

موضوع



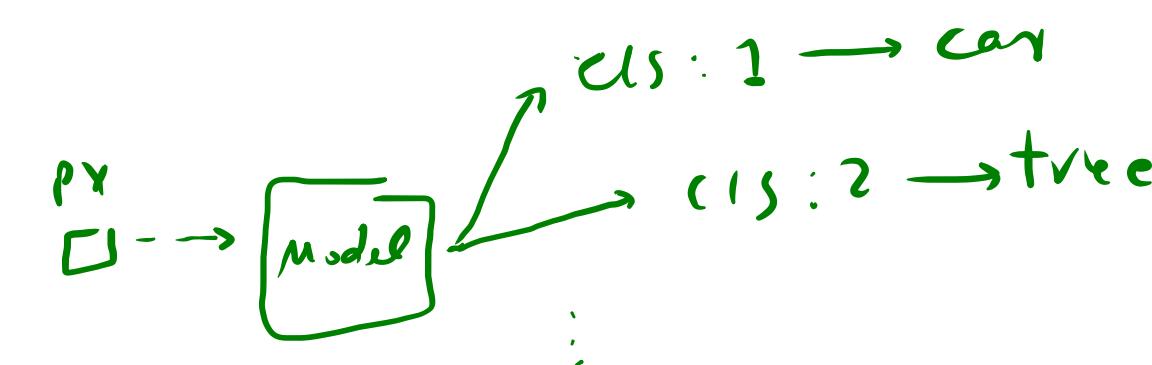
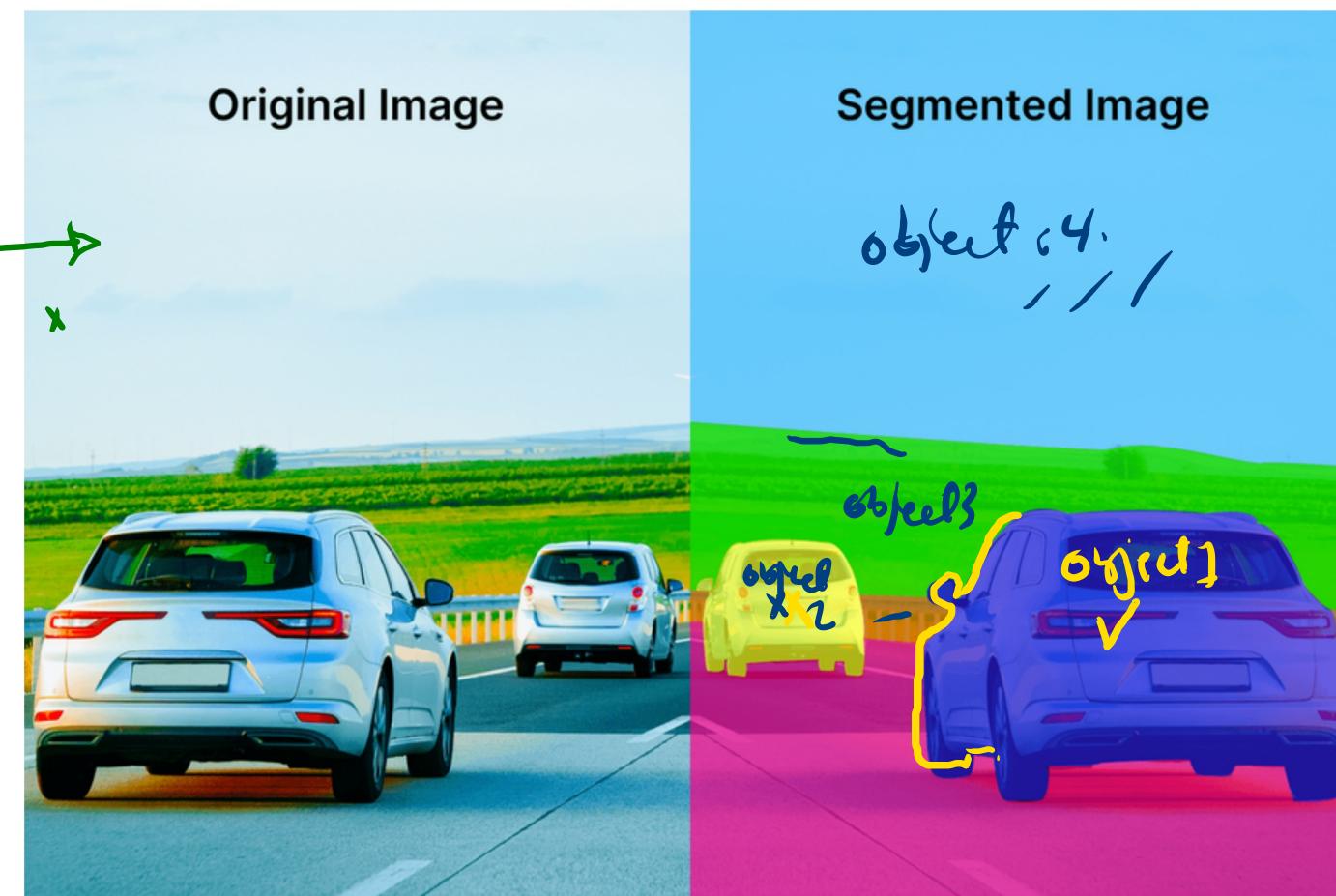
Image Segmentation

ناحیه بندی تصاویری

[pixel classification .]

ناحیه بندی تصویر یک تکنیک بینایی کامپیوتر است که یک تصویر دیجیتال را به گروههای مجزا از پیکسل‌ها - بخش‌های تصویر - تقسیم می‌کند تا تشخیص شیء و وظایف مرتبط را ممکن سازد.

