



UNIVERSITÀ
DEGLI STUDI
FIRENZE

DIEF

Dipartimento di
Ingegneria Industriale

Neural Networks for Semantic Segmentation

Machine Learning Course

Antonio Castellucci
Michela Crulli

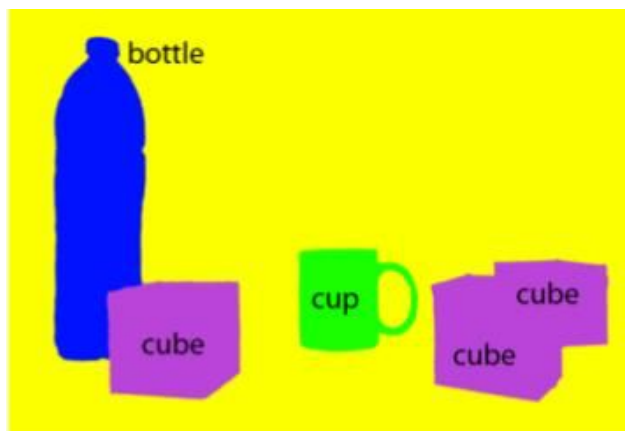
- Introduction
 - Semantic segmentation
 - Applications
 - Neural networks for semantic segmentation
- DUC + HDC
 - Approach
 - Results
- Deeplab V3+
 - Approach
 - Results
- ICNet
 - Approach
 - Results
- NAS
 - Introduction
 - DARTS
 - Auto-DeepLab
- Conclusions

- Semantic Segmentation

- The goal of semantic image segmentation is to label *each pixel* of an image with a corresponding class of what is being represented where:

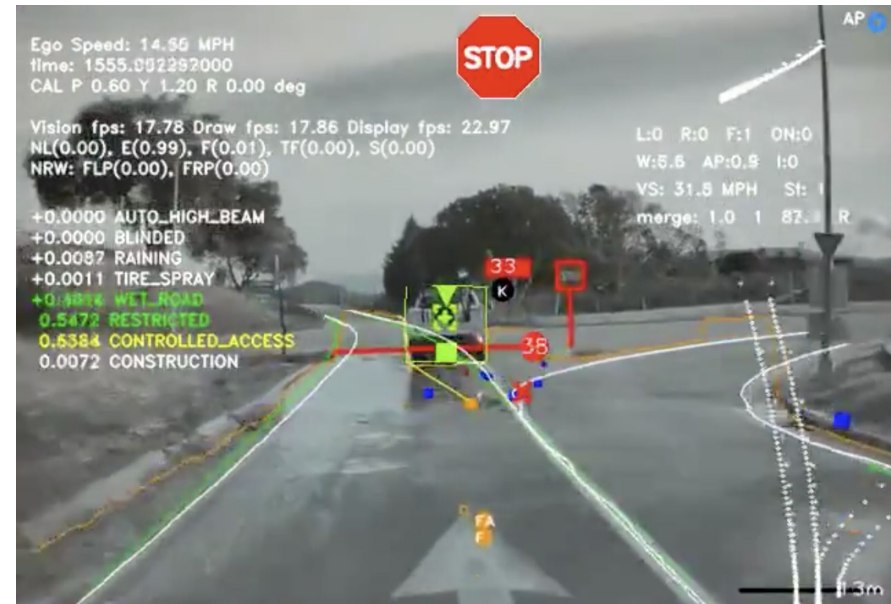
$$L = \{l_1, l_2, \dots, l_k\} \quad X = \{x_1, x_2, \dots, x_N\}$$

X is a **2D** image of $W \times H = N$ pixels x



• Applications

- GeoSensing
- Autonomous driving
- Facial Segmentation
- Precision Agriculture



- Neural Networks for semantic segmentation

Before:

- Features hand-written
- Random Forest, Boosting
- Prediction of only central pixel of a patch

With gpu's growth:

- Neural Networks
- Convolutional Neural Networks



- Neural Networks for semantic segmentation

3 key components:

- A fully convolutional network (FCN)
- Conditional Random Fields (CRF)
- Dilated convolution (Atrous convolution)

Common improvements:

- Deeper FCN models
- More powerful CRFs

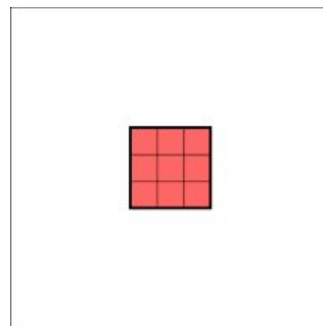


Understanding Convolution for Semantic Segmentation

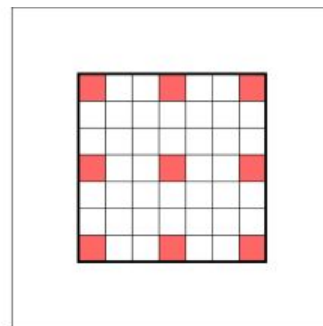
(Panqu Wang, Pengfei Chen, Ye Yuan, Ding Liu, Zehua Huang, Xiaodi Hou, Garrison Cottrel, 2018)

- Possible solution for denser feature maps:

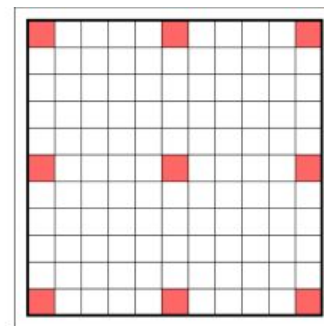
- Atrous convolution:



Kernel: 3x3
Dilation rate: 1



Kernel: 3x3
Dilation rate: 3

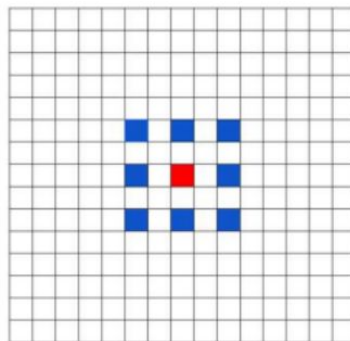


Kernel: 3x3
Dilation rate: 5

$$y[i] = \sum_k x[i + r \cdot k] w[k]$$

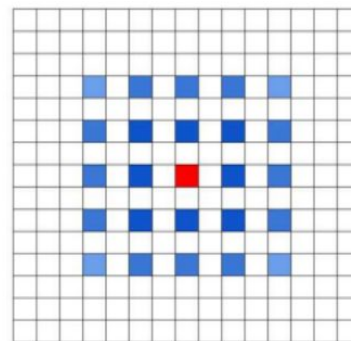
- Problem:

