

1 Contribution & Proposed Method

Performance Drop by Domain Shift

Training: $\mathcal{D}^{S_1}, \mathcal{D}^{S_2}, \dots, \mathcal{D}^{S_n}$
Inference: $\mathcal{D}^{T_1}, \mathcal{D}^{T_2}, \dots, \mathcal{D}^{T_n}$

Domain Generalization (DG)

- + No target data/knowledge required
- + Generalization across multiple unknown domains

Previous

DG Method for Segmentation

→ PREVIOUS approaches require complex methods or modules

→ Ours only relies on fine-tuning a strong vision-language pre-trained backbone

Our **key contributions** are:

- By **simple fine-tuning WITHOUT any additional modules or methods for segmentation AND detection** our approach shows similar or stronger performance than previous works.
- New SOTA performance on Cityscapes, ACDC and synthetic-to-real** benchmarks and proposing new real-to-real and synthetic-to-real evaluation scheme

2 Comparison of Pre-Trainings

Pre-Training					mIoU in %				
Method		Data	Sup.	Self-Sup.	\mathcal{D}^{CS}_{val}	\mathcal{D}^{BDD}_{val}	\mathcal{D}^{MV}_{val}	\mathcal{D}^{ACDC}_{val}	DG mean
\mathcal{D}^S : GTA5	SegFormer [105]	ImgNet-1K	✓	✗	46.6	45.6	50.1	36.4	44.7
	Supervised*	ImgNet-21K	✓	✗	49.3	47.0	52.2	43.5	48.0
	MoCov3* [13]	ImgNet-1K	✗	✓	49.7	46.2	52.4	39.1	46.9
	DeiT3* [94]	ImgNet-21K	✓	✓	53.7	52.6	59.3	49.0	53.7
	SAM* [53]	SA-1B	✓	✗	53.2	50.3	58.8	45.5	52.0
	LDM [77]	LAION-5B	✓	✗	49.2	-	-	-	-
	CLIP \blacklozenge [73]	WIT	✓	✗	53.2	49.8	57.1	45.0	51.2
	CLIP \circ [73]	WIT	✓	✗	55.6	52.5	59.9	51.5	54.9
EVA-02-CLIP \blacklozenge [90]	Merged-2B	✓	✗	55.2	51.3	57.4	47.3	52.8	
EVA-02-CLIP \circ [90]	Merged-2B	✓	✗	65.3	58.3	66.0	62.6	63.1	
\mathcal{D}^S : SYNTHIA	SegFormer [105]	ImgNet-1K	✓	✗	41.1	36.2	42.4	32.6	38.1
	Supervised*	ImgNet-21K	✓	✗	44.3	37.1	43.1	34.8	39.8
	MoCov3* [13]	ImgNet-1K	✗	✓	40.2	35.4	41.5	31.7	37.2
	DeiT3* [94]	ImgNet-21K	✓	✓	47.8	39.1	45.4	34.7	41.8
	SAM* [53]	SA-1B	✓	✗	51.6	40.4	50.1	40.1	45.6
	CLIP \blacklozenge [73]	WIT	✓	✗	46.1	41.8	45.8	35.1	42.2
	CLIP \circ [73]	WIT	✓	✗	51.1	44.7	50.6	40.7	46.8
	EVA-02-CLIP \blacklozenge [90]	Merged-2B	✓	✗	48.3	42.6	46.4	37.1	43.6
EVA-02-CLIP \circ [90]	Merged-2B	✓	✗	56.8	51.9	55.1	48.5	53.1	

+9.4%

+7.5%

→ Vision-Language pre-training performs significantly better than vision-only and benefits from large-scale text-image datasets

3 DG Benchmarks

DG Method	Enc. Params	mIoU (%) on				
		\mathcal{D}_{val}^{CS}	\mathcal{D}_{val}^{BDD}	\mathcal{D}_{val}^{MV}	DG mean	
Vision DG	Baseline	81.4M	46.6	45.6	50.1	47.4
	ReVT	81.4M	50.0	48.0	52.8	50.3
	SHADE	81.4M	53.3	48.2	55.0	52.2
	IBAFformer	81.4M	56.3	49.8	58.3	54.8
	CMFormer	197M	55.3	49.9	60.1	55.1
	HRDA	81.4M	57.4	49.1	61.2	55.9
Vision-Language DG	DGinStyle	81.4M	58.6	52.3	62.5	57.8
	DIDEX	81.4M	62.0	54.3	62.0	59.7
	PromptFormer	-	52.0	-	-	-
	CLOUDS	198M	60.2	57.4	67.0	61.5
	Rein	307M	65.3	60.5	64.9	63.6
	VLTseg	304M	65.3	58.3	66.0	63.2