<https://www.ibm.com/think/topics/image-segmentation>

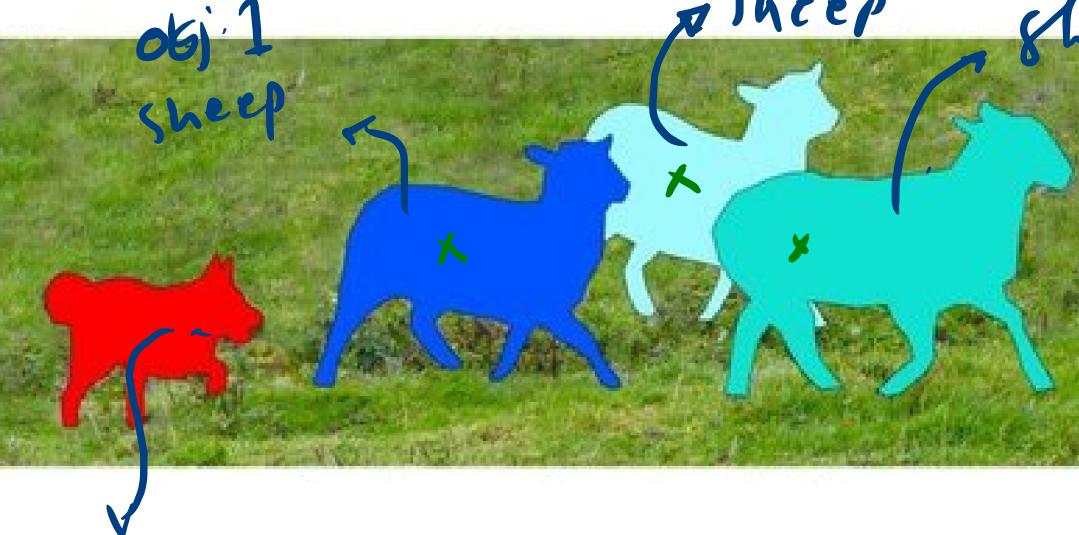


Object Detection

original



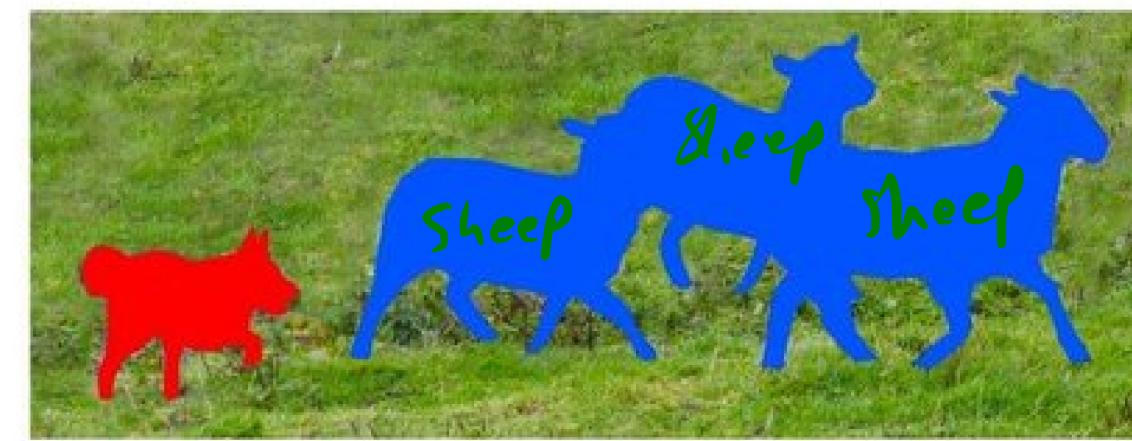
{
cls
object

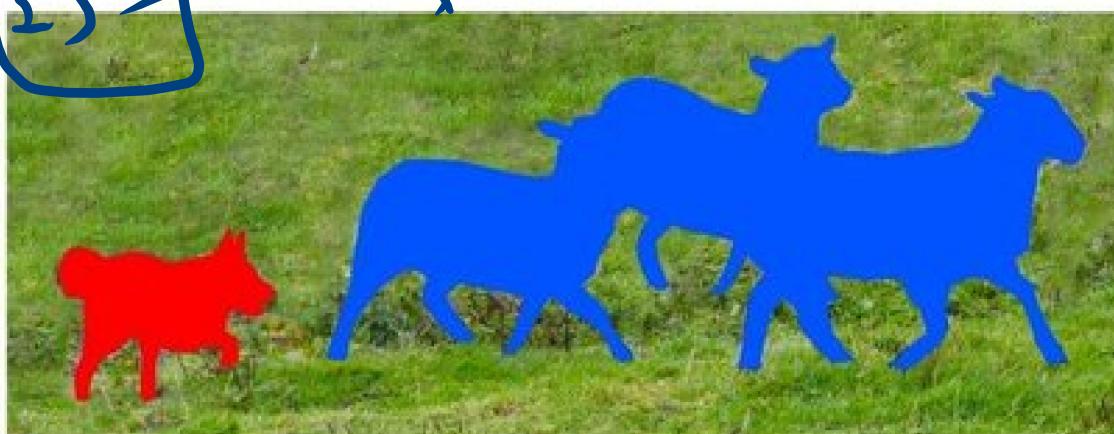
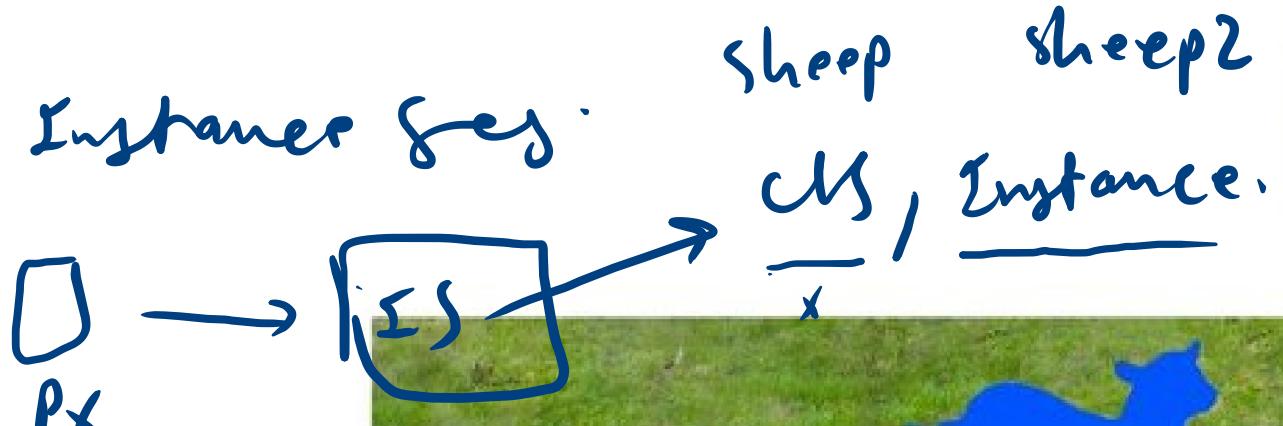
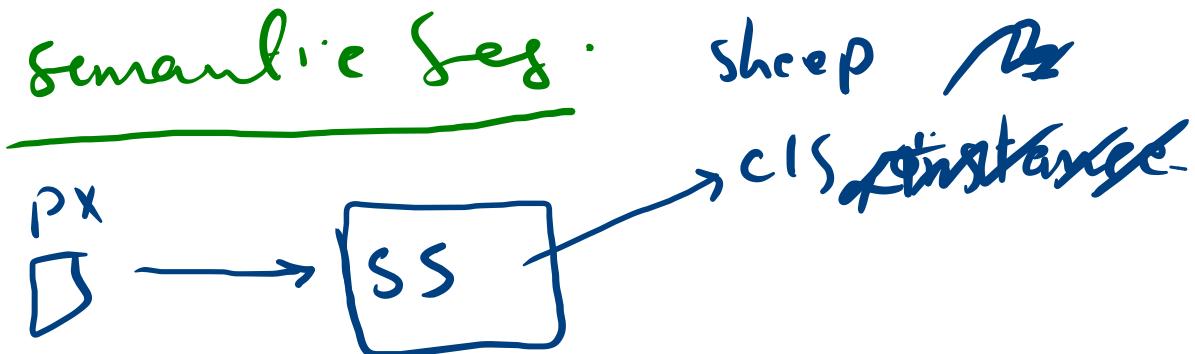


