Introduction

- Most deepfake detectors use ConvNets as feature extractors.
- The community is hesitant to use large ViTs:
 - Need extensive data to train.
 - Suboptimal generalization with small- and medium-size datasets.
- Recently, ViTs trained with SSL strageries demonstrated remarkable performances:
- Robust feature extractors \rightarrow Applicable for deepfake detection?
- Explicit semantic segmentation \rightarrow Explanation & localization?
- Strong transfer learning across various downstream tasks \rightarrow Better generalizability?

Methodology

Overview:

- Approach 1: Generalized form of the recent DFD approaches.
- Approach 2: Leveraging the strong transferability & segmentation of SSL ViTs.

Comparative study:

- ConvNets vs. ViTs.
- Supervised vs. unsupervised.
- Frozen backbone with adaptor vs. partial fine-tuning.



Patch embedder Pre-transformer block Transformer block 1 Transformer block n - k Transformer block n Post-transforme block

> -----Real / Fake _____

Final linear

Approach 2 Partial fine-tuning

Backbone	Are
EfficientNetV2 Large	Co
DeiT III L/16-LayerScale	Tra
EVA-02-CLIP-L/14	Tra
MĀE VĪT-L/16	Tra
DINO (various versions)	Tra
DINOv2 (various versions)	Tra



Model	k	FF++	FF++	FF++	FF++	FF++	FF++	DFD	DFD	Vid-
		Real	DF	F2F	FS	NT	FSh	Real	Fake	TIMIT
Approach 1										
EfficientNetV2 Large	4	60.90	84.70	83.40	79.50	75.30	83.50	61.00	85.50	59.75
DeiT III L/16-LayerScale	4	38.80	88.00	87.10	81.60	82.20	82.70	37.60	92.00	93.45
EVA-02-CLIP-L/14	4	50.40	97.40	90.80	93.20	80.80	90.30	44.50	97.20	83.30
MĀĒ VĪT-L/16	4	47.20	97.50	93.70	94.20	91.00	97.70	54.20	98.40	91.75
DINOv2 ViT-L/14-Reg	4	72.40	94.60	87.00	89.30	76.50	89.40	69.40	98.30	96.80
Approach 2										
EfficientNetV2 Large	1	65.00	94.60	89.00	89.30	84.10	90.50	22.90	91.50	99.20
DeiT III L/16-LayerScale	1	56.90	97.00	89.50	89.90	84.10	86.50	33.40	95.00	96.35
EVA-02-CLIP-L/14	1	53.50	97.50	92.70	92.40	83.30	92.40	62.40	98.80	96.55
MĀĒ VĪT-L/16	1	68.60	98.10	91.70	94.20	85.70	91.90	64.50	97.20	96.00
DINOv2 ViT-L/14-Reg	1	75.60	97.20	92.80	94.80	81.60	93.40	30.80	99.60	99.75
MAE ViT-L/16	15	80.50	99.20	94.10	93.50	90.30	95.80	77.60	99.70	99.70
DINOv2 ViT-L/14-Reg	11	85.10	98.60	94.30	95.40	89.70	97.20	67.50	99.60	99.85

Results on unse	eer	datase	t									
Model	k	Threshold	Real	Repaint P2	Repaint LDM	LaMa	Pluralistic	Acc.	TPR	TNR	EER	HTER
Approach 1			1									
EfficientNetV2 Large	4	0.6355	47.89	52.89	55.78	49.11	56.67	52.47	47.89	53.61	49.22	49.25
DeiT III L/16-LayerScale	4	0.9983	52.89	52.67	56.44	44.44	52.67	51.82	52.89	51.56	47.83	47.78
EVA-02-CLIP-L/14	4	0.0737	63.44	43.00	52.44	31.22	51.22	48.27	63.44	44.47	45.75	46.04
MAE VIT-L/16	4	0.9385	56.00	53.22	62.78	19.33	64.67	51.20	56.00	50.00	47.31	47.00
DINOv2 ViT-L/14-Reg	4	0.0759	60.44	53.00	66.89	43.11	70.89	58.87	60.44	58.47	40.67	40.54
Approach 2			1									
EfficientNetV2 Large	1	0.5479	63.00	50.22	53.89	65.78	64.89	59.56	63.00	58.69	39.58	39.15
DeiT III L/16-LayerScale	1	0.9999	56.00	58.56	69.56	39.56	69.56	58.64	56.00	59.31	42.56	42.35
EVA-02-CLIP-L/14	1	0.9999	45.44	71.11	83.44	12.22	82.11	58.87	45.44	62.22	45.20	46.17
MAE ViT-L/16	1	0.1769	65.44	47.11	58.22	13.67	71.00	51.09	65.44	47.50	44.22	43.53
DINOv2 ViT-L/14-Reg	1	0.9980	50.78	70.22	78.22	65.00	86.78	70.20	50.78	75.06	36.28	37.08
MĀĒ ViT-L/16	15	0.8948	69.89	50.78	68.89	22.56	76.44	57.71	69.89	54.67	37.56	37.72
DINOv2 ViT-L/14-Reg	11	0.7418	70.22	53.22	73.22	93.00	74.56	72.84	70.22	73.50	27.61	28.14

Inter-University Research Institute Corporation Research Organization of Information and Systems National Institute of Informatics





72.30

74.90

81.50

86.10

84.20

75.10

76.90

85.10

challenging problem!

hhuy@nii.ac.jp

71.30

62.70

77.30

92.40

99.70

99.74

99.92

142M Not used

Linear	(2 layers)
65.69	76.98
73.10	80.10
77.06	77.21
74.16	77.99
78.40	81.83

81.07

-P2 16, 571. Repaint-Latent 26,200 21,000 ---- DINOv2 ViT-L/14-Reg **---** MAE ViT-L/16 Deepfakes 1 Deepfakes 19.77 19.96 16.77



16.51

DINOv2 ViT-L/14-Reg Most of the behaviors are similar to

<u>augmentations → Less robus</u>t. 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 However, fgeneralizability is still a

2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20





Conclusion & Future Work

Conclusion:

- Backbone: SSL ViTs are better than supervised ConvNets
- An appropriate SSL training strategy is needed (DINOv2 is better than MAE).
- Partial fine-tuning the backbone is better than using adaptors with a frozen backbone.

Future work:

- Improve deepfake localization with the self-attention mechanism (without grond truth).
- Explore SSL on unlabeled deepfake datasets.
- Improve generalizability.

