



Fully Test-time Adaptation for Tabular Data

Zhi Zhou*, Kun-Yang Yu*, Lan-Zhe Guo[†], Yu-Feng Li[✉]

{zhouz, yuky, guolz, liyf}@lamda.nju.edu.cn

* These authors contributed equally ✉ Corresponding author

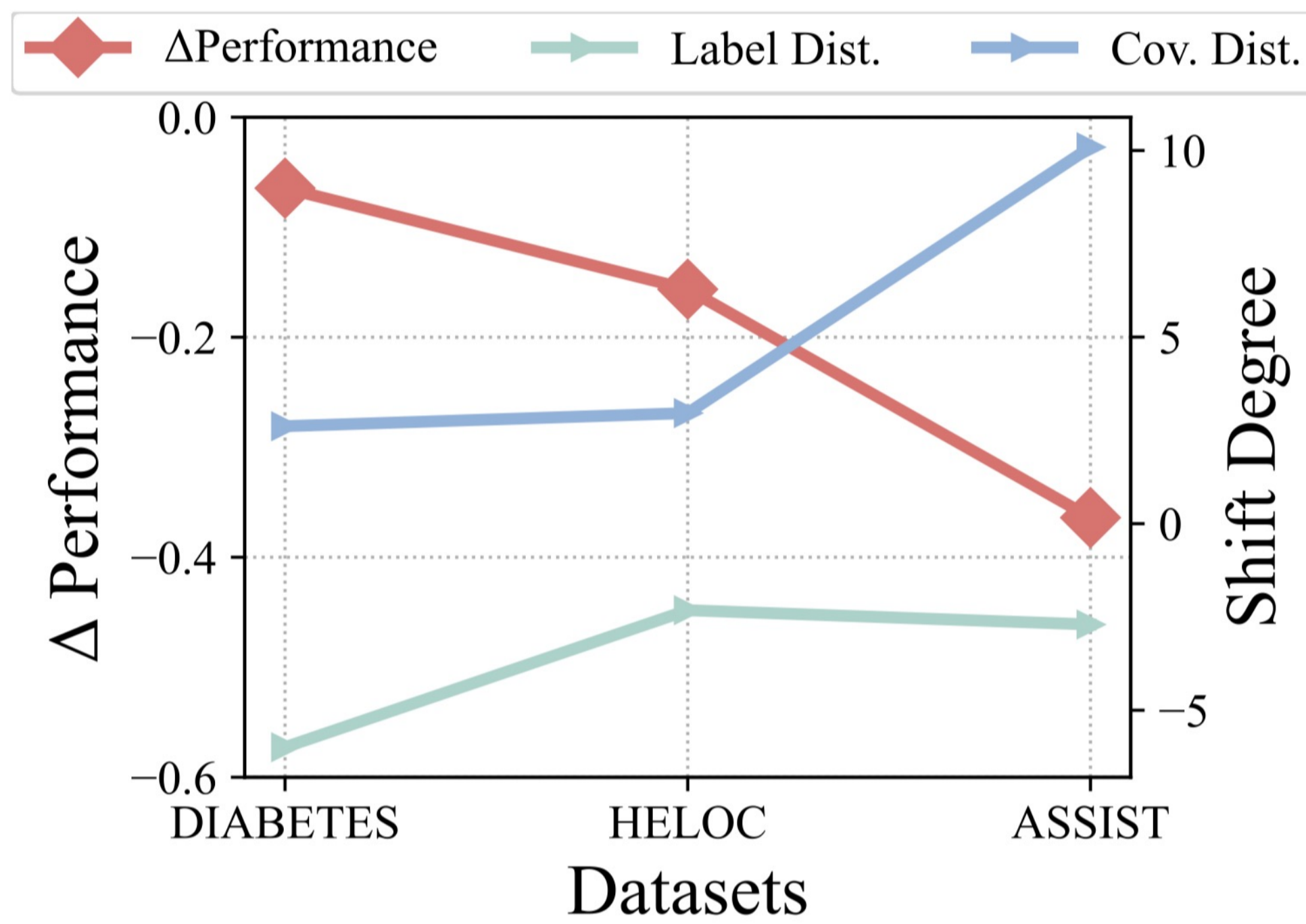


TL; DR We design a fully test-time adaptation for tabular data that addresses covariate and label distribution shifts.

