# Towards Fair Knowledge Distillation using Student Feedback

Abhinav Java[1], Surgan Jandial[1], Chirag Agarwal[2]

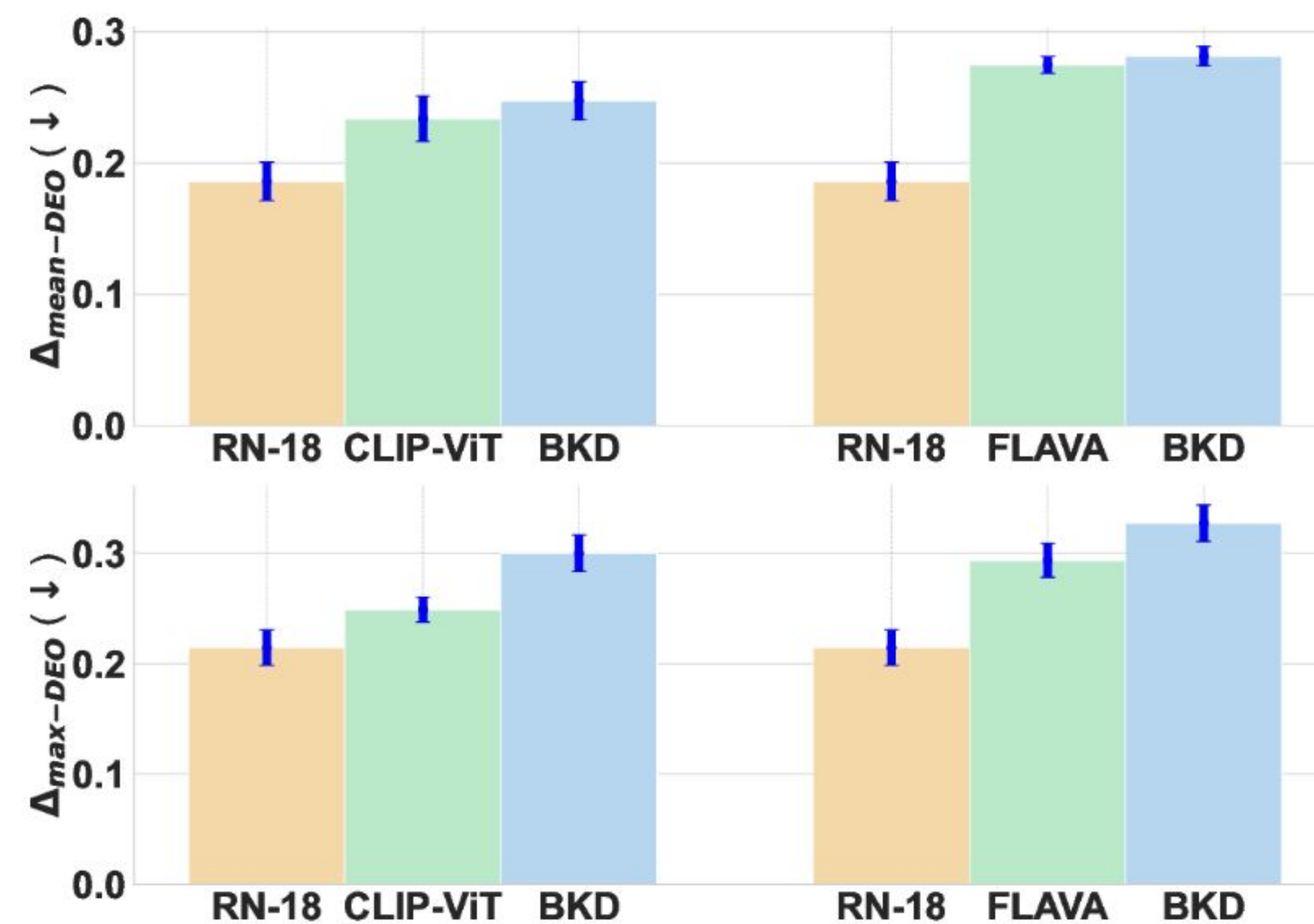[1]MDSR Labs, Adobe, India        [2]Harvard University

