



# FastFCN

C.Feng@AMC Oct 01, 2019

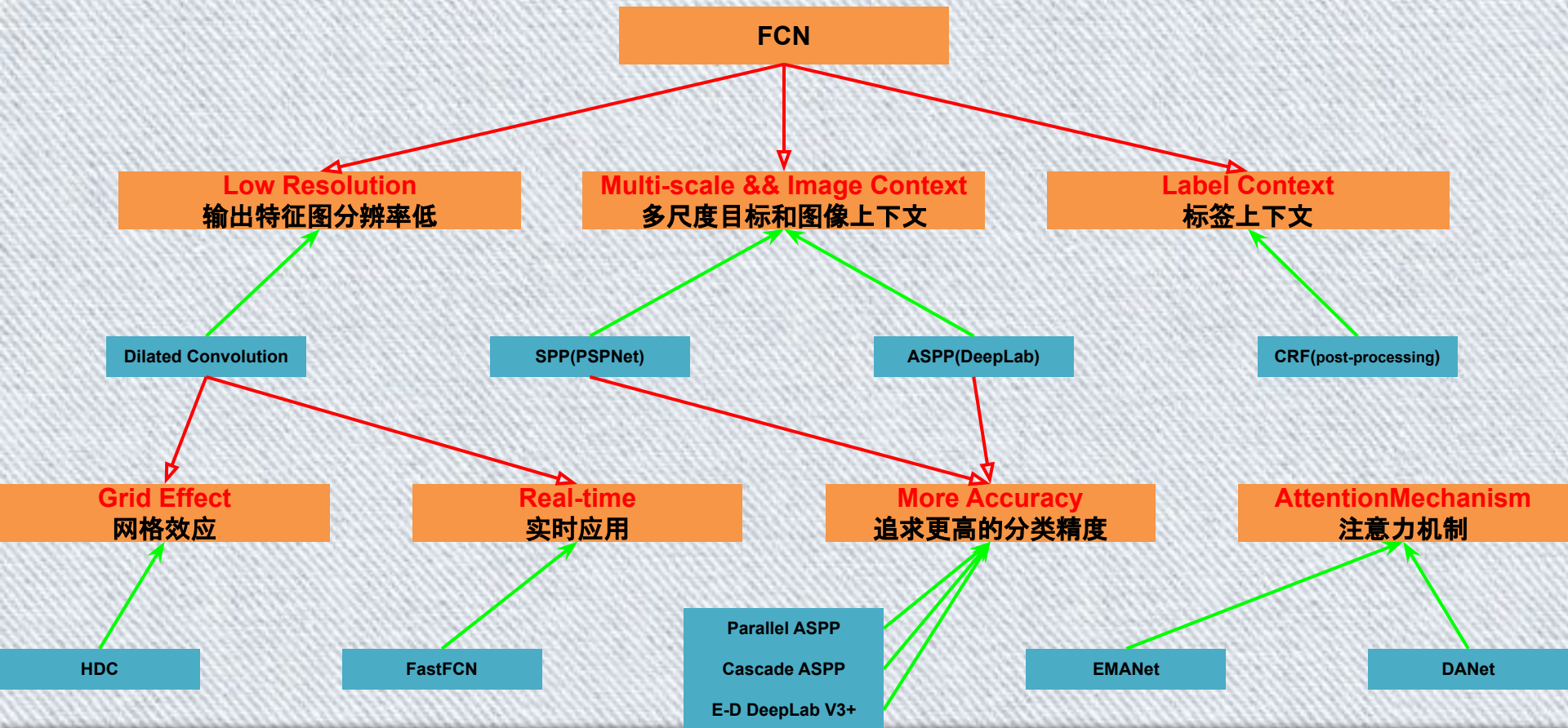
---

Wu, H., Zhang, J., Huang, K., Liang, K., & Yu, Y. (2019).

FastFCN: Rethinking Dilated Convolution in the Backbone for Semantic Segmentation.

