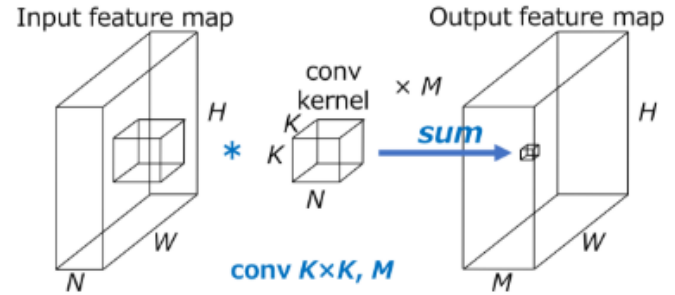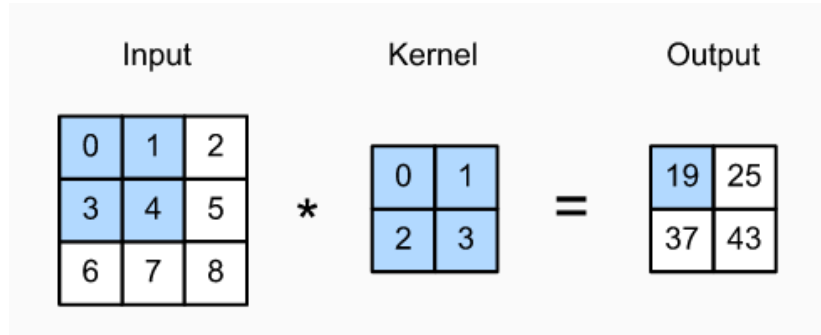# *Pixel-Adaptive Convolutional Neural Networks (CVPR 2019)*
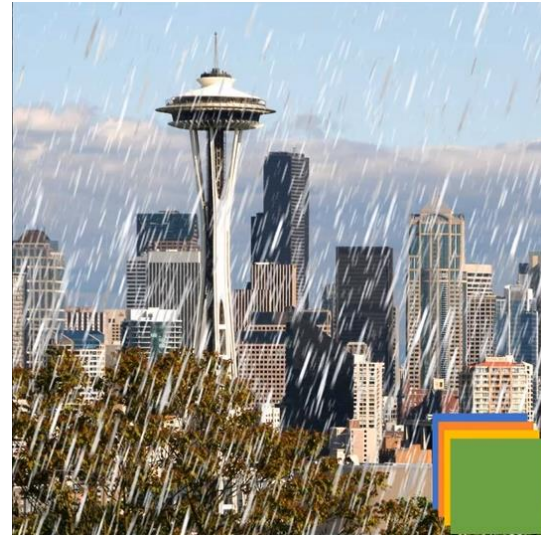
Hang Su[1], Varun Jampani[2], Deqing Sun[2], Orazio Gallo[2], Erik Learned-Miller[1], and Jan Kautz[2]