## Fairness & Knowledge Distillation: How can we incorporate student feedback to perform Bias Aware Distillation?

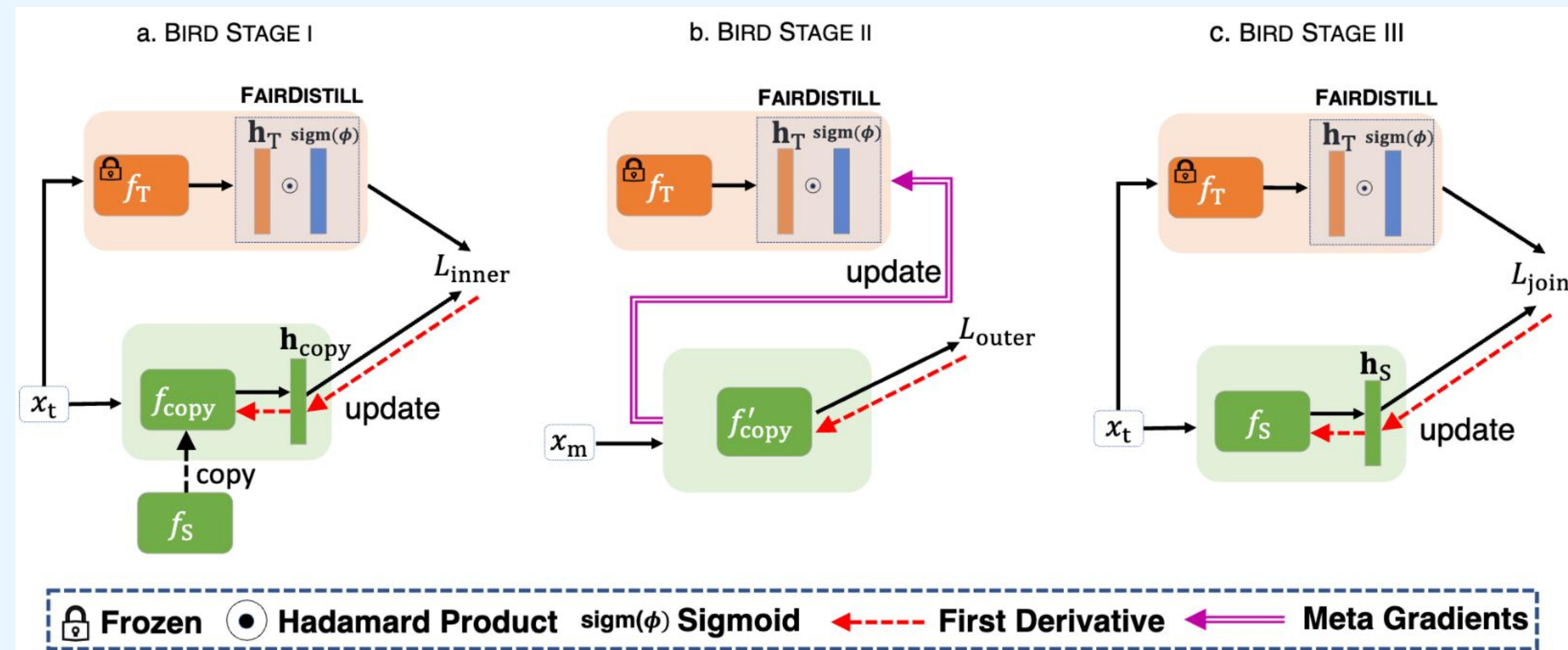### Student mimics the fairness properties of the teacher

Fairness scores for baseline teacher (CLIP-ViT/B32, FLAVA), baseline student (ResNet-18), and distilled student models using base KD (BKD)

### Overview of our BIRD 🐦 framework

**BIRD** learns bias-aware representations from the teacher $f_T$ by training the *FAIRDISTILL* operator using a meta-learning framework:
a. In Stage I, **BIRD** updates a copy of the student model with $L_{inner}$,
b. in Stage II, the updated model $f'_C$ is used to train $\phi$ with bias-feedback information in the form of meta-gradients from $L_{outer}$, and
c. in Stage III, the student model $f_S$ is distills unbiased representations using *FAIRDISTILL* (from Stage II).