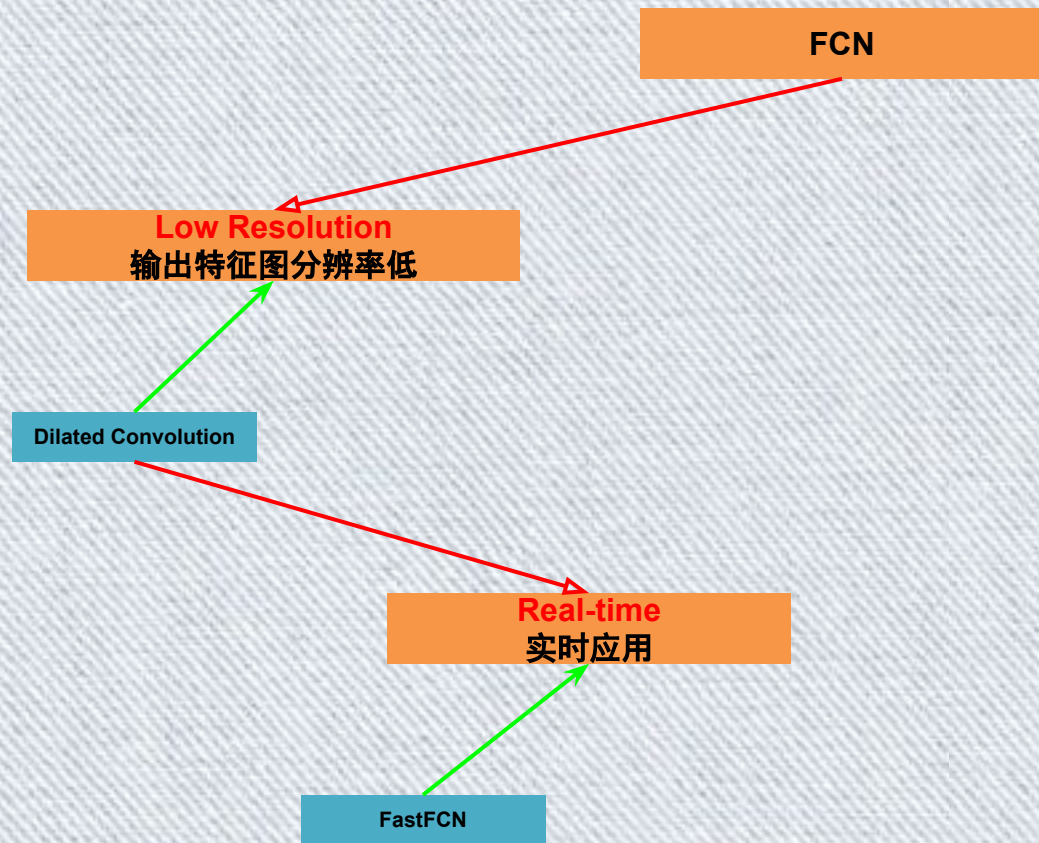
ArXiv, abs/1903.11816.