Problem Setting

Distribution shifts in testing data render tabular models ineffective. Test-time adaptation offers a potential solution, but ...

Observation 1: Covariate distribution and label distribution shifts in tabular data hinder the performance of FTTA methods.



▲ Figure 1: The label and covariate distribution shifts between training and testing in tabular data degrade the model performance. The shift degree is taken logarithm for aesthetic purposes.

Observation 2: Typical augmentation used in test-time adaptation is ineffective for tabular data.

Method	DIABETE	HELOC	ASSIST
Non-Adaptation	60.82 ± 0.22	54.37 ± 5.35	55.86 ± 3.81
σ = 0.2	60.46 ± 0.20	46.40 ± 3.08	54.89 ± 1.88
σ = 0.4	59.18 ± 0.42	43.36 ± 0.25	54.86 ± 3.00
σ = 0.6	57.73 ± 0.64	43.06 ± 0.07	54.51 ± 2.26
σ = 0.8	56.19 ± 0.83	43.07 ± 0.03	53.79 ± 3.80
σ = 1.0	54.74 ± 0.77	43.09 ± 0.01	54.23 ± 3.56

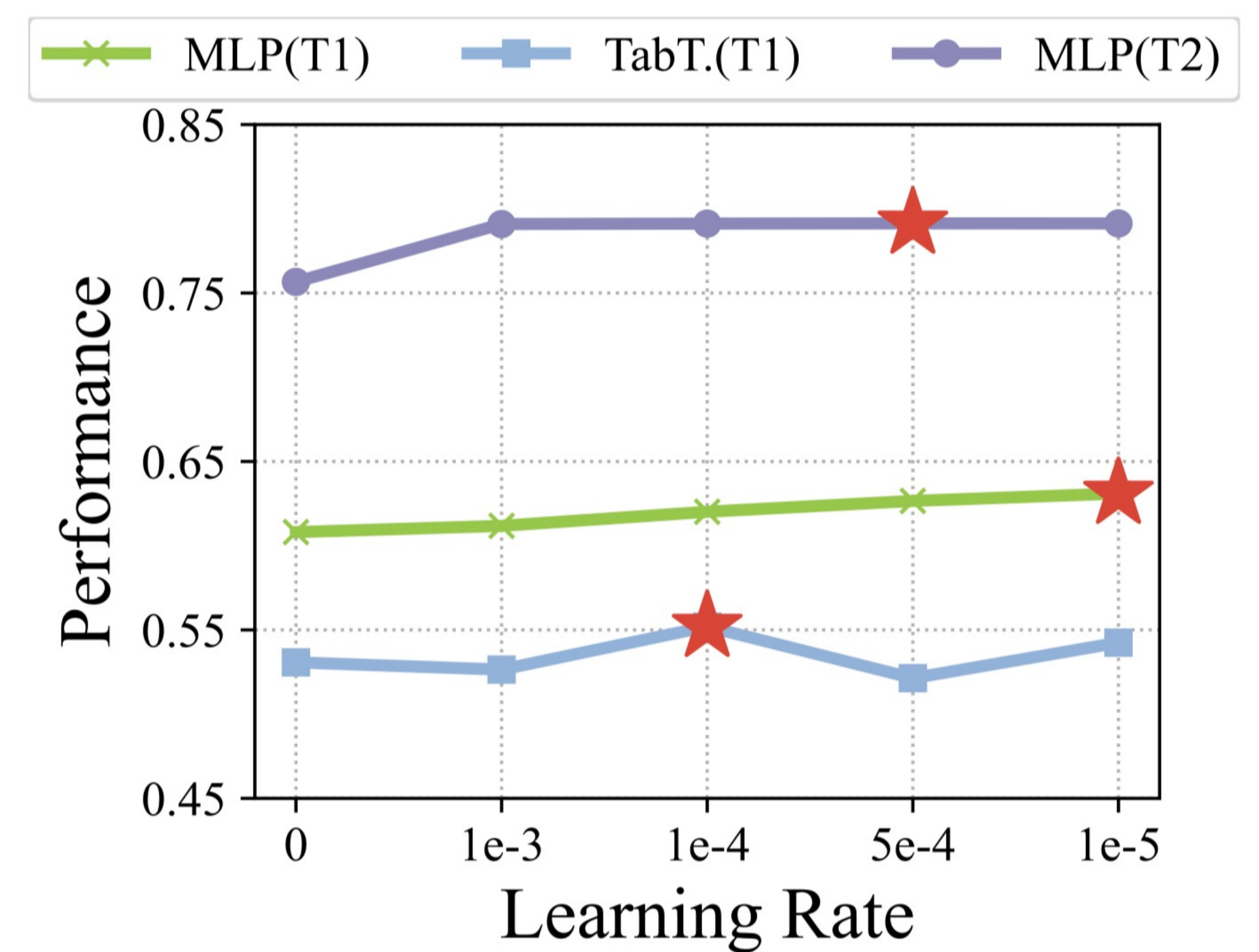
▲ Table 1: Performance of the non-adaptation baseline and CoTTA method with different augmentation strengths σ using the MLP model. The best performance is in bold.