deg
object 1

تفاوت سگمنتیشن در این دو تصویر چیست؟

2 seg.





Semantic Segmentation

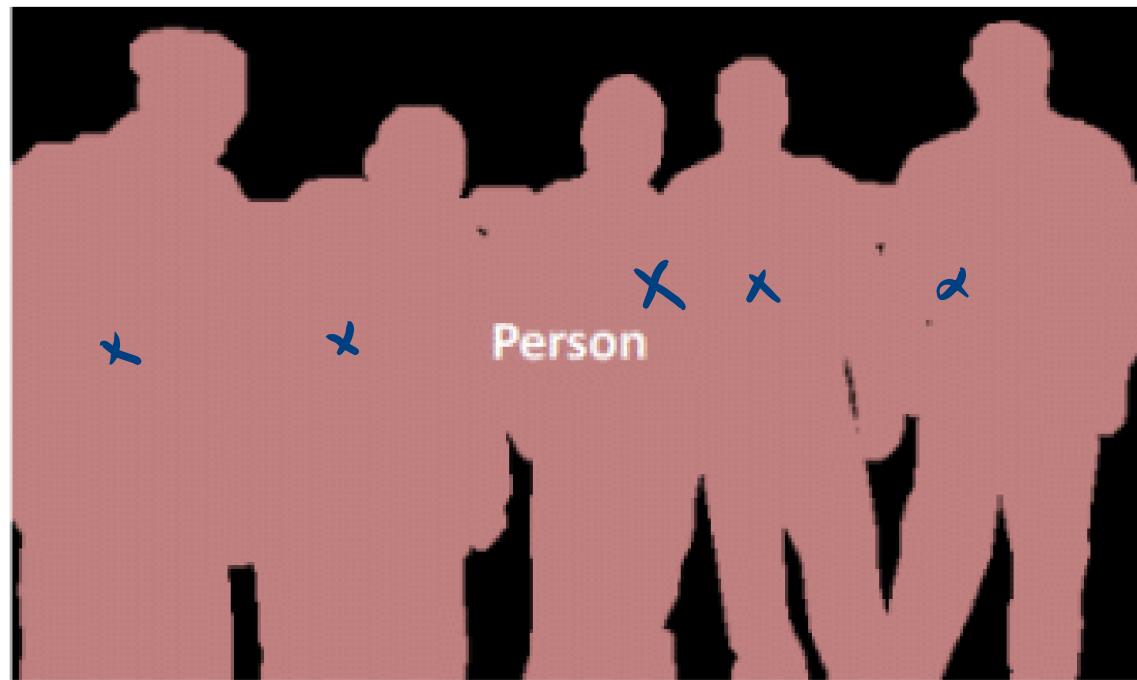
هر پیکسل تصویر را به یک کلاس معنایی اختصاص می‌دهد!



Instance Segmentation

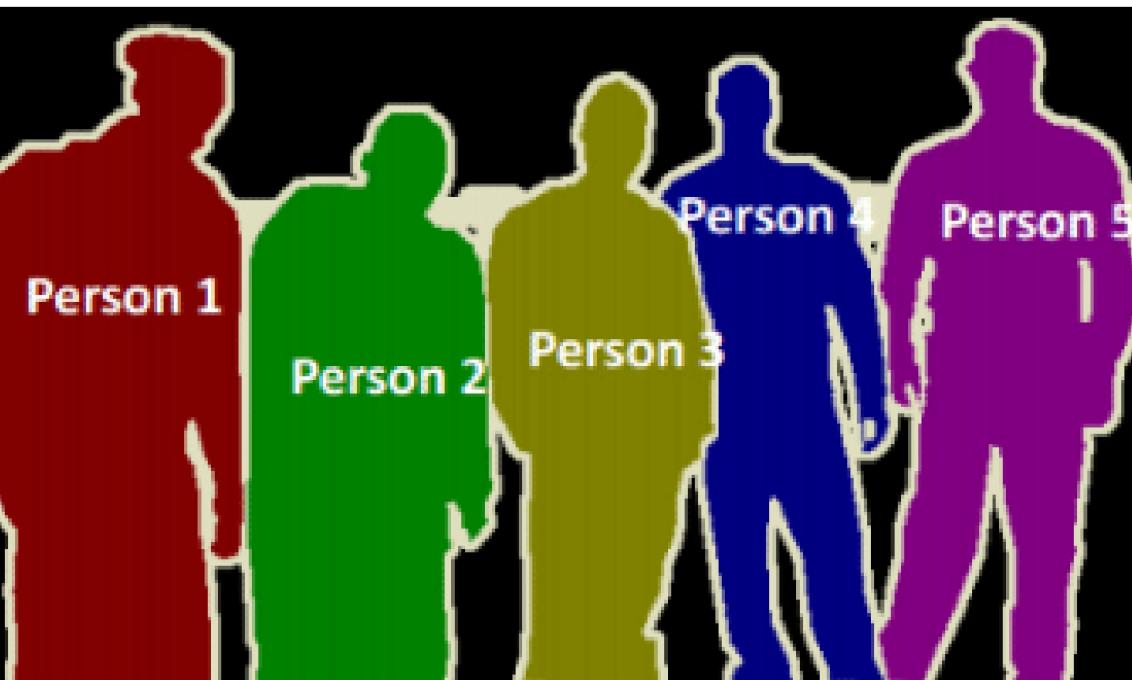
بین نمونه‌های مختلف از یک کلاس تمایز قائل می‌شود!

