



# تمرین سری دوم

## یادگیری ژرف

امیر حسین محمدی

۹۹۲۰۱۰۸۱



## مسئله ۱ . Network Design

(آ) فرض کنید که شبکه ای با ساختار زیر داریم، خروجی هر لایه و تعداد پارامترهای آن را محاسبه کنید.

$$\text{Input} = 224 * 224 * 3 \rightarrow 1$$

$$1 \rightarrow \text{Conv}[3 * 3 * 64, \text{stride} = 1, \text{pad} = 1] \rightarrow 2$$

$$2 \rightarrow \text{ReLU} \rightarrow 3$$

$$3 \rightarrow \text{Pooling}[2 * 2 * 32, \text{stride} = 2] \rightarrow 4$$

$$4 \rightarrow \text{FC}[10 - \text{class}] \rightarrow 5$$

\* اعداد ۱ تا ۴ برای مشخص کردن خروجی و ورودی های لایه ها استفاده شده اند.

(ب) در صورتیکه بخواهیم فقط با یک لایه ی تمام متصل (FC) تصویری با ابعاد ورودی و خروجی مشابه داشته باشیم، تعداد پارامترهای لایه ی تمام متصل را محاسبه کنید

(ج) تفاوت Deformable Convolution با Standard Convolution در چیست و چه مزایایی نسبت به حالت استاندارد دارد؟



$$Input = 224 * 224 * 3 \rightarrow 1$$

$$1 \rightarrow Conv[3 * 3 * 64, stride = 1, pad = 1] \rightarrow 2$$

تعداد پارامترهای تولید شده بعد از لایه ی  
کانوولوشنی داریم:

parameters=((shape of width of filter\*shape of height  
filter\*number of filters in the previous layer+1)\*number of  
filters)



$$Parameters = (3 * 3 * 3 + 1) * 64 = 1792$$

ابعاد خروجی لایه ی کانوولوشنی داریم:

$$output = \frac{input - kernel_{size} + 2 * padding}{stride} + 1$$



$$output = \frac{224 - 3 + 2 * 1}{1} + 1 = 224$$



$$output = (224, 224, 64)$$



$$۲ \rightarrow ReLU \rightarrow ۳$$

تعداد پارامترهای تولید شده

ابعاد خروجی

$Parameters = 0$

$output = (224, 224, 64)$

توابع فعالیت صرفاً سطح بالاتری از ویژگی‌ها را تولید می‌کنند و پارامتر اضافه‌تری تولید نمی‌کنند

ابعاد تغییری نمی‌کنند و زیرا صرفاً به ازای هر ورودی مجدداً یک خروجی تولید می‌شود



$$3 \rightarrow \text{Pooling}[2 * 2 * 32, \text{stride} = 2] \rightarrow 4$$

تعداد پارامترهای تولید شده بعد از لایه ی

pooling

$$\text{Parameters} = 0$$

لایه pooling صرفاً یک تابع است که ویژگی های را استخراج می کند و بنابراین پارامتری تولید نمی کند

ابعاد خروجی بعد از لایه ی pooling

$$\text{output} = \frac{\text{input} - \text{kernel}_{\text{size}}}{\text{stride}} + 1$$

$$\text{output} = \frac{224 - 2}{2} + 1 = 112$$

$$\text{output} = (112, 112, 64)$$



$$x \rightarrow FC[10 - class] \rightarrow o$$

ابعاد خروجی

$$output = (10)$$

تعداد پارامترهای تولید شده با ازای لایه  
های کاملاً متصل

$$((\text{current layer neurons } c * \text{previous layer neurons } p) + 1 * c)$$

$$Parameters = (112 * 112 * 64) * 10 + 10 = 8028170$$

با توجه به  
فرضیات مسئله



(ب) در صورتیکه بخواهیم فقط با یک لایه ی تمام متصل (FC) تصویری با ابعاد ورودی و خروجی مشابه داشته باشیم، تعداد پارامترهای لایه ی تمام متصل را محاسبه کنید



$$Input = 224 * 224 * 3 \rightarrow 1$$

$$1 \rightarrow FC[10 - class] \rightarrow 2$$

تعداد پارامترهای تولید شده

ابعاد خروجی

$$output = (10)$$

$$((\text{current layer neurons } c * \text{previous layer neurons } p) + 1 * c)$$

$$Parameters = (224 * 224 * 3) * 10 + 10 = 1505290$$



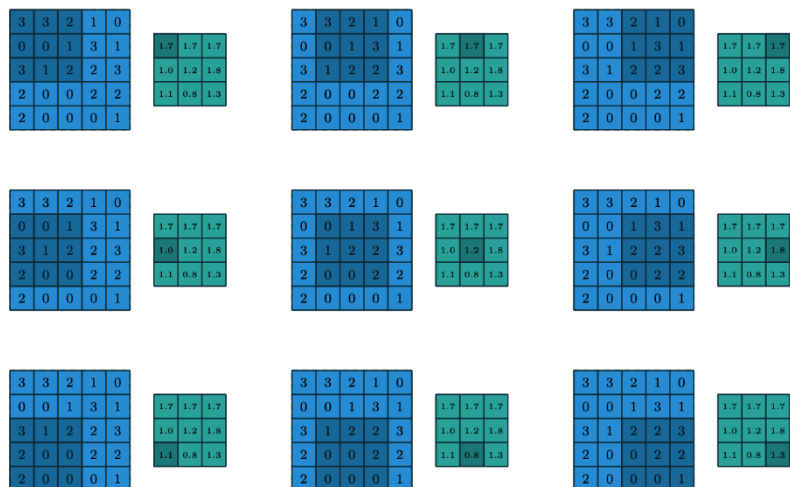


(ج) تفاوت Deformable Convolution با Standard Convolution در چیست و چه مزایایی نسبت به حالت استاندارد دارد؟



# Standard Convolution

در کانولوشن استاندارد ما یک کرنل با ابعاد ثابت داریم که با گام های ثابت حرکت می کند و ویژگی های تصاویر را استخراج می کند. در واقع نقاط کانولوشنی ما به صورت منظم قرار دارند. که رابطه ی آن به صورت زیر تعریف می شود.



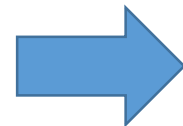
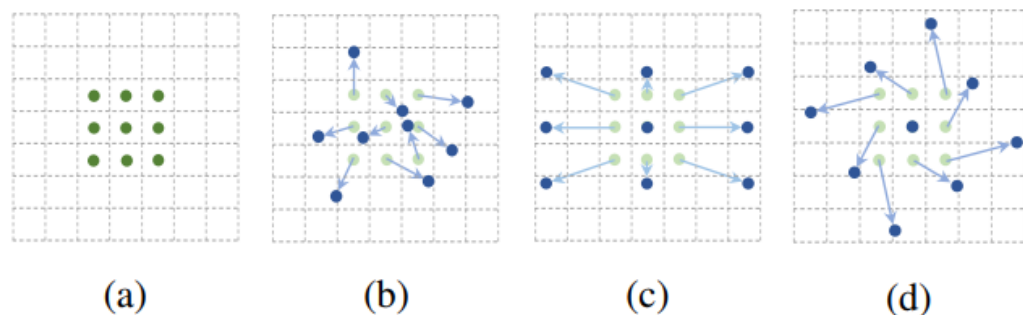
$$y(p_0) = \sum_{p_n \in \mathcal{R}} w(p_n) \cdot x(p_0 + p_n)$$

این مدل استخراج ویژگی گاهی به دلیل ویژگی ها و اسکیل ها متفاوت **object** ها در تصاویر به خوبی عمل نمی کند. زیرا همانطور که گفته شد ابعاد و تعداد گام ها در این مدل یکسان است و از طرفی به دلیل اینکه در این مدل از کانولوشن کرنل های ما شکل مربعی دارند لزوماً **object** های با اشکال گوناگون به خوبی شناسایی نمی شوند (زیرا ما تصاویر و **object** های با اسکیل های متفاوت داریم). بنابراین ما نیاز داریم تا کانولوشن هایی اعمال کنیم که شکل های با عرض و ارتفاع های مختلف و ... را شناسایی کنیم.