Observation 3: Existing FTTA methods degrades when dealing with tabular data

Method	DIABETE	HELOC	ASSIST
Non-Adaptation	60.82 ± 0.22	54.37 ± 5.35	55.86 ± 3.81
Optimize Parameters	61.34 ± 0.33	54.35 ± 5.38	50.87 ± 0.32
Optimize Predictions	61.47 ± 0.35	43.10 ± 0.00	45.12 ± 0.18

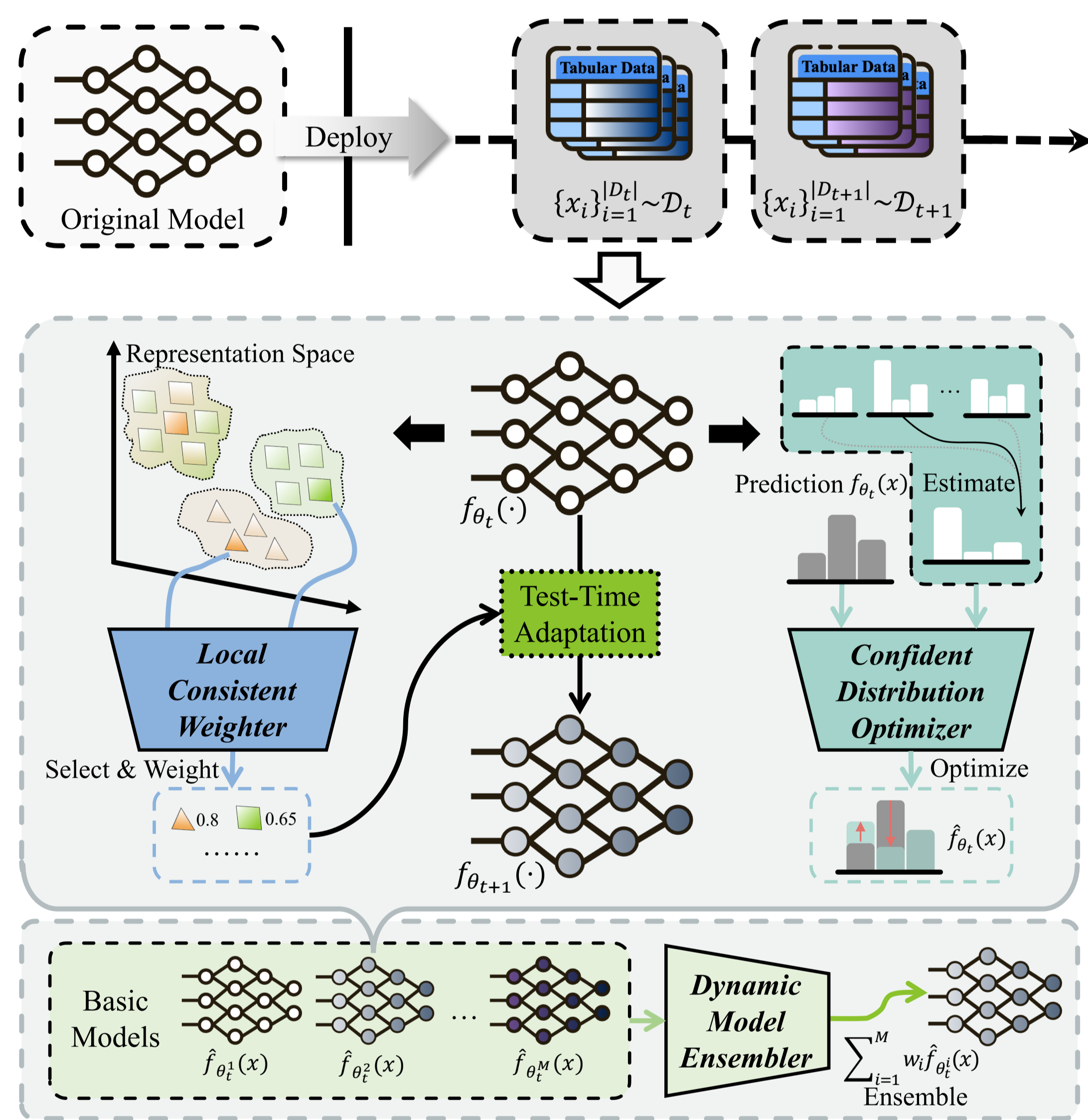
▲ Table 2: Performance of the non-adaptation baseline and two representative FTTA methods using an MLP backbone model. Degraded performance is underlined.