- Gridding Problem

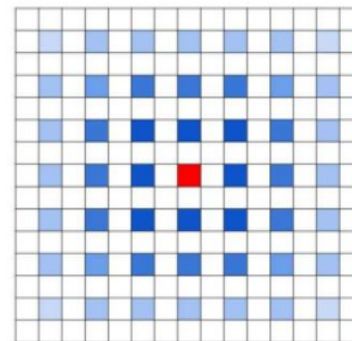


K = 3

r = 2

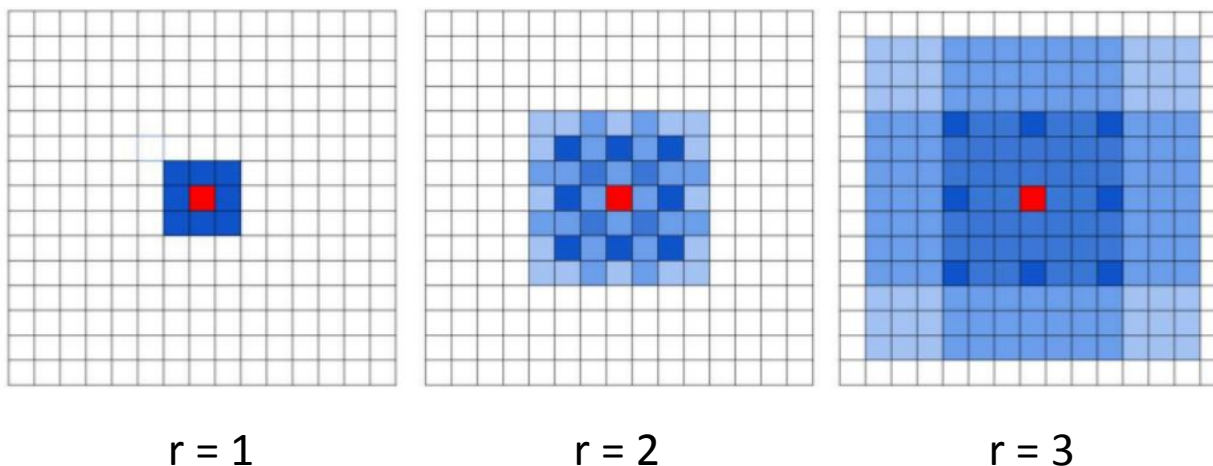


r = 2



r = 2

- Solution: Hybrid Dilated Convolution (HDC)
 - different dilatation rate r

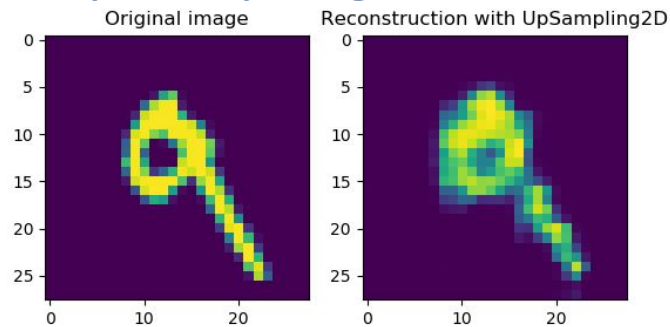


- maximum distance between two nonzero values

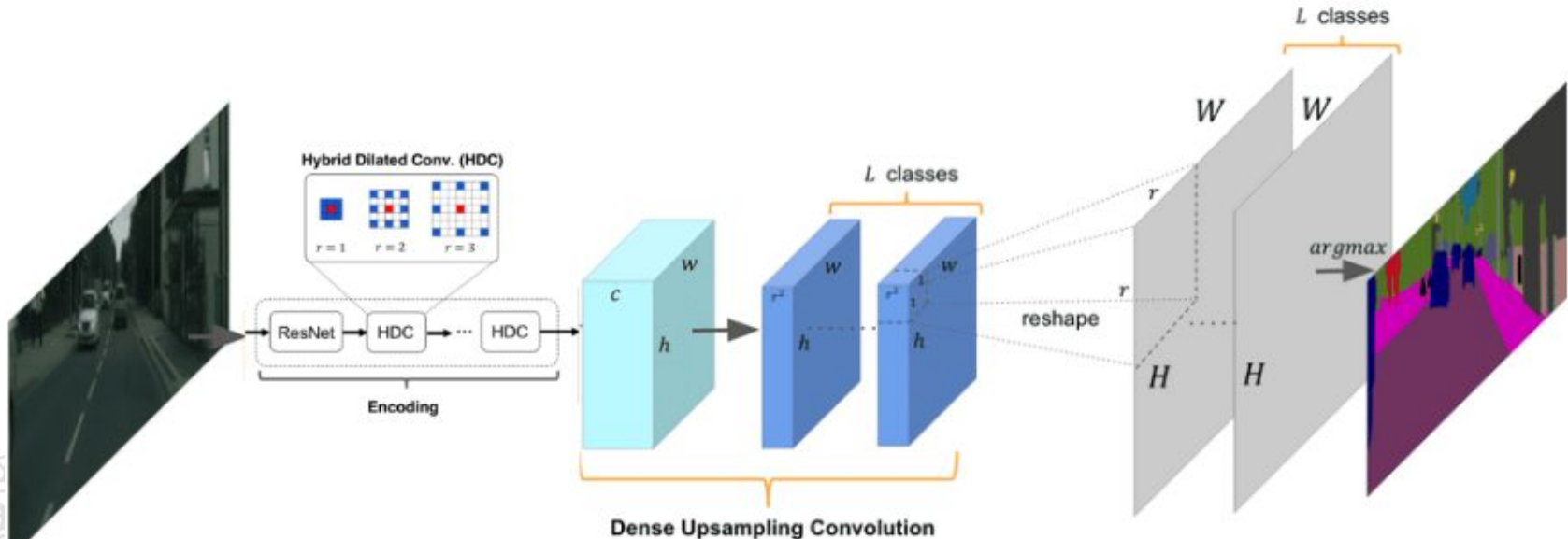
$$M_i = \max[M_{i+1} - 2r_i, M_{i+1} - 2(M_{i+1} - r_i), r_i]$$

$$M_n = r_n, \quad M_2 \leq K$$

- Problem: Bilinear Upsampling

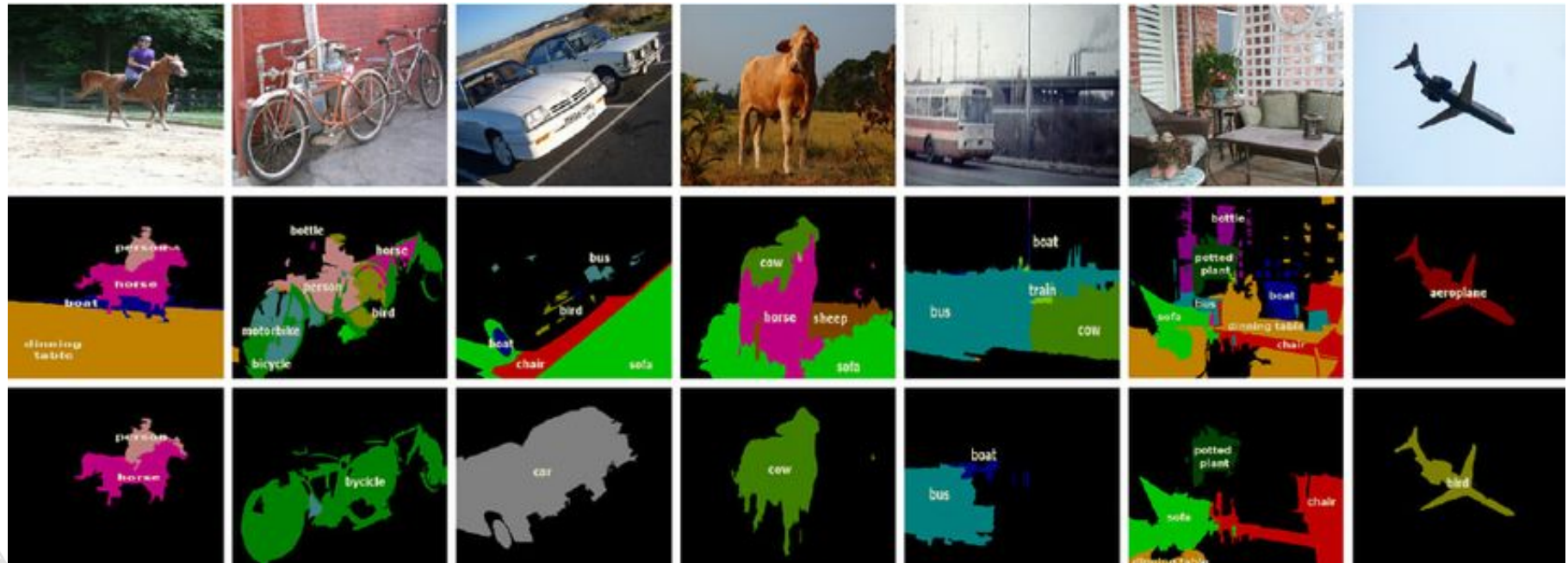


- Solution: Dense Upsampling Convolution (DUC)



- Dataset: PASCAL VOC2012

- 20 foreground object classes
- 1 background class
- 1456 training images, 1449 validation images and 1456 test images
- Usually, the training set is augmented to have 10582 images



- Dataset: Cityscapes

- Street scenes from 50 different cities
- High quality pixel-level annotations of 5 000 frame
- 20 000 weakly annotated frames
- 19 classes



- Dataset: Cityscapes

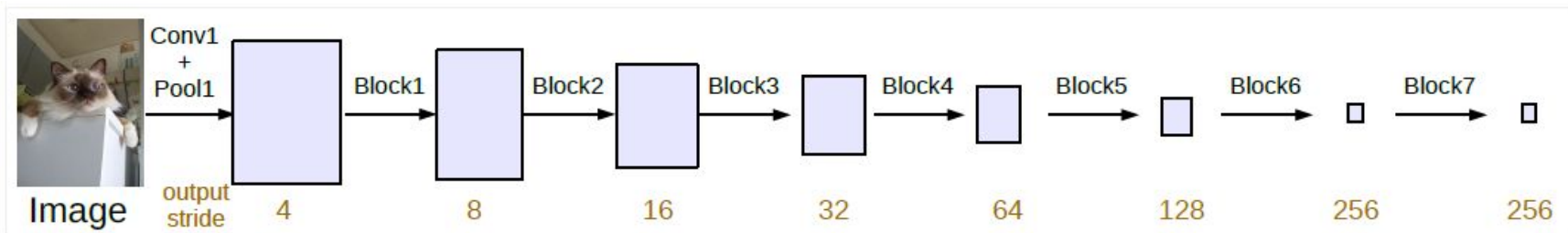
Method	mxNet	mIoU
Custom DUC+HDC	0.11.0	81.9
Original DUC+HDC	1.5.1	80.1

Encoder-Decoder with Atrous Separable Convolution for Semantic Image Segmentation