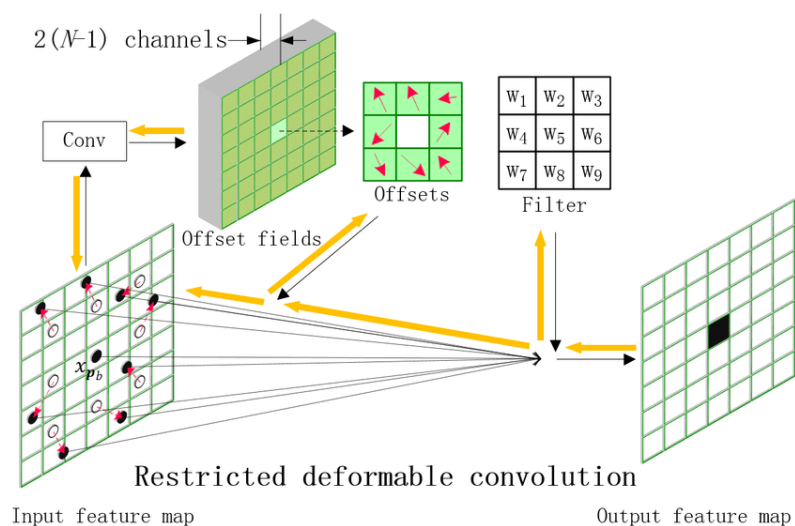
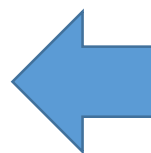


# Deformable Convolution

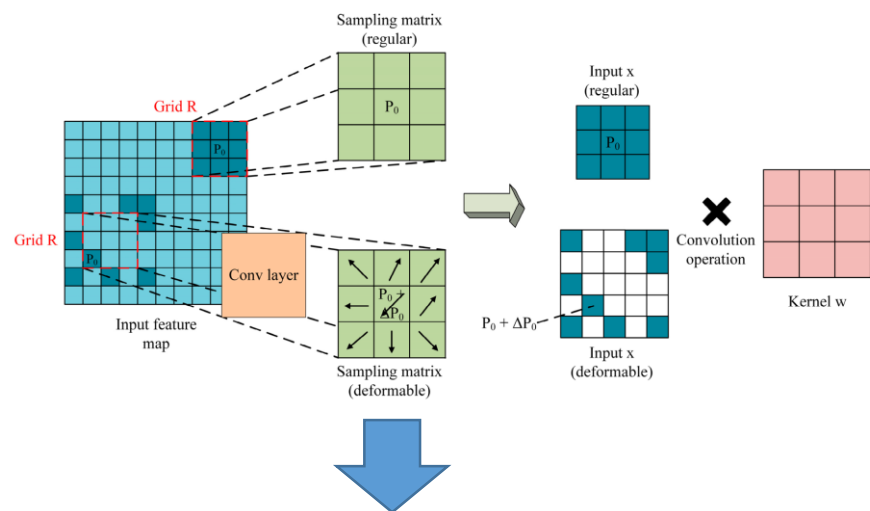
با توجه به مشکلات مطرح شده در کانولوشن های استاندارد، ایده استفاده از کانولوشن های شکل پذیر ایجاد شد به گونه ای که بتوانند برخلاف شکل مربعی و ثابتی که در فرم استاندارد وجود داشت **Object** های با شکل های پیچیده تر را نیز شناسایی کنند. برای این کار کافی بود تا به معادله ی محاسبه ی کانولوشن که در کانولوشن استاندارد به آن پرداخته شد یک پارامتر قابل یادگیری اضافه کنیم که در هنگام آموزش شبکه به روزرسانی می شود و شکل کانولوشن را تغییر داده و به نوعی کانولوشن شکل پذیر ایجاد می کند که رابطه ی آن به شکل زیر تعریف می شود.



$$y(\mathbf{p}_0) = \sum_{\mathbf{p}_n \in \mathcal{R}} w(\mathbf{p}_n) \cdot x(\mathbf{p}_0 + \mathbf{p}_n + \Delta \mathbf{p}_n).$$



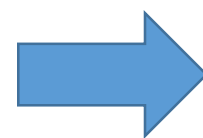
همانطور که می بینیم بر اساس پارامتر  $\Delta P_n$  که به معادله ی کانولوشن اضافه کردیم کانولوشن های ما همانند شکل های بالا شکل پذیر شده اند مثلاً در ابتدای آموزش شکل کانولوشن همانند کانولوشن استاندارد است (همانطور که در شکل a می بینیم) اما در هنگام یادگیری شکل کانولوشن تغییر می کند و مثلاً در شکل c همانطور که می بینیم شکل کانولوشن در راستای محور ایکس رشد می کند و مثلاً می توانند اشکال عریض مثل اتوبوس را شناسایی کند. تفاوت کانولوشن استاندارد و کانولوشن شکل پذیر در شکل مقابل مشخص است.



در شکل مقابل مثال هایی از کانولوشن های شکل پذیر را می بینیم که ابعاد و شکل کانولوشن متناسب با Object های موجود در تصویر تغییر کرده است. نتایج نشا می دهد که استفاده از کانولوشن های شکل پذیر در مسائلی که شناسایی object ها اهمیت دارند بسیار مفید و موثر خواهند بود. همچنین استفاده از کانولوشن های شکل پذیر تعداد پارامتر ها ویژگی های زیادی را به شبکه تحمیل نمی کند. همچنین در مسائل Object detection بسیار بهتر عمل میکند. همچنین در مسائلی که تصاویر اسکیل های زیادی دارند بهتر عمل می کنند. نمونه ای از برترها و مقایسه های بین کانولوشن استاندارد و کانولوشن شکل پذیر (deformable) بر اساس بررسی شبکه های مختلف در زیر آورده شده است.



Test Accuracy	Regular CNN	Deformable CNN
Regular MNIST	98.74%	97.27%
Scaled MNIST	57.01%	92.55%



مقایسه ای از عملکرد کانولوشن های استاندارد و کانولوشن های شکل پذیر بر روی مجموعه دادگان MNIST که اسکیل شده اند و نشده اند.



deformation modules	DeepLab mIoU@V / @C	class-aware RPN mAP@0.5 / @0.7	Faster R-CNN mAP@0.5 / @0.7	R-FCN mAP@0.5 / @0.7
atrous convolution (2,2,2) (default)	69.7 / 70.4	68.0 / 44.9	78.1 / 62.1	80.0 / 61.8
atrous convolution (4,4,4)	73.1 / 71.9	72.8 / 53.1	78.6 / 63.1	80.5 / 63.0
atrous convolution (6,6,6)	73.6 / 72.7	73.6 / 55.2	78.5 / 62.3	80.2 / 63.5
atrous convolution (8,8,8)	73.2 / 72.4	73.2 / 55.1	77.8 / 61.8	80.3 / 63.2
deformable convolution	<b>75.3 / 75.2</b>	<b>74.5 / 57.2</b>	78.6 / 63.3	81.4 / 64.7
deformable RoI pooling	N.A	N.A	78.3 / 66.6	81.2 / 65.0
deformable convolution & RoI pooling	N.A	N.A	<b>79.3 / 66.9</b>	<b>82.6 / 68.5</b>

ارزیابی کانوولوشن های شکل پذیر با روش های دیگر مثل کانوولوشن های atrous بر روی شبکه ی ResNet. همانطور که مشخص است کانوولوشن های Deformable به دقت بهتری رسیده اند.



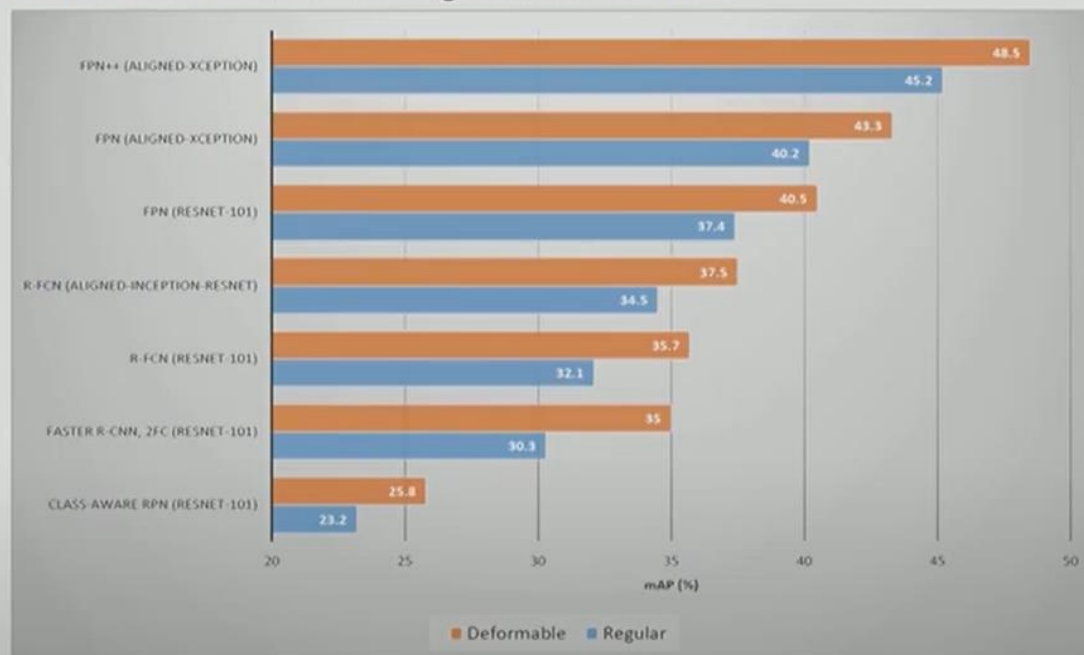
method	# params	net. forward (sec)	runtime (sec)
DeepLab@C	46.0 M	0.610	0.650
<b>Ours</b>	46.1 M	0.656	0.696
DeepLab@V	46.0 M	0.084	0.094
<b>Ours</b>	46.1 M	0.088	0.098
class-aware RPN	46.0 M	0.142	0.323
<b>Ours</b>	46.1 M	0.152	0.334
Faster R-CNN	58.3 M	0.147	0.190
<b>Ours</b>	59.9 M	0.192	0.234
R-FCN	47.1 M	0.143	0.170
<b>Ours</b>	49.5 M	0.169	0.193

همانطور که مشخص است کانوولوشن های Deformable بر اساس تعداد پارمترها و زمان اجرا و پیچیدگی محاسباتی سربار محاسباتی و حافظه ای چشم گیری به شبکه تحمیل نمی کند.



