# Rethinking Learnable Tree Filter for Generic Feature Transform





#### NeurIPS 2020

宋林 (Lin Song) stevengrove@stu.xjtu.edu.cn



Conventional methods for visual context modeling

**Local-based:** increase the receptive region of convolutional kernels e.g., Dilated Convolution, Deformable Convolution, ASPP, PSP Global-based: modeling the pairwise relations based on the visual attention mechanism e.g., Non-Local, CCNet, LatentGNN





Conventional methods for visual context modeling

However, there is still a conflict between long-range dependencies modeling and object details preserving.





LTF-VI: Learnable Tree Filter for Structure-preserving Feature Transform

LTF-VI uses the minimum spanning tree generated by the low-level guided features, which can retain the structural details with linear complexity w.r.t. vertex number





#### LTF-VI: Learnable Tree Filter for Structure-preserving Feature Transform

#### LTF-V1 is a differentiable and plug-and-play module





#### LTF-VI: Learnable Tree Filter for Structure-preserving Feature Transform

### LTF-V1 can model higher-order relations, which is crucial for feature discriminability





LTF-VI: Learnable Tree Filter for Structure-preserving Feature Transform

Nevertheless, the geometric constraint and the non-differentiable spanning tree construction in the LTF-V1 module are found to be its Achilles' heels, which impede the usage for more generic feature transform





#### **Reformulation of LTF-VI**

The statistical expectation of the sampling from input features under a specific distribution:

$$y_{oldsymbol{i}} = \mathbb{E}_{oldsymbol{h}_{oldsymbol{i}} \sim P_{oldsymbol{\mathcal{G}}_{T}}}[x_{oldsymbol{h}_{oldsymbol{i}}}] = \sum_{orall j \in \mathcal{V}} P_{oldsymbol{\mathcal{G}}_{T}}(h_{oldsymbol{i}} = j) x_{oldsymbol{j}}$$

The distribution is modeled by Markov Random Field:  $P_{\mathcal{G}_T}(\mathbf{H}|\mathbf{O}=X) = rac{1}{Z} \prod_{orall i \in \mathcal{V}} \phi_i(h_i, x_i) \prod_{orall (i,j) \in \mathcal{E}} \psi_{i,j}(h_i, h_j)$ 

The distribution for LTF-VI can be considered as a specific Markov Random Field:

 $\phi_i(h_i, x_i) \equiv 1$ 

$$\psi_{m{i},m{j}}(h_{m{i}},h_{m{j}}) \coloneqq egin{cases} \delta(h_{m{i}}-h_{m{j}}) & h_{m{i}} 
otin \mathrm{Desc}_{\mathcal{G}_T}(m{i}) \ exp(-\omega_{m{i},m{j}})\delta(h_{m{i}}-h_{m{j}}) & h_{m{i}} \in \mathrm{Desc}_{\mathcal{G}_T}(m{i}) \end{cases}$$

#### NeurIPS 2020

j, j



#### **Reformulation of LTF-V1**

The statistical expectation of the sampling from input features under a specific distribution:

$$y_{oldsymbol{i}} = \mathbb{E}_{h_{oldsymbol{i}} \sim P_{\mathcal{G}_{T}}}[x_{h_{oldsymbol{i}}}] = \sum_{orall j \in \mathcal{V}} P_{\mathcal{G}_{T}}(h_{oldsymbol{i}} = j) x_{oldsymbol{j}}$$

The distribution is modeled by Markov Random Field:  $P_{\mathcal{G}_T}(\mathbf{H}|\mathbf{O}=X) = rac{1}{Z} \prod_{orall i \in \mathcal{V}} \phi_i(h_i, x_i) \prod_{orall (i,i) \in \mathcal{E}} \psi_{i,j}(h_i, h_j)$ 

The distribution for LTF-VI can be considered as a specific Markov Random Field:  $\phi_i(h_i, x_i) \equiv 1$  Constant modeling and geometric constraint force the LTF-V1 to focus on the nearby region, leading to the difficulty of long-range interactions