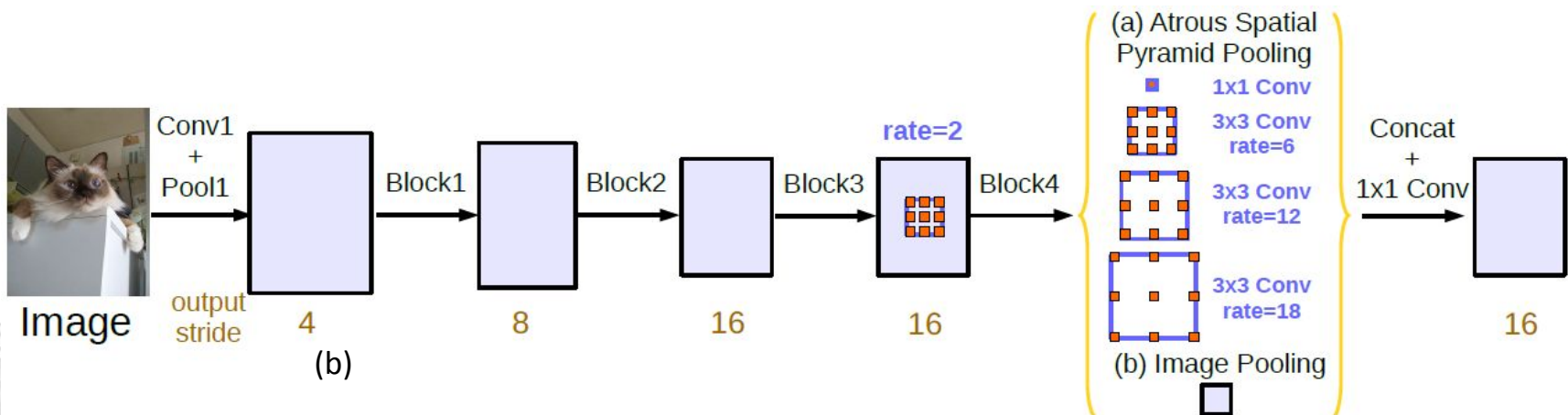
(Liang- Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, Hartwig Adam, 2018)

• Deeplab V3

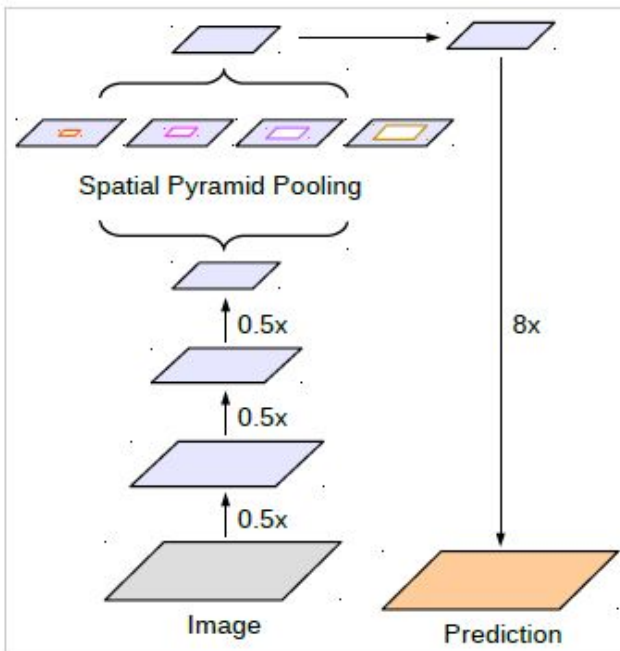
- Going deeper without atrous convolution: location/spatial information is lost at the deeper layers.



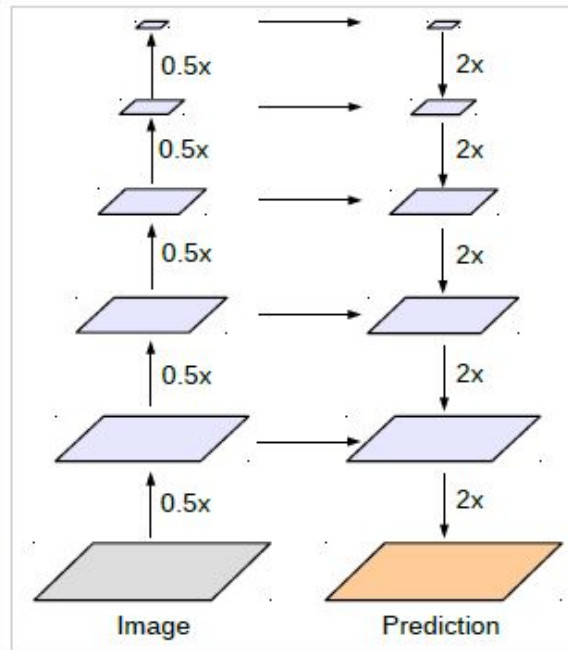
- Going deeper with atrous convolution: information lost in decoding phase.



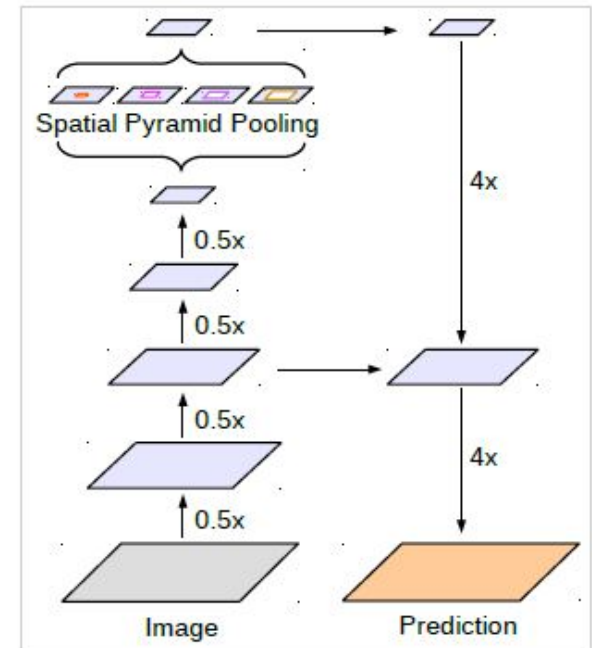
- Deeplab V3+



(a) Spatial Pyramid Pooling

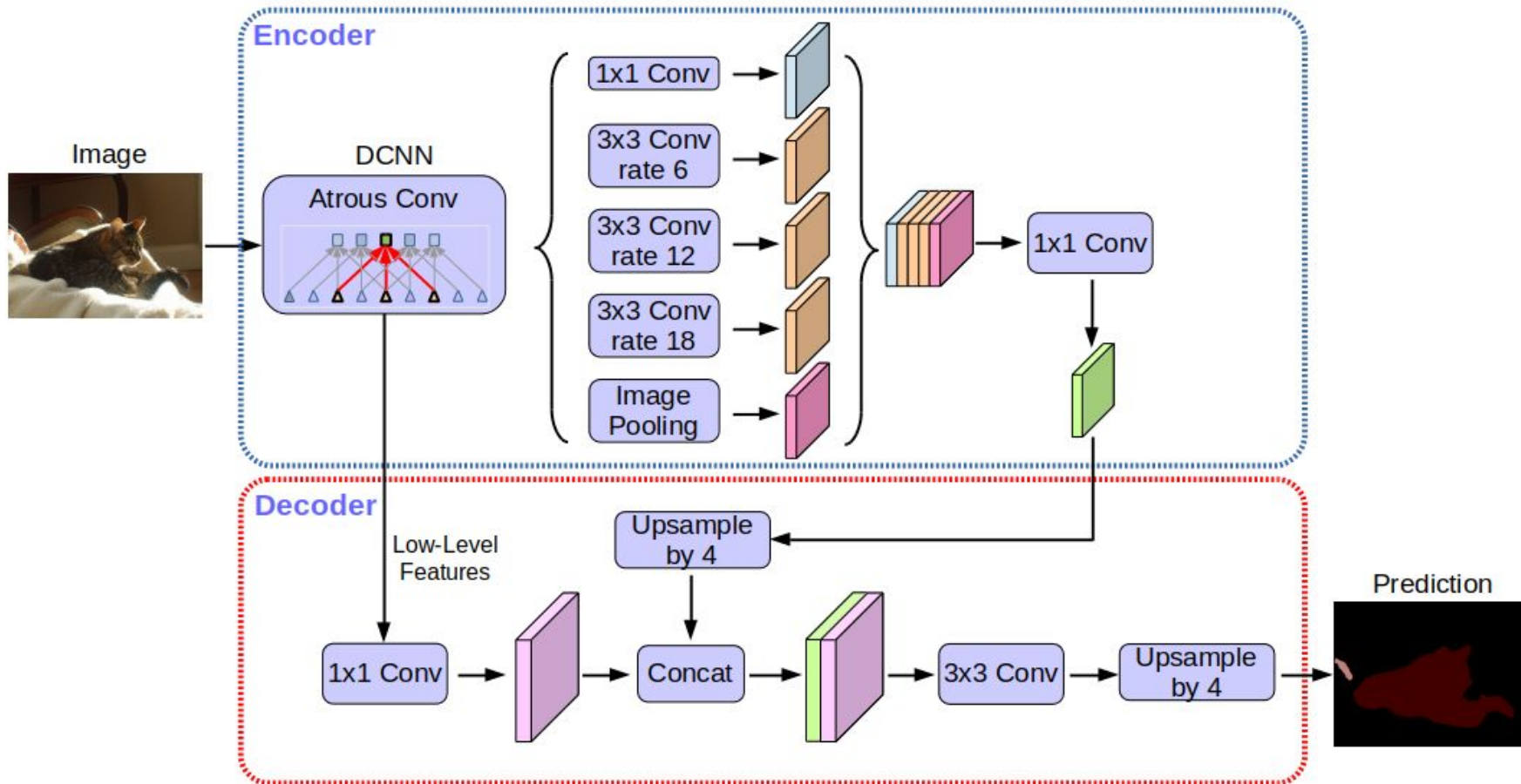


(b) Encoder-Decoder



(c) Encoder-Decoder with Atrous Con

- Deeplab V3+



- Dataset: Pascal VOC 2012

Method	Backbone	OS	mIoU
Custom Tensorflow	ResNet 101	16	77.31
Original Tensorflow	ResNet 101	16	78.85
Custom Pytorch	ResNet 101	16	78.31
Custom Pytorch	Mobilnet	16	70.89