# Object Detection on COCO

- Deformable ConvNets v.s. regular ConvNets



همانطور که مشخص است در مسائل تشخیص شی هم، کانولوشن های شکل پذیر نتایج بهتری نسبت به کانولوشن های استاندارد دارند.





## مسئله ۲ . Transposed Convolution

کانولوشن بین ورودی  $X$  و فیلتر  $W$  به صورت زیر است:

$$X = \begin{bmatrix} x_{(.,.)} & x_{(.,1)} & x_{(.,2)} \\ x_{(1,.)} & x_{(1,1)} & x_{(1,2)} \\ x_{(2,.)} & x_{(2,1)} & x_{(2,2)} \end{bmatrix} \quad W = \begin{bmatrix} w_{(.,.)} & w_{(.,1)} \\ w_{(1,.)} & w_{(1,1)} \end{bmatrix}$$

می توان عملیات کانولوشن را به صورت ضرب ماتریسی نوشت که ورودی و خروجی را به صورت یک بردار و فیلتر را به صورت یک ماتریس در نظر گرفت. ورودی  $X$  را به صورت بردار زیر نمایش می دهیم:





$$X = \begin{bmatrix} x_{(.,0)} & x_{(.,1)} & \dots & x_{(r,r)} \end{bmatrix}^T$$

(آ) عملیات کانولوشن بالا با  $S = 1$  را به صورت ضرب ماتریسی  $Y = AX$  بنویسید.

(ب) با استفاده از نمایش ماتریسی بالا می توان گرادیان پس انتشار نسبت به ورودی  $X$ ، را به صورت  $\frac{\partial L}{\partial X} = A^T \frac{\partial L}{\partial Y}$  نمایش داد.

عملیات transposed convolution را می توان مشابه عملیات گرادیان پس انتشار نسبت به ورودی در نظر گرفت و می توان این عملیات را به صورت کانولوشن مستقیم با ماتریس  $A^T$  به عنوان فیلتر دانست.

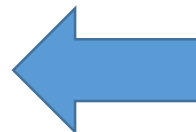
فرض کنید  $X$  خروجی یک عملیات کانولوشن با فیلتر  $W$  و  $\text{stride} = 2$  است. خروجی transposed convolution را با ورودی های زیر محاسبه کنید.

$$X = \begin{bmatrix} 2 & 1 \\ 3 & 6 \end{bmatrix} \quad W = \begin{bmatrix} 1 & 2 \\ 4 & 1 \end{bmatrix}$$



❖ طبق فرضیات مسئله داریم:

$$X = \begin{bmatrix} x_{(\cdot,\cdot)} & x_{(\cdot,1)} & x_{(\cdot,2)} \\ x_{(1,\cdot)} & x_{(1,1)} & x_{(1,2)} \\ x_{(2,\cdot)} & x_{(2,1)} & x_{(2,2)} \end{bmatrix} \quad W = \begin{bmatrix} w_{(\cdot,\cdot)} & w_{(\cdot,1)} \\ w_{(1,\cdot)} & w_{(1,1)} \end{bmatrix}$$



❖ می توانیم ماتریس ورودی را به صورت زیر تعریف کنیم:

$$X = [x_{0,0}, x_{0,1}, x_{0,2}, x_{1,0}, x_{1,1}, x_{1,2}, x_{2,0}, x_{2,1}, x_{2,2}]^T_{1 \times 9}$$



همانطور که می بینیم ابعاد ماتریس  $X$  به صورت ۱ در ۹ است. برای اینکه ابعاد ماتریس  $w$  یا همان کرنل را بدست آوریم، می دانیم در تبدیل عملیات کانوولوشن به ضرب ماتریسی اگر ماتریس ورودی ( $X$ ) در  $n$  باشد و ابعاد ماتریس کرنل ( $w$ ) در  $m$  در  $m$  باشد آنگاه برای اینکه بتوانیم ماتریس  $X$  مسئله ما که ابعاد ۱ در ۹ دارد بتواند با ما ماتریس  $w$  ضرب شود ابعاد ماتریس  $A$  که در خود به نوعی اطلاعات ماتریس  $w$  را دارد به ابعاد  $n^2$  در  $(n - m + 1)^2$  به تغییر کند یعنی ما باید یک ماتریس به نام  $A$  با ابعاد ۴ در ۹ داریم. بنابراین باید این چهار مقدار ماتریس  $w$  به نوعی در ماتریس جدید که با ابعاد ۴ در ۹ وجود دارد قرار گیرد که بتواند دقیقاً عملیات کانوولوشن را شبیه سازی کند اما به صورت ضرب ماتریسی، همچنین با توجه به  $\text{stride}=1$  باید مقادیر وزن ها طوری در ماتریس  $A$  قرار گیرند که این شرط نیز لحاظ شود، بنابراین ماتریس  $A$  با ابعاد ۴ در ۹ به صورت زیر تعریف می شود:

$$A = \begin{bmatrix} w_{0,0} & w_{0,1} & 0 & w_{1,0} & w_{1,1} & 0 & 0 & 0 & 0 \\ 0 & w_{0,0} & w_{0,1} & 0 & w_{1,0} & w_{1,1} & 0 & 0 & 0 \\ 0 & 0 & 0 & w_{0,0} & w_{0,1} & 0 & w_{1,0} & w_{1,1} & 0 \\ 0 & 0 & 0 & 0 & w_{0,0} & w_{0,1} & 0 & w_{1,0} & w_{1,1} \end{bmatrix}_{4 \times 9}$$

ماتریس بالا بر اساس توضیح مسئله مثل  $\text{stride}=1$  و ابعاد ماتریس کانوولوشن ایجاد شده است.



خواسته ی  
مسئله

$$\diamond y = AX$$

$$A = \begin{bmatrix} w_{0,0} & w_{0,1} & 0 & w_{1,0} & w_{1,1} & 0 & 0 & 0 & 0 \\ 0 & w_{0,0} & w_{0,1} & 0 & w_{1,0} & w_{1,1} & 0 & 0 & 0 \\ 0 & 0 & 0 & w_{0,0} & w_{0,1} & 0 & w_{1,0} & w_{1,1} & 0 \\ 0 & 0 & 0 & 0 & w_{0,0} & w_{0,1} & 0 & w_{1,0} & w_{1,1} \end{bmatrix}_{4 \times 9}$$

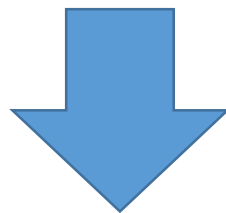
$$X = [x_{0,0}, x_{0,1}, x_{0,2}, x_{1,0}, x_{1,1}, x_{1,2}, x_{2,0}, x_{2,1}, x_{2,2}]^T_{1 \times 9}$$

$$y = \begin{bmatrix} w_{0,0} & w_{0,1} & 0 & w_{1,0} & w_{1,1} & 0 & 0 & 0 & 0 \\ 0 & w_{0,0} & w_{0,1} & 0 & w_{1,0} & w_{1,1} & 0 & 0 & 0 \\ 0 & 0 & 0 & w_{0,0} & w_{0,1} & 0 & w_{1,0} & w_{1,1} & 0 \\ 0 & 0 & 0 & 0 & w_{0,0} & w_{0,1} & 0 & w_{1,0} & w_{1,1} \end{bmatrix}_{4 \times 9} * \begin{bmatrix} x_{0,0} \\ x_{0,1} \\ x_{0,2} \\ x_{1,0} \\ x_{1,1} \\ x_{1,2} \\ x_{2,0} \\ x_{2,1} \\ x_{2,2} \end{bmatrix}_{9 \times 1} = \begin{bmatrix} w_{0,0} * x_{0,0} + w_{0,1} * x_{0,1} + w_{1,0} * x_{1,0} + w_{1,1} * x_{1,1} \\ w_{0,0} * x_{0,1} + w_{0,1} * x_{0,2} + w_{1,0} * x_{1,1} + w_{1,1} * x_{1,2} \\ w_{0,0} * x_{1,0} + w_{0,1} * x_{1,1} + w_{1,0} * x_{2,0} + w_{1,1} * x_{2,1} \\ w_{0,0} * x_{1,1} + w_{0,1} * x_{1,2} + w_{1,0} * x_{2,1} + w_{1,1} * x_{2,2} \end{bmatrix}_{4 \times 1}$$

❖ خروجی ما در حال حاضر متناسب با ابعاد خروجی کانولوشن نیست (به خاطر تغییر سایزهایی که دادیم) بنابراین با خروجی ما ریشپ شود (۲ در ۲) تا دقیقا خروجی حاصل ضرب بالا با خروجی کانولوشن برابر شود.



