# MoETTA: Test-Time Adaptation Under Mixed Distribution Shifts with MoE-LayerNorm
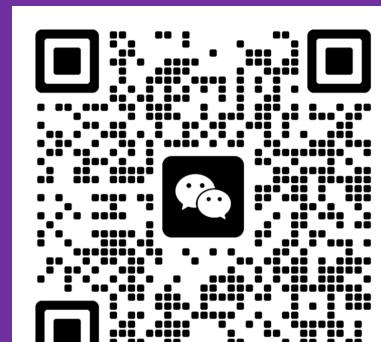
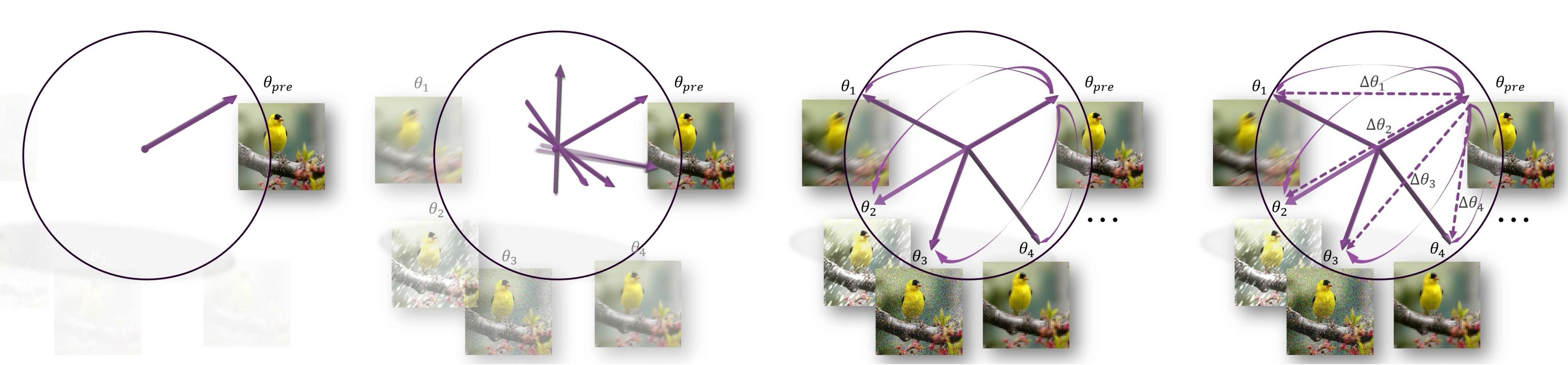Xiao Fan[1,2*], Jingyan Jiang[3†], Zhaoru Chen[3], Fanding Huang[2], Xiao Chen[2], Qinting Jiang[2], Bowen Zhang[3], Xing Tang[3], Zhi Wang[2]