Observation 4: Adaptation is sensitive to both tasks and models for tabular data.



▲ Figure 2: The performance of FTAT with different learning rates. The optimal value differs across backbones and tasks. A red star marks the highest point of each line.

FTAT Method



▲ Figure 3: The overall illustration of FTAT

Confident Distribution Optimizer

Optimize model predictions using the equation:

$$\hat{f}_{\theta_{t+1}}(x_k) = \frac{f_{\theta_{t+1}}(x_k) \circ \hat{P}_t}{P_0}$$

where P_0 is training label distribution and \hat{P}_t is estimated label distribution at timestamp t .

1. Biased estimation \hat{P}_t is computed $\hat{P}_t = \frac{\sum_{x_i=1}^{D_t} \mathbb{I}[\text{Entropy}(\hat{f}_{\theta_t}(x_i)) < \epsilon] \cdot \hat{f}_{\theta_t}(x_i)}{\sum_{x_i=1}^{D_t} \mathbb{I}[\text{Entropy}(\hat{f}_{\theta_t}(x_i)) < \epsilon]}$ on low-entropy samples:
2. Debias using covariate matrix, $\sum_{x_i=1}^{D_t} \mathbb{I}[\arg \max_j \hat{f}_{\theta_t}(x_i)_j = k] \cdot \hat{f}_{\theta_t}(x_i)$ where the k -th row is defined as: $\sum_{x_i=1}^{D_t} \mathbb{I}[\arg \max_j \hat{f}_{\theta_t}(x_i)_j = k]$
3. Estimated \hat{P}_t is tracked smoothly throughout the whole time: $\hat{P}_t = \text{Norm}(\hat{P}_{t-1} - \alpha \cdot \hat{C}_t^{-1} \hat{P}_t)$

Local Consistent Weighter

The test-time adaptation is optimized based on the entropy minimization objective with sample weight $\mathcal{W}(x_k, D_t, \theta_t)$:

$$\left[\max \hat{f}_{\theta_t}(x_k) - \min \hat{f}_{\theta_t}(x_k) \right] \cdot \mathcal{I}(x_k, D_t, \theta_t)$$

where A_1 indicates the margin of prediction and A_2 is the consistency indicator of prediction.

1. The neighbor set of the sample is defined as follows: $N(x_k, D_t) = \{x | \text{Dist}(x, x_k) < \text{Dist}_t, x \in D_t\}$
2. The threshold is the average distance between samples in the data batch at each timestamp t : $\text{Dist}_t = \frac{2}{|D_t|(|D_t|-1)} \sum_{i=1}^{|D_t|} \sum_{j=i+1}^{|D_t|} \text{Dist}(x_i, x_j)$
3. The consistency is defined as the consistency between its prediction and the predictions of its neighbors: $\mathcal{I}(x_k, D_t, \theta_t) = \begin{cases} 1, & \left| \hat{f}_{\theta_t}(x_k) - \frac{\sum_{x \in N(x_k, D_t)} \hat{f}_{\theta_t}(x)}{|N(x_k, D_t)|} \right| < \beta \\ 0, & \text{Otherwise.} \end{cases}$

Dynamic Model Ensemble

The final predictions are defined as the dynamic ensemble of adapted models using different learning rates:

$$\sum_{i=1}^M w_i \cdot \hat{f}_{\theta_t^i}(x)$$

1. The weight is proportional to the loss $R_t^i(D_t)$: $w_i \propto 1 - R_t^i(D_t)$
2. The weight satisfied the constraint: $\sum_{i=1}^M w_i = 1$

