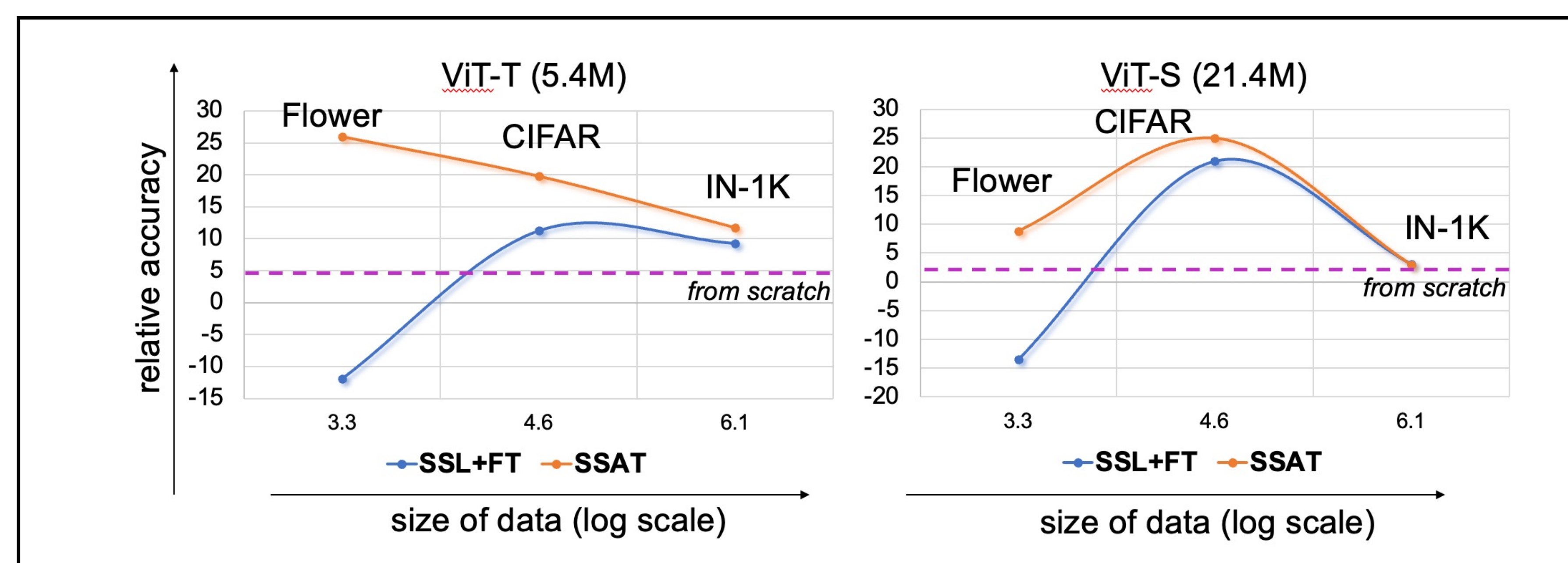




UNC Charlotte, DTU, BITS Pilani Hyderabad, IIIT Guwahati, Stony Brook University

## Motivation

- Training transformers require **large-scale data**
- **What about small-scale data distribution?**
  - Use pre-trained models
- **What about domains with Limited data?**
  - For example – **Medical data**



Relative classification accuracy on three datasets with different sizes: (i) *Oxford Flower* (2K samples), (ii) *CIFAR* (50K samples), and (iii) *ImageNet-1K* (IN-1K, 1.2M samples). **Self-supervised Auxiliary Task** (SSAT) consistently outperforms others on all three datasets with two backbones. On the other hand, given the same Self-supervised Learning (SSL) method, SSL+ Fine-tuning (FT) achieves a compromised performance than SSAT, especially on the tiny *Oxford Flower* dataset (even worse than training from scratch).