- Dataset: Cityscapes

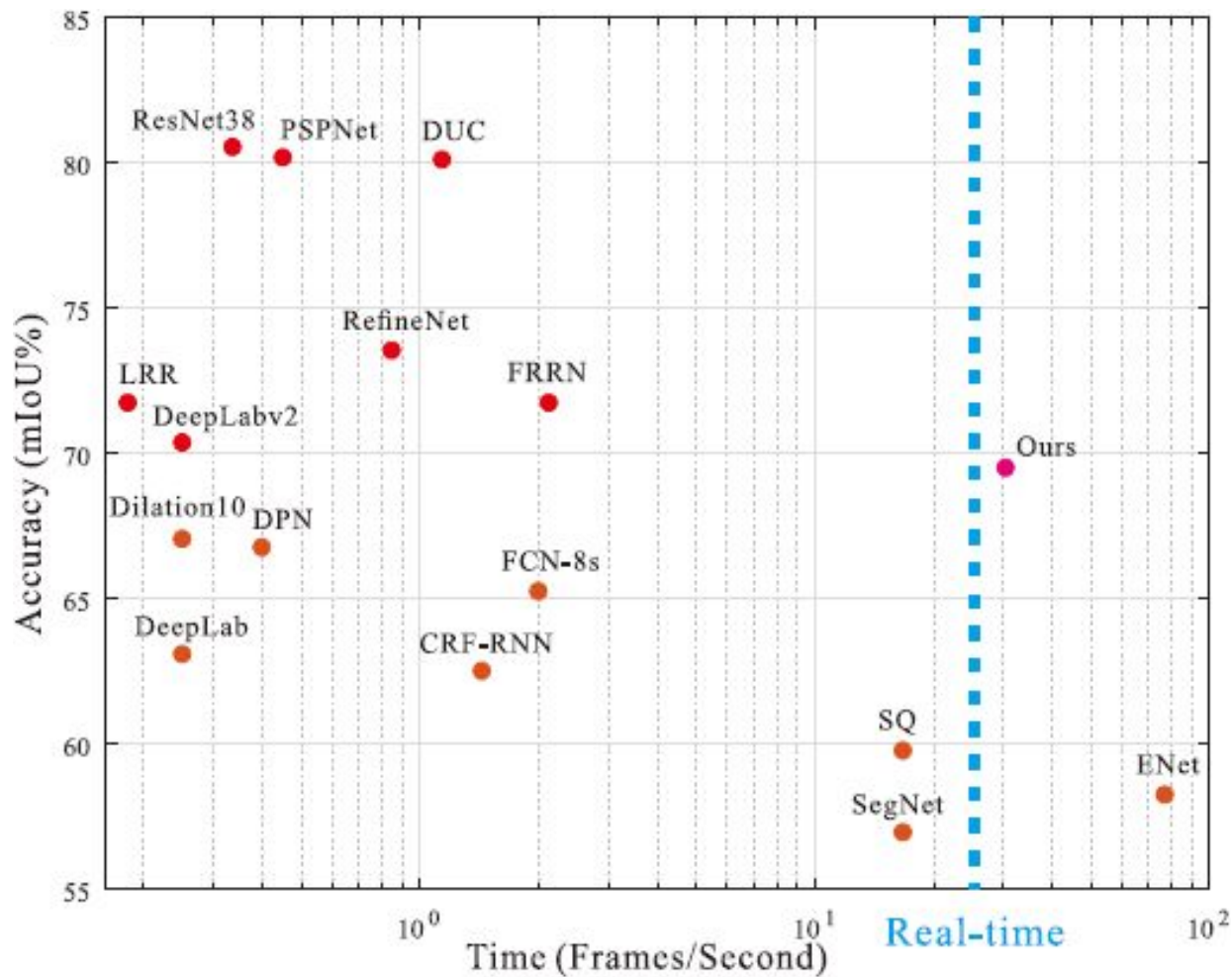
Method	Backbone	OS	mIoU
Custom Pytorch	Mobilnet	16	72.1
Original Tensorflow	ResNet 101	16	81.9

ICNet for Real-Time Semantic Segmentation on High-Resolution Images

(Hengshuang Zhao, Xiaojuan Qi, Xiaoyong Shen, Jianping Shi, Jiaya Jia, 2018)