# 分割领域多种问题及其解决方案(Q&A)



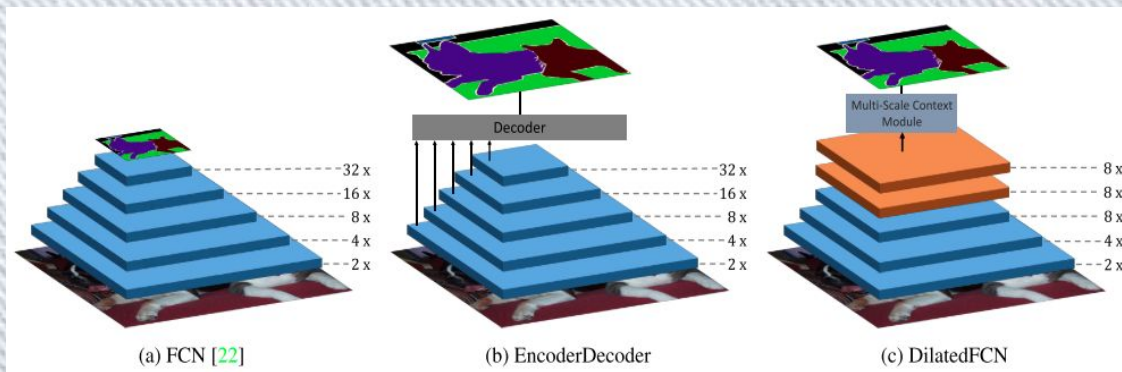


# 空洞卷积计算实时性解决方案(Q&A)





# Why



DeepLab 方法移除FCN最后两层下采样操作并引入扩张卷积来保持特征图感受野不变,后跟一个多尺度的语义模块从而得到最终效果,如图(c)所示,其中扩张卷积在保持最终特征图的分辨率方面作用明显,大大提升了编解码语义分割方法的分割精度,然而扩展卷积大大增加了计算复杂度和内存占用,限制了其在实时问题上的应用.





