

# Pruning Parameterization with Bi-level Optimization for Efficient Semantic Segmentation on the Edge

WED-PM-290

Changdi Yang\*<sup>1</sup> Pu Zhao\*<sup>1</sup> Yanyu Li<sup>1</sup> Wei Niu<sup>2</sup> Jiexiong Guan<sup>2</sup> Hao Tang<sup>3</sup> Minghai Qin<sup>1</sup> Bin Ren<sup>2</sup> Xue Lin<sup>1</sup> Yanzhi Wang<sup>1</sup>  
<sup>1</sup>Northeastern University <sup>2</sup>College of William & Mary <sup>3</sup>CVL, ETH Zurich



# Preview

- We aim to derive ViTs with fewer computations and fast inference speed
- We propose a pruning parameterization method to formulate the pruning problem of semantic segmentation.
- We adopt a bi-level optimization method to solve this problem with the help of implicit gradients.

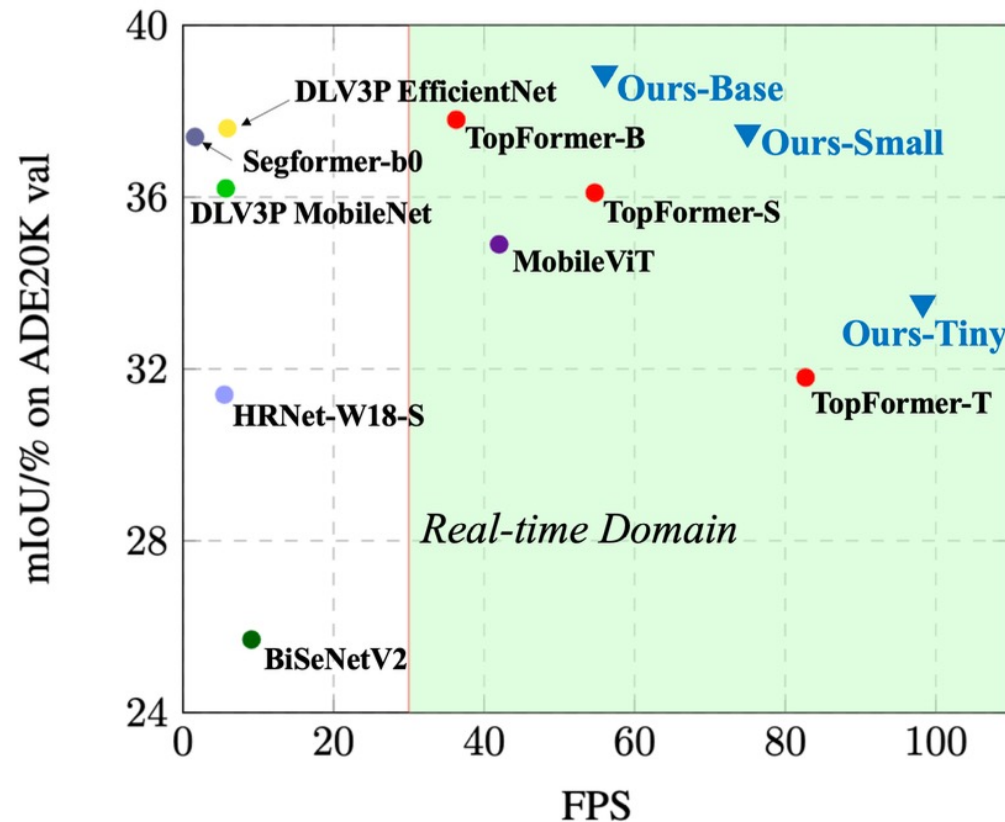


Figure 1. Comparison of accuracy versus FPS on ADE20K.



# Introduction

---

- **Background**

- Segmentation models usually have heavy computational cost.
- Many light-weight models are designed for desktop GPUs.
- NAS and pruning methods cost huge memory and training time.

- **Goal**

- Derive ViTs with fast inference speed on edge devices.