بنابراین پس از reshape  
خواهیم داشت:



$$y = AX = \begin{bmatrix} w_{0,0} * x_{0,0} + w_{0,1} * x_{0,1} + w_{1,0} * x_{1,0} + w_{1,1} * x_{1,1} & w_{0,0} * x_{0,1} + w_{0,1} * x_{0,2} + w_{1,0} * x_{1,1} + w_{1,1} * x_{1,2} \\ w_{0,0} * x_{1,0} + w_{0,1} * x_{1,1} + w_{1,0} * x_{2,0} + w_{1,1} * x_{2,1} & w_{0,0} * x_{1,1} + w_{0,1} * x_{1,2} + w_{1,0} * x_{2,1} + w_{1,1} * x_{2,2} \end{bmatrix}_{2 \times 2}$$



(ب) با استفاده از نمایش ماتریسی بالا می توان گرادیان پس انتشار نسبت به ورودی  $X$ ، را به صورت  $\frac{\partial L}{\partial X} = A^T \frac{\partial L}{\partial Y}$  نمایش داد.

عملیات transposed convolution را می توان مشابه عملیات گرادیان پس انتشار نسبت به ورودی در نظر گرفت و می توان این عملیات را به صورت کانولوشن مستقیم با ماتریس  $A^T$  به عنوان فیلتر دانست.

فرض کنید  $X$  خروجی یک عملیات کانولوشن با فیلتر  $W$  و  $\text{stride} = 2$  است. خروجی transposed convolution را با ورودی های زیر محاسبه کنید.

$$X = \begin{bmatrix} 2 & 1 \\ 3 & 6 \end{bmatrix} \quad W = \begin{bmatrix} 1 & 2 \\ 4 & 1 \end{bmatrix}$$

❖ با توجه به فرضیات مسئله و نتایج قسمت الف داریم:

$$x = [2 \quad 1 \quad 3 \quad 6]_{1 \times 4}$$

$$W = \begin{bmatrix} 1 & 2 & 0 & 0 & 4 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 2 & 0 & 0 & 4 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 2 & 0 & 0 & 4 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 2 & 0 & 0 & 4 & 1 \end{bmatrix}_{4 \times 16}$$

با توجه به  $\text{stride} = 2$  و ابعاد ماتریس های  $X$  و  $W$

❖ به دلیل اینکه طبق فرضیات  $X$  خروجی مسئله است بنابراین باید  $W$  را بر اساس سائز ورودی مسئله طراحی کنیم برای بدست آوردن سائز ورودی مسئله داریم:

$$\frac{\text{input} - 2 * \text{kernel}_{\text{size}}}{\text{stride}} + 1 = 2$$

بنابراین ابعاد ورودی مسئله ۴ در ۴ بوده است

بنابراین همانطور که در قسمت قبل نیز اشاره شده برای بدست آوردن ابعاد  $W$  با ابعاد آنرا به گونه ای در نظر بگیریم که مقدار خروجی آن ۱۶ باشد تا با reshape کردن بتوانیم ویژگی هایی با ابعاد سائز ورودی تولید کنیم. بنابراین تعداد ستون های ماتریس  $W$  ۱۶ و تعداد سطرهای آن ۴ در نظر گرفته می شود (تا بتوانیم ضرب ماتریسی انجام دهیم). همچنین به دلیل اینکه  $\text{stride} = 2$  در نظر گرفته شده بین هر دو عدد غیر از صفر دو صفر فاصله می اندازیم (برای ایجاد stride)



$$[2 \ 1 \ 3 \ 6]_{1 \times 4} * \begin{bmatrix} 1 & 2 & 0 & 0 & 4 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 2 & 0 & 0 & 4 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 2 & 0 & 0 & 4 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 2 & 0 & 0 & 4 & 1 \end{bmatrix}_{4 \times 16} = [2, 4, 1, 2, 8, 2, 4, 1, 3, 6, 6, 12, 12, 3, 24, 6]_{1 \times 16}$$

همانطور که گفتیم باید متناسب با  
ابعاد ورودی بردار حاصل را reshape  
کنیم



$$\text{Transposed Convolution} = \begin{bmatrix} 2 & 4 & 1 & 2 \\ 8 & 2 & 4 & 1 \\ 3 & 6 & 6 & 12 \\ 12 & 3 & 24 & 6 \end{bmatrix}_{4 \times 4}$$



### مسئله‌ی ۳. Object Detection

یکی از کاربردهای شبکه‌های کانولوشن، آشکارسازی اشیاء است. مقالات مرتبط را مطالعه کرده و به سوالات زیر پاسخ دهید.

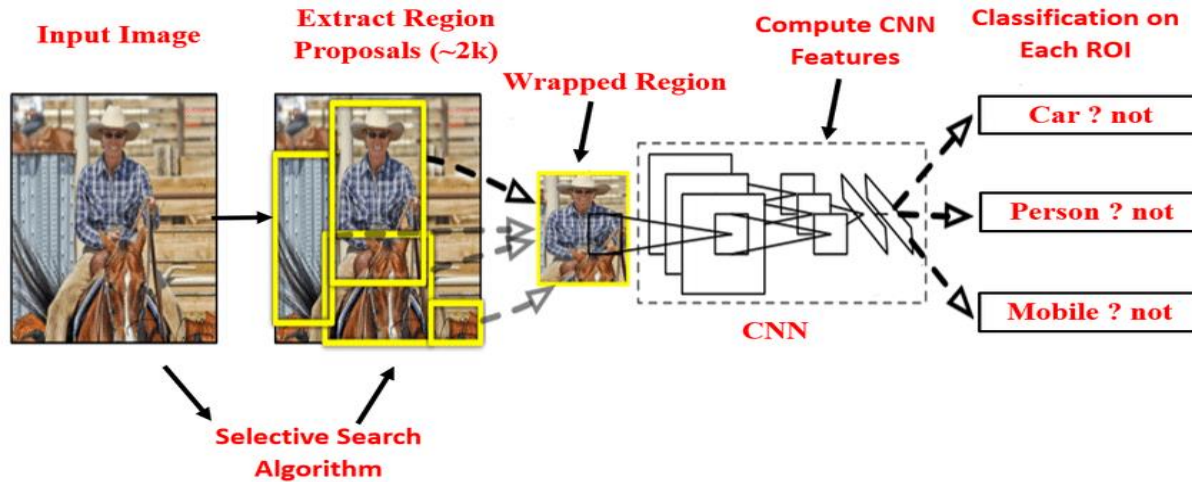
(آ) روش‌های YOLO و RCNN و RFCN را بررسی کرده و با یکدیگر مقایسه کنید.

(ب) در آشکارسازی‌های real-time استفاده از کدام یک را پیشنهاد می‌کنید.

(ج) لایه Spatial Pyramid Pooling را مختصراً توضیح داده و مزایای استفاده از آن را، در آشکارسازی تصاویر بیان کنید.



روش RCNN یک روش برای شناسایی Object های موجود در یک تصویر است (Object detection). بر اساس این روش یک تصویر به نواحی مختلف تقسیم شده (Region proposal) و از یک الگوریتم به اسم selective search برای ایجاد این نواحی استفاده می کند. این نواحی که تقریباً طبق گفته مقاله ۲۰۰۰ ناحیه متفاوت ROI (regions of interest) هستند به شبکه ی کانوولوشنی مورد نظر به عنوان ورودی داده می شود مثل Alex net از این شبکه با حجم قابل توجهی از داده های ورودی مواجه می شود برای همین معمولاً دیتاست های معمول برای object detection حجم کمی دارند مثلاً coco ۹۰۰۰ داده دارد. همانطور که در شکل زیر می بینیم، همچنین تابع هدف آن به صورت زیر تعریف می شود:



$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i N_{cls}(p_i, p_i^*) + \lambda \frac{1}{N_{reg}} \sum_i p_i^* L_{reg}(t_i, t_i^*)$$