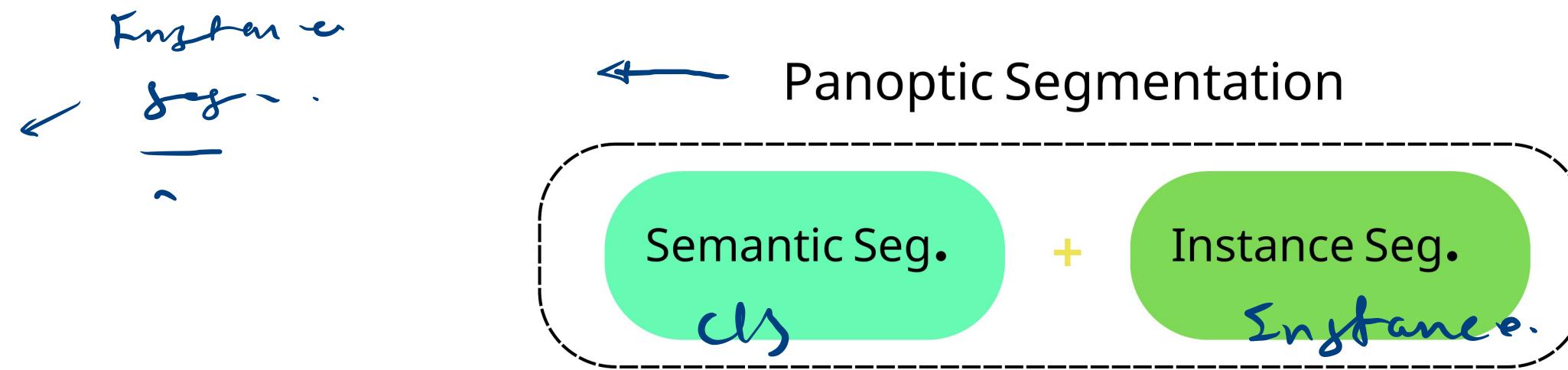
Semantic



نوع سگمنتیشن دو تصویر زیر ؟

Instance Seg.