VERITAS  
VIRTUS



# Problem Formulation

---

- **Previous methods**

- Depend on the magnitudes of model weight.
- Adopt the in-differentiable sorting operations.
- Leads to inconsistent performance and additional overhead.

- **Our method**

- Uses a soft mask to indicate whether to prune.
- Get rid of sorting operations.

VERITAS  
VIRTUS



# Problem Formulation

- **Soft mask construction**

- Adopt channel pruning to search for a suitable width for each convolution (CONV) layer.
- Insert a depth-wise CONV layer following each CONV layer.
- $s_l$  can serve as a soft mask or pruning indicator

$$\mathbf{a}_l = \mathbf{s}_l \odot (\mathbf{w}_l \odot \mathbf{a}_{l-1}).$$

LVX  
VERITAS  
VIRTUS



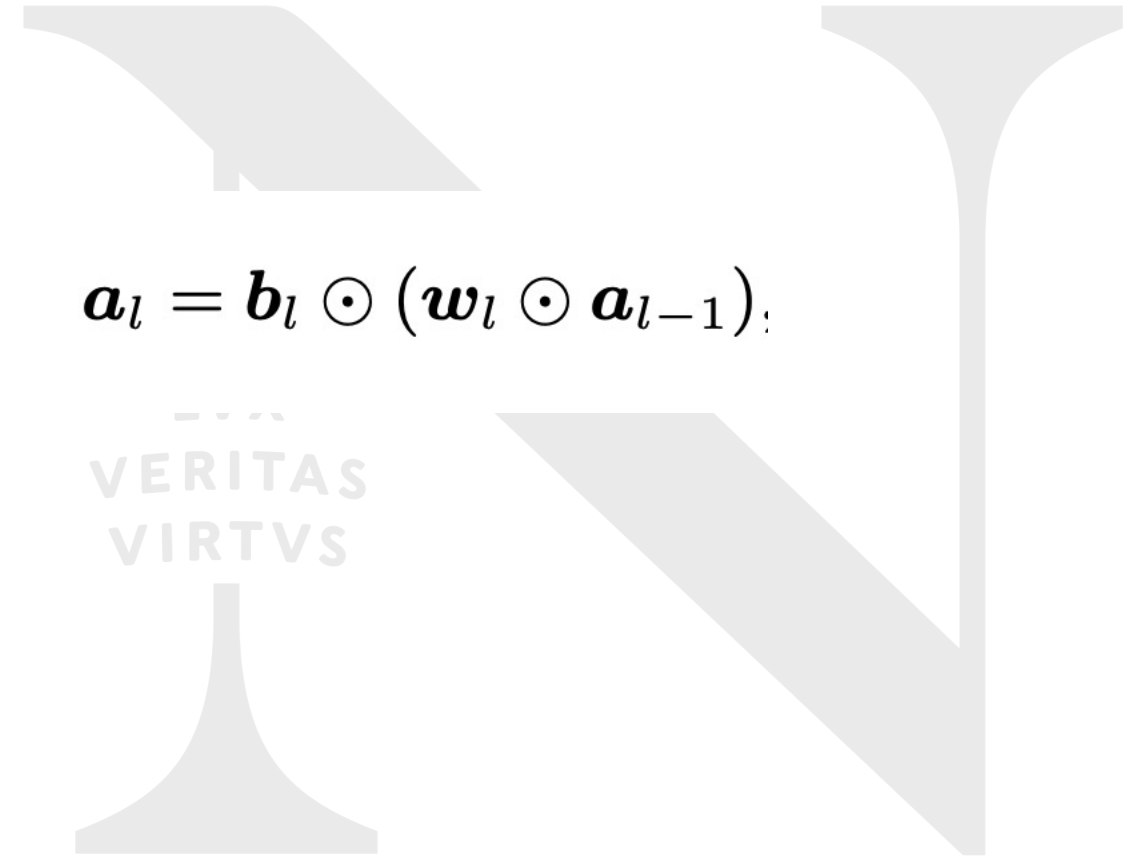
# Problem Formulation

- **Forward and backward propagation**

- We adopt a threshold  $\tau$ , and the forward pass with the mask is represented as:

$$\mathbf{b}_l = \begin{cases} 1, & \mathbf{s}_l > \tau \\ 0, & \mathbf{s}_l \leq \tau \end{cases} \quad (\text{element-wise}),$$

$$\mathbf{a}_l = \mathbf{b}_l \odot (\mathbf{w}_l \odot \mathbf{a}_{l-1});$$



# Problem Formulation

- **Training Loss**
  - **Cross-Entropy loss and regularization loss**

$$\mathcal{L}_m(\mathbf{w}, \mathbf{s}) = \mathcal{L}(\mathbf{w}, \mathbf{s}) + \beta \cdot \mathcal{L}_{reg}(\mathbf{s})$$

$$\mathcal{L}_{reg} = \left| \sum_l o'_l \times i_l \times t_l \times t'_l \times k^2 - \mathcal{C} \right|^2$$

VERITAS  
VIRTUS



# Problem Formulation

- Bi-Level optimization

$$\begin{aligned} \min_{\mathbf{s}} \quad & \mathcal{L}_m(\mathbf{w}^*, \mathbf{s}), \\ \text{s.t.} \quad & \mathbf{w}^* = \arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \mathbf{s}) + \frac{1}{2\lambda} \|\mathbf{w}\|_2^2. \end{aligned}$$

VIRTUS





# Proposed method

- **Bi-Level optimization method**

- In the first step, we update the weights  $w$  with a few training steps for a fixed mask  $s$ .
- Next in the second step, we update  $s$  with implicit gradients.

$$\frac{d\mathcal{L}_m(\mathbf{w}^*, \mathbf{s})}{d\mathbf{s}} = \frac{d\mathbf{w}^*}{d\mathbf{s}} \nabla_{\mathbf{w}} \mathcal{L}_m(\mathbf{w}^*, \mathbf{s}) + \nabla_{\mathbf{s}} \mathcal{L}_m(\mathbf{w}^*, \mathbf{s})$$

$$\frac{d\mathcal{L}_m(\mathbf{w}^*, \mathbf{s})}{d\mathbf{s}} = \nabla_{\mathbf{s}} \mathcal{L}_m(\mathbf{w}^*, \mathbf{s}) - \lambda \mathbf{b} \cdot [\nabla_{\mathbf{w}} \mathcal{L}_m(\mathbf{w}^*, \mathbf{s})]^2$$



# Problem Formulation

- **Difference with other pruning works**
  - Decouple the pruning policy from model parameter magnitudes.
  - Our method does not have such a constraint that pruned weights should be zero.
  - Our mask method is more straightforward and effective.

$$\mathbf{b}_l = \begin{cases} 1, & \mathbf{s}_l > \tau \\ 0, & \mathbf{s}_l \leq \tau \end{cases} \quad (\text{element-wise}), \quad \mathbf{a}_l = \mathbf{b}_l \odot (\mathbf{w}_l \odot \mathbf{a}_{l-1});$$