# How

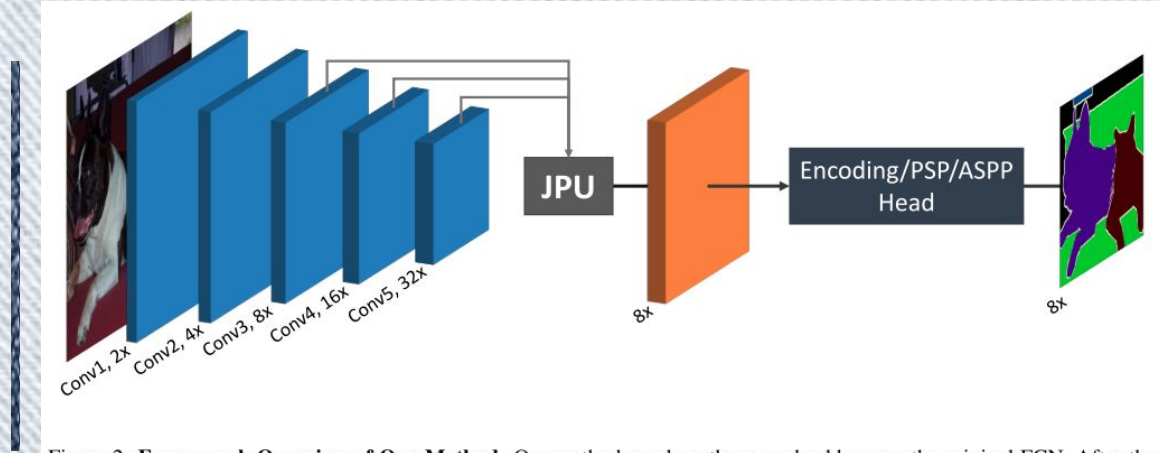
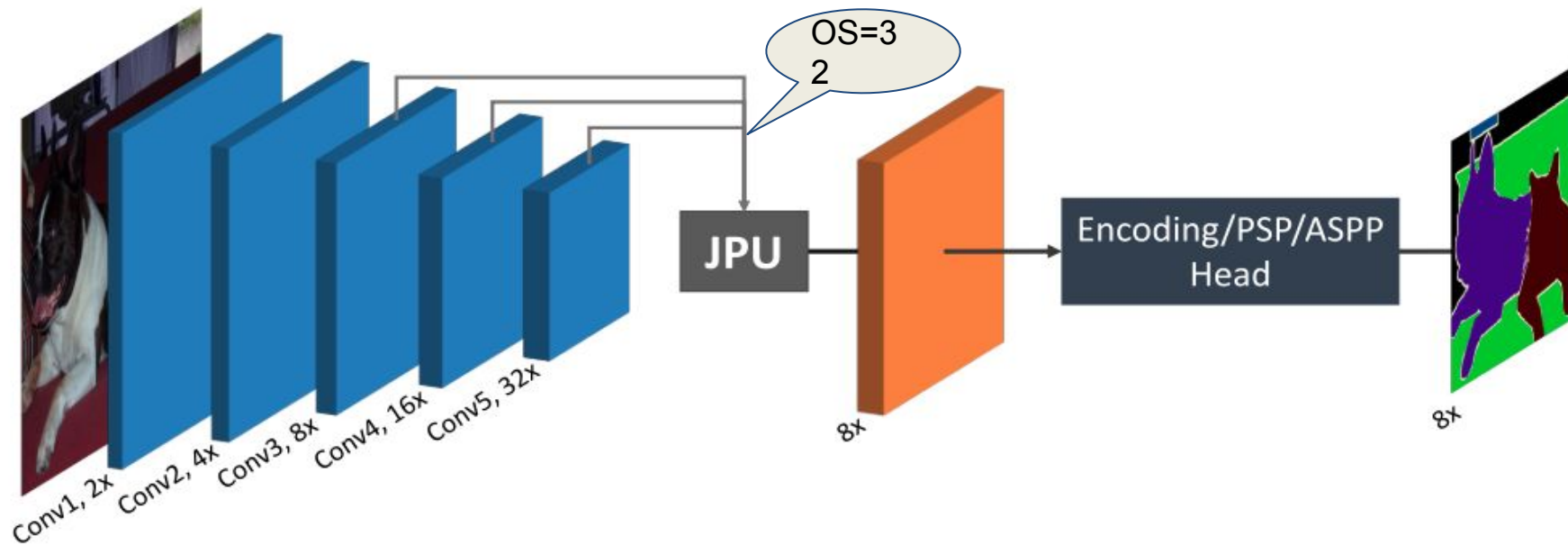