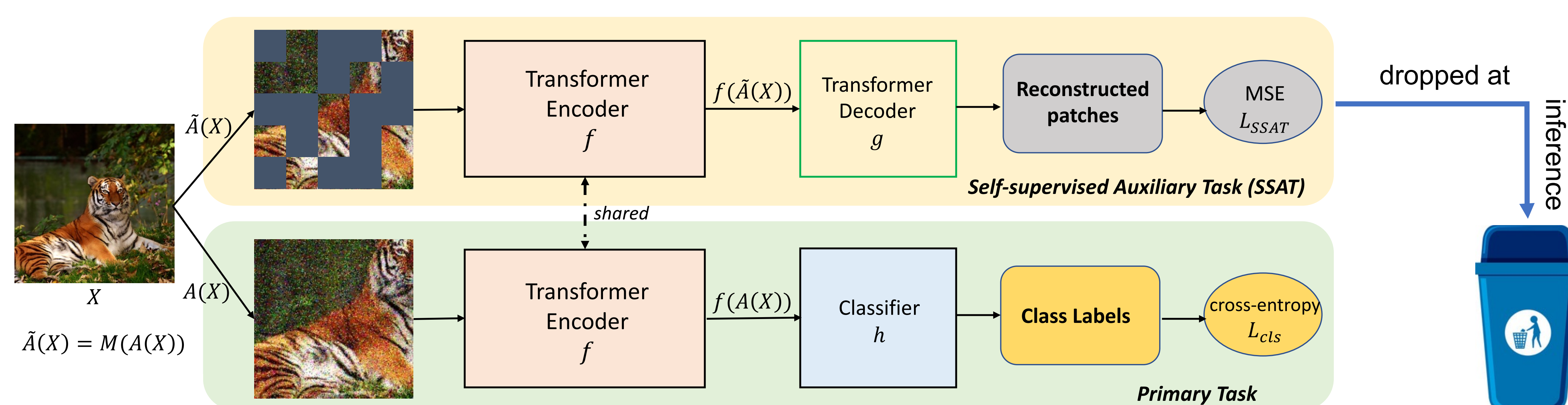
## Straightforward implementation

### Training SSAT

- Inflate the input batch of images with two set of Augmentations
- Two streams -> **Primary + Auxiliary**
- **Primary** stream – Performing classification of the image
- **Auxiliary** Stream – Reconstructing the input image from 25% of the input
- $Total Loss = \lambda Loss_{Primary} + (1 - \lambda) Loss_{Auxiliary}$

## Self-supervised Auxiliary Task (SSAT)

### Which SSL task?



SSAT improves ViT performance while reducing carbon footprints

## Experiments



SSAT trained with outperforms on 12 Computer Vision tasks

Top-1 classification accuracy (%) of different ViT variants with and without SSAT on *CIFAR-10*, *CIFAR-100*, *Flowers102*, and *SVHN* datasets. All models were trained for 100 epochs.

Does SSAT promotes overfitting?

Method	IN-1K	CIFAR-100-p	IN-1K-p
ViT-T	65.0	25.1	48.3
+SSAT	<b>72.7</b>	<b>37.6</b>	<b>59.6</b>
ViT-S	74.2	22.5	62.7
+SSAT	<b>76.4</b>	<b>43.9</b>	<b>64.5</b>