Experiments

Method	MLP			TabTransformer			FT-Transformer		
	Acc.	Balanced Acc.	F1	Acc.	Balanced Acc.	F1	Acc.	Balanced Acc.	F1
Non-Adaptation	62.45	64.61	60.59	60.86	63.08	58.32	59.69	62.49	54.29
TENT	58.43	61.63	50.97	58.32	61.40	51.73	55.41	55.41	36.34
EATA	61.43	63.69	60.11	60.33	62.36	60.09	56.04	58.95	44.62
LAME	59.48	62.32	58.47	59.15	62.50	58.39	58.90	61.98	51.86
CoTTA	61.59	63.78	60.57	60.37	62.82	59.75	59.64	62.43	53.41
ODS	59.18	62.15	57.83	59.22	62.02	58.46	59.05	61.70	51.41
SAR	61.16	63.49	59.18	60.30	62.77	59.72	59.27	62.11	57.04
FTAT	66.77	64.96	72.00	66.14	64.40	69.03	64.01	62.54	69.56

▲ Table 3: Average performance of each method on 6 datasets using different backbones.

Method	HELOC			ANES			Health Ins.		
	Acc.	Balanced Acc.	F1	Acc.	Balanced Acc.	F1	Acc.	Balanced Acc.	F1
Non-Adaptation	54.37 ± 5.35	58.25 ± 3.56	40.02 ± 16.8	79.11 ± 0.31	75.66 ± 0.46	84.24 ± 0.16	65.79 ± 0.63	70.68 ± 0.44	66.21 ± 0.90
TENT	54.35 ± 5.38	58.24 ± 3.58	39.95 ± 16.9	78.07 ± 0.35	74.09 ± 0.65	83.76 ± 0.13	64.30 ± 0.70	69.79 ± 0.47	63.87 ± 1.06
EATA	54.37 ± 5.35	58.25 ± 3.56	40.02 ± 16.8	78.13 ± 0.30	74.20 ± 0.59	83.79 ± 0.10	65.78 ± 0.63	70.68 ± 0.44	66.21 ± 0.90
LAME	43.10 ± 0.00	50.00 ± 0.00	30.10 ± 0.00	63.50 ± 0.00	54.60 ± 0.00	46.80 ± 0.00	63.44 ± 1.69	69.14 ± 1.09	62.61 ± 2.69
CoTTA	54.36 ± 5.35	58.25 ± 3.56	40.03 ± 16.8	78.13 ± 0.30	74.20 ± 0.59	83.79 ± 0.10	65.79 ± 0.63	70.68 ± 0.44	66.21 ± 0.90
ODS	43.10 ± 0.00	50.00 ± 0.00	30.10 ± 0.00	63.50 ± 0.00	54.60 ± 0.00	46.80 ± 0.00	63.45 ± 1.68	69.14 ± 1.07	62.62 ± 2.68
SAR	52.32 ± 6.05	56.74 ± 3.99	33.16 ± 0.90	78.13 ± 0.30	74.20 ± 0.59	83.79 ± 0.10	65.79 ± 0.63	70.68 ± 0.44	66.21 ± 0.90
FTAT	64.09 ± 1.14	63.64 ± 0.93	67.80 ± 2.71	80.09 ± 0.23	79.12 ± 0.20	83.42 ± 0.25	72.42 ± 0.20	65.30 ± 0.15	80.83 ± 0.23

▲ Table 4: Performance of FTAT approach and comparison methods on 6 datasets using MLP.