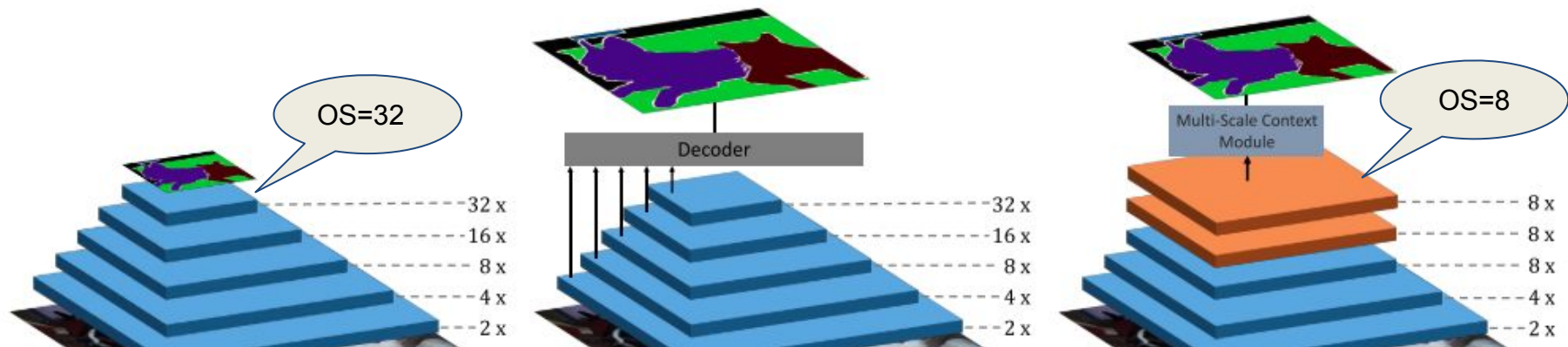
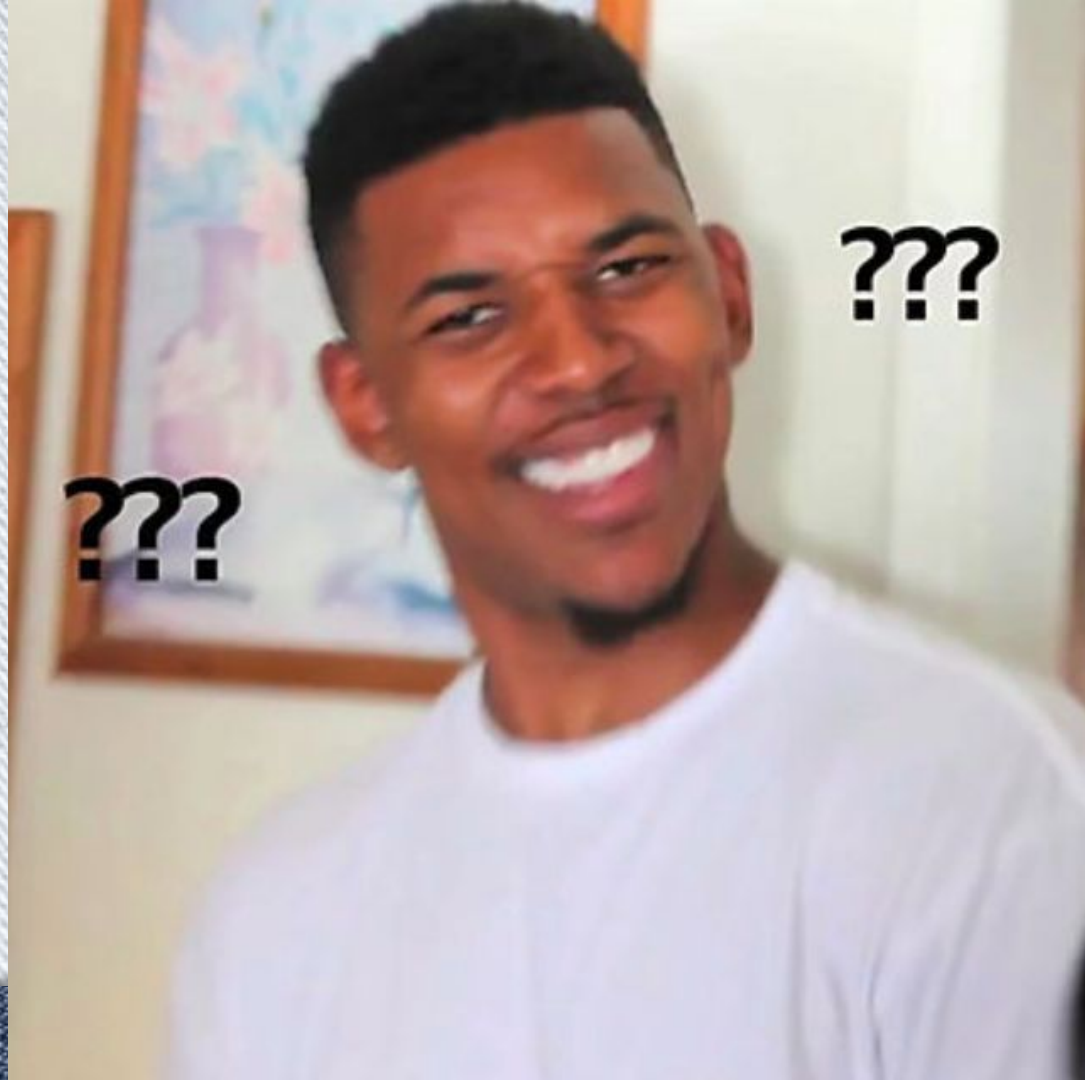