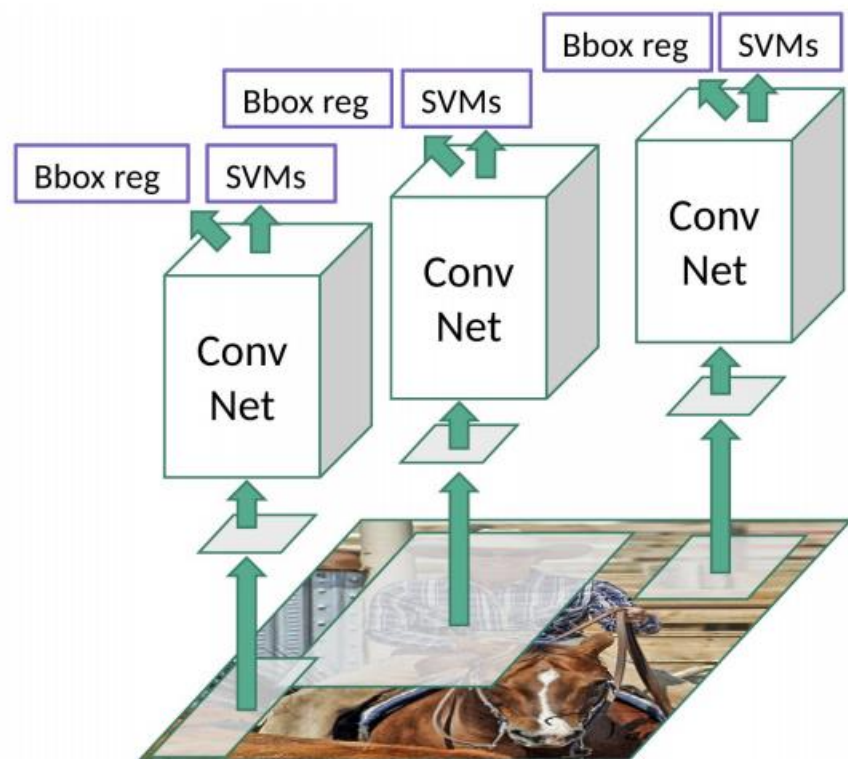
که در این فرمول  $P_i$  احتمال پیش بینی شده است،  $P_i^*$  احتمال واقعی است که در صورتی که object وجود داشته باشد مقدار ۱ و در صورتی که object وجود نداشته باشد مقدار ۰ به خود می گیرد.  $t_i$  مقدار مختصات جعبه های محاطی است و  $t_i^*$  مقدار واقعی جعبه های محاطی است.  $L_{cls}$  مقدار لاگ خطای دسته بندی است و  $L_{reg}$  خطای رگرسیون است.





این روش بر اساس یک سری جعبه های محاط شده **bounding box** استفاده می کند و بر اساس مدل رگرسیونی مکانی که به عنوان مختصات جعبه محاطی پیش بینی کرده را با جعبه محاطی واقعی مقایسه می کند و مقدار تابع هزینه را محاسبه می کند. برای تعیین میزان همپوشانی مقدار جعبه محاطی تخمین زده شده با جعبه محاطی واقعی از معیار **intersection-over-union(IoU)** استفاده می شود. این الگوریتم از سه بخش اصلی تشکیل شده است که عبارت است از (همانطور که از شکل مقابل مشخص است):

- شبکه CNN برای استخراج ویژگی
- **SVM classifier** برای **object detection**
- مدل رگرسیونی برای تعیین کردن مختصات جعبه های محاطی



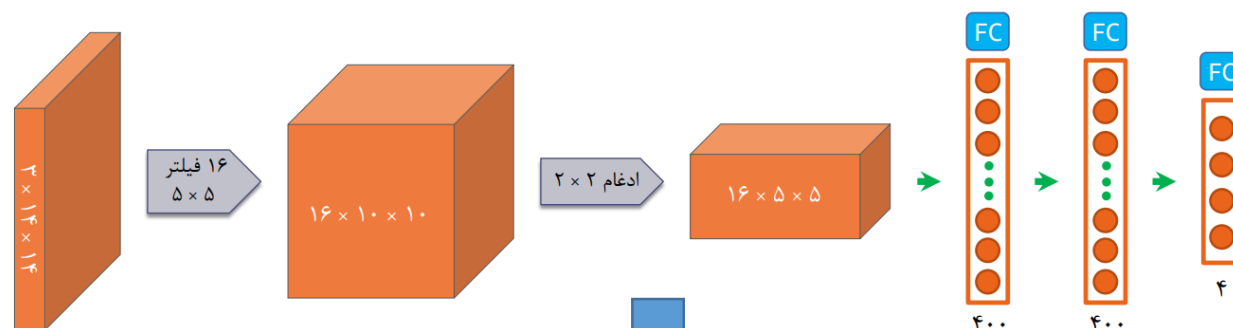
بنابراین این الگوریتم شناسایی **Object** بسیار کند است (زیرا هر یک از موارد بالا بسیار زمان بر است) بنابراین برای زمانی که می خواهیم الگوریتم به صورت **Real-time** کار کند نمی توانیم از این الگوریتم استفاده کنیم. در حال حاضر ورژن سریع تری از این الگوریتم وجود دارد **faster RCNN** و یا **fast RCNN**



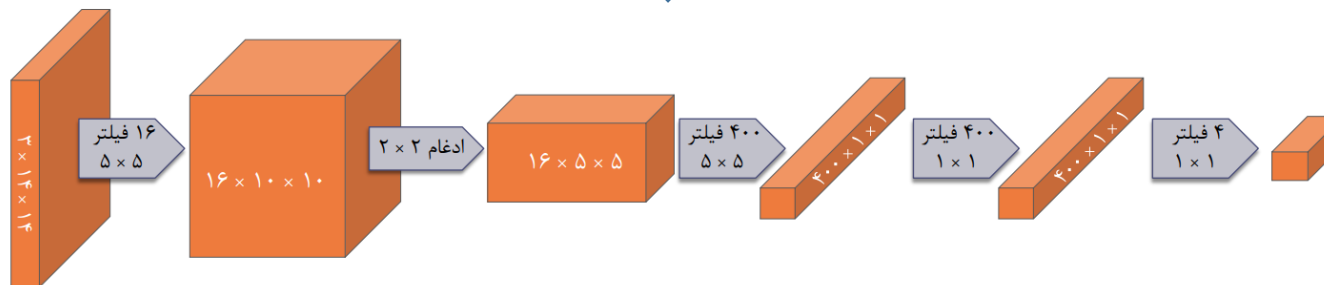
## توضیحی در مورد R-FCN

برخلاف RCCN که از نظر محاسباتی بسیار پیچیده و از نظر زمان بسیار کند است الگوریتم R-FCN بسیار بهینه و با دقت است. همانطور که اسم این الگوریتم مشخص است region based fully convolutional network، این الگوریتم از لایه های کاملاً کانولوشنی استفاده می کند. در واقع ما تمامی لایه های fully connect را نیز به صورت لایه های کانولوشنی پیاده سازی می کنیم (زیرا لایه های کانولوشنی سربار محاسباتی زیادی دارند) همانطور که در شکل زیر می بینیم:

شبکه ی کانولوشنی  
با لایه های  
fully  
connected



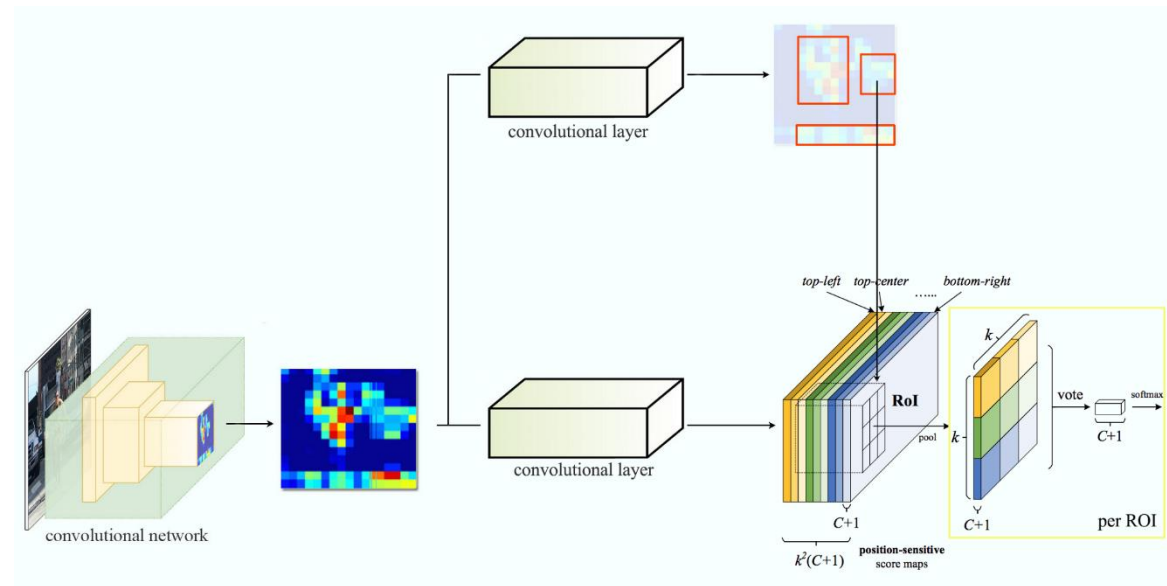
شبکه ی کانولوشنی  
فقط با لایه های  
کانولوشنی