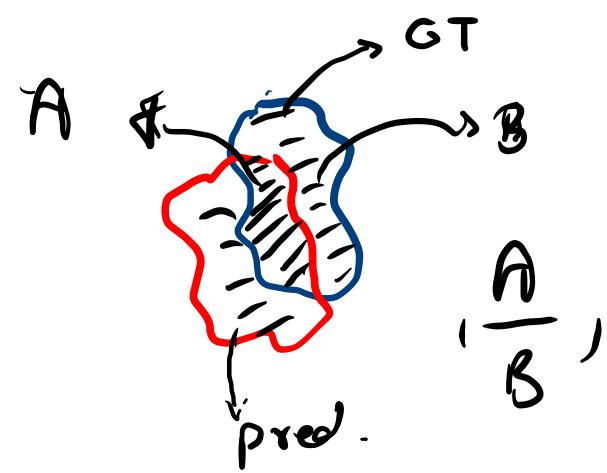
x

Semantic Segmentation

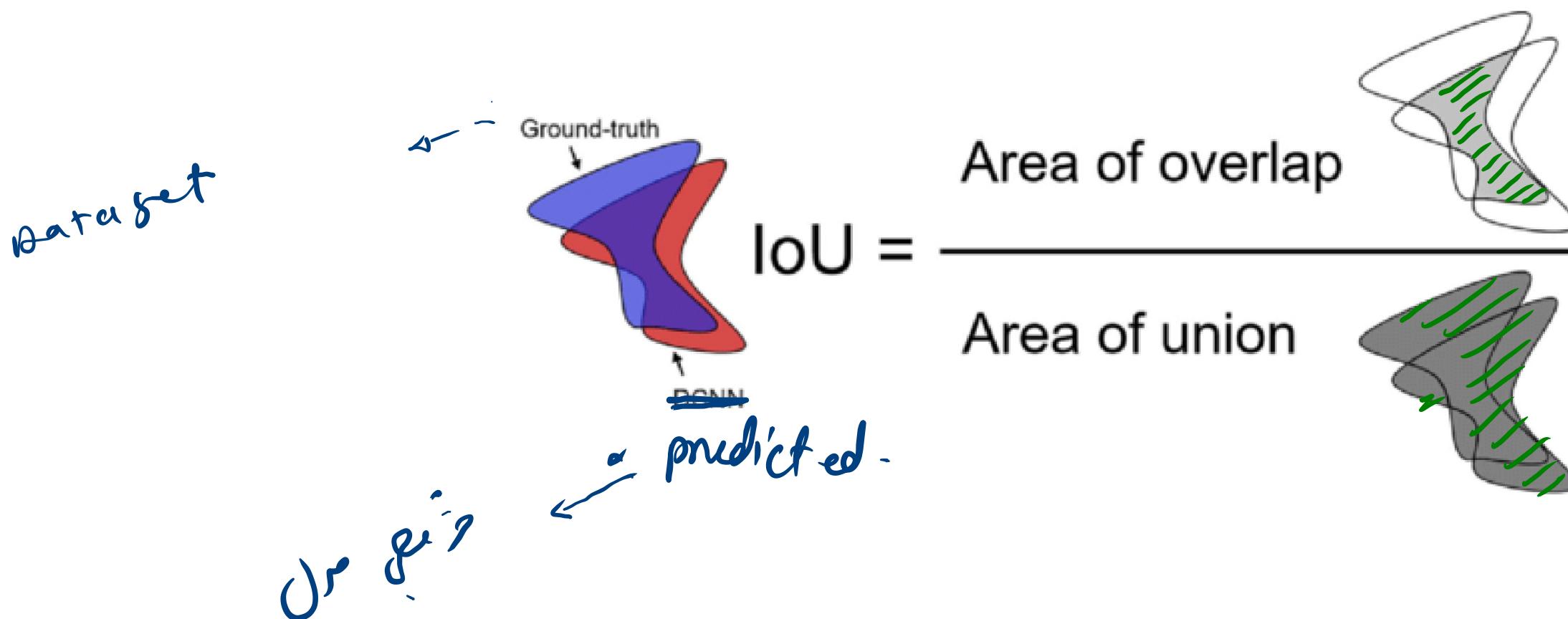


موضوع

ناحیه بندی معنایی

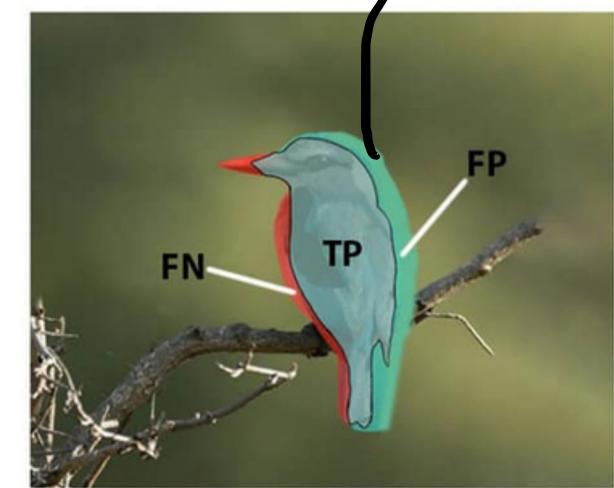
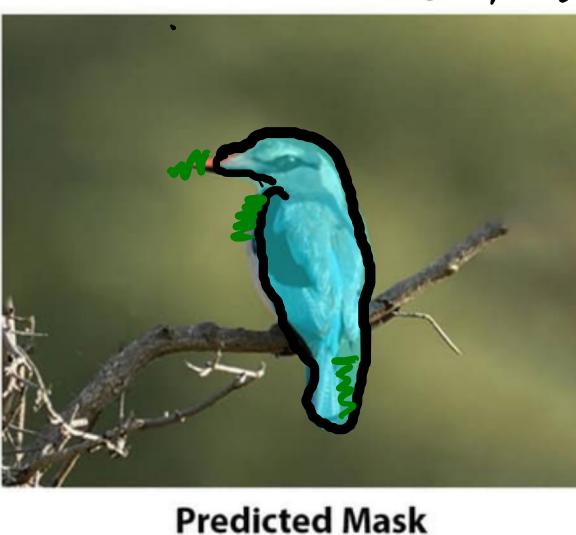


✓ $\frac{\text{Area of overlap}}{\text{Area of union}}$ IoU (Intersection over Union)



$$IoU = \frac{TP}{(TP + FP + FN)}$$

IoU (Intersection over Union)



مُنْتَهِيَةٌ مُنْتَهِيَةٌ مُنْتَهِيَةٌ

— سُلْطَانٌ صَفَرَ رَأَرَنَ رَأَنَ حَدَرَ GT

— سُلْطَانٌ صَفَرَ رَأَرَنَ رَأَنَ حَدَرَ FP

— سُلْطَانٌ صَفَرَ رَأَرَنَ رَأَنَ حَدَرَ TP

mIoU (mean Intersection over Union)

سیویلیتی
Semantic
Seg.

$$\text{IoU}_c = \frac{\text{TP}_c}{\text{TP}_c + \text{FP}_c + \text{FN}_c},$$

$$\text{mean_IoU} = \frac{1}{C} \sum_c \text{IoU}_c.$$

$$\begin{aligned} \Sigma \text{IoU}_p &\rightarrow \text{person} : .98 \\ \Sigma \text{IoU}_b &\rightarrow \text{ball} : .7 \\ \Sigma \text{IoU}_c &\rightarrow \text{car} : .82 \end{aligned}$$

$$\text{mIoU} = \frac{.98 + .7 + .82}{3} = \frac{2.5}{3} = .83$$

Pixel Accuracy (PA)

برابر با

$$PA = \frac{\text{تعداد پیکسل هایی که درست کласیفای شده اند}}{\text{تعداد کل پیکسل ها}}$$

تعداد پیکسل هایی که درست کласیفای شده اند!

تعداد کل پیکسل ها

mean Pixel Accuracy (mPA)

$$mPA = \frac{\text{مجموع تعداد پیکسل هایی که درست کласیفای شده اند! (روی همه کلاس ها)}}{\text{تعداد کل پیکسل های همه کلاس ها}}$$

TABLE 1. Popular semantic segmentation datasets.

Dataset	Training	Validation	Testing	Classes	Resolution
MS COCO [207]	118K	5000	-	80	vary in size
Pascal VOC 2012 [9]	1464	1449	-	20	vary in size
Pascal Context [10]	4998	5105	9637	59	vary in size
Pascal-Part [211]	1716	-	1817	6	vary in size
Cityscapes [11]	2975	500	1525	19	2048 × 1024
GTA5 [52]	-	-	-	19	1914 × 1052
KITTI [208]	-	-	-	3	vary in size
CamVid [12]	367	101	233	11	960 × 720
ADE20K [13]	20K	2000	-	150	vary in size
SYNTHIA [55]	-	-	-	13	1280 × 960
Stanford Background [213]	-	-	-	8	320 × 240

↓ ↓

Heavy and Lightweight Deep Learning Models for Semantic Segmentation: A Survey

CRISTINA CĂRUNTA^①, ALINA CĂRUNTA^②, AND CĂLIN-ADRIAN POPA^①, (Member, IEEE)

^①Department of Computers and Information Technology, Politehnica University of Timișoara, 300223 Timișoara, Romania

^②Department of Computer Science, West University of Timișoara, 300223 Timișoara, Romania

Corresponding author: Călin-Adrian Popa (calin.popa@cs.upt.ro)

14 January 2025

Heavy and Lightweight Deep Learning Models for Semantic Segmentation: A Survey

CRISTINA CĂRUNTA^①, ALINA CĂRUNTA^②, AND CĂLIN-ADRIAN POPA^①, (Member, IEEE)

^①Department of Computers and Information Technology, Politehnica University of Timișoara, 300223 Timișoara, Romania

^②Department of Computer Science, West University of Timișoara, 300223 Timișoara, Romania

Corresponding author: Călin-Adrian Popa (calin.popa@cs.upt.ro)

14 January 2025

Heavy

TABLE 2. Performance of complex semantic segmentation models ("T/V" represents test/validation set).