Figure 2: **Framework Overview of Our Method.** Our method employs the same backbone as the original FCN. After the backbone, a novel upsampling module named Joint Pyramid Upsampling (JPU) is proposed, which takes the last three feature maps as the inputs and generates a high-resolution feature map. A multi-scale/global context module is then employed to produce the final label map. Best viewed in color.





???

???



# What is Joint Upsampling

**联合上采样 Joint Upsampling** 给定一个低分辨率的目标图像和一个高分辨率的指导图像，联合上采样旨在从指导图像中转换出结构和细节来生成高分辨率的目标图像。一般而言，低分辨率的目标图像  $y_l$  由低分辨率的指导图像  $x_l$  通过变换  $f(\cdot)$  生成，也就是说  $y_l = f(x_l)$ 。给定  $x_l$  和  $y_l$ ，我们需要获得变换  $\hat{f}(\cdot)$  来趋近  $f(\cdot)$ ，其中  $\hat{f}(\cdot)$  的计算复杂度远低于  $f(\cdot)$ 。举例而言，如果  $f(\cdot)$  是多层感知机 MLP，那么  $\hat{f}(\cdot)$  可以被简化为线性变换。高分辨率的目标图像  $y_h$  就通过  $\hat{f}(\cdot)$  作用在高分辨率的指导图像  $x_h$  来获得。严格来说，给定  $x_l, y_l$  和  $x_h$ ，联合上采样定义如式 7:

$$y_h = \hat{f}(x_h), \text{ where } \hat{f}(\cdot) = \underset{\hat{h}(\cdot) \in \mathcal{H}}{\operatorname{argmax}} \|y_l - h(x_l)\| \quad (7)$$

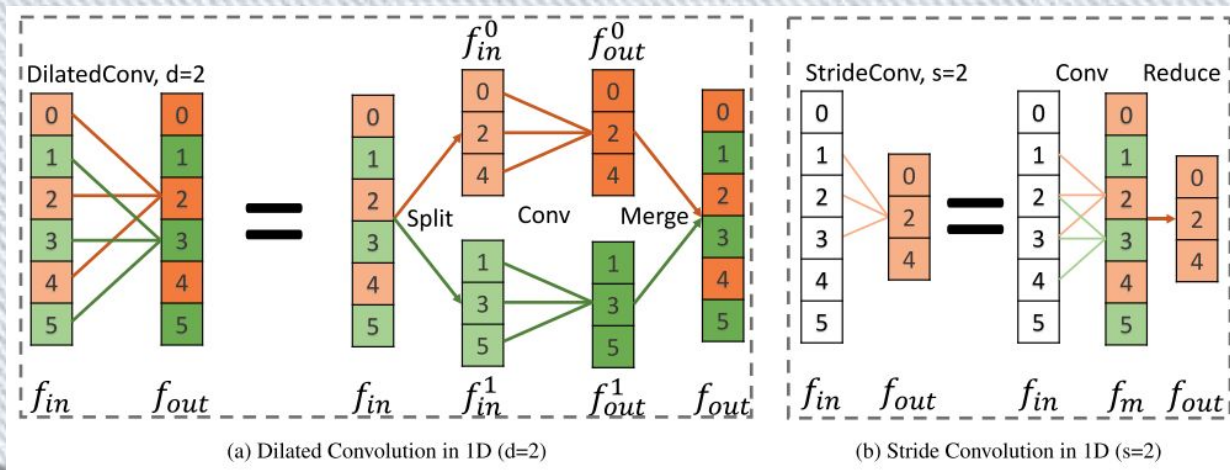
其中  $\mathcal{H}$  是所有可能的变换函数的集合， $\|\cdot\|$  是预定义的距离度量。



# The Essence of Convolution

如图(a) 所示,  $\text{DilatedConv} = \text{Split}(S) + \text{Conv}(\text{Cr}) + \text{Merge}(M)$ . 其中 Split 就是将输入按照序号奇偶不同分为两列, Conv 就是一般的卷积, Merge 就是交错融合.

如图(b) 所示,  $\text{StrideConv} = \text{Conv}(\text{Cr}) + \text{Reduce}(R)$ . 首先对输入  $f_{in}$  进行一次卷积得到中间特征  $f_m$ , 再删除基数位置数值得到最后的结果.