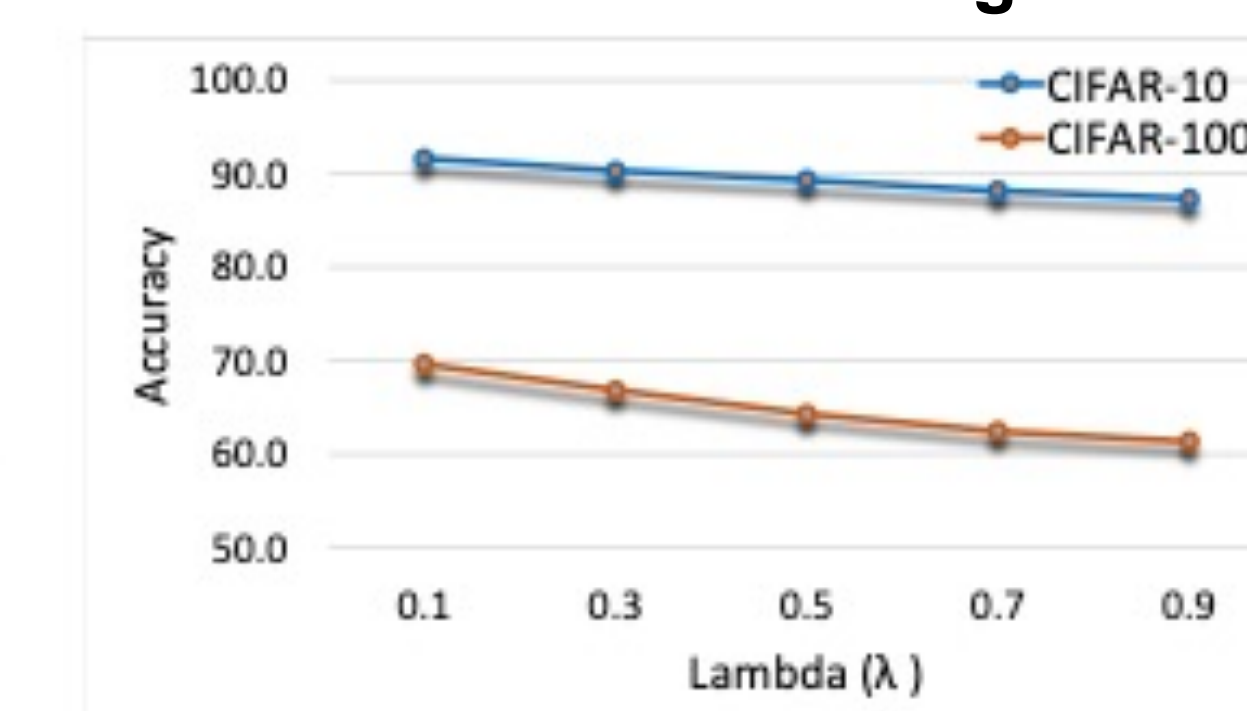
Comparison of SSAT with SSL+FT

Method	GFLOPs	CIFAR-10	CIFAR-100	IN-1K	Train time	Kg CO <sub>2</sub> eq.
Scratch	1.26	79.47	55.11	65.0	60	5.96
SSL+FT	0.43+1.26	85.33	60.43	70.09	55	5.46
SSL+FT	0.43+1.26	86.48	63.28	71.1	82	8.15
SSL+FT	0.43+1.26	88.72	67.53	74.07	104	10.33
Ours	1.67	<b>91.65</b>	<b>69.64</b>	<b>72.69</b>	78	7.55

Appropriate SSL for SSAT

SSAT (SSL)	CIFAR-10	CIFAR-100
SimCLR (Chen et al., ICLR 20)	55.21	36.49
DINO (Caron et al., CVPR 20)	80.07	60.6
MAE (He et al., CVPR 20)	<b>91.65</b>	<b>69.64</b>