### Problem Formulation

Given a dataset $\mathbf{D}^{train}$ and a biased teacher model $f_T$ optimized for predictive performance on $\mathbf{D}^{train}$, we aim to learn a student model $f_S$ whose representations do not reflect any undesirable discriminatory biases (i.e., they are fair) and achieve high predictive performance (i.e., they are accurate).

| Model | Method | AUROC ($\uparrow$) | $\Delta_{\text{mean-DEO}}(\downarrow)$ | $\Delta_{\text{max-DEO}}(\downarrow)$ |
|---|---|---|---|---|
| Flava | Baseline | $84.43\pm0.12$ | $27.48\pm0.64$ | $29.37\pm1.53$ |
| | BKD | $84.42\pm0.11$ | $27.39\pm0.58$ | $29.36\pm1.41$ |
| | FitNet | $84.47\pm0.10$ | $26.59\pm0.62$ | $28.56\pm0.68$ |
| | AD | $84.35\pm0.05$ | $10.54\pm0.80$ | $12.93\pm0.79$ |
| | MFD | $84.45\pm0.11$ | $26.64\pm0.62$ | $28.63\pm0.68$ |
| | **BIRD** | $85.48\pm0.02$ | $\mathbf{2.53}\pm0.17$ | $\mathbf{4.12}\pm0.59$ |
| CLIP-ViT/32 | Baseline | $87.01\pm0.26$ | $23.38\pm1.72$ | $24.91\pm1.15$ |
| | BKD | $87.07\pm0.26$ | $23.26\pm1.67$ | $24.62\pm1.14$ |
| | FitNet | $87.17\pm0.13$ | $22.84\pm1.03$ | $24.17\pm1.22$ |
| | AD | $88.20\pm0.17$ | $17.02\pm1.03$ | $17.82\pm0.97$ |
| | MFD | $87.22\pm0.11$ | $21.99\pm0.70$ | $23.70\pm1.58$ |
| | **BIRD** | $88.55\pm0.03$ | $\mathbf{3.44}\pm0.92$ | $\mathbf{5.19}\pm1.06$ |
| CLIP-R50 | Baseline | $87.72\pm0.06$ | $21.11\pm0.30$ | $21.97\pm0.41$ |
| | BKD | $87.72\pm0.06$ | $21.10\pm0.40$ | $22.07\pm0.41$ |
| | FitNet | $87.54\pm0.14$ | $22.01\pm1.05$ | $23.30\pm1.15$ |
| | AD | $88.51\pm0.02$ | $5.33\pm0.19$ | $7.93\pm0.22$ |
| | MFD | $87.49\pm0.12$ | $22.56\pm0.56$ | $23.52\pm0.33$ |
| | **BIRD** | $87.93\pm0.01$ | $\mathbf{2.65}\pm0.29$ | $\mathbf{4.49}\pm0.48$ |

### BIRD improves fairness of knowledge distillation

Shown is the comparative performance of **BIRD** on CelebA Dataset (Left) for three foundation models and on CIFAR10-S dataset (Bottom) for ResNet18➜ResNet18. Note that all results indicate avg. performance across five independent runs. Arrows ($\uparrow\downarrow$) indicate the direction of desired performance. *BIRD* retains the predictive power (AUROC) of the baseline while improving fairness criterion (shaded).

| Method | AUROC ($\uparrow$) | F1-score ($\uparrow$) | $\Delta_{\text{mean-DEO}}(\downarrow)$ | $\Delta_{\text{max-DEO}}(\downarrow)$ |
|---|---|---|---|---|
| Baseline | $98.91\pm0.02$ | $88.34\pm0.17$ | $26.26\pm0.70$ | $47.94\pm1.94$ |
| BKD | $98.95\pm0.02$ | $88.90\pm0.13$ | $25.30\pm0.63$ | $46.92\pm2.16$ |
| FitNet | $98.89\pm0.01$ | $88.15\pm0.08$ | $26.55\pm0.66$ | $48.86\pm1.85$ |
| AT | $98.99\pm0.02$ | $88.95\pm0.12$ | $25.16\pm0.33$ | $46.08\pm2.27$ |
| AD | $98.44\pm0.11$ | $85.98\pm0.43$ | $\mathbf{16.20}\pm1.18$ | $\mathbf{31.94}\pm3.89$ |
| MFD | $98.93\pm0.03$ | $88.32\pm0.10$ | $27.27\pm0.34$ | $49.16\pm1.62$ |
| **BIRD** | $\mathbf{99.12}\pm0.02$ | $\mathbf{89.45}\pm0.14$ | $19.77\pm0.37$ | $38.26\pm1.73$ |

java.abhinav99@gmail.com