Method	Year	Backbone	VOC 2012	mIoU Cityscapes	ADE20K
FCN [8]	2015	VGG-16	62.2 T	-	-
DeconvNet [16]	2015	VGG-16	72.5 T	-	-
ParseNet [15]	2016	VGG-16	69.8 T	-	-
RefineNet [90]	2017	ResNet-152	83.4 T	-	40.7 V
		ResNet-101	-	73.6 T	-
PSPNet [30]	2017	ResNet-101	85.4 T	80.2 T	-
		ResNet-269	-	-	44.94 V
DeepLabV3+ [84]	2018	Xception	89.0 T	82.1 T	-
DenseASPP [96]	2018	DenseNet-161	-	80.6 T	-
EncNet [92]	2018	ResNet-101	85.9 T	-	44.65 V
DANet [1]	2019	ResNet-101	82.6 T	81.5 T	-
FDNet [193]	2019	-	84.2 T	-	-
CFNet [205]	2019	ResNet-101	87.2 T	79.60 T	44.89 V
CCNet [87]	2020	ResNet-101	-	81.9 T	45.76 V
SANet [88]	2020	ResNet-101	83.2 T	-	-
CDN [191]	2020	-	85.7 T	80.5 T	55.27 T
HRNet [24]	2021	HRNetV2-W48	-	82.5 T	-
SETR [98]	2021	ViT-L	-	81.64 T 82.15 V	50.28 V
TrSeg [99]	2021	DRN [202]	-	79.9 T	-
Segmenter [100]	2021	ViT-L	-	81.3 V	53.63 V
PVT [102]	2021	PVT-Large	-	-	44.8 V
Swin [35]	2021	Swin-L	-	-	53.5 V
			-	-	62.8 T
FLANet [89]	2022	HRNet-W48	-	83.6 T	46.99 V
		ResNet-150	88.5 T	-	-
SegNeXt [97]	2022	MSCAN-L	90.6 T	-	-
ConvNeXt [42]	2022	ConvNeXt-XL	-	-	54.0 V
PVT V2 [103]	2022	PVT V2-B5	-	-	48.7 V
HRViT [104]	2022	HRViT-b3	-	83.16 V	50.20 V
P2T [105]	2022	P2T-Large	-	-	49.4 V
FsaNet [22]	2023	ResNet-101	-	83.05 T	44.10 V
InternImage [94]	2023	InternImage-H	-	86.1 T 87.0 V	62.9 V
ONE-PEACE [38]	2023	-	-	-	63.0 V
SERNet-Former [222]	2024	Efficient-ResNet	-	87.35 V	-

Heavy and Lightweight Deep Learning Models for Semantic Segmentation: A Survey

CRISTINA CĂRUNTA^①, ALINA CĂRUNTA^②, AND CĂLIN-ADRIAN POPA^①, (Member, IEEE)

^①Department of Computers and Information Technology, Politehnica University of Timișoara, 300223 Timișoara, Romania

^②Department of Computer Science, West University of Timișoara, 300223 Timișoara, Romania

Corresponding author: Călin-Adrian Popa (calin.popa@cs.upt.ro)

14 January 2025

TABLE 2. Performance of complex semantic segmentation models. ("T/V" represents test/validation set).

Method	Year	Backbone	VOC 2012	mIoU Cityscapes	ADE20K
FCN [8]	2015	VGG-16	62.2 T	-	-
DeconvNet [16]	2015	VGG-16	72.5 T	-	-
ParseNet [15]	2016	VGG-16	69.8 T	-	-
RefineNet [90]	2017	ResNet-152	83.4 T	-	40.7 V
		ResNet-101	-	73.6 T	-
PSPNet [30]	2017	ResNet-101	85.4 T	80.2 T	-
		ResNet-269	-	-	44.94 V
DeepLabV3+ [84]	2018	Xception	89.0 T	82.1 T	-
DenseASPP [96]	2018	DenseNet-161	-	80.6 T	-
EncNet [92]	2018	ResNet-101	85.9 T	-	44.65 V
DANet [1]	2019	ResNet-101	82.6 T	81.5 T	-
FDNet [193]	2019	-	84.2 T	-	-
CFNet [205]	2019	ResNet-101	87.2 T	79.60 T	44.89 V
CCNet [87]	2020	ResNet-101	-	81.9 T	45.76 V
SANet [88]	2020	ResNet-101	83.2 T	-	-
CDN [191]	2020	-	85.7 T	80.5 T	55.27 T
HRNet [24]	2021	HRNetV2-W48	-	82.5 T	-
SETR [98]	2021	ViT-L	-	81.64 T 82.15 V	50.28 V
TrSeg [99]	2021	DRN [202]	-	79.9 T	-
Segmenter [100]	2021	ViT-L	-	81.3 V	53.63 V
PVT [102]	2021	PVT-Large	-	-	44.8 V
Swin [35]	2021	Swin-L	-	-	53.5 V 62.8 T
FLANet [89]	2022	HRNet-W48	-	83.6 T	46.99 V
		ResNet-150	88.5 T	-	-
SegNeXt [97]	2022	MSCAN-L	90.6 T	-	-
ConvNeXt [42]	2022	ConvNeXt-XL	-	-	54.0 V
PVT V2 [103]	2022	PVT V2-B5	-	-	48.7 V
HRViT [104]	2022	HRViT-b3	-	83.16 V	50.20 V
P2T [105]	2022	P2T-Large	-	-	49.4 V
FsaNet [22]	2023	ResNet-101	-	83.05 T	44.10 V
InternImage [94]	2023	InternImage-H	-	86.1 T 87.0 V	62.9 V
ONE-PEACE [38]	2023	-	-	-	63.0 V
SERNet-Former [222]	2024	Efficient-ResNet	-	87.35 V	-