 $\psi_{m{i},m{j}}(h_{m{i}},h_{m{j}}) \coloneqq egin{cases} \delta(h_{m{i}}-h_{m{j}}) & h_{m{i}} 
otin \mathrm{Desc}_{\mathcal{G}_T}(m{i},m{j}) \ \exp(-\omega_{m{i},m{j}})\delta(h_{m{i}}-h_{m{j}}) & h_{m{i}} \in \mathrm{Desc}_{\mathcal{G}_T}(m{i},m{j}) \end{cases}$ 



#### Our Method

#### Proposed LTF-V2

We use data-dependent modeling for unary term:  $\phi_{i}(h_{i}, x_{i}) \coloneqq \begin{cases} f(x_{i}) & h_{i} = i \\ \exp(-\beta) & h_{i} \neq i \end{cases}$   $\psi_{i,j}(h_{i}, h_{j}) \coloneqq \begin{cases} \delta(h_{i} - h_{j}) & h_{i} \notin \text{Desc}_{\mathcal{G}_{T}}(i, j) \\ \exp(-\omega_{i,j})\delta(h_{i} - h_{j}) & h_{i} \in \text{Desc}_{\mathcal{G}_{T}}(i, j) \end{cases}$ 



#### **Our Method**

#### Proposed LTF-V2

We use data-dependent modeling for unary term:  $\phi_{i}(h_{i}, x_{i}) \coloneqq \begin{cases} f(x_{i}) & h_{i} = i \\ \exp(-\beta) & h_{i} \neq i \end{cases}$   $\psi_{i,j}(h_{i}, h_{j}) \coloneqq \begin{cases} \delta(h_{i} - h_{j}) & h_{i} \notin \operatorname{Desc}_{\mathcal{G}_{T}}(i, j) \\ \exp(-\omega_{i,j})\delta(h_{i} - h_{j}) & h_{i} \in \operatorname{Desc}_{\mathcal{G}_{T}}(i, j) \end{cases}$ 

We can derive the closed-form solution by using belief propagation algorithm:

$$y_{oldsymbol{i}} = rac{1}{z_{oldsymbol{i}}} \sum_{orall x_{oldsymbol{j} \in oldsymbol{X}}} S_{\mathcal{G}_{T}}(E_{oldsymbol{j},oldsymbol{i}}) \exp(-eta)^{|E_{oldsymbol{j},oldsymbol{i}}|} f(x_{oldsymbol{j}}) x_{oldsymbol{j}}$$

Learnable Tree Filter V2 (LTF-V2): relaxes the geometric constraint and enables efficient long-range interactions



### **Our Method**

### Framework of Learnable Tree Filter V2 Module The learnable spanning tree module enables fully end-to-end training:



**Algorithm 1:** Close random spanning tree **Input:** A 4-connected graph  $\mathcal{G}$ . **Output:** Random spanning tree  $\mathcal{G}_T$ . 1  $\mathcal{G}_T \leftarrow \emptyset$ . 2 for  $e \in E(\mathcal{G})$  do  $| l(e) \leftarrow e.$ ▲ 3 4 while |V(G)| > 1 do for  $v_i \in V(\mathcal{G})$  do 5  $\begin{bmatrix} e_i \sim E_{\mathcal{G}}(v_i). \\ \mathcal{G}_T \leftarrow \mathcal{G}_T \cup \{l(e_i)\}. \end{bmatrix}$ 6 7  $\operatorname{Contract}(E(\mathcal{G})).$ 8  $\operatorname{Flatten}(\mathcal{G}).$ 9 10 return  $\mathcal{G}_T$ .