- Quality Vs Speed

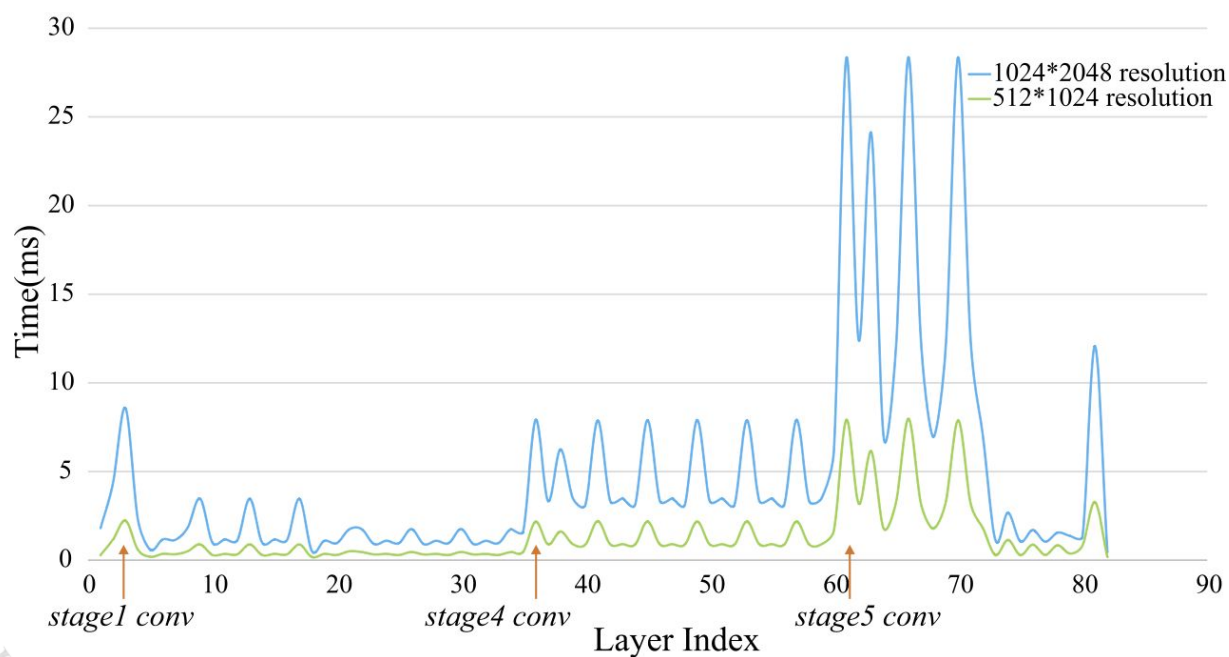


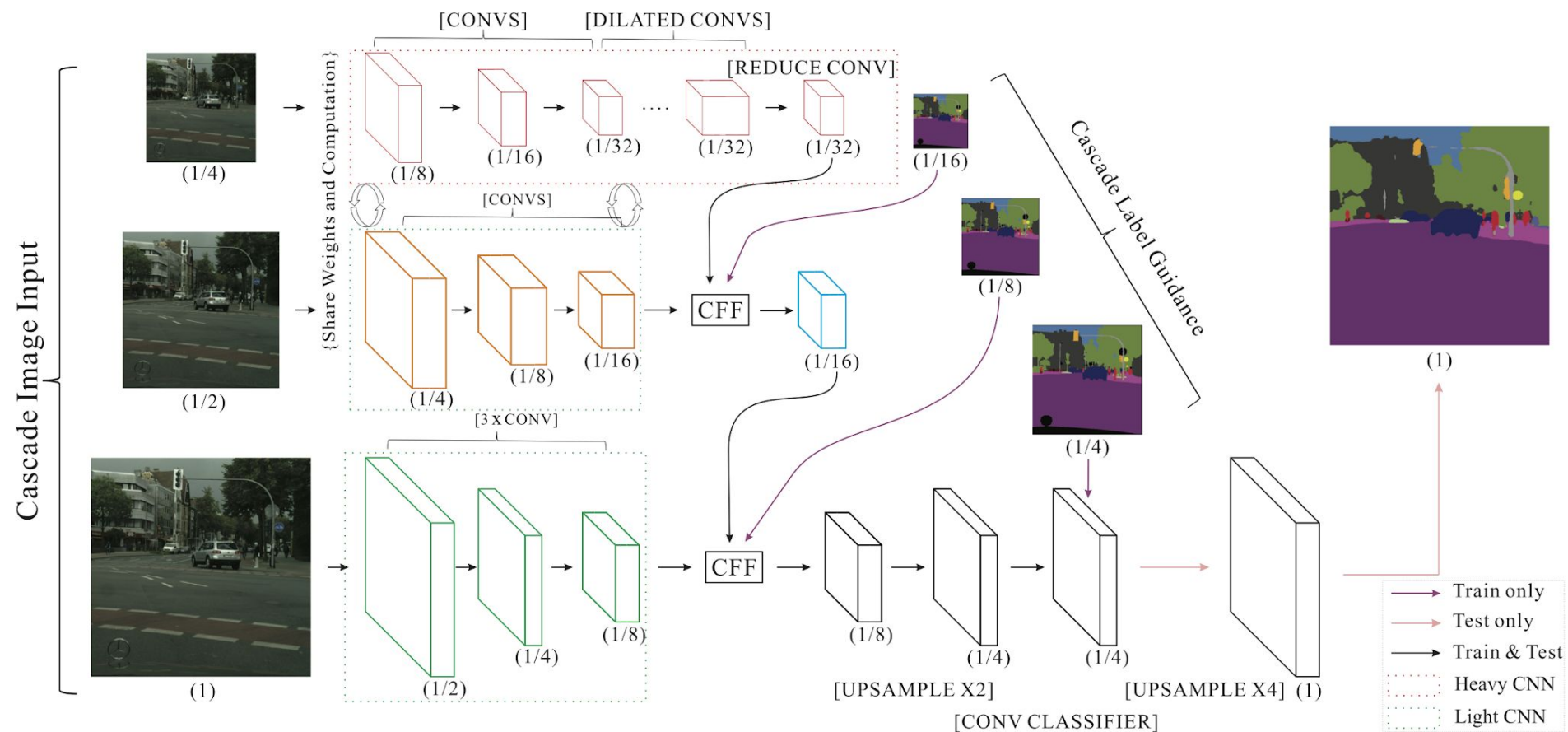
- Convolution Speed Analysis (PSPNet50)

$$\Phi: V \rightarrow U$$

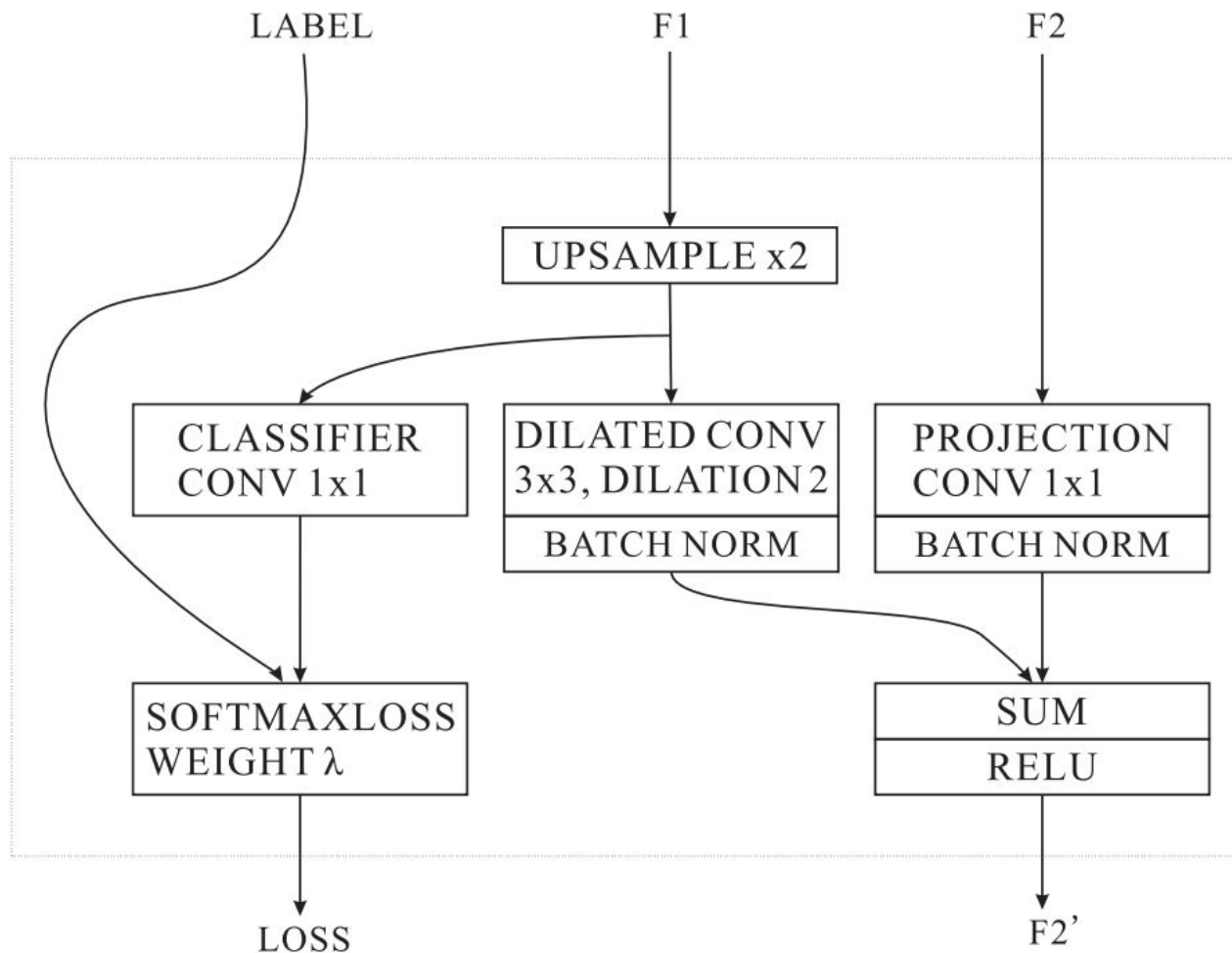
$$V \in \mathbb{R}^{c \times h \times w} \quad U \in \mathbb{R}^{c' \times h' \times w'} \quad h' = h/s \quad w' = w/s$$

$$O(\Phi) \approx c'ck^2hw/s^2$$



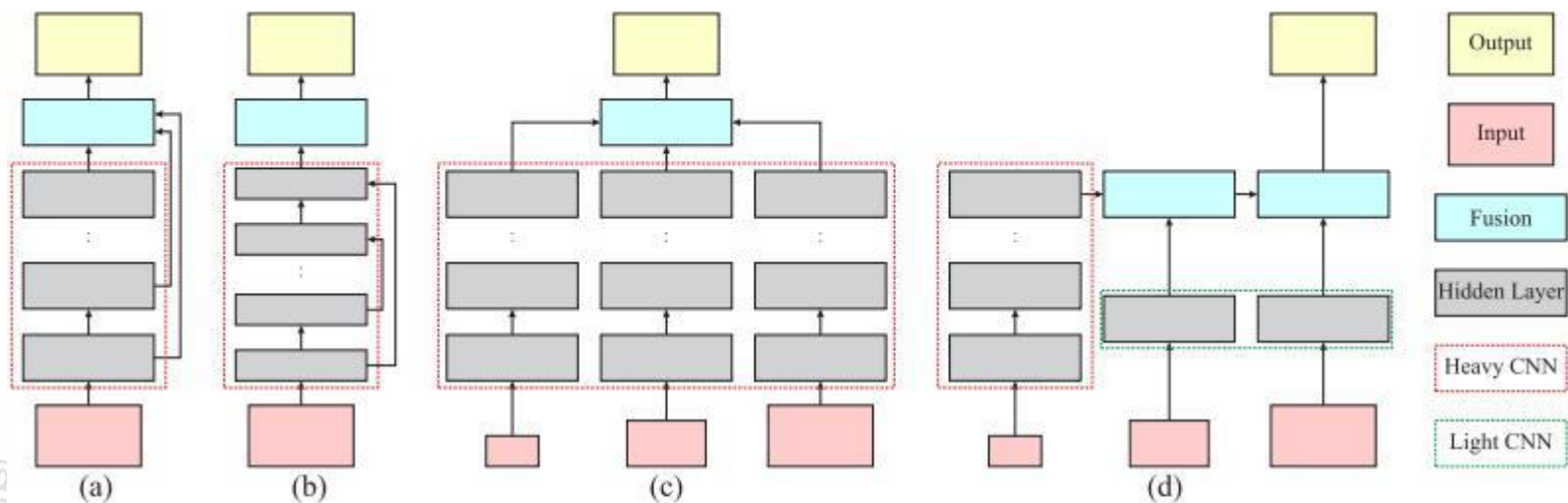


- Cascade Feature Fusion (CFF)