Ablation for loss scaling factor



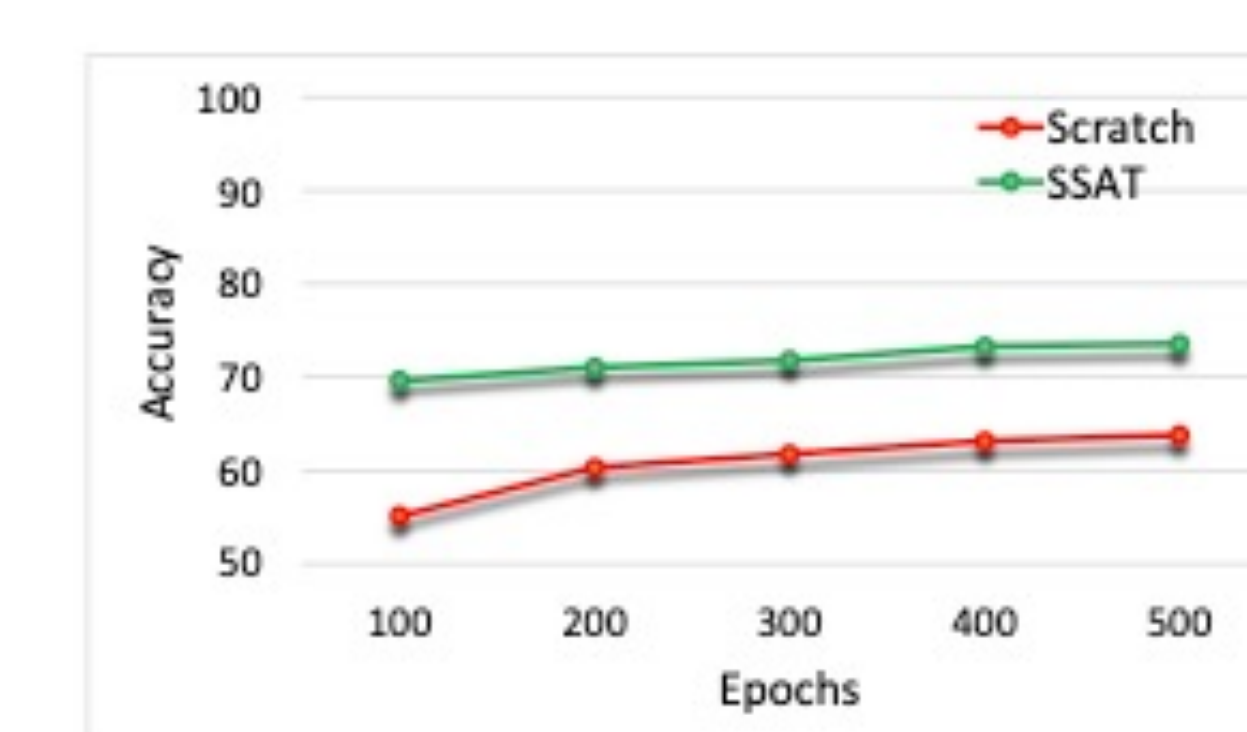
Method	# params. (M)	CIFAR-10	CIFAR-100	Flowers102	SVHN
ViT-T (Touvron et al., ICML 21)	5.4	79.47	55.11	45.41	92.04
+SSAT	5.8	<b>91.65 (+12.18)</b>	<b>69.64 (+14.53)</b>	<b>57.2 (+11.79)</b>	<b>97.52 (+5.48)</b>
ViT-S (Touvron et al., ICML 21)	21.4	79.93	54.08	56.17	94.45
+SSAT	21.8	<b>94.05 (+14.12)</b>	<b>73.37 (+19.29)</b>	<b>61.15 (+4.98)</b>	<b>97.87 (+3.42)</b>
CVT-13 (Wu et al., ICCV 21)	20	89.02	73.50	54.29	91.47
+SSAT	20.3	<b>95.93 (+6.91)</b>	<b>75.16 (+1.66)</b>	<b>68.82 (+14.53)</b>	<b>97 (+5.53)</b>
Swin-T (Liu et al., CVPR 21)	29	59.47	53.28	34.51	71.60
+SSAT	29.3	<b>83.12 (+23.65)</b>	<b>60.68 (+7.4)</b>	<b>54.72 (+20.21)</b>	<b>85.83 (+14.23)</b>
ResNet-50 (He et al., CVPR 16)	25.6	91.78	72.80	46.92	96.45