# Experiments & Results

- **Datasets**

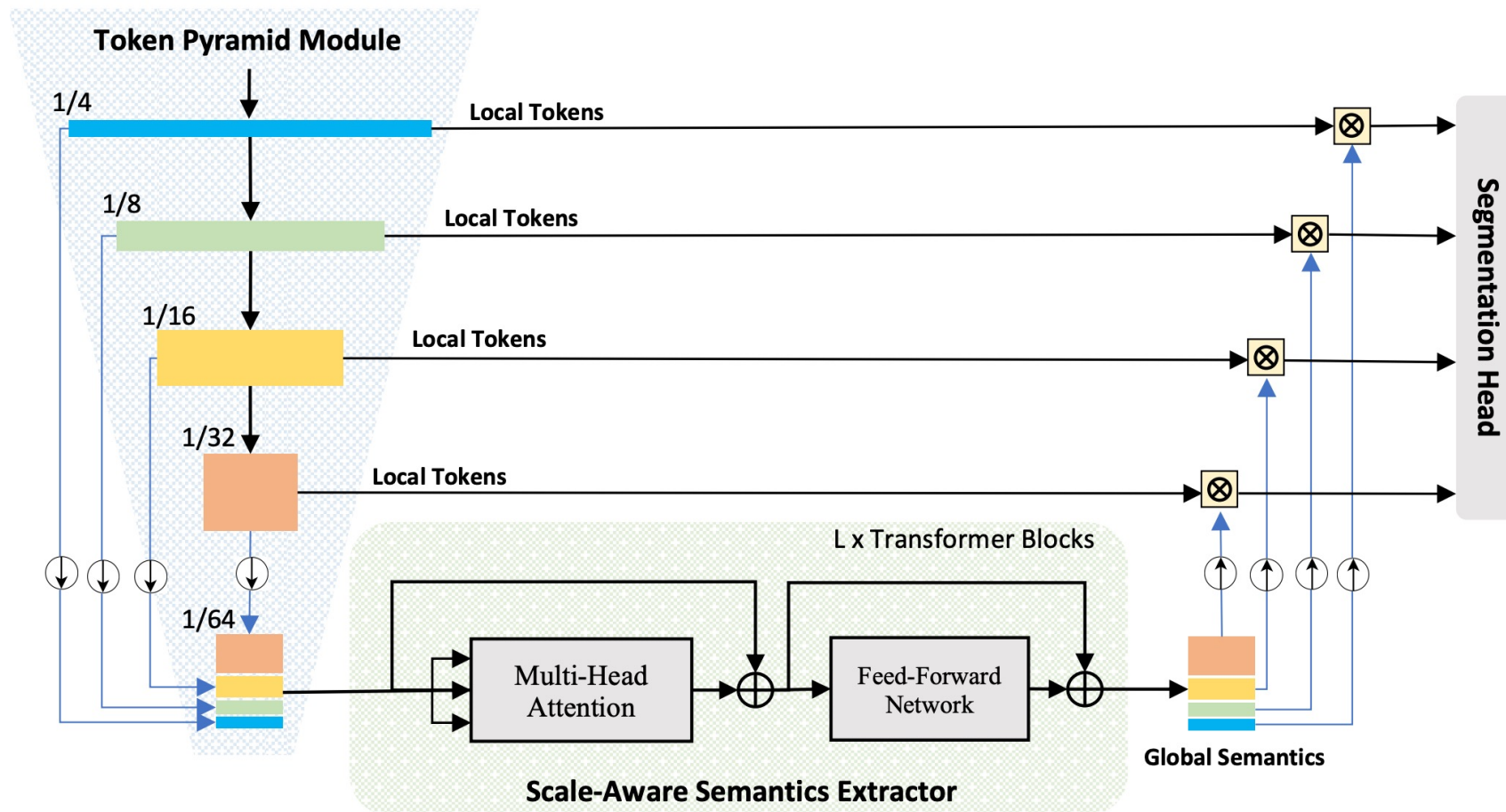
- **ADE20K:** 25k training, 2k validation with 150 label classes.
- **Cityscapes:** 2,975 training, 500 validation, and 1,525 for testing with 19 label classes.
- **Pascal VOC 2012:** 1,464 training and 1,449 images validation.

LVX  
VERITAS  
VIRTVS



# Experiments & Results

- **Tested model**
  - TopFormer



# Experiments & Results

Table 1. Comparison of our searched model and prior arts on the ADE20K val dataset. We compare with popular handcraft baselines in the first segment, NAS-based models in the second segment, pruning-based methods in the third segment and lightweight ViT-based models in the fourth segment. We measure the FPS on the Qualcomm Adreno 660 GPU of the Samsung Galaxy S21 mobile phone. Some FPS results are not available due to unsupported operations on the mobile device.