- Cascade Label Guidance (CLG)

$$\mathcal{L} = - \sum_{t=1}^{\mathcal{T}} \lambda_t \frac{1}{y_t x_t} \sum_{y=1}^{y_t} \sum_{x=1}^{x_t} \log \frac{e^{\mathcal{F}_{\hat{n}, y, x}^t}}{\sum_{n=1}^{\mathcal{N}} e^{\mathcal{F}_{n, y, x}^t}}$$



- Dataset: Cityscapes
 - Experiments on gtx 1070 max-q

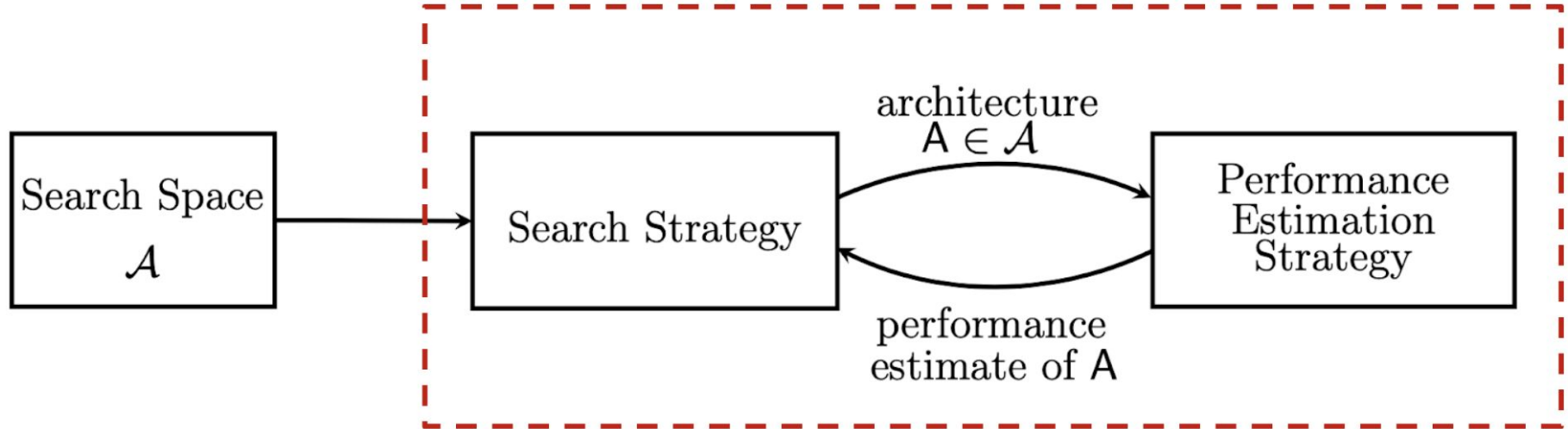
Method	Model	DR	mIoU	Time(ms)	Frame(fps)
ICNetC	30k_bn	no	67,35	66	15,14
ICNetC	90k_bn	no	80,08	69	14,33

- Experiments on TitanX Maxwell

Method	Model	DR	mIoU	Time(ms)	Frame(fps)
ICNetO	30k	no	67,7	33	30,3
ICNetO	90k	no	70,6	33	30,3



One-shot approach:
learning model architecture parameters and weights together

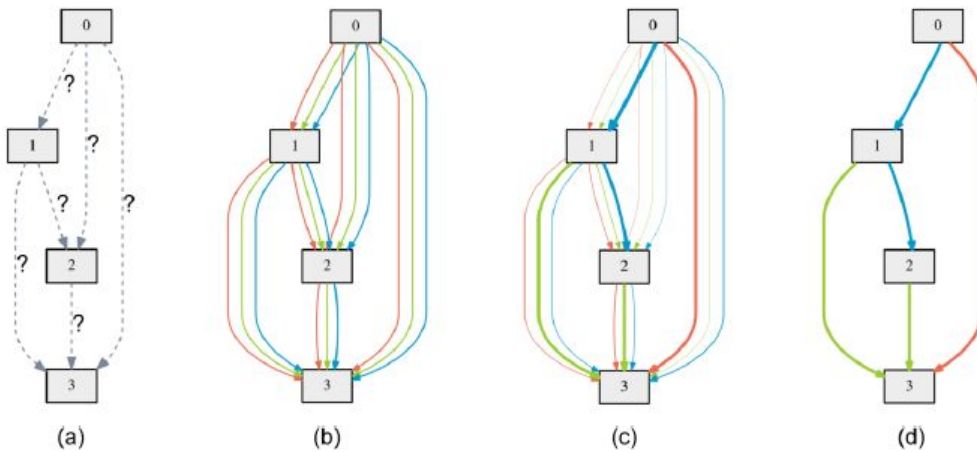


DARTS: Differentiable Architecture Search

(Hanxiao Liu, Karen Simonyan, Yiming Yang 2019)



Cell Architecture

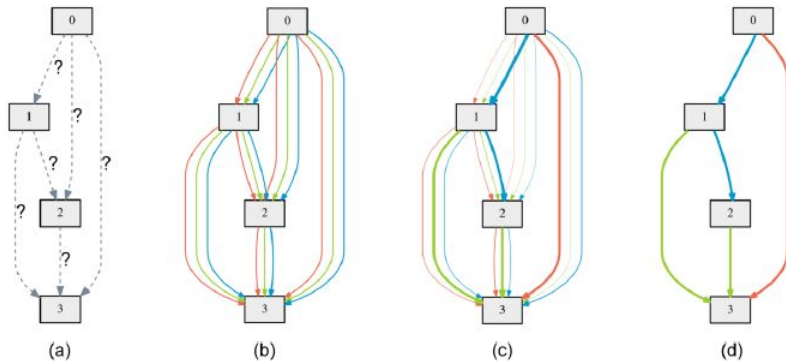


- 3×3 depthwise-separable conv
- 5×5 depthwise-separable conv
- 3×3 atrous conv with rate 2
- 5×5 atrous conv with rate 2
- 3×3 average pooling
- 3×3 max pooling
- skip connection
- no connection (zero)

Each intermediate node is computed based on all of its predecessors:

$$x^{(j)} = \sum_{i < j} o^{(i,j)}(x^{(i)})$$

DARTS: Continuous Relaxation and Optimization



Softmax over all possible operations:

$$\bar{o}^{(i,j)}(x) = \sum_{o \in \mathcal{O}} \frac{\exp(\alpha_o^{(i,j)})}{\sum_{o' \in \mathcal{O}} \exp(\alpha_{o'}^{(i,j)})} o(x)$$

A discrete architecture can be obtained by replacing each mixed operation $\bar{o}^{(i,j)}$ with the most likely operation:

$$o^{(i,j)} = \operatorname{argmax}_{o \in \mathcal{O}} \alpha_o^{(i,j)}$$

$$\min_{\alpha} \quad \mathcal{L}_{val}(w^*(\alpha), \alpha)$$

$$\text{s.t.} \quad w^*(\alpha) = \operatorname{argmin}_w \mathcal{L}_{train}(w, \alpha)$$

- First Order Approximation

$$\begin{aligned} \min_{\alpha} \quad & \mathcal{L}_{val}(w^*(\alpha), \alpha) \\ \text{s.t.} \quad & w^*(\alpha) = \operatorname{argmin}_w \mathcal{L}_{train}(w, \alpha) \end{aligned}$$

$$\nabla_{\alpha} \mathcal{L}_{val}(w^*(\alpha), \alpha) \approx \nabla_{\alpha} \mathcal{L}_{val}(w - \xi \nabla_w \mathcal{L}_{train}(w, \alpha), \alpha) \quad (1)$$

Algorithm 1: DARTS – Differentiable Architecture Search

Create a mixed operation $\bar{o}^{(i,j)}$ parametrized by $\alpha^{(i,j)}$ for each edge (i, j)

while *not converged* **do**

1. Update weights w by descending $\nabla_w \mathcal{L}_{train}(w, \alpha)$
2. Update architecture α by descending $\nabla_{\alpha} \mathcal{L}_{val}(w - \xi \nabla_w \mathcal{L}_{train}(w, \alpha), \alpha)$

Replace $\bar{o}^{(i,j)}$ with $o^{(i,j)} = \operatorname{argmax}_{o \in \mathcal{O}} \alpha_o^{(i,j)}$ for each edge (i, j)