* Work completed during an internship at Tsinghua University.   † Corresponding author.
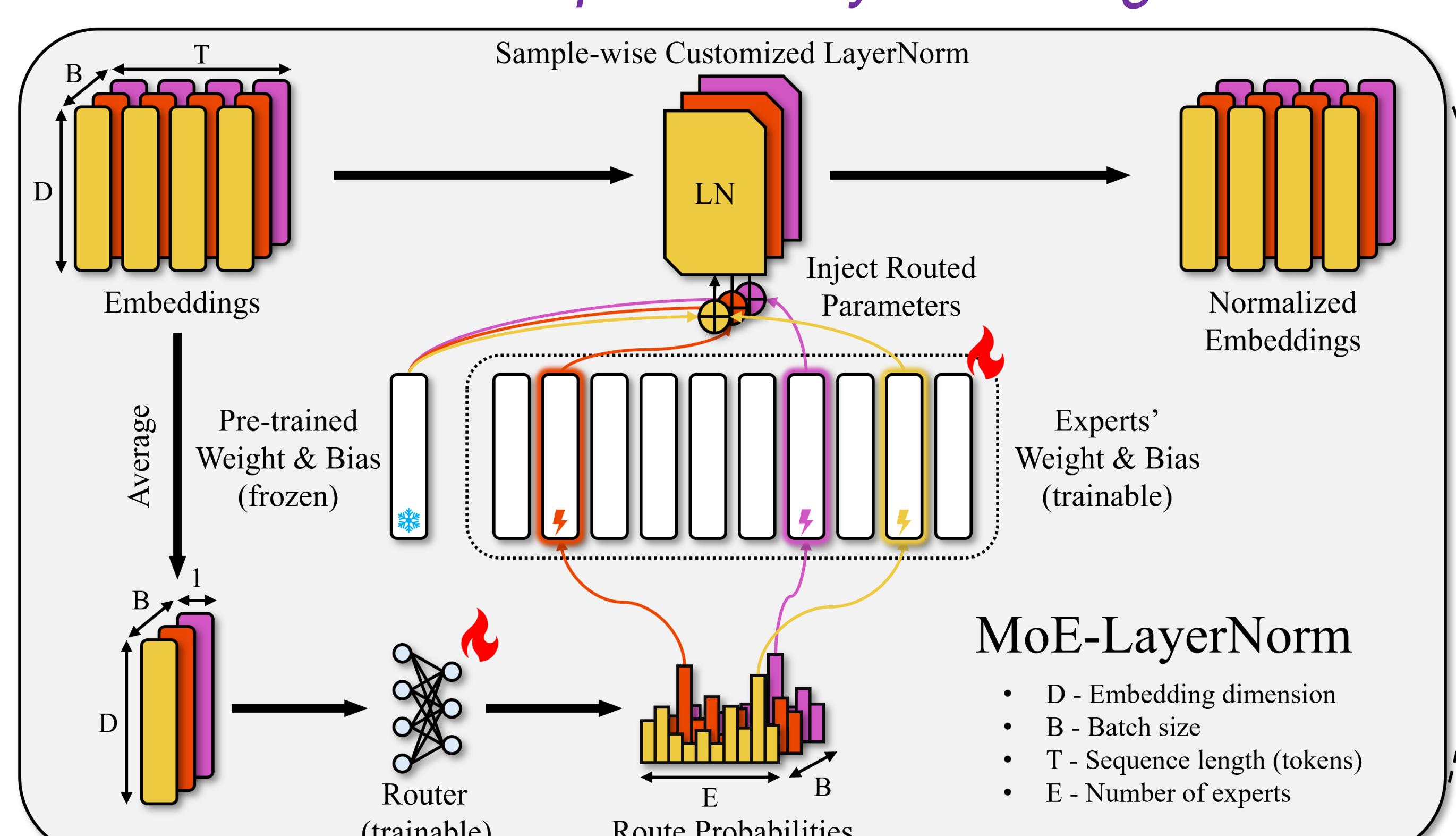
WeChat

