ Label Prediction

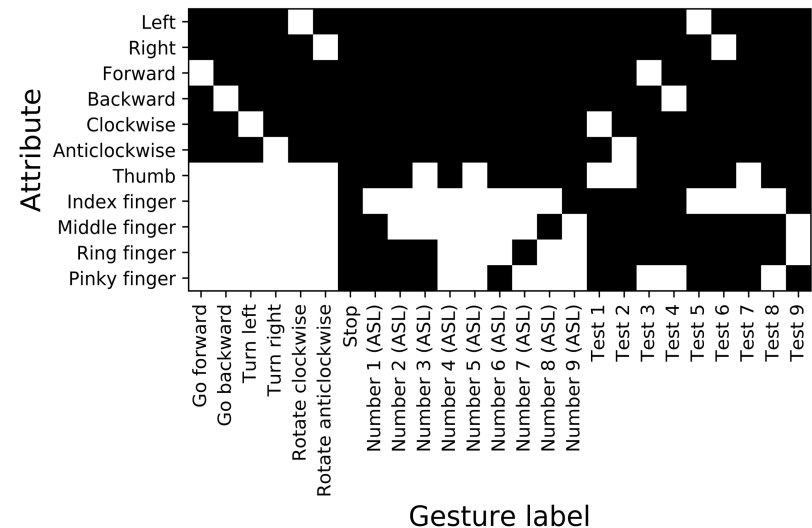
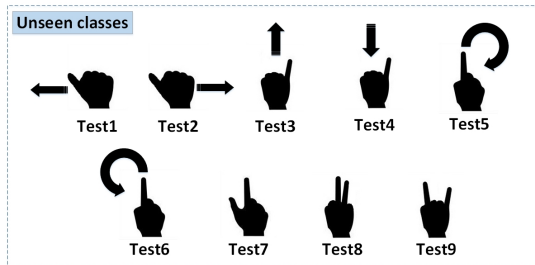
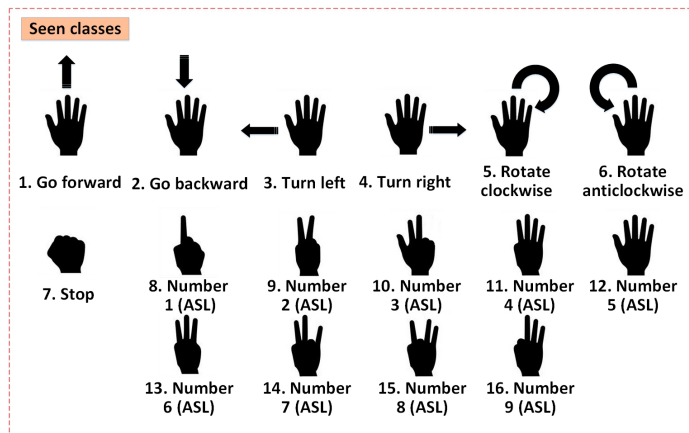
- Comparing the minimum distance in the prototype space $d_m(x)$ with the thresholds $Th(x)$.
- Seen categories: PBD result $\varepsilon(x)$
- Unseen categories: SAE result $\varepsilon_u(x)$

$$label(x) = \begin{cases} \varepsilon(x), d_m(x) \leq Th(x) \\ \varepsilon_u(x), d_m(x) > Th(x) \end{cases}$$

- **Introduction**
- **Methods**
- **Results**
- **Conclusion**

Dataset

- 16 seen gestures and 9 unseen gestures
- 11 attributes including hand movement and finger bending states



□ Experimental Results

■ State-of-the-art Comparisons

- Zero-shot gesture recognition method: ESZSL¹
- Generalized zero-shot object recognition method: CADA-VAE² and f-CLSWGAN³

Methods	Acc_s	Acc_u	H
ESZSL [15]	77.81%	13.89%	23.57%
CADA-VAE [11]	80.00%	53.89%	64.40%
f-CLSWGAN [12]	79.79%	55.00%	65.08%
End-to-End Framework (Ours)	89.06%	58.33%	70.49%

1. Madapana, Naveen, and Juan Wachs. Zsgl: zero shot gestural learning.

2. Schonfeld, Edgar, et al. Generalized zero-and few-shot learning via aligned variational autoencoders.

3. Xian, Yongqin, et al. Feature generating networks for zero-shot learning.

□ Experimental Results

■ Ablation Analysis

- The traditional SAE¹ without the prototype-based detector
- The framework with a fixed threshold
- The framework where two branches are trained separately

Methods	Acc_s	Acc_u	H	Test Time
BLSTM+SAE [6]	91.88%	15.00%	25.79%	0.023s
End-to-End Framework (Fixed Threshold)	84.69%	50.56%	63.31%	0.022s
PBD+SAE	90.63%	57.22%	70.15%	0.026s
End-to-End Framework	89.06%	58.33%	70.49%	0.022s

1. Kodirov, Elyor, Tao Xiang, and Shaogang Gong. Semantic autoencoder for zero-shot learning.

- **Introduction**
- **Methods**
- **Results**
- **Conclusion**

- **A prototype-based GZSL framework for hand gesture recognition**
 - **An end-to-end framework with two branches**
 - **A novel hand gesture dataset**
 - **Comprehensive experiments demonstrate the effectiveness of our proposed approach**

Thanks for your attention!

- Jinting Wu
- Institute of Automation, Chinese Academy of Sciences;
University of Chinese Academy of Sciences
- E-mail: wujinting2016@ia.ac.cn