# Analysis of Output

$$y_d = x \rightarrow C_r \rightarrow \underbrace{C_d \rightarrow \dots \rightarrow C_d}_n$$

$$= x \rightarrow C_r \rightarrow \underbrace{SC_r M \rightarrow \dots \rightarrow SC_r M}_n = x \rightarrow C_r \rightarrow S \rightarrow \underbrace{C_r \rightarrow \dots \rightarrow C_r}_n$$

$$\rightarrow M \\ = y_m \rightarrow S \rightarrow C_r^n \rightarrow M \\ = \{y_m^0, y_m^1\} \rightarrow C_r^n \rightarrow M$$

$$y_s = x \rightarrow C_s \rightarrow \underbrace{C_r \rightarrow \dots \rightarrow C_r}_n \\ = x \rightarrow C_r \rightarrow R \rightarrow \underbrace{C_r \rightarrow \dots \rightarrow C_r}_n$$

$$= y_m \rightarrow R \rightarrow C_r^n \\ = y_m^0 \rightarrow C_r^n$$

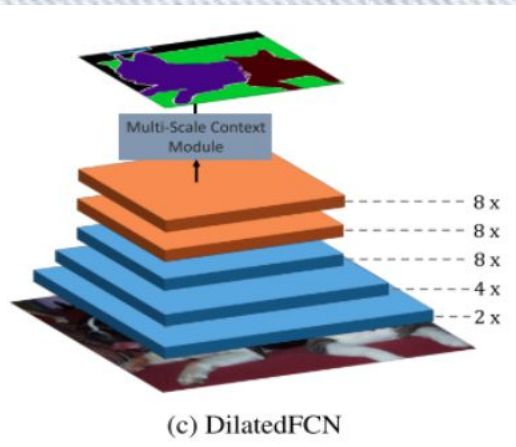
其中的,  $C_r$ ,  $C_d$  和  $C_s$  分别表示的是常规卷积、扩展卷积与步幅卷积,  $S$ ,  $M$  和  $R$  表示的分离、聚合和减去操作, 其中的  $C_{nr}$  表示  $n$  层常规卷积操作. 通过观察上面的两个公式可以发现不论是步幅卷积还是扩展卷积都可以通过常规卷积得到.

之所以比较二者的不同是因为二者展示的是一个优化的问题, 出  $x, y_s$ , 利用二者之间的联系可以得到.

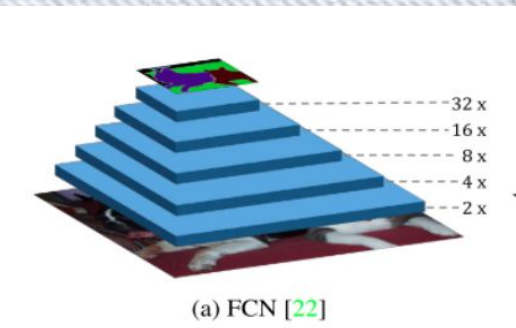
然后作者用 CNN 模型来近似这个优化过程, 也就出现了 JPU (Joined Pyramid Upsampling).

$$y = \{y_m^0, y_m^1\} \rightarrow \hat{h} \rightarrow M \\ \text{where } \hat{h} = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \|y_s - h(y_m^0)\|,$$

$$y_m = x \rightarrow C_r$$



基于Dilated Convolution的Backbone



基于Stride Convolution的Backbone



# Joined Pyramid Upsampling (JPU) Module

如上述公式阐明, 联合上采样过程在梯度下降过程中的收敛会花费很多时间, 于是本文通过设计CNN 模块来解决这个问题. 具体的, 首先在给定 $x$  的条件下生成 $y_m$ , 然后 $y_{0m}$ 和 $y_{1m}$ 的特征聚合在一起学习一个映射 $\hat{h}$ , 最后通过卷积模块得到最终的预测 $y$ , 由此设计的JPU 模块如图所示.

首先每一个输入特征图通过一个普通卷积模块(阶段a): 通过 $x$  生成 $y_m$ , 将 $f_m$  转换为具有缩减尺寸的嵌入空间, 有助于融合和降低计算复杂度.