Category	Method	Backbone	Parameters	GMACs	FPS	val mIoU (%)
Human Design	PSPNet [71]	ResNet50-D8	49.1M	178.8	0.4	41.1
	DeepLabV3+ [6]	EfficientNet	17.1M	26.9	5.9	37.6
	BiSeNetV2 [66]	N.A.	3.34M	12.4	9.1	25.7
	SFNet [40]	ResNet-50	—	75.7	—	42.8
	HRNet-W18-S [60]	HRNet-W18-S	4.0M	10.2	5.5	31.4
NAS	HR-NAS-A [15]	Searched	2.5M	1.4	—	33.2
	HR-NAS-B [15]	Searched	3.9M	2.2	—	34.9
	NASViT [24]	Searched	—	2.5	—	37.9
	HRViT-b1 [26]	Searched	8.2M	14.6	—	45.9
Prune	EagleEye [38]	N.A.	3.4M	1.2	59.2	34.3
	DMCP [28]	N.A.	3.3M	1.2	63.8	33.9
	ResPep [16]	N.A.	3.3M	1.2	64.9	35.0
	CHEX [32]	N.A.	3.3M	1.2	64.2	35.2
ViT	Segmenter [57]	Searched	6.7M	4.6	4.1	39.9
	MobileViT [48]	Searched	3.9M	2.2	41.8	34.9
	SegFormer-B0 [63]	MiT-B0	3.8M	8.4	2.6	37.4
	TopFormer-Base [68]	N.A.	5.1M	1.8	36.3	37.8
	TopFormer-Small [68]	N.A.	3.1M	1.2	54.7	36.1
	TopFormer-Tiny [68]	N.A.	1.4M	0.6	82.7	31.8
Ours	<b>Ours-Base</b>	Searched	<b>3.7M</b>	<b>1.8</b>	<b>56.5</b>	<b>38.9</b>
	<b>Ours-Small</b>	Searched	<b>3.3M</b>	<b>1.2</b>	<b>75.2</b>	<b>37.5</b>
	<b>Ours-Tiny</b>	Searched	<b>1.3M</b>	<b>0.7</b>	<b>98.0</b>	<b>33.5</b>

# Experiments & Results

Table 2. Our search results on the Cityscapes val dataset. We compare with popular handcraft baselines, NAS-based models, pruning based methods and lightweight ViT-based models. We measure the FPS on the Qualcomm Adreno 660 GPU of the Samsung Galaxy S21 mobile phone. Some FPS results are not available due to unsupported operations on the mobile device.

Category	Method	Backbone	Resolution	#params	GMACs	FPS	mIoU%
Human Design	ENet [50]	N.A.	$512 \times 1024$	354.9K	5.9	—	58.5
	PSPNet [71]	ResNet101	$1024 \times 2048$	68.07M	525.0	0.2	78.8
	BiSeNetV2 [66]	N.A.	$512 \times 1024$	3.34M	24.6	5.0	73.4
	DeepLabV3+ [6]	MBv2	$512 \times 1024$	2.26M	9.5	9.4	69.0
	STDC1-Seg50 [21]	STDC1	$512 \times 1024$	12.05M	31.1	4.3	72.2
	STDC2-Seg50 [21]	STDC2	$512 \times 1024$	16.08M	44.3	3.7	74.2
NAS	Auto-DeepLab-S [45]	N.A.	$1024 \times 2048$	10.15M	333.3	—	79.7
	FasterSeg [7]	N.A.	$1024 \times 2048$	—	28.2	—	73.1
	DCNAS [69]	N.A.	$1024 \times 2048$	—	294.6	—	85.0
Prune	EagleEye [38]	N.A.	$512 \times 1024$	3.6M	2.4	27.6	69.6
	DMCP [28]	N.A.	$512 \times 1024$	3.5M	2.4	32.0	70.3
	ResPep [16]	N.A.	$512 \times 1024$	3.5M	2.4	28.2	71.3
	CHEX [32]	N.A.	$512 \times 1024$	3.4M	2.4	35.5	71.7
ViT	HRViT-b1 [26]	N.A.	$512 \times 1024$	8.1M	28.2	—	81.6
	SegFormer-B0 [63]	MiT-B0	$512 \times 1024$	3.8M	17.7	0.9	71.9
	TopFormer-Base [68]	N.A.	$512 \times 1024$	5.1M	2.7	21.5	70.6
	TopFormer-Tiny [68]	N.A.	$512 \times 1024$	1.4M	1.2	42.8	66.1
Ours	<b>Ours-Base</b>	Searched	$512 \times 1024$	<b>3.7M</b>	<b>3.6</b>	<b>30.8</b>	<b>74.7</b>
	<b>Ours-Small</b>	Searched	$512 \times 1024$	<b>3.3M</b>	<b>2.4</b>	<b>38.7</b>	<b>73.6</b>
	<b>Ours-Tiny</b>	Searched	$512 \times 1024$	<b>1.3M</b>	<b>1.4</b>	<b>52.6</b>	<b>71.5</b>