Method	Year	Backbone	GPU	Resolution	FPS	Params (M)	mIoU
ENet [113]	2016	-	Titan X	640 × 360	135.4		
				1280 × 720	46.8	0.37	
				1920 × 1080	21.6		
LinkNet [115]	2017	ResNet-18	Titan X	640 × 360	65.8		
				1280 × 720	18.7	11.5	76.4 V
ERFNet [116]	2018	-	Titan X	640 × 360	83.3		
				1280 × 720	24.4	-	69.7 T
ContextNet [125]	2018	-	Titan X	1920 × 1080	11.4		
				2048 × 1024	18.3	0.85	65.9 V
BiSeNet [128]	2018	Xception-39 ResNet-18	Titan XP	2048 × 1024	105.8	5.8	
					65.5	49.0	74.7 T
ShelfNet [118]	2019	ResNet-18	GTX 1080Ti	2048 × 1024	36.9	-	74.8 T
SwiftNetRN-18 [126]	2019	ResNet-18	GTX 1080Ti	2048 × 1024	39.9	11.8	75.5 T
LEDNet [143]	2019	ResNet	GTX 1080Ti	1024 × 512	71	0.94	70.6 T
FDDWNet [120]	2020	-	RTX 2080Ti	1024 × 512	60	0.80	71.5 T
LDPNet [160]	2020	-	GTX 1080Ti	1024 × 512	87	0.8	71.1 T
CARNet [144]	2020	ResNet-152	Titan X	1024 × 512	29.4	63.3	75.2 T
ERPNet [121]	2021	ResNet-18 MobileNet V3	GTX 2080	896 × 896	66.1	10.98	75.3 T
				780 × 780	109.5	1.46	73.6 T
DDRNet [139]	2021	ResNet	GTX 2080Ti	2048 × 1024	101.6	5.7	77.4 T
BiSeNet V2 [129]	2021	-	GTX 1080Ti	2048 × 1024	156	-	72.6 T
SegFormer [39]	2021	MiT-B5	-	1024 × 1024	2.5	84.7	84.0 V
RTFormer [164]	2022	-	RTX 2080Ti	2048 × 1024	110	4.8	76.3 V
RELAXNet [153]	2022	-	GTX 2080Ti	1024 × 512	64	1.9	74.8 T
JPNNet [152]	2022	-	GTX 1080Ti	1024 × 512	109.9	3.49	71.62 T
PP-LiteSeg [203]	2022	STDC1 [131] STDC2 [131]	GTX 1080Ti	1024 × 512	273.6		72.0 T
				1536 × 768	102.6	-	73.1 V
BFMNet [140]	2023	ResNet-18	RTX 2080Ti	1536 × 768	63.7		77.7 T
				2048 × 1024	31.4	22.3	78.9 T
MLFNet [141]	2023	ResNet-18 ResNet-34	Titan XP	1024 × 512	95.1	9.90	71.0 T
					72.2	13.03	72.1 T
BiSeNet V3 [135]	2023	STDC2	GTX 1080Ti	1536 × 768	93.8	-	79.0 T
LETNet [166]	2023	-	RTX 3090	1024 × 512	150	0.95	72.8 T
AFFormer [168]	2023	-	V100 NVIDIA	2048 × 1024	22	3.0	78.7 V
PIDNet [33]	2023	ResNet	RTX 3090	2048 × 1024	93.2	7.6	78.6 T
DSNet [221]	2024	-	RTX 3090	2048 × 1024	31.1	36.9	80.6 T
					37.6	6.8	80.4 V

Heavy and Lightweight Deep Learning Models for Semantic Segmentation: A Survey

CRISTINA CĂRUNTA^①, ALINA CĂRUNTA^②, AND CĂLIN-ADRIAN POPA^①, (Member, IEEE)

^①Department of Computers and Information Technology, Politehnica University of Timișoara, 300223 Timișoara, Romania

^②Department of Computer Science, West University of Timișoara, 300223 Timișoara, Romania

Corresponding author: Călin-Adrian Popa (calin.popa@cs.upt.ro)

14 January 2025

TABLE 3. Performance of real-time models on Cityscapes dataset. ("T/V" represents test/validation set; 'FPS' denotes frames per second; 'Params (M)' is the number of learnable parameters).

Method	Year	Backbone	GPU	Resolution	FPS	Params (M)	mIoU
ENet [113]	2016	-	Titan X	640 × 360	135.4		
				1280 × 720	46.8	0.37	58.3 T
				1920 × 1080	21.6		
LinkNet [115]	2017	ResNet-18	Titan X	640 × 360	65.8		
				1280 × 720	18.7	11.5	76.4 V
ERFNet [116]	2018	-	Titan X	1920 × 1080	8.5		
				640 × 360	83.3		
ContextNet [125]	2018	-	Titan X	1280 × 720	24.4	-	69.7 T
				1920 × 1080	11.4		
BiSeNet [128]	2018	Xception-39 ResNet-18	Titan XP	2048 × 1024	18.3	0.85	65.9 V
				2048 × 1024	105.8	5.8	68.4 T
ShelfNet [118]	2019	ResNet-18	GTX 1080Ti	2048 × 1024	36.9	-	74.8 T
SwiftNetRN-18 [126]	2019	ResNet-18	GTX 1080Ti	2048 × 1024	39.9	11.8	75.5 T
LEDNet [143]	2019	ResNet	GTX 1080Ti	1024 × 512	71	0.94	70.6 T
FDDWNet [120]	2020	-	RTX 2080Ti	1024 × 512	60	0.80	71.5 T
LDPNet [160]	2020	-	GTX 1080Ti	1024 × 512	87	0.8	71.1 T
CARNet [144]	2020	ResNet-152	Titan X	1024 × 512	29.4	63.3	75.2 T
ERPNet [121]	2021	ResNet-18 MobileNet V3	GTX 2080	896 × 896	66.1	10.98	75.3 T
				780 × 780	109.5	1.46	73.6 T
DDRNet [139]	2021	ResNet	GTX 2080Ti	2048 × 1024	101.6	5.7	77.4 T
BiSeNet V2 [129]	2021	-	GTX 1080Ti	2048 × 1024	156	-	72.6 T
SegFormer [39]	2021	MiT-B5	-	1024 × 1024	2.5	84.7	84.0 V
RTFormer [164]	2022	-	RTX 2080Ti	2048 × 1024	110	4.8	76.3 V
RELAXNet [153]	2022	-	GTX 2080Ti	1024 × 512	64	1.9	74.8 T
JPNANet [152]	2022	-	GTX 1080Ti	1024 × 512	109.9	3.49	71.62 T
PP-LiteSeg [203]	2022	STDC1 [131] STDC2 [131]	GTX 1080Ti	1024 × 512	273.6		72.0 T
				1536 × 768	102.6		73.1 V 77.5 T 78.2 V
BFMNet [140]	2023	ResNet-18	RTX 2080Ti	1536 × 768	63.7		77.7 T
				2048 × 1024	31.4	22.3	78.9 T
MLFNet [141]	2023	ResNet-18 ResNet-34	Titan XP	1024 × 512	95.1	9.90	71.0 T
				72.2	13.03		72.1 T
BiSeNet V3 [135]	2023	STDC2	GTX 1080Ti	1536 × 768	93.8	-	79.0 T
LETNet [166]	2023	-	RTX 3090	1024 × 512	150	0.95	72.8 T
AFFormer [168]	2023	-	V100 NVIDIA	2048 × 1024	22	3.0	78.7 V
PIDNet [33]	2023	ResNet	RTX 3090	2048 × 1024	93.2	7.6	78.6 T
DSNet [221]	2024	-	RTX 3090	2048 × 1024	37.6	6.8	80.4 V