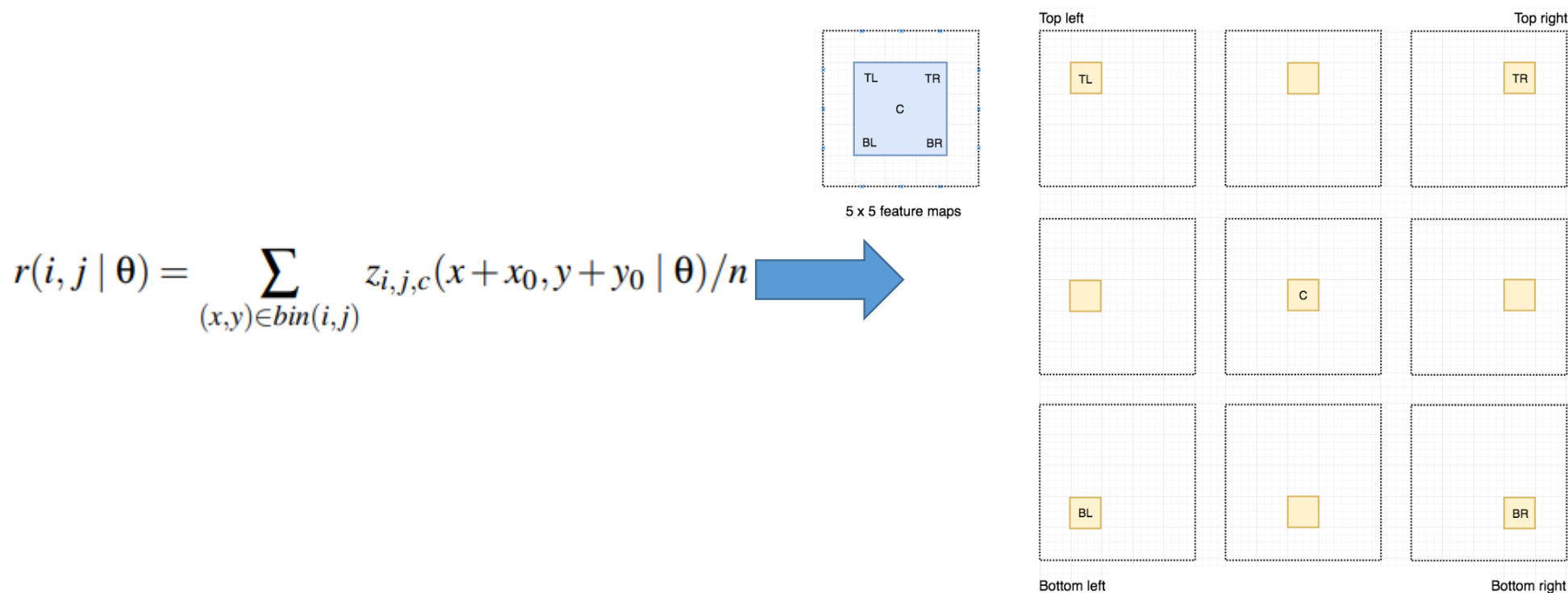
این الگوریتم به دنبال ایجاد یک مصالحه بین translation-invariance در مسائل دسته بندی و translation-variance در تشخیص تصاویر با استفاده از position-sensitive score map و position-sensitive ROI pooling است.

بر اساس این الگوریتم ما به جای استفاده از نواحی مختلف استفاده می کنیم همانطور که در شکل زیر می بینم با استفاده از یک شبکه کانولوشنی و بر اساس یک region proposal network (RPN) شبیه RCNN، ROI های تصویر را استخراج می کنیم. که بخش های مختلف یک تصویر را پوشش می دهد. در واقع این الگوریتم به دنبال استفاده از امکان انجام محاسبات به صورت share هستیم و این مسئله به ما کمک می کند که بخش های مختلف یک تصویر را به صورت موازی و همزمان پردازش کنیم همانطور که در زیر میبینیم:



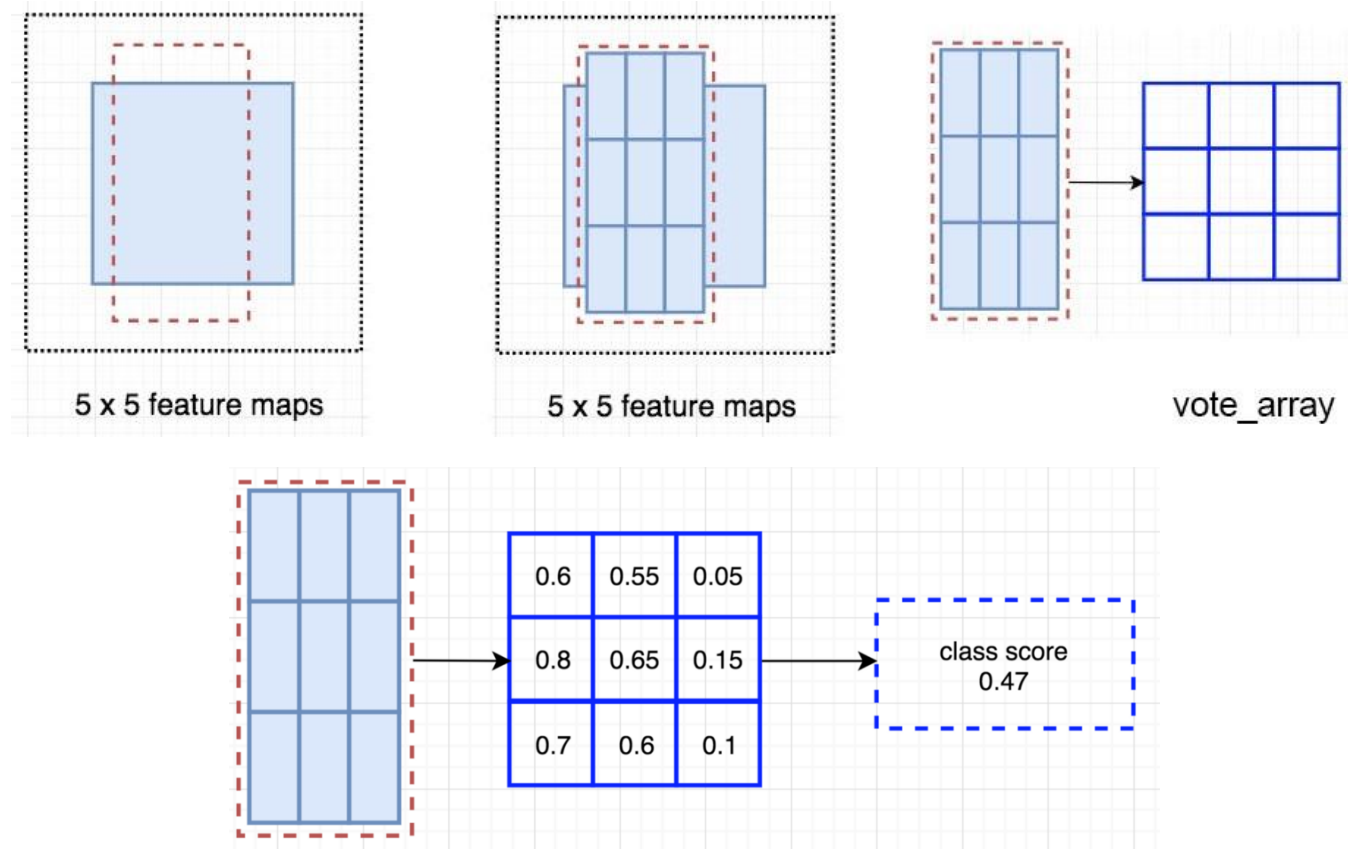


همانطور که در شکل زیر می بینیم علاوه بر لایه ی کانولوشنی لایه کانولوشنی برای استخراج ROI یک لایه ی کانولوشنی دیگر نیز وجود دارد که ویژگی های استخراج شده از آن توسط یک  $k \times k$  در  $k$  به بخش های مختلف تقسیم می شود و احتمال وجود هر قسمت از آن Object در آن قسمت از  $grid$  محاسبه می شود (score map) همانطور که در شکل زیر می بینیم. در واقع به این ویژگی های استخراج شده به ازای هر خانه از  $grid$ ،  $position-sensitive$  score maps می گویند. بنابراین ما به ازای هر کلاس این  $k \times k$  در  $k$  را داریم بنابراین کل ویژگی های استخراج شده برابر می شود با  $K \times K \times (c+1)$  است که  $c$  تعداد کلاس ها و ۱ برای حالتی است که object وجود ندارد (background).