接着所有生成的特征图进行上采样和融合生成 $y_c$ (阶段b):  $y_c$  通过多尺度语义模块, 其中空洞率为1 的平行连接用于捕获 $y_{0m}$ 和 $y_m$  其他部分的关系( $y_{0m}$ 和 $y_{1m}$ 之间的关系), 其他空洞率的平行连接用于学习将 $y_{0m}$ 转化为 $y_s$  的映射 $\hat{h}$ , 映射关系如下图. 所以JPU 可以在多层特征图上提取多尺度文本信息, 与ASPP(在单层上提取多尺度信息) 不同.

阶段对提取到的特征进行一次Conv, 最后得到分割图(阶段c).

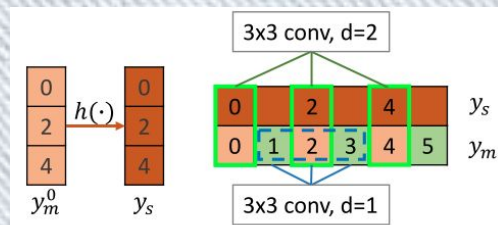
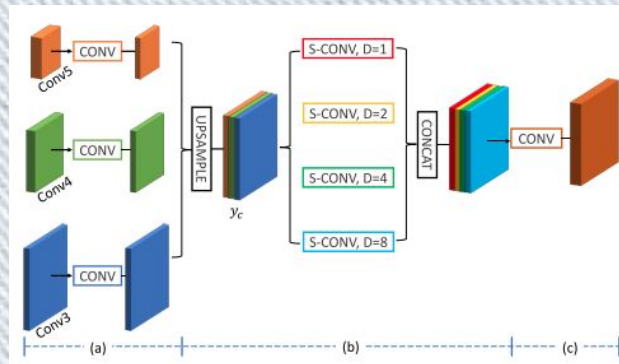


Figure 5: The convolution with dilation rate 1 focuses on  $y_m^0$  and the rest part of  $y_m$ , and the convolution with dilation rate 2 aims at  $y_m^0$  and  $y_s$ . Best viewed in color.

# Back 2

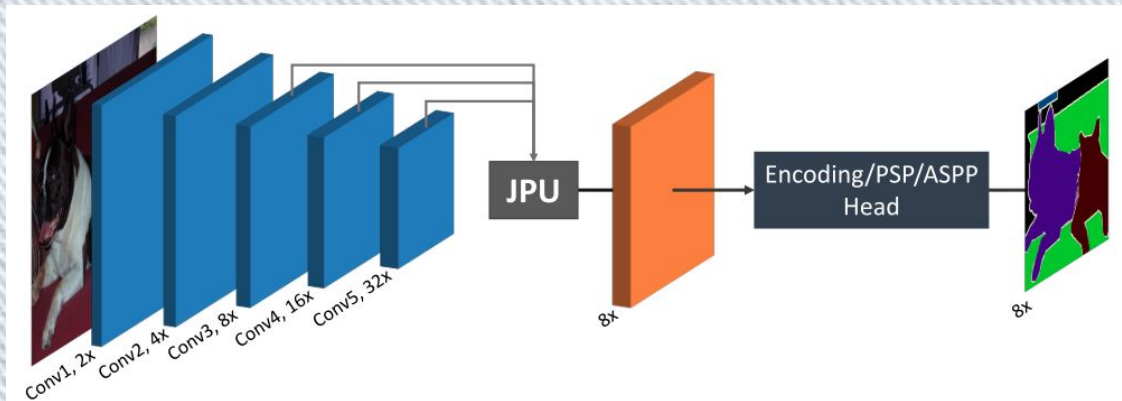


Figure 2: **Framework Overview of Our Method.** Our method employs the same backbone as the original FCN. After the backbone, a novel upsampling module named Joint Pyramid Upsampling (JPU) is proposed, which takes the last three feature maps as the inputs and generates a high-resolution feature map. A multi-scale/global context module is then employed to produce the final label map. Best viewed in color.





莱特兄弟-鸟-飞机

郎之万-海豚-声纳

斯帕拉捷-蝙蝠-雷达

仿生学-控制论

利奥那多·达·芬奇-鸟-飞行器