# Path Representation Learning of Mixture of Experts Based on Contrastive Learning



東京工業大学

Masahiro Kada<sup>1</sup>, Ryota Yoshihashi<sup>1</sup>, Rei Kawakami<sup>1</sup>, Satoshi Ikehata<sup>1,3</sup> , Ikuro Sato<sup>1,2</sup>

<sup>1</sup>Tokyo Institute of Technology, <sup>2</sup>Denso IT Laboratory, <sup>3</sup>National Institute of Informatics

Recognition and Learning Algorithm Laboratory

### **1. Deep Mixture of Experts (MoE)**

WhatA type of new NN architecturescontaining Router and Expert blocks.



**Advantage** Scalability: Can add experts while keeping inference time constant.

**Dis-**<br/>advantageVulnerability: Output can change abruptlydue to the change in feedforward path<br/>even for small input deformation.

## **3. Proposed Method**

#### Basic idea

Making the router outputs (paths) closer to each other between a pair of images augmented from the same image.



Making the router outputs (paths) closer in each pair of corresponding image patches.

### **2. Contribution**

- Proposed Path Representation Learning that aims to robustify router output (`path') under input deformations.
- Confirmed that local path representations of transformer-based MoE become more insensitive to input deformations.
- Improved the performance over the previous method.

### **3. Previous Method**

[C. Riquelme+, NeurIPS2021] proposed Vision Transformer-based MoE.



The Router selects an expert in a patch-wise manner by taking the arg-max of softmax in each



JSD: Jensen Shannon Divergence  $JSD(p,q) = \frac{1}{2} \sum \left( p \log \frac{p}{M} + q \log \frac{q}{M} \right) \text{, where } M = \frac{p+q}{2}$ 

### **4. Preliminary Experiment**

#### Qualitative results

Previous method

Inconsistencies of expert selections in corresponding patches

0-	6	5	5	4	5	5	5	5	4	2	5	5	6	7	1
50 -	8	2	8	2	2	2	8	5	2	2	2	2	7	6	İ
50	2	2	2	2	8	3	2	2	2	2	2	2	6	5	i
LOO -	4	1	8	8	8	8	1	3	2	8	2	2	6	5	i
L50 -	5	8	8	8	2	2	2	2	8	2	2	2	6	5	İ
	4	2	2	2	2	2	2	2	2	2	2	2	2	2	İ
200 -	4	2	2	2	2	2	2	2	2	2	2	2	2	2	
50 -	4	2	8	8	2	2	2	2	8	8	2	2	2	2	
	4	2	6	8	6	2	2	2	8	8	6	2	2	1	
800 -	8	2	6	6	6	4	2	3	6	8	6	2	2	8	
	3	2	3	3	3	4	2	3	3	3	3	2	2	3	

6	7	7	5	5	7	7	7	7	5	5	6
7	6	7	5	7	7	6	7	7	5	5	6
6	5	5	5	5	6	6	3	3	3	3	3
6	5	5	5	5	6	6	3	3	3	3	3
6	5	5	5	4	2	6	3	4	3	3	3
2	2	5	5	4	6	4	6	4	4	4	4
2	2	2	2	2	2	1	8	4	3	4	4
2	2	1	8	8	4	4	6	4	4	4	4
2	1	2	2	2	4	4	3	4	3	_4	4
2	8	2	2	1	4	4	3	3	3	3	3
2	3	2	3	3	2	2	3	3	3	3	3

75.91%



#### **Observation**

Object labels tend to be *roughly* associated with paticular experts after supervised training.



#### Our method

More consistent expert selections in corresponding patches



#### Classification Accuracy Outperformed

NeurIPS2021] slightly.

[C. Riquelme+,

	Accuracy							
Previous method	75.84%							

Ours

Dataset: ImageNet-1K

This work is supported by DENSO IT LAB Recognition and Learning Algorithm Collaborative Research Chair (Tokyo Tech.).