[1]UMass Amherst    [2]NVIDIA

# *standard convolution*



Input

| 0 | 1 | 2 |
|---|---|---|
| 3 | 4 | 5 |
| 6 | 7 | 8 |

*

Kernel

| 0 | 1 |
|---|---|
| 2 | 3 |

=

Output

| 19 | 25 |
|----|----|
| 37 | 43 |

Input feature map

$H$

$W$

$N$

conv kernel

$\times M$

$K$

$K$

$N$

*

**sum**

**conv K×K, M**
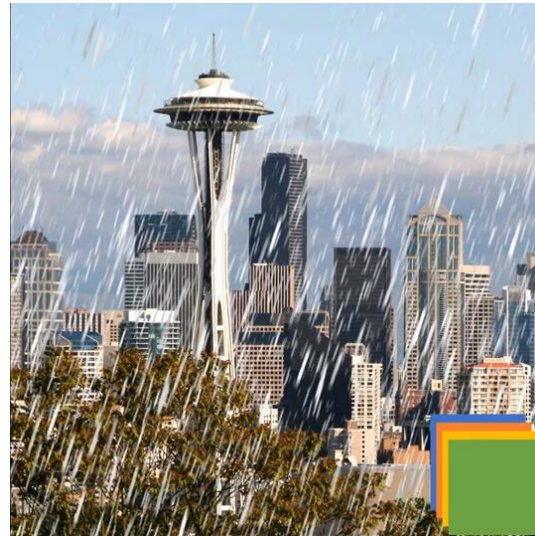
Output feature map

$H$

$W$

$M$

# *Convolution is spatially-shared*

# Convolution is content-agnostic
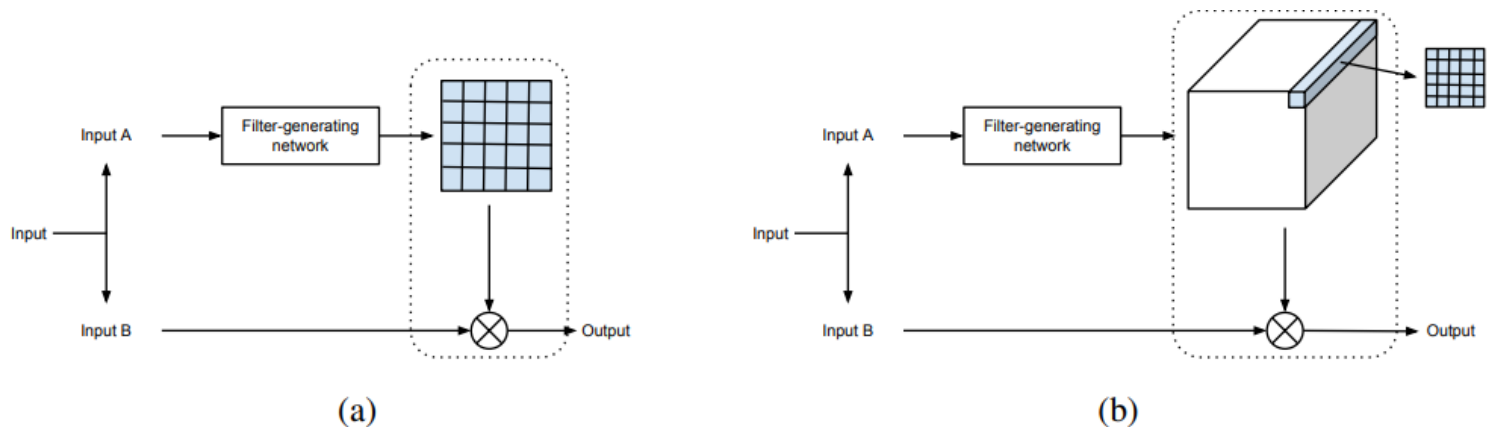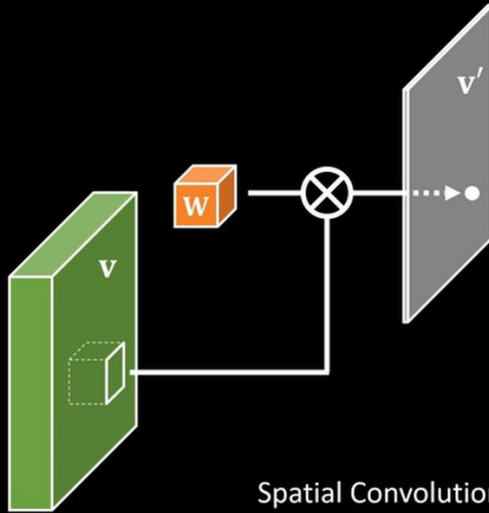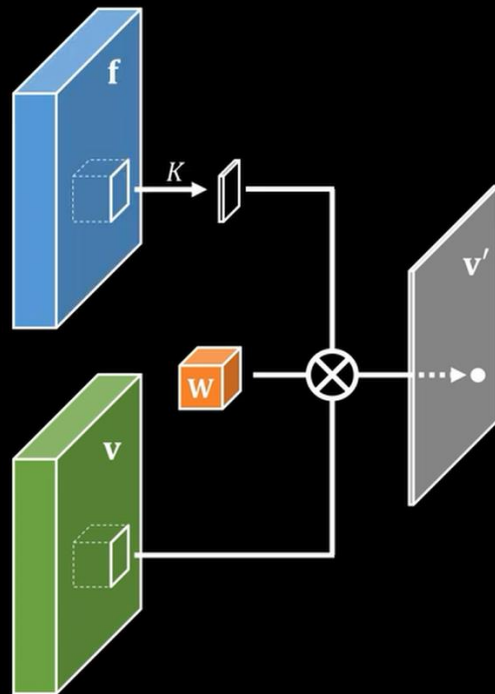
# Dynamic Filter Networks  (NIPS 2016)



Figure 2: *Left:* Dynamic convolution: the filter-generating network produces a single filter that is applied convolutionally on $I_B$. *Right:* Dynamic local filtering: each location is filtered with a location-specific dynamically generated filter.