Heavy and Lightweight Deep Learning Models for Semantic Segmentation: A Survey

CRISTINA CĂRUNTA^①, ALINA CĂRUNTA^②, AND CĂLIN-ADRIAN POPA^①, (Member, IEEE)

^①Department of Computers and Information Technology, Politehnica University of Timișoara, 300223 Timișoara, Romania

^②Department of Computer Science, West University of Timișoara, 300223 Timișoara, Romania

Corresponding author: Călin-Adrian Popa (calin.popa@cs.upt.ro)

14 January 2025

mIoU > 70

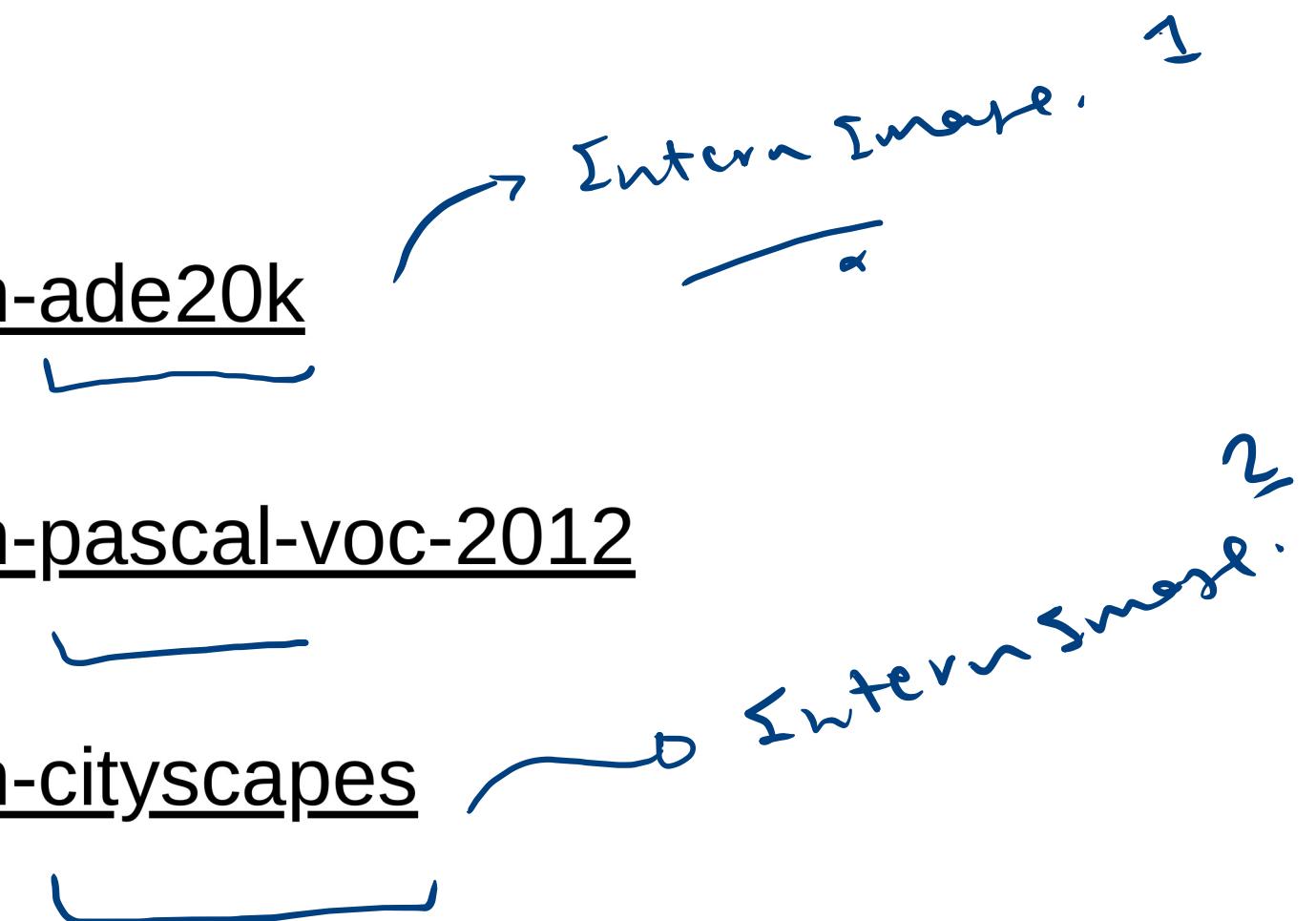
FPS > 100

TABLE 3. Performance of real-time models on Cityscapes dataset. ("T/V" represents test/validation set; 'FPS' denotes frames per second; 'Params (M)' is the number of learnable parameters).

<https://paperswithcode.com/sota/semantic-segmentation-on-ade20k>

<https://paperswithcode.com/sota/semantic-segmentation-on-pascal-voc-2012>

<https://paperswithcode.com/sota/semantic-segmentation-on-cityscapes>



Semantic
yes.

InternImage: Exploring Large-Scale Vision Foundation Models with Deformable Convolutions

Wenhai Wang^{1*}, Jifeng Dai^{2,1*}, Zhe Chen^{3,1*}, Zhenhang Huang^{1*}, Zhiqi Li^{3,1*}, Xizhou Zhu^{4*}, Xiaowei Hu¹, Tong Lu³, Lewei Lu⁴, Hongsheng Li⁵, Xiaogang Wang^{4,5}, Yu Qiao^{1✉}

¹Shanghai AI Laboratory ²Tsinghua University

³Nanjing University ⁴SenseTime Research ⁵The Chinese University of Hong Kong

<https://github.com/OpenGVLab/InternImage>

The End