#### Visualization on Instance Segmentation/Object Detection



Visualization of ground-truth (left), vanilla mask-rcnn (mid) and learnable tree filter (right) on COCO val set



#### Visualization on Semantic Segmentation



Visulaization of learnable tree filter (middle row) and ground-truth (bottom row) on VOC2012 val set



#### Ablation studies for instance-aware tasks

resource consumption on COCO dataset

Model	Stage	$AP_{box}$	$\mathrm{AP}^{50}_{box}$	$AP_{box}^{75}$	$AP_{seg}$	$AP_{seg}^{50}$	$AP_{seg}^{75}$	#FLOPs	#Params
ResNet-50 (1x)	-	38.8	58.7	42.4	35.2	55.6	37.6	279.4B	44.4M
+Non-Local [11] +CCNet [12] +LatentGNN [14] +GCNet [15]	4 345 345 All	39.5 40.1 40.6 40.7	59.6 60.4 61.3 61.0	42.7 44.1 44.5 44.2	35.6 36.0 36.6 36.7	56.7 57.4 58.1 58.1	37.6 38.4 39.2 39.2	+10.67B +16.62B +3.59B +0.35B	+2.09M +6.88M +1.07M +10.0M
+LTF-V1	345	40.0	60.4	43.7	36.1	57.5	38.4	+0.31B	+0.06M
+LTF-V2	3 4 5 345	40.1 40.6 40.2 <b>41.2</b>	59.9 61.0 60.5 <b>61.6</b>	43.9 44.4 43.6 <b>45.2</b>	36.0 36.6 36.1 <b>37.0</b>	57.1 58.2 57.5 <b>58.4</b>	38.3 39.0 38.4 <b>39.5</b>	+0.43B +0.26B +0.17B +0.68B	+0.02M +0.04M +0.08M +0.14M

#### NeurIPS 2020

#### Compared with previous work, our method (LTF-V2) achieves higher performance with much less

Ablation study on COCO2017 val set



#### Ablation studies for semantic segmentation

achieves superior performance on Cityscapes dataset

				Results on Cityscapes test set (fine only)					
Ablatio	n study on Citys	capes val set (fin	e only)	Model	Backbone	mIoU (%)			
Model	MS&Flip	mIoU (%)	mAcc (%)	PSPNet [9]	ResNet-101	78.4			
ResNet-50	X	72.9	95.5	DFN [39]	ResNet-101	79.3			
+LTF-V1	×	75.9	95.8	DenseASPP [50]	DenseNet-161	80.6			
+LTF-V2	×	77.4	<b>96.0</b>	LTF-V1 [16]	ResNet-101	80.8			
ResNet-50		75 5	95.9	CCNet [12]	ResNet-101	81.4			
+LTF-V1	· _	77.2	96.0	DANet [51]	ResNet-101	81.5			
+LTF-V2	✓	<b>78.9</b>	<b>96.1</b>	SPNet [52]	ResNet-101	82.0			
				Ours (LTF-V2)	ResNet-101	82.1			

#### NeurIPS 2020

## Compared with the state-of-the-art work, our method (LTF-V2) based on a simple FPN architecture



#### **Empirical Runtime**

#### For GPU devices, we parallelize the algorithm along with batches, channels, and nodes of the same depth.



#### NeurIPS 2020

### Learnable Tree Filter is easy to use. You only need to add 2 lines to your PyTorch code.



### Future Work

#### Potential Applications

- Replacing the attention module of transformer for natural language processing
- Efficiently modeling higher-order relations, e.g., solving maze problem
- Enhancing sequential representation for video analysis



#### NeurIPS 2020

ner for natural language processing e.g., solving maze problem <mark>eo analysis</mark>







![](_page_18_Picture_0.jpeg)

For any questions, please contact

stevengrove@stu.xjtu.edu.cn

![](_page_18_Picture_4.jpeg)

https://github.com/Megvii-BaseDetection/TreeFilter-Torch https://github.com/StevenGrove/LearnableTreeFilterV2

#### NeurIPS 2020

## Thanks

#### Source code is available

![](_page_18_Picture_9.jpeg)