**v** input features

**v′** output features
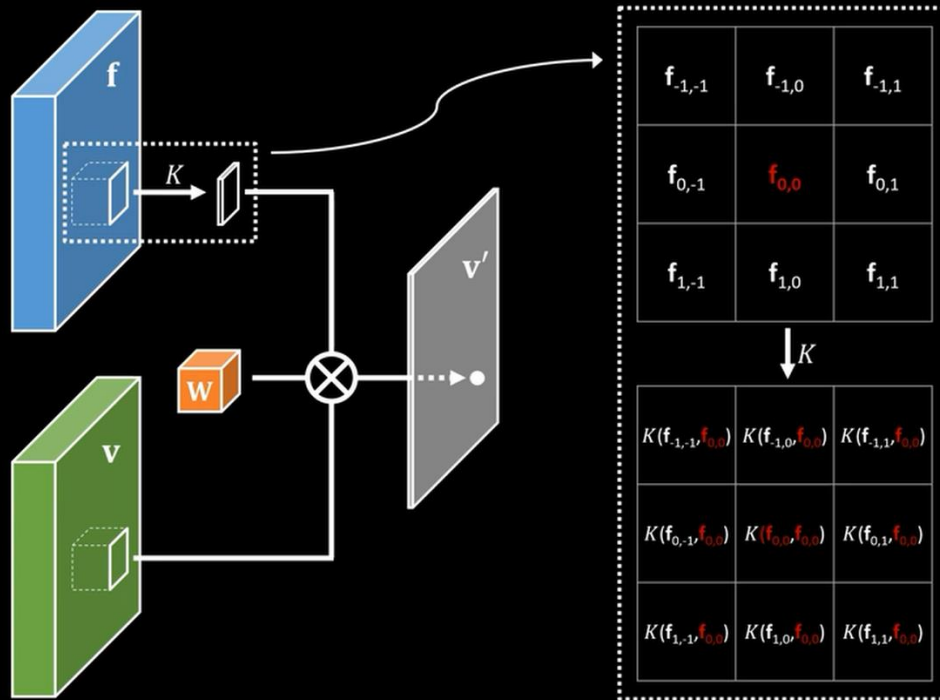
p $(x, y)$ coordinates

**W** filter weights

Spatial Convolution $\qquad \mathbf{v}_i' = \sum_{j \in \Omega(i)} \mathbf{W}[\mathrm{p}_i - \mathrm{p}_j]\mathbf{v}_j + \mathbf{b}$

**v**   input features

**v′**  output features

p   $(x, y)$ coordinates

**W**  filter weights

**f**   adapting features

v input features

v′ output features

p $(x, y)$ coordinates

W filter weights

f adapting features

$K$ adapting kernel

f

$K$

v′

W

v

$f_{-1,-1}$ | $f_{-1,0}$ | $f_{-1,1}$
$f_{0,-1}$ | $f_{0,0}$ | $f_{0,1}$
$f_{1,-1}$ | $f_{1,0}$ | $f_{1,1}$