Performance of SSAT based ViT on Medical and DomainNet datasets

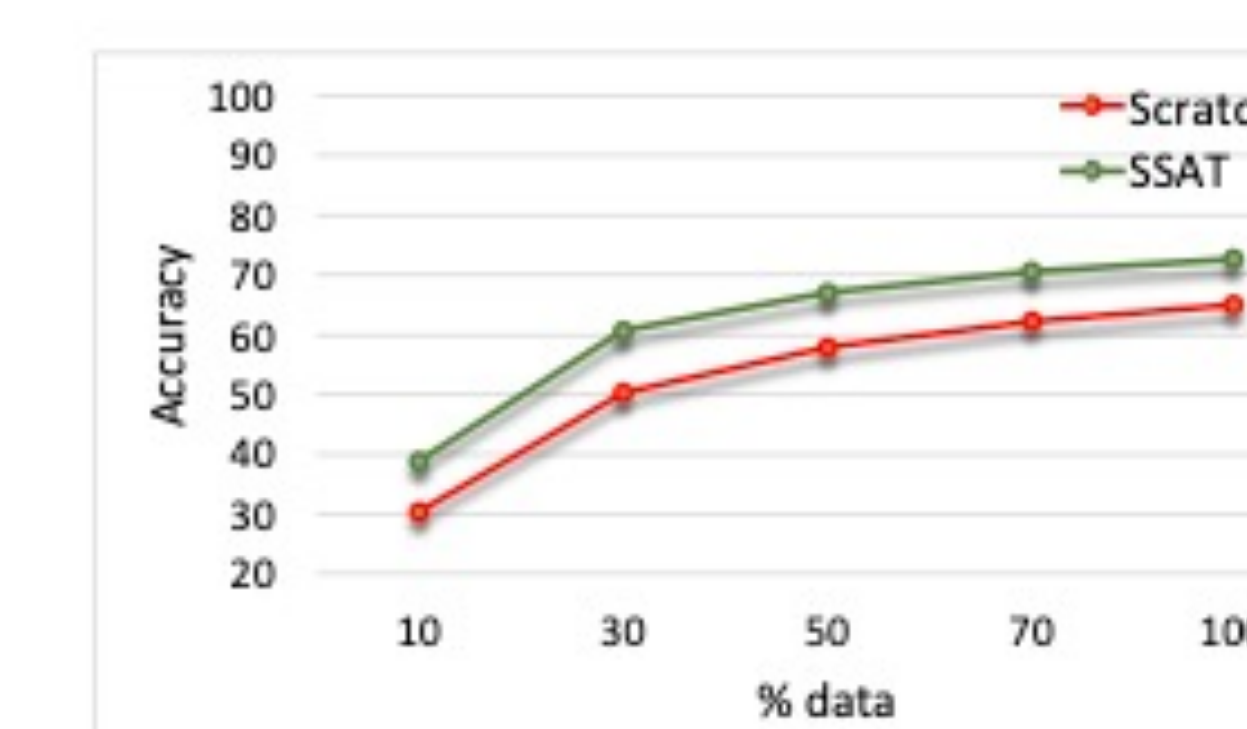
Method	Chaoyang	PMNIST
Scratch	77.37	90.22
ViT-T	78.78	91.99
IN-1K pretrained + FT	<b>82.52</b>	<b>93.11</b>
Scratch + SSAT	<b>82.52</b>	<b>93.11</b>
Scratch	80.04	91.19
ViT-S	80.18	92.63
IN-1K pretrained + FT	<b>81.25</b>	<b>93.27</b>
Scratch + SSAT	<b>81.25</b>	<b>93.27</b>

Method	ClipArt	Infograph	Sketch
ViT-T	29.66	11.77	18.95
+SSAT	<b>47.95</b>	<b>16.37</b>	<b>46.22</b>
CVT-13	60.34	19.39	56.98
+L <sub>drloc</sub> (Liu et al., NeurIPS 21)	<b>60.64</b>	<b>20.05</b>	<b>57.56</b>
+SSAT	<b>60.66</b>	<b>21.27</b>	<b>57.71</b>

ViT vs. ViT+SSAT for longer training epochs



ViT vs. ViT+SSAT for different fraction of data.