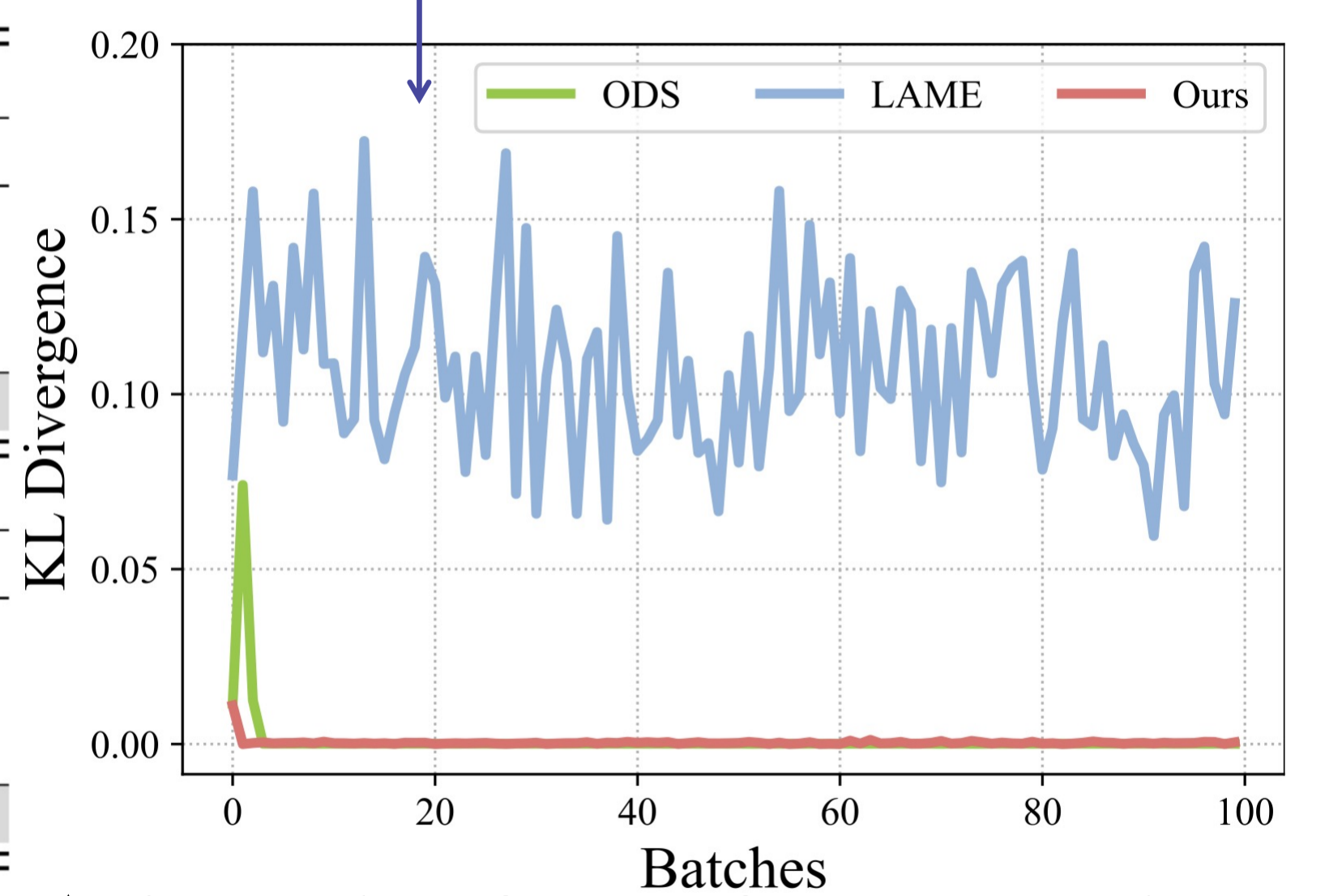
RQ1: Does FTAT method outperform comparisons across different datasets using different backbones?

RQ2: Does each component contribute to performance in the FTAT method?

RQ3: Does the FTAT method accurately estimate the label distributions?

DIABETE			
Method	Acc.	Balanced Acc.	F1
Non-Adaptation	60.81 ± 0.21	60.59 ± 0.23	51.18 ± 1.69
FTAT w/o CDO	60.85 ± 0.22	60.61 ± 0.24	51.26 ± 1.69
FTAT w/o LCW	61.43 ± 0.16	61.28 ± 0.20	55.61 ± 1.67
FTAT	61.66 ± 0.30	61.54 ± 0.28	59.27 ± 0.96

▲ Table 5: Ablation study. The performance of the FTAT approach using MLP backbone when removing different components.



▲ Figure 4: The performance of LAME, ODS, and FTAT in estimating label distribution.

- ✓ If you are interested in this paper, feel free to contact Zhi Zhou and Kun-Yang Yu (zhouz@lamda.nju.edu.cn, yuky@lamda.nju.edu.cn).
- ✓ To obtain more details of our paper, please visit the project homepage (<https://wnjxyk.github.io/FTTA>).
- ✓ This research was supported by the NSFC (Grant No. 624B2068, 62176118, and 62306133), the Key Program of Jiangsu Science Foundation (BK20243012), and the Fundamental Research Funds for the Central Universities (022114380023).