$\downarrow K$

$K(f_{-1,-1}, f_{0,0})$ | $K(f_{-1,0}, f_{0,0})$ | $K(f_{-1,1}, f_{0,0})$
$K(f_{0,-1}, f_{0,0})$ | $K(f_{0,0}, f_{0,0})$ | $K(f_{0,1}, f_{0,0})$
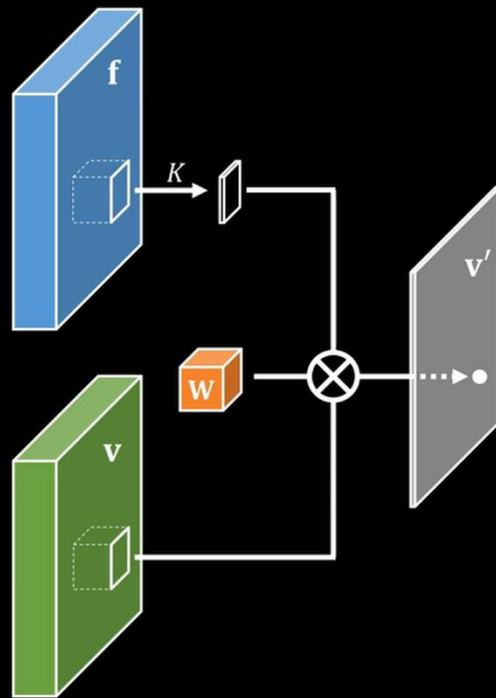$K(f_{1,-1}, f_{0,0})$ | $K(f_{1,0}, f_{0,0})$ | $K(f_{1,1}, f_{0,0})$

# Pixel Adaptive Convolution (PAC)

$$\mathbf{v}'_i = \sum_{j \in \Omega(i)} K(\mathbf{f}_i, \mathbf{f}_j)\mathbf{W}[\mathrm{p}_i - \mathrm{p}_j]\mathbf{v}_j + \mathbf{b}$$

Spatial Convolution
$$\mathbf{v}'_i = \sum_{j \in \Omega(i)} \mathbf{W}[\mathrm{p}_i - \mathrm{p}_j]\mathbf{v}_j + \mathbf{b}$$
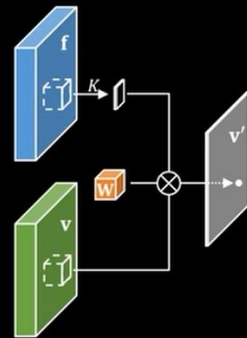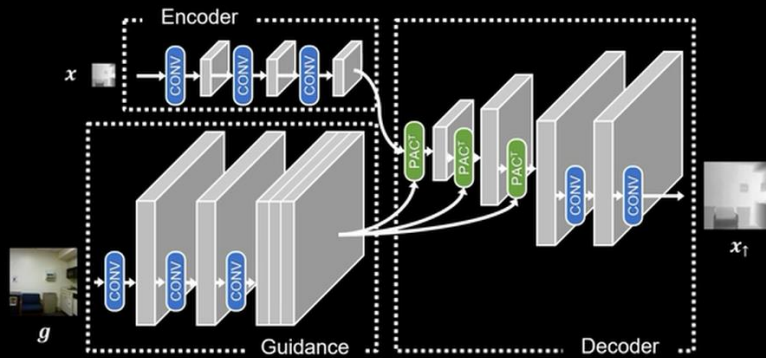


9

$$\mathbf{v}'_i = \sum_{j \in \Omega(i)} K(\mathbf{f}_i, \mathbf{f}_j) \mathbf{W}[\mathrm{p}_i - \mathrm{p}_j] \mathbf{v}_j + \mathbf{b}$$
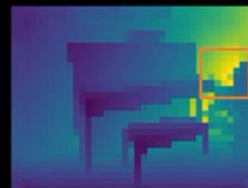
- Spatial convolution $\quad K(\mathbf{f}_i, \mathbf{f}_j) = 1$

- Bilateral filtering $\quad \mathbf{f} = (r, g, b), K(\mathbf{f}_i, \mathbf{f}_j) = \exp(-\frac{1}{2\alpha_1} \|\mathbf{f}_i - \mathbf{f}_j\|^2)$

  $$\mathbf{W}[\mathrm{p}_i - \mathrm{p}_j] = \exp(-\frac{1}{2\alpha_2} \|\mathrm{p}_i - \mathrm{p}_j\|^2)$$

- Average pooling $\quad K(\mathbf{f}_i, \mathbf{f}_j) = 1, \mathbf{W}[\mathrm{p}_i - \mathrm{p}_j] = \frac{1}{F^2}$

- Detail-preserving pooling [1] $\quad K(\mathbf{f}_i, \mathbf{f}_j) = \alpha + (\|\mathbf{f}_i - \mathbf{f}_j\|^2 + \epsilon^2)^\lambda$
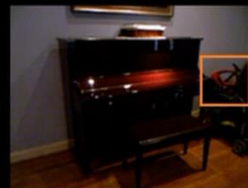
PAC generalizes many existing filtering techniques

[1] F. Saeedan, et al. Detail Preserving Pooling in Deep Networks. CVPR '18.