- Second Order Approximation

Applying chain rule to the approximate architecture gradient (eq. 1) yields:

$$\nabla_{\alpha} \mathcal{L}_{val}(w', \alpha) - \xi \nabla_{\alpha, w}^2 \mathcal{L}_{train}(w, \alpha) \nabla_{w'} \mathcal{L}_{val}(w', \alpha) \quad (2)$$

where: $w' = w - \xi \nabla_w \mathcal{L}_{train}(w, \alpha)$

Let ϵ be a small scalar and $w^{\pm} = w \pm \epsilon \nabla_{w'} \mathcal{L}_{val}(w', \alpha)$

Then:

$$\nabla_{\alpha, w}^2 \mathcal{L}_{train}(w, \alpha) \nabla_{w'} \mathcal{L}_{val}(w', \alpha) \approx \frac{\nabla_{\alpha} \mathcal{L}_{train}(w^+, \alpha) - \nabla_{\alpha} \mathcal{L}_{train}(w^-, \alpha)}{2\epsilon}$$

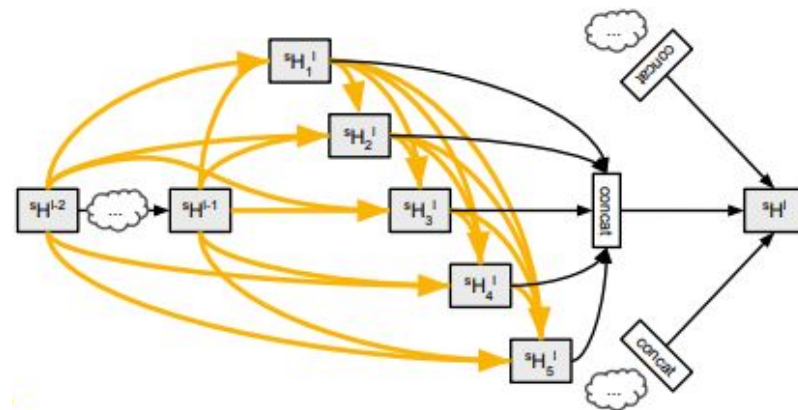
$$O(|\alpha||w|) \text{ to } O(|\alpha| + |w|)$$

Auto-DeepLab: Hierarchical Neural Architecture Search for Semantic Image Segmentation

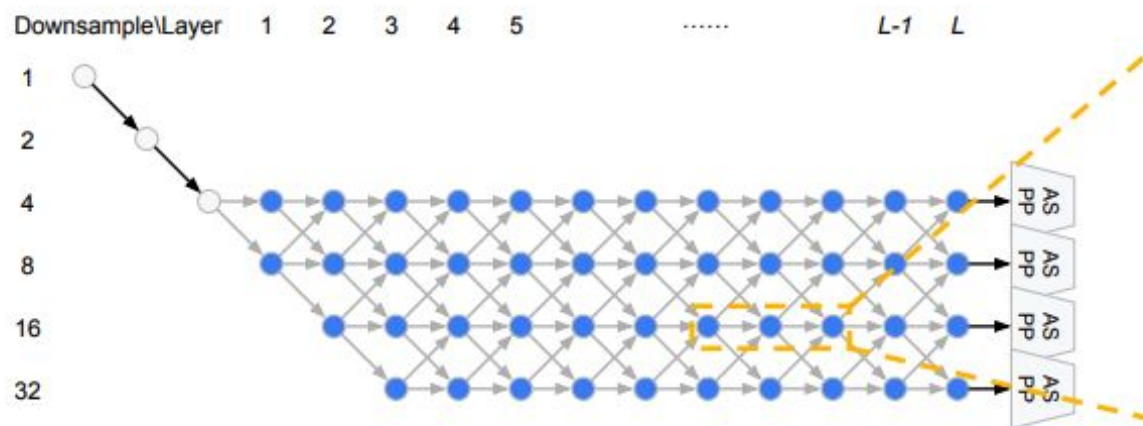
(Chenxi Liu, Liang-Chieh Chen, Florian Schroff, Hartwig
Adam, Wei Hua, Alan Yuille, Li Fei-Fei
Johns Hopkins University - Google - Stanford University)



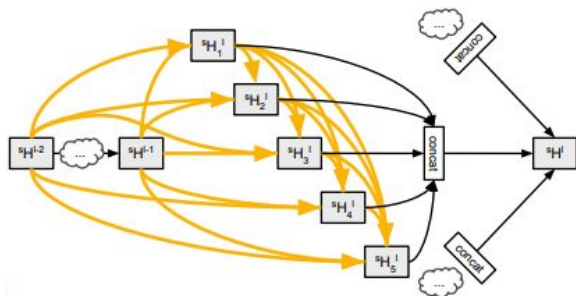
- Cell Architecture:



- Network Architecture:



• Cell Architecture



- a *cell* is a directed acyclic graph consisting of B blocks.
- a *block* i in cell l may be specified using a 5-tuple $(l_1', l_2', O_1', O_2', C)$, where $l_1', l_2' \in l_i'$, $O_1', O_2' \in O$ and $C \in C$.
- The cell's output tensor H_l is simply the concatenation of the blocks' output tensors H_1^l, \dots, H_B^l in this order.
- The set of possible input tensors, l_i' , consists of the output of the previous cell H^{l-1} , the output of the previous-previous cell H^{l-2} , and previous blocks' output in the current cell $\{H_1^l, \dots, H_{i-1}^l\}$.

- Cell Architecture

Every block's output tensor depends on:

i.
$$H_i^l = \sum_{H_j^l \in \mathcal{I}_i^l} O_{j \rightarrow i}(H_j^l)$$

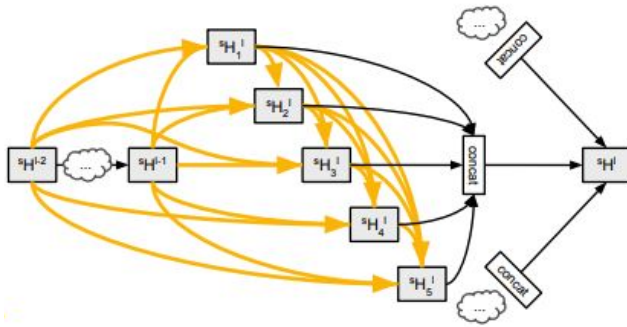
ii.
$$\bar{O}_{j \rightarrow i}(H_j^l) = \sum_{O^k \in \mathcal{O}} \alpha_{j \rightarrow i}^k O^k(H_j^l)$$

where
$$\sum_{k=1}^{|\mathcal{O}|} \alpha_{j \rightarrow i}^k = 1 \quad \forall i, j$$

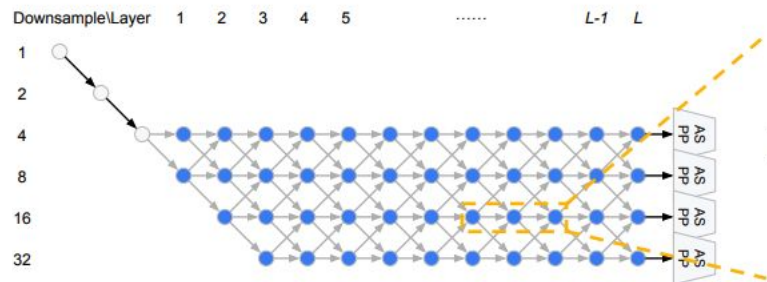
$$\alpha_{j \rightarrow i}^k \geq 0 \quad \forall i, j, k$$

Together with Eq. (i) and Eq. (ii), the cell level update may be summarized as:

$$H^l = \text{Cell}(H^{l-1}, H^{l-2}; \alpha)$$



- Network Architecture



- Each *layer* l will have at most 4 hidden state $\{^4H^l, ^8H^l, ^{16}H^l, ^{32}H^l\}$, with the upper left superscript indicating the spatial resolution.

- The network level update is:

$$\begin{aligned} {}^sH^l = & \beta_{\frac{s}{2} \rightarrow s}^l \text{Cell}(\frac{s}{2} H^{l-1}, {}^sH^{l-2}; \alpha) \\ & + \beta_{s \rightarrow s}^l \text{Cell}({}^sH^{l-1}, {}^sH^{l-2}; \alpha) \\ & + \beta_{2s \rightarrow s}^l \text{Cell}(2s H^{l-1}, {}^sH^{l-2}; \alpha) \end{aligned}$$

where $s = 4, 8, 16, 32$ and $l = 1, 2, \dots, L$.

- The scalars β are normalized such that:

$$\begin{aligned} \beta_{s \rightarrow \frac{s}{2}}^l + \beta_{s \rightarrow s}^l + \beta_{s \rightarrow 2s}^l &= 1 & \forall s, l \\ \beta_{s \rightarrow \frac{s}{2}}^l \geq 0 \quad \beta_{s \rightarrow s}^l \geq 0 \quad \beta_{s \rightarrow 2s}^l \geq 0 & & \forall s, l \end{aligned}$$

- Optimization

- Apply first-order approximation and partition the training data into two disjoint sets *trainA* and *trainB*.
- The optimization alternates between:
 1. Update network weights w by $\nabla_w \mathcal{L}_{trainA}(w, \alpha, \beta)$
 2. Update architecture α, β by $\nabla_{\alpha, \beta} \mathcal{L}_{trainB}(w, \alpha, \beta)$

where the loss function L is the cross entropy calculated on the semantic segmentation mini-batch.