SOTA-Level

→Competitive to SOTA and significantly better than other, more complex works

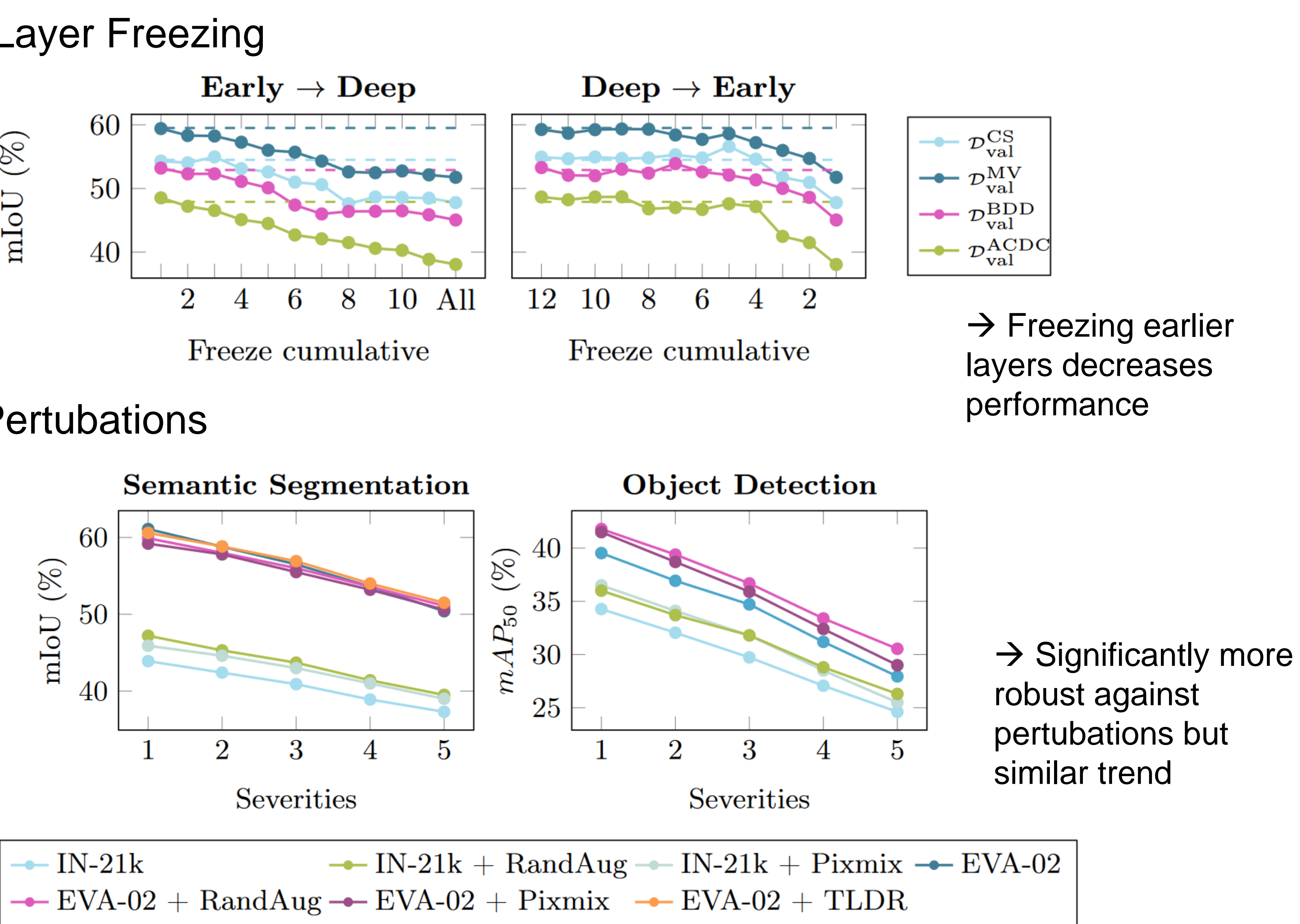
Method	Encoder	Init	Day Clear	Night Clear	Dusk Rainy	Night Rainy	Day Foggy	Average
S-DGOD [103]	R-101	IN1k	56.1	36.6	28.2	16.6	33.5	28.7
G-NAS [104]	R-101	IN1k	58.4	45.0	35.1	17.4	36.4	33.5
PDDOC [58]	R-101	IN1k	53.6	38.5	33.7	19.2	39.1	32.6
CLIP The GAP [97]	R-101	CLIP	51.3	36.9	32.3	18.7	38.5	31.6
PODA [26]	R-101	CLIP	-	43.4	40.2	20.5	44.4	37.1
VLTDet	R-101	CLIP	60.5	44.6	38.4	22.1	42.3	36.9
VLTDet	ViT-L-14	EVA02-CLIP	56.6	44.4	43.6	26.6	41.8	39.1

→ On par with SOTA for ResNet-101 for object detection, SOTA with EVA-backbone

→ Significantly enhanced real-to-real generalization

→ Enhanced in-domain performance

4 Analysis



5 Real-to-Real DG

Cityscapes (\mathcal{D}^{CS}_{train}) → Cityscapes (\mathcal{D}^{CS}_{test})					Cityscapes (\mathcal{D}^{CS}_{train}) → ACDC ($\mathcal{D}^{ACDC}_{test}$)		
Method	Rank	Params	Iter.	mIoU in %	Method	UDA/DG	mIoU in %
ViT-Adapter-L [14]	5	571M	80k	85.2	HRDA [41]	UDA	67.96
InverseForm [7]	3	-	-	85.6	CISS [80]	UDA	69.55
HS3 [6]	4	-	-	85.8	PromptFormer [32]	DG	62.0
InternImage [99]	2	1.2B	80k	86.1	Rein [102]	DG	77.56
VLTSeg	1	304M	40k	86.4	VLTSeg (1024 ²)	DG	77.91

→ New SOTA

→ New SOTA

→ +0.35%

Paper

Project Page