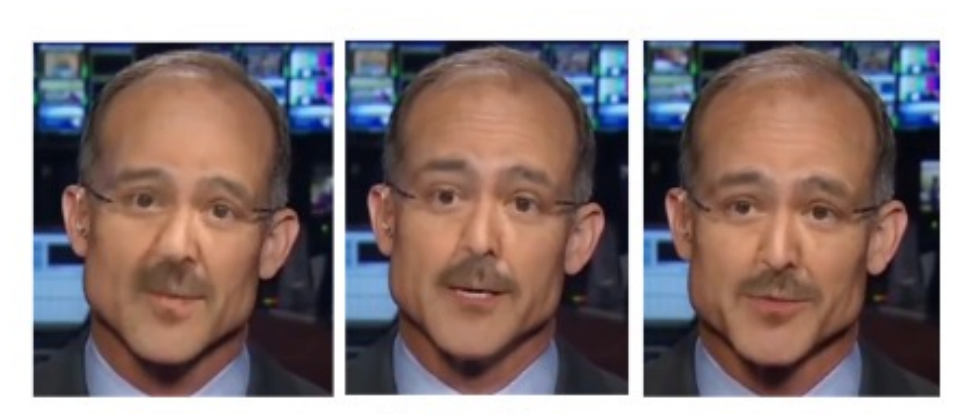
## SSAT for Video DeepFake Detection



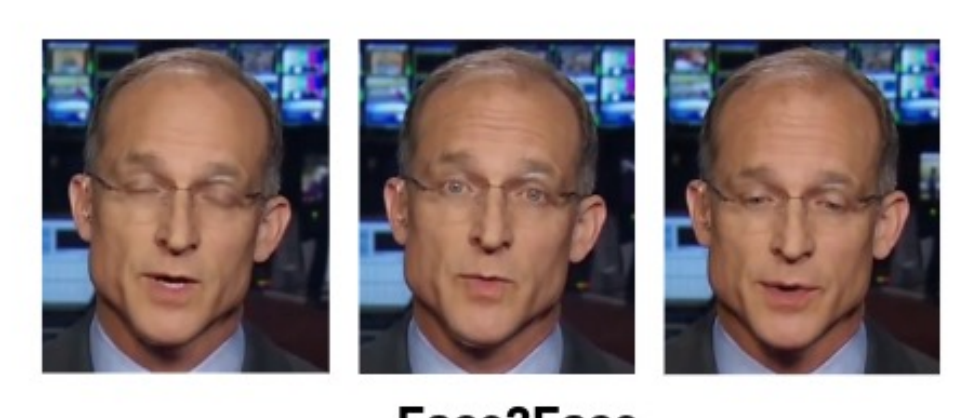
### Manipulations



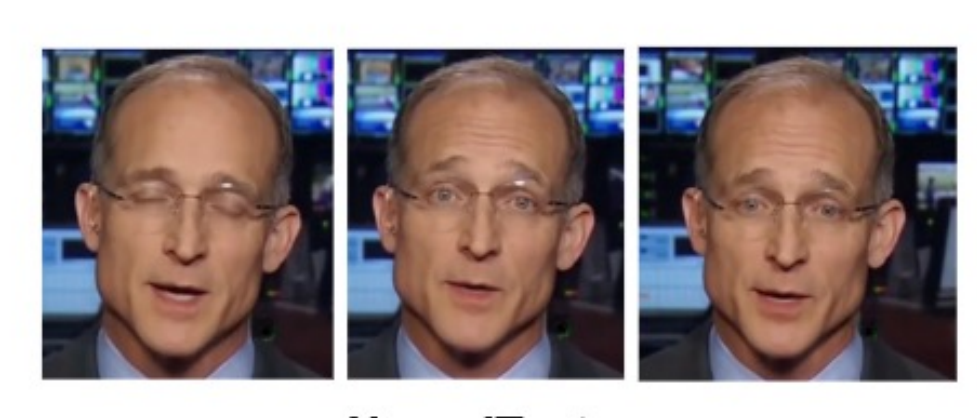