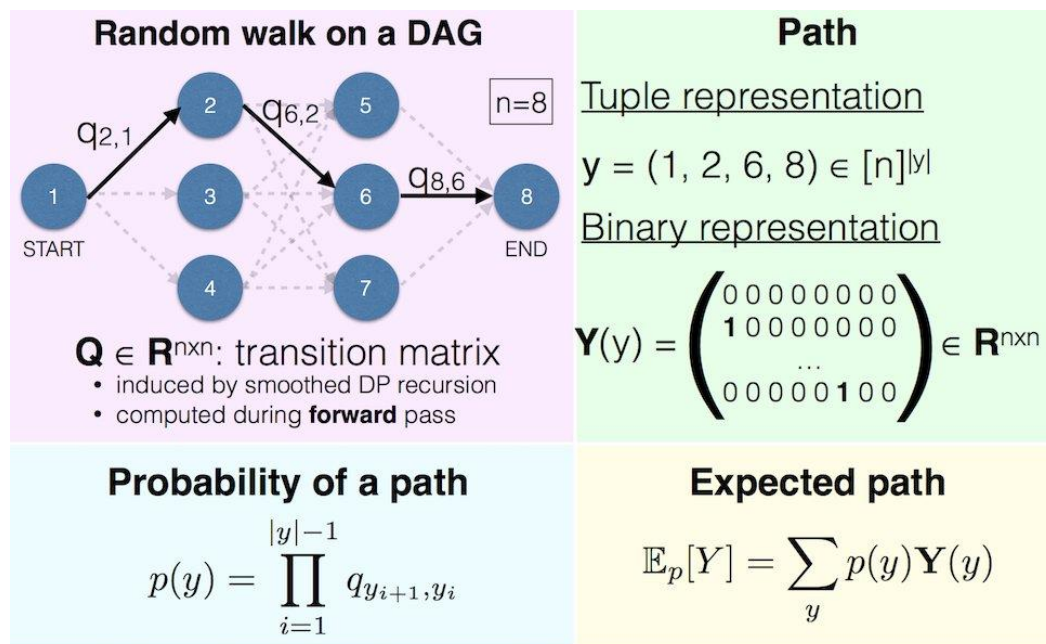
• Decoding Discrete Architecture

- Cell architecture

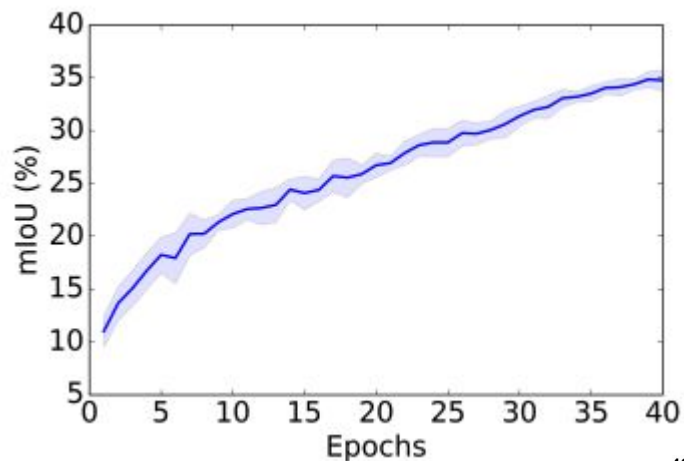
Decoding by first retaining the 2 strongest predecessors for each block:

$$\max_{k, O^k \neq \text{zero}} \alpha_{j \rightarrow i}^k$$

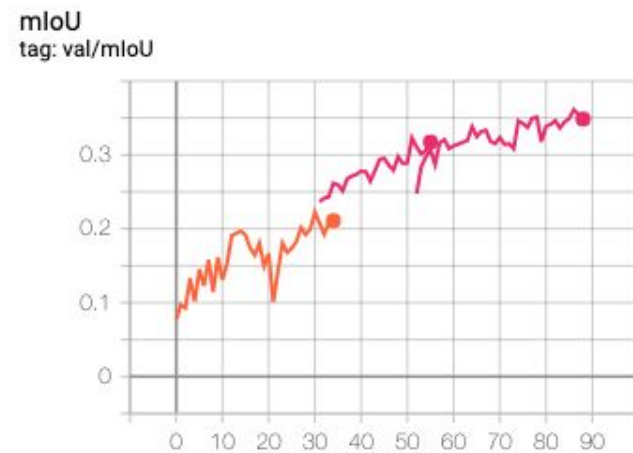
- Network Architecture



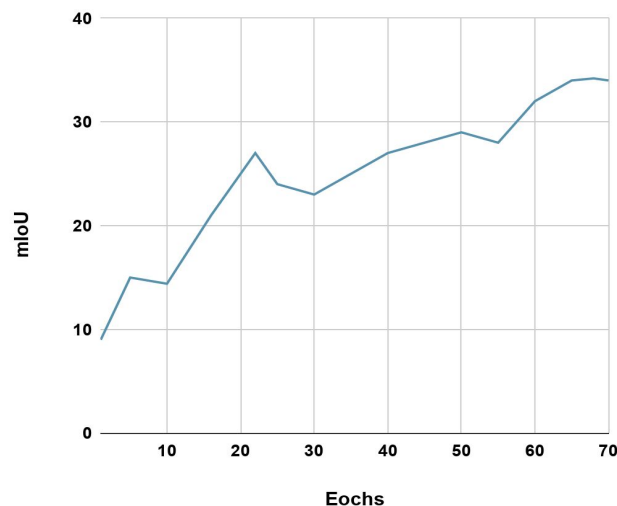
- Search results



filter_multiplier 8, crop_size 512



filter_multiplier 4, crop_size 224



filter_multiplier 8, crop_size 224

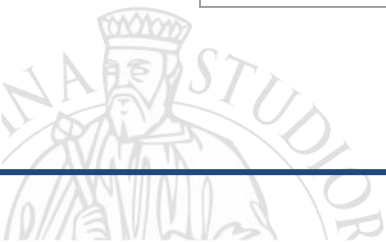
- Retrain results

- Cityscapes

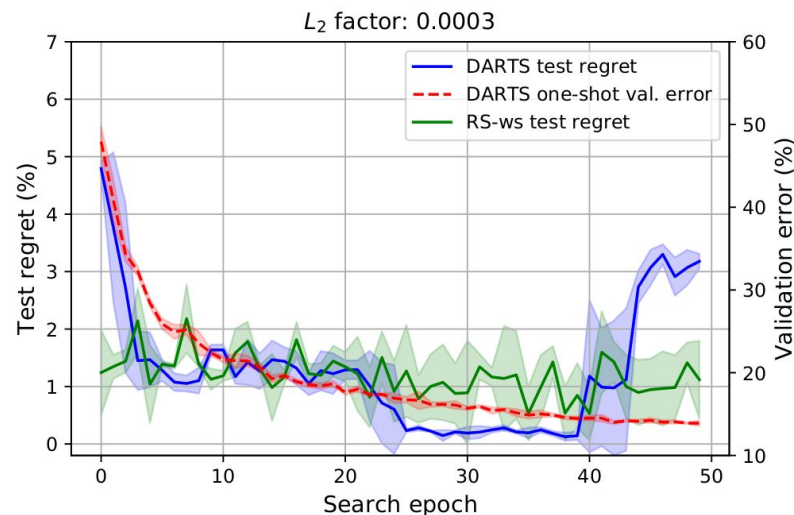
Method	F	<i>Crop-Size</i>	mIoU(%)
Auto-DeepLab-Half	32	256	69.41
Auto-DeepLab-S	20	512	79.74
Auto-DeepLab-M	32	512	80.04
Auto-DeepLab-L	48	512	80.33

- Multi-Scale + Coarse

Auto-DeepLab-L	48	512	82.1
----------------	----	-----	------



- Understanding and Robustifying Differentiable Architecture Search (Arber Zela, Thomas Elsken, Tonmoy Saikia, Yassine Marrakchi, Thomas Brox & Frank Hutter, ICLR 2020)



- a simple heuristic:

Let $\underline{\lambda}_{\max}^{\alpha}(i)$ denote the value of λ_{\max}^{α} smoothed over $k = 5$ epochs around i ;

then, we stop if $\lambda_{\max}^{\alpha}(i - k) / \underline{\lambda}_{\max}^{\alpha} < 0.75$ and return the architecture from epoch $i - k$.

- Results comparison
 - Cityscapes

Method	ImageNet	mIoU(%)
HDC+DUC	✓	81.9 (+1,8)
DeepLabv3+	✓	81.9 (-9,8)
ICNet	✓	80.08 (+9,48)
AutoDeepLab		82.1 (-12,69)

Thanks for the attention!