Homepage

## Motivation

### Illustration of Adaptation Process Under Mixed Distribution Shifts



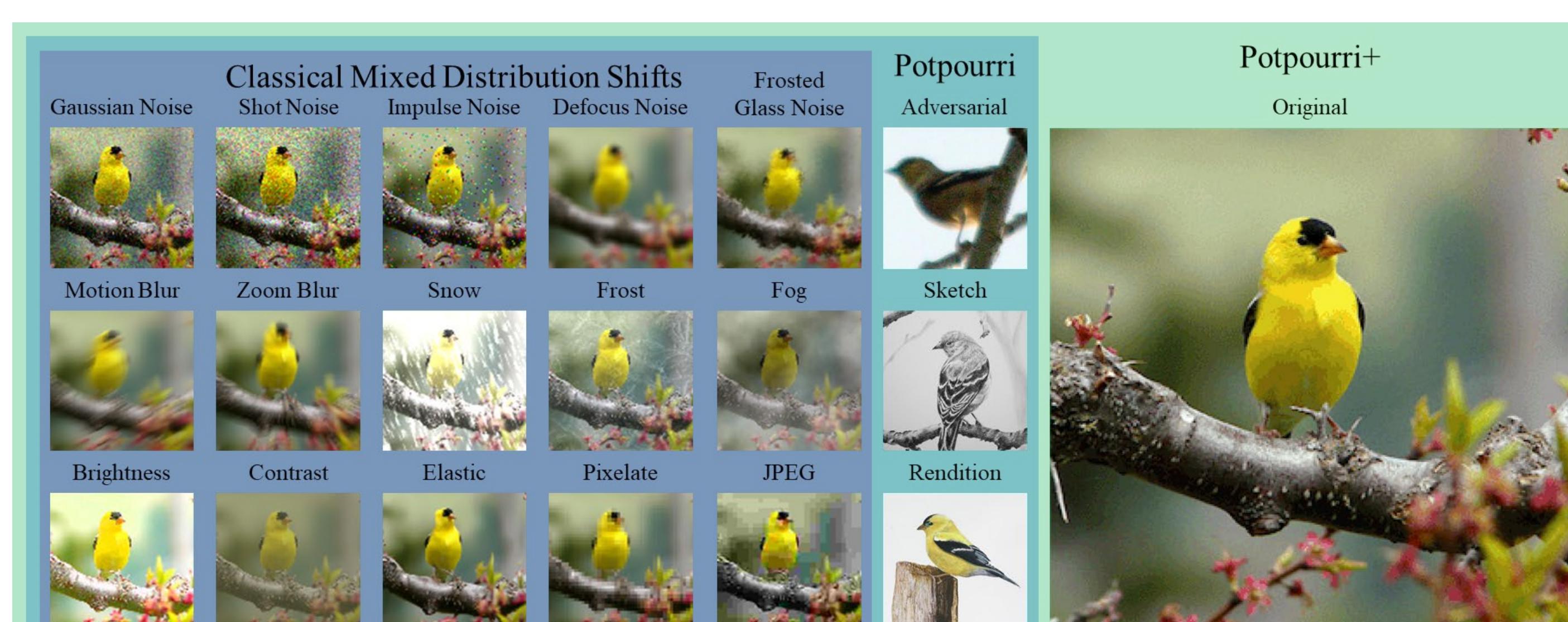*Domain-Specific gradient signals can be inconsistent or even conflicting.*

*Leveraging a diverse set of experts to represent multiple adaptation solutions within one model is particularly advantageous for mixed distribution shifts.*

### Theory

**Expectation of cosine similarity between adaptation directions of different domains.**

Let $\theta_1, \theta_2, \theta \in \mathbb{R}^d$ be i.i.d. $\mathcal{N}(\mathbf{0}, \boldsymbol{I})$, then $\lim_{d \to \infty} \mathbb{E}[\cos\langle\theta_1 - \theta, \theta_2 - \theta\rangle] = 0.5$.

### Empirical Evidence

| Model | Tent | SAR | EATA | DeYO |
|---|---|---|---|---|
| ViT-B/16 | 0.69 | 0.64 | 0.66 | 0.71 |
| ViT-L/16 | 0.34 | 0.37 | 0.28 | 0.26 |