### Deepfakes



### FaceSwap



### Face2Face



### NeuralTextures

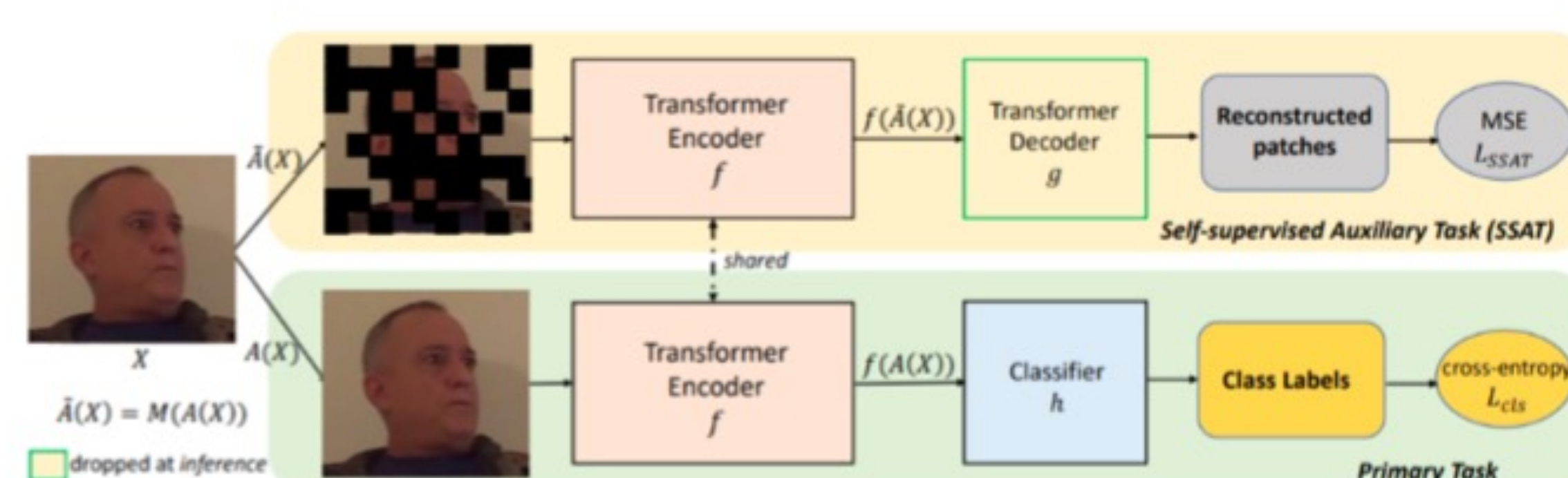
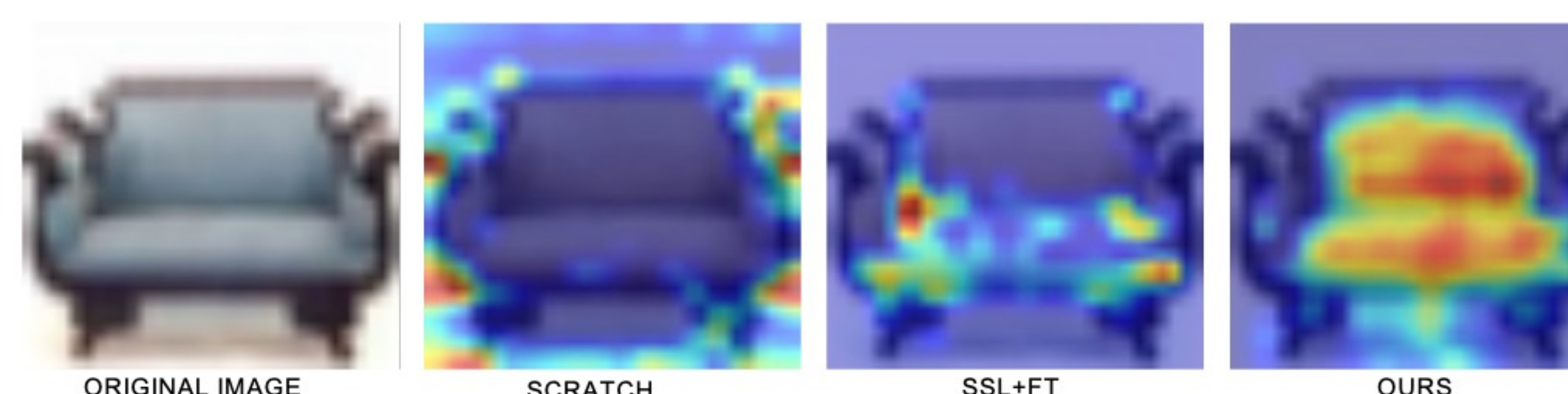