# Experiments & Results

Table 3. Comparison of search cost on the Cityscapes val dataset.

Method	GPU Days	GMACs	mIoU
Auto-DeepLab [45]	3	695.0	82.1
GAS [44]	6.7	-	73.5
FasterSeg [7]	2	28.2	73.1
Fast-NAS [49]	8	435.7	78.9
SparseMask [62]	4.2	36.4	68.6
DCNAS [69]	5.6	294.6	85.0
LDP [34]	4.3	—	75.8
<b>Without implicit gradients</b>	<b>1.1</b>	<b>2.4</b>	<b>71.9</b>
<b>Ours-Small</b>	<b>1.3</b>	<b>2.4</b>	<b>73.6</b>



# Experiments & Results

Table 4. Results on the PASCAL VOC 2012 test dataset. We compare our results with popular CNN-based models and lightweight ViT-based models.

Method	#params	GMACs	mIoU%	FPS
EfficientNet-B7 [59]	66.0M	194.0	85.2	0.1
EMANet [41]	10.0M	43.1	80.1	2.5
PSANet [2]	18.5M	56.3	78.5	1.4
DeepLabV3+ R101 [6]	43.9M	58.5	77.4	2.2
DeepLabV3+ R50 [6]	24.9M	37.8	76.3	3.1
DeepLabV3+ MBv2 [6]	2.3M	5.7	70.5	5.1
TopFormer-B [68]	5.1M	1.8	71.0	36.8
TopFormer-S [68]	3.1M	1.2	69.8	55.2
TopFormer-T [68]	1.4M	0.6	65.7	81.5
MobileViT-XXS [48]	1.9M	1.7	73.6	43.8
<b>Ours-Base</b>	<b>3.7M</b>	<b>1.8</b>	<b>74.3</b>	<b>56.8</b>
<b>Ours-Small</b>	<b>3.3M</b>	<b>1.2</b>	<b>73.4</b>	<b>75.0</b>
<b>Ours-Tiny</b>	<b>1.3M</b>	<b>0.7</b>	<b>70.5</b>	<b>97.6</b>





# Experiments & Results

Table 5. Our searched results for DeepLabV3+ with MobileNetV2 backbone on Cityscapes val. The input resolution is  $512 \times 1024$ .

Method	#params	GMACs	mIoU%	FPS
BiSeNetV2 [66]	3.34M	24.6	73.4	5.0
STDC1-Seg50 [21]	12.05M	31.1	72.2	4.3
STDC2-Seg50 [21]	16.08M	44.3	74.2	3.7
DeepLabV3+ [6]	2.26M	9.5	69.0	9.4
<b>Ours-Base-DeepLab</b>	<b>1.21M</b>	<b>7.6</b>	<b>70.9</b>	<b>22.3</b>
<b>Ours-Small-DeepLab</b>	<b>569.0K</b>	<b>4.3</b>	<b>70.2</b>	<b>28.1</b>



# Experiments & Results



(a) Input

(b) Ours-Base

(c) TopFormer-Base

Figure A2. Visualization results on samples of Cityscapes validation dataset.



# Ablation Study

Table A1. Results on Cityscapes with different  $\beta$ .

$\beta$	0.001	0.01	0.1	1.0
mIoU	73.1	73.6	72.8	71.2

Table A2. Results on Cityscapes with different  $\lambda$ .

$\lambda$	0.01	0.1	1.0
mIoU	73.1	73.6	70.8



Thanks