## Methodology



MoE-LayerNorm
- D - Embedding dimension
- B - Batch size
- T - Sequence length (tokens)
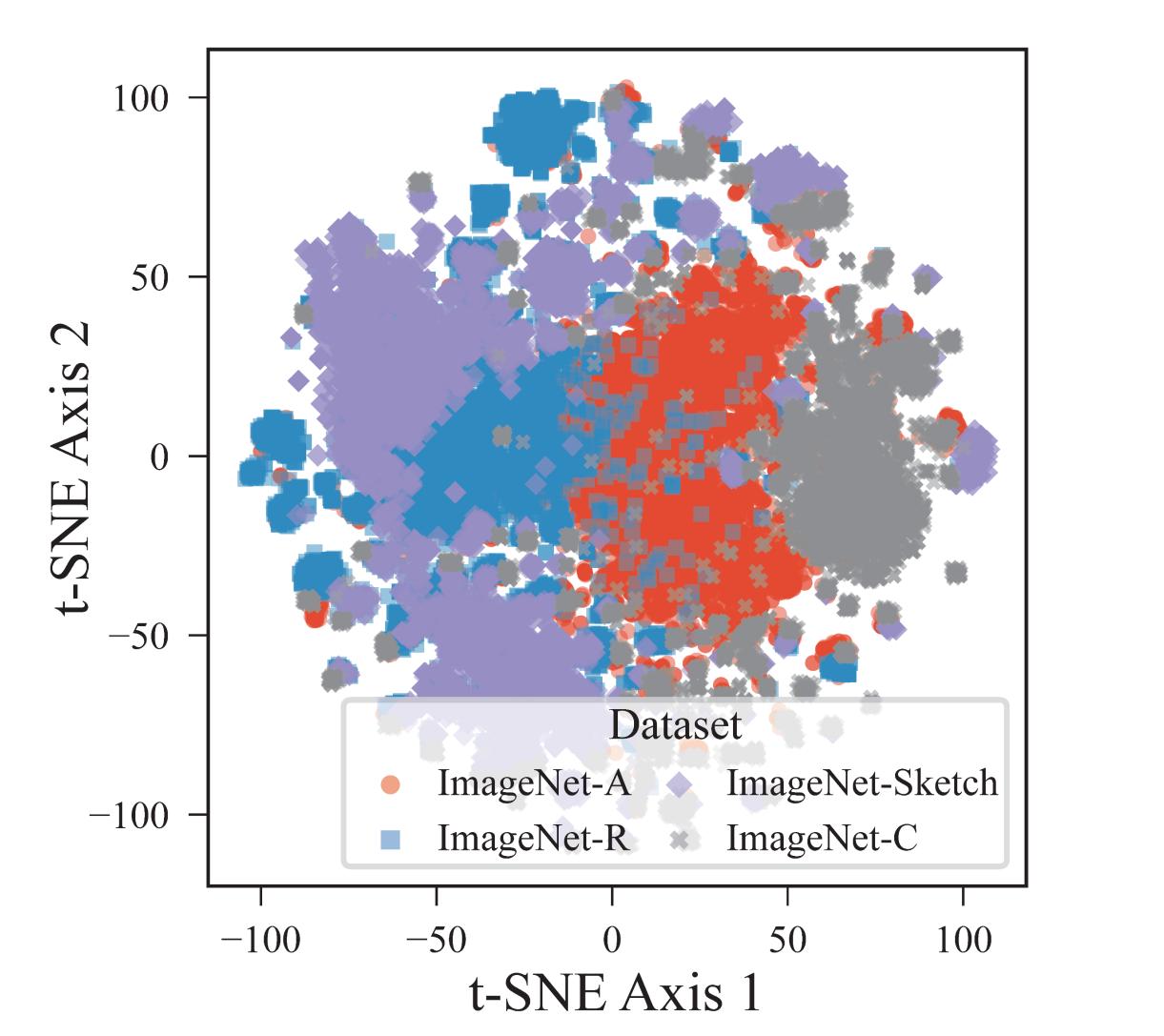- E - Number of experts

### Overall Loss Function

$$\frac{1}{|\mathcal{S}_t|} \sum_{x \in \mathcal{S}_t} \exp[\mathrm{E}_0 - \mathrm{Ent}(x)]\mathrm{Ent}(x) + \alpha_t \sum_{i=1}^{M} \mathcal{L}_{\text{load balancing}}^i$$

**Sample selection**
$$\mathcal{S}_t = \{x | \mathrm{Ent}(x) < E_{\max}^t \wedge x \in \mathcal{B}_t\}$$

**Dynamic threshold**
$$E_{\max}^t = \begin{cases} E_{\text{avg}}^0, & t = 0 \\ E_{\max}^{t-1} \times \frac{E_{\text{avg}}^t}{E_{\text{avg}}^{t-1}}, & t \geq 1. \end{cases}$$

**Trade-off coefficient**
$$\alpha_t = \begin{cases} \lambda \times E_{\text{avg}}^0, & t = 0 \\ \alpha_{t-1} \times \frac{E_{\text{avg}}^t}{E_{\text{avg}}^t}, & t \geq 1. \end{cases}$$

**Differentiable load balancing loss**
$$\mathcal{L}_{\text{load balancing}} = N \times \sum_{i=1}^{N} \boldsymbol{F}_i \times \boldsymbol{P}_i,$$
$$\boldsymbol{P}_i = \frac{1}{|\mathcal{B}_t|} \sum_{x \in \mathcal{B}_t} \boldsymbol{p}_i(x),$$
$$\boldsymbol{F}_i = \frac{1}{|\mathcal{B}_t|} \sum_{x \in \mathcal{B}_t} \mathbb{1}_{\{\arg\max_k \boldsymbol{p}_k(x)=i\}}.$$

## Benchmark

*Propose two more benchmarks*
- Potpourri: More OOD samples
- Potpourri+: Include ID samples





ViT CLS token t-SNE projection.