Joint upsampling network with PAC

Encoder · Guidance · Decoder

$x$ · $g$ · $x_\uparrow$

$x$ (low-res input)

$g$ (guidance)

$x_\uparrow$ (Ours)

bilinear baseline

Joint depth upsampling (16x)    Joint optical flow upsampling (16x)

# PAC for efficient CRF inference



Input · · · Unaries → Softmax → | PAC dilation=16 / PAC dilation=64 → ⊕ → Softmax | · · · → | PAC dilation=16 / PAC dilation=64 → ⊕ → Softmax | → Prediction

$MF_1$  $MF_T$  PAC-CRF

dilation=16  dilation=64

| Input | GT | Unary | Full-CRF [1] | Ours |

[1] Philipp Krähenbühl and Vladlen Koltun. Efficient Inference in Fully Connected CRFs with Gaussian Edge Potentials. NIPS '11.
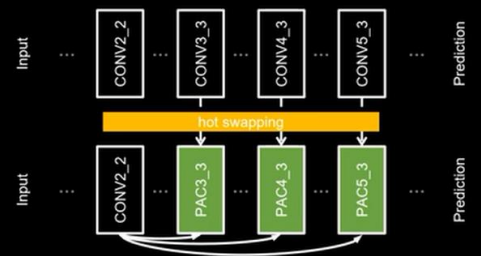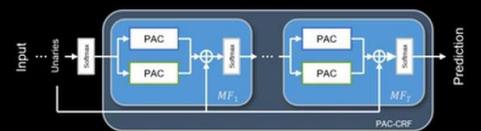
12

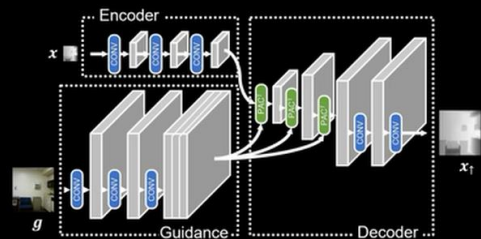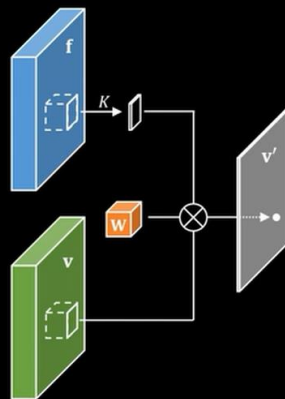# Layer "hot-swapping" with PAC



| Method | CRF | mIoU | Runtime |
|---------|:---:|:-----:|:--------:|
| FCN | | 67.20 | 39ms |
| FCN | ✓ | 69.82 | 117ms |
| PAC-FCN | | 69.18 | 41ms |
| PAC-FCN | ✓ | 71.34 | 118ms |

*evaluated on Pascal VOC 2012 test set

# Summary

- Pixel Adaptive Convolution Sec.3:
  - Content-adaptive
  - Generalizes several existing filtering techniques
- Three use cases:
  - Joint upsampling networks Sec. 4
  - Efficient CRF inference Sec. 5
  - Network layer hot-swapping Sec. 6

```python
in_ch, g_ch = 16, 8                    # channel sizes of input and guidance
stride, f, b, h, w = 5, 2, 64, 64      # stride, filter size, batch size, input height and width
input = torch.rand(b, in_ch, h, w)
guide = torch.rand(b, g_ch, h, w)      # guidance feature


pool = nn.AvgPool2d(f, stride)
out_pool = pool(input)                 # standard spatial convolution


pacpool = PacPool2d(f, stride)
out_pac = pacpool(input, guide)        # PAC
out_pac = pacpool(input, None, guide_k) # alternative interface
                                       # guide_k is pre-computed 'K'
                                       # of shape [b, g_ch, f, f, h, w]. packernel2d can be
                                       # used for its creation.
```