Figure 11. Mask Autoencoder as a Self-Supervised Auxiliary Task for deepfake detection.

Cross training evaluation and zero-shot transfer results of DeepFake detection on *FaceForensics++* with SSAT.

Method	cross-training evaluation				zero-shot transfer			
	Deepfakes	Face2Face	FaceSwap	NeuralTextures	Deepfakes	Face2Face	FaceSwap	NeuralTextures
Scratch	84.48	79.21	56.63	<b>82.08</b>	-	-	-	-
Cross-efficient-vit (Coccomini et al., ICIAP 22)	82.67	69.89	79.93	64.87	-	-	-	-
DFDC winner (Seferbekov, DFDC 20 winner)	96.43	73.93	86.07	58.57	88.57	57.50	80.36	54.64
VideoMAE SSL (0.95)	82.67	64.16	58.42	63.44	86.28	49.82	69.18	51.97
VideoMAE SSL (0.75)	78.34	65.59	57.35	61.65	82.67	48.39	65.23	51.97
VideoMAE (0.95) + SSAT	92.42	79.21	89.61	81.36	<b>92.42</b>	<b>61.65</b>	<b>92.83</b>	<b>62.37</b>
VideoMAE (0.75) + SSAT	<b>96.75</b>	<b>80.65</b>	<b>91.40</b>	72.76	87.73	60.57	88.17	61.65

## GradCAM visualizations



State-of-the-art comparison

Method	# enc. params.	epochs	CIFAR-10	CIFAR-100
CVT-13+L <sub>drloc</sub> (Liu et al., NeurIPS 21)	-	-	90.30	74.51
CVT-13+ SSAT	20M	100	<b>95.93</b>	<b>75.16</b>
ViT (scratch)	-	-	93.58	73.81
SL-ViT (Lee et al., WACV 23)	2.8M	300	94.53	76.92
ViT <sup>†</sup> (SSL+FT) (Gani et al., BMVC 22)	-	-	94.2	76.08
ViT + SSAT	-	-	<b>95.1</b>	<b>77.8</b>
DeiT-Ti+L <sub>guidance</sub> (Li et al., ECCV 22)	6M	300	-	78.15
DeiT-Ti+L <sub>guidance</sub> + SSAT	-	-	-	<b>79.46</b>

## Conclusion

- We presented a very **simple** method of using **SSL** to train ViTs on domains with **Limited data**
- **SSAT** – Jointly optimize the **primary task** with SSL as an **auxiliary task**
- Effectiveness validated on **10 image classification datasets + 2 video datasets**
- If you plan to use **ViTs on a small training distribution**, consider using **SSAT!**