### Robustness to Mixed Distribution Shifts

| Model | Setting | Noadapt | Tent | EATA | CoTTA | SAR | DeYO | MGTTA | BECoTTA | Ours |
|---|---|---|---|---|---|---|---|---|---|---|
| ViT-B/16 | Classical | 55.52 | 63.20$_{0.08}$ | 64.28$_{0.09}$ | 60.53$_{0.57}$ | 60.76$_{0.04}$ | 63.97$_{0.04}$ | 66.20$_{0.01}$ | 61.57$_{0.08}$ | **67.20$_{0.03}$** |
| | Potpourri | 54.18 | 60.99$_{0.05}$ | 61.99$_{0.11}$ | 59.67$_{1.21}$ | 58.71$_{0.03}$ | 61.66$_{0.02}$ | 62.98$_{0.26}$ | 59.08$_{0.86}$ | **65.12$_{0.08}$** |
| | Potpourri+ | 55.92 | 62.28$_{0.03}$ | 63.17$_{0.06}$ | 59.26$_{0.68}$ | 59.99$_{0.07}$ | 62.90$_{0.03}$ | 64.35$_{0.07}$ | 58.87$_{3.23}$ | **66.15$_{0.06}$** |
| ConvNeXt-B (CNN Arch.) | Classical | 54.81 | 58.88$_{0.06}$ | 64.50$_{0.06}$ | 59.65$_{0.04}$ | 61.67$_{2.65}$ | 64.32$_{0.03}$ | - | 50.16$_{7.96}$ | **67.40$_{0.02}$** |
| | Potpourri | 53.91 | 58.23$_{0.05}$ | 62.69$_{0.07}$ | 58.57$_{0.00}$ | 61.16$_{0.41}$ | 62.46$_{0.07}$ | - | 28.28$_{20.54}$ | **65.70$_{0.05}$** |
| | Potpourri+ | 55.69 | 59.69$_{0.04}$ | 63.94$_{0.07}$ | 60.02$_{0.02}$ | 62.72$_{0.10}$ | 63.57$_{0.06}$ | - | 48.92$_{9.35}$ | **66.68$_{0.07}$** |

## Experiment

### Computation efficiency comparison

| Method | #Act. params per sample | #Fwd | #Bwd | Used time |
|---|---|---|---|---|
| Noadapt | 0 | 100% | 0% | 100% |
| Tent | 0.04M | 100% | 100% | 226% |
| EATA | 0.04M | 100% | 80% | 239% |
| SAR | 0.03M | 199% | 175% | 440% |
| DeYO | 0.04M | 196% | 53% | 317% |
| CoTTA | 86.42M | 199% | 100% | 798% |
| MGTTA | 2.80M | 100% | 100% | 227% |
| BECoTTA | 0.13M | 100% | 86% | 334% |
| Ours | 0.23M | 100% | 76% | 247% |

### Ablation analyses

| | Classical | Pot. | Pot.+ | Avg. |
|---|---|---|---|---|
| **Full method** | **67.25** | **65.14** | **66.21** | **66.20** |
| **Loss Components** | | | | |
| w/o Sample selection | 67.04 | 64.01 | 57.61 | 62.89 |
| w/o Entropy re-weight | 62.86 | 60.51 | 61.79 | 61.72 |
| w/o $\mathcal{L}_{\text{load balancing}}$ | 26.27 | 16.27 | 21.29 | 21.28 |
| **MoE Architecture** | | | | |
| w/o Grad to router | 65.17 | 62.80 | 63.92 | 63.96 |
| w/o Sample-wise router | 28.69 | 28.60 | 24.96 | 27.42 |
| w/o MoE-LayerNorm | 22.38 | 17.94 | 26.93 | 22.42 |
| w/o Layer-wise router | 17.40 | 27.09 | 15.18 | 19.89 |

## Evolution and Final-State Statistics of Expert Parameter Similarity Across MoE-LayerNorms



**Hyper-Parameter Sensitivity**