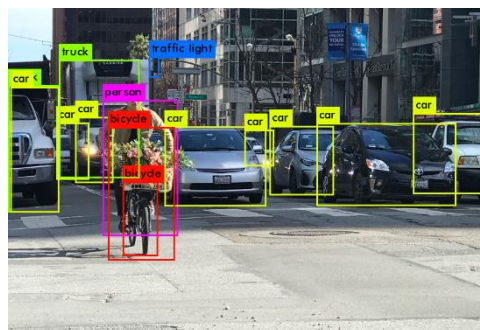
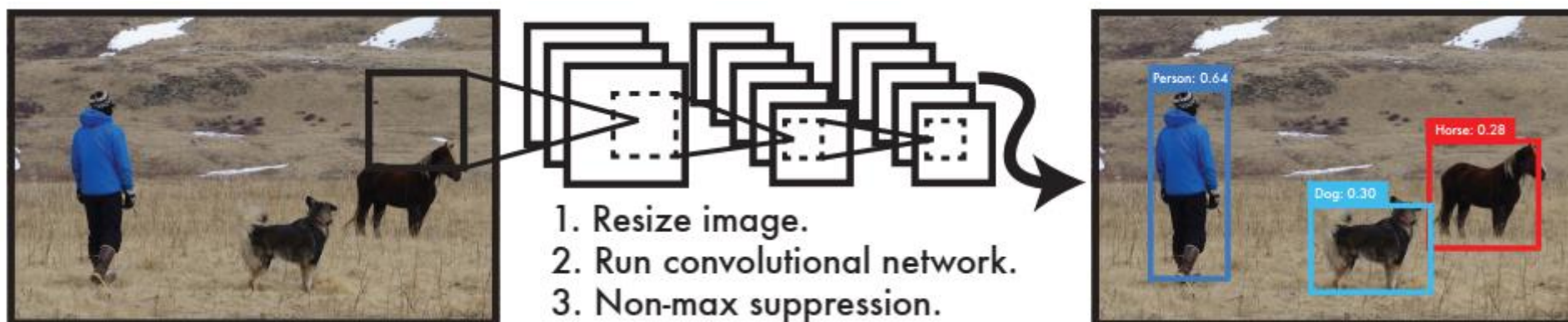
عملیات بعدی ترکیب ROI های بدست آمده با position-sensitive score maps است که به این مرحله position-sensitive ROI-pool می گویند. که به ازای هر class ما یک جدول vote داریم که با استفاده از یک میانگین و سپس soft max احتمال هر کلاس را بدست می آوریم.



R-FCN با کاهش میزان سربار مورد نیاز برای هر ROI سرعت را بهبود می بخشد. region-based feature maps مستقل از ROI هستند و می توانند خارج از هر ROI محاسبه شوند. بنابراین محاسبات ساده تر است و بنابراین R-FCN سریعتر از سریعتر R-CNN است.



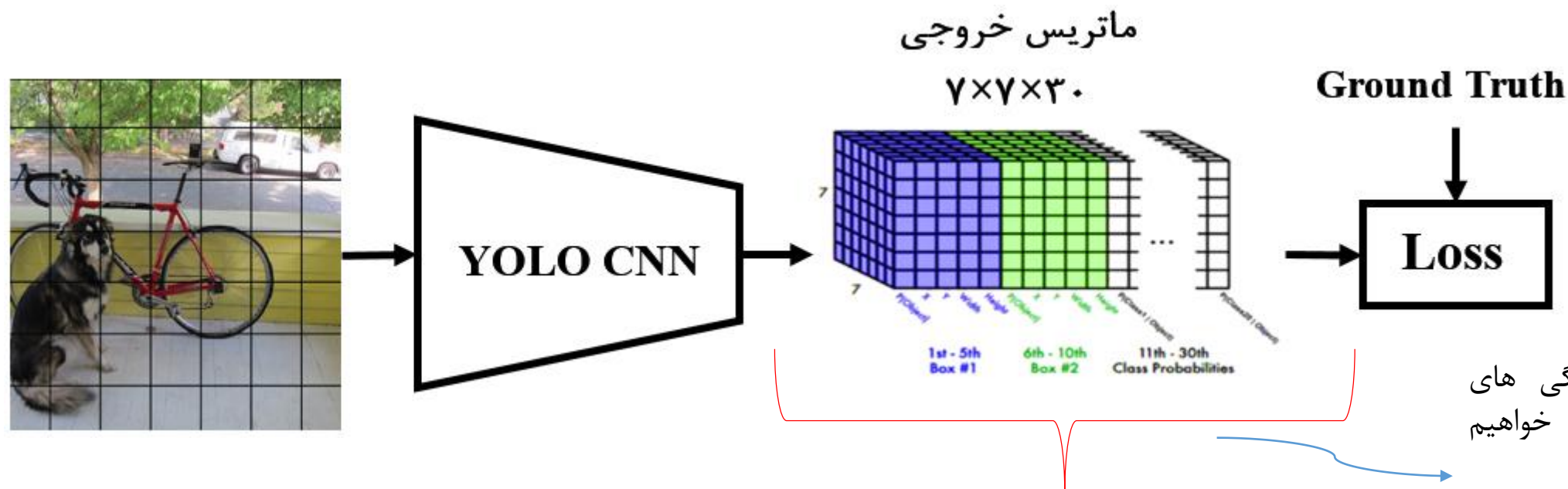
مخفف عبارت You Only Look Once، به معنای “شما فقط یکبار به تصویر نگاه می‌کنید” هست. درواقع، این عبارت به همان قابلیت سیستم بینایی انسان اشاره دارد که با یک نگاه عمل تشخیص اشیا را انجام می‌دهد. بر خلاف روش‌های موجود در تشخیص اشیا مثل RCNN و RFCN الگوریتم YOLO مسئله‌ی تشخیص تصاویر را مثل یک مسئله‌ی رگرسیون در نظر می‌گیرد و مستقیماً از پیکسل‌های تصویر به مختصات باکس و احتمال کلاس‌ها می‌رسد. همانطور که در شکل زیر نگاه می‌کنیم.







در الگوریتم YOLO تصویر ورودی با ابعاد  $448 \times 448 \times 3$  به یک Grid یا شبکه  $S \times S$  تقسیم بندی می شود. این تصویر به شبکه YOLO داده می شود. خروجی شبکه کانولوشنی، ماتریسی به ابعاد  $S \times S \times 30$  خواهد بود. هریک از درایه های ماتریس  $S \times S$  خروجی معادل با یک سلول در شبکه  $S \times S$  ورودی است (به ورودی و خروجی در شکل ۵ دقت کنید). خروجی  $S \times S \times 30$  شامل مختصات باکس ها و احتمال هاست. اگر در فرآیند آموزش Train باشیم، خروجی  $S \times S \times 30$  به همراه باکس های واقعی (Ground Truth) به تابع هزینه داده می شود. مقدار  $S$  در یولو نسخه ۱، برابر با ۷ در نظر گرفته شده است. اگر در فرآیند آزمایش Test باشیم، خروجی  $S \times S \times 30$  به الگوریتم حذف غیر حداکثرها Non-maximum Suppression داده می شود (برای این که ممکن است چندین باکس به صورت هم پوشانی شده روی یک شی قرار بگیرند) تا باکس های ضعیف از بین بروند و تنها باکس های درست در خروجی نمایش داده شوند. در ادامه، در مورد طراحی شبکه YOLO، نحوه آموزش شبکه، تابع هزینه، آزمایش شبکه و غیره توضیح خواهیم داد.



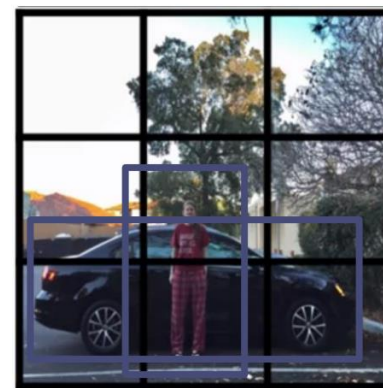
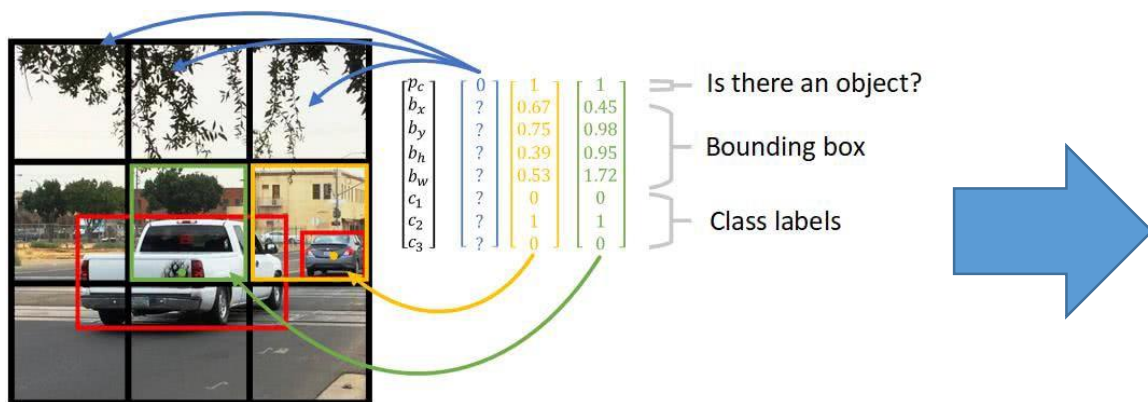


در الگوریتم YOLO ما به ازای هر قسمت از تصویر که Grid بندی شده است یک سری بردار به عنوان  $y$  در نظر می گیریم که به صورت زیر تعریف می شود.

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ C_1 \\ C_2 \\ C_3 \end{bmatrix}$$

همانطور که گفته شد YOLO مثل یک مدل رگرسیونی عمل می کند بنابراین به ازای هر قسمت از grid این برداری را محاسبه می کند که شامل احتمال حضور یک object در آن قسمت است  $p$ . مشخصات جعبه محاطی  $b_x, b_y, b_h, b_w$  و احتمال تعلق به هر کلاس  $c_1$  تا  $c_n$

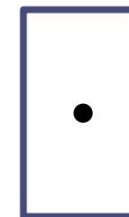
همچنین در این الگوریتم از یک ویژگی به اسم کنگرد استفاده می شود که این لنگر ها به نوعی میانگینی از ابعاد باکس متناظر برای object های متفاوت است. مثلاً لنگر مربوط به اتوبوس باکسی شبیه مستطیل دارد. به ازای لنگر های مختلف ما برداری شبیه بردار بالا را تعریف می کنیم.



$$y = \begin{bmatrix} y_{a1} \\ y_{a2} \end{bmatrix}$$

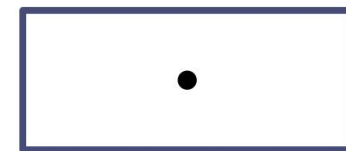
به ازای هر خانه

جعبه لنگر ۱



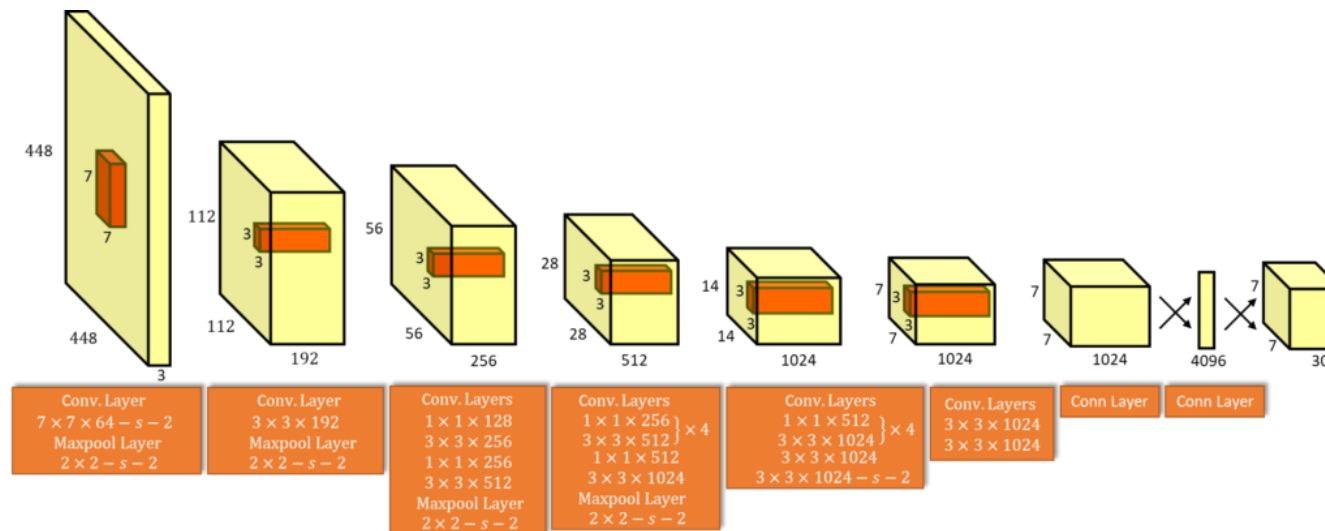
$$y_{a1} = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ C_1 \\ C_2 \\ C_3 \end{bmatrix}$$

جعبه لنگر ۲



$$y_{a2} = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ C_1 \\ C_2 \\ C_3 \end{bmatrix}$$





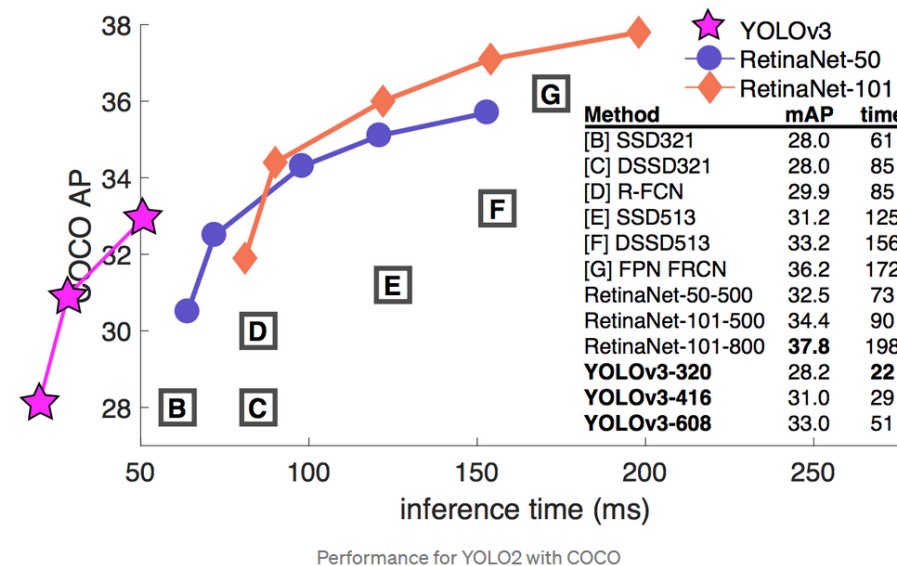
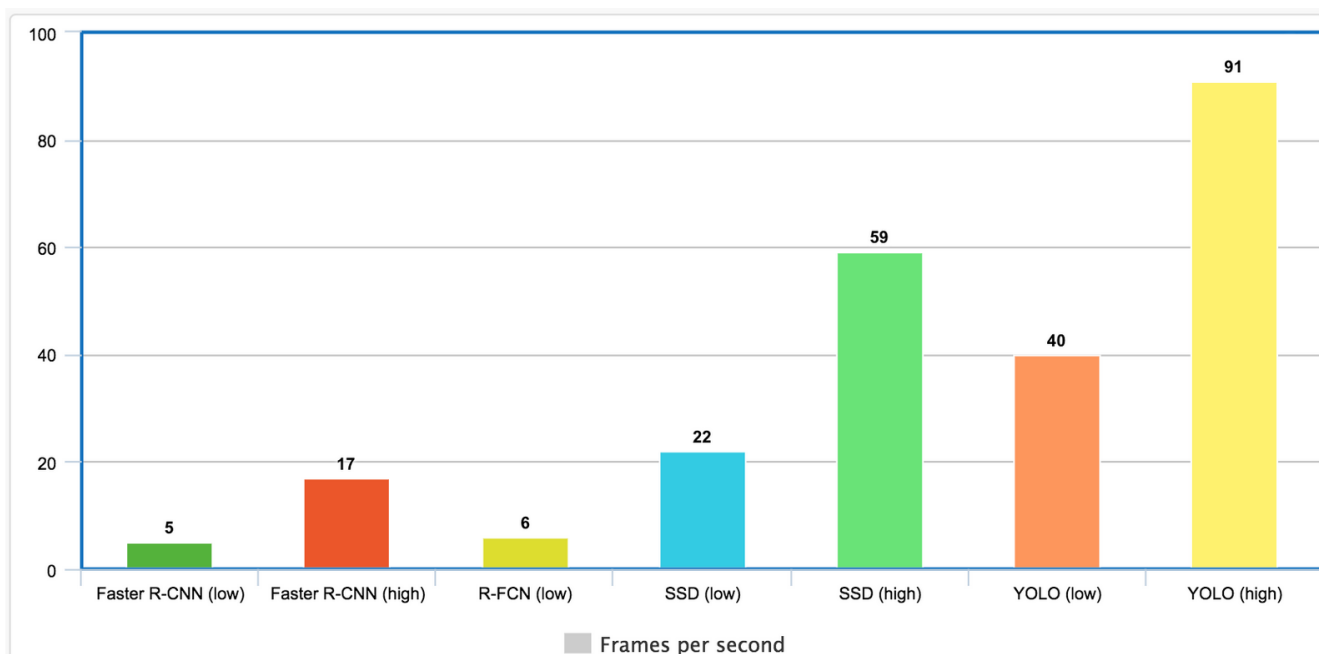
شامل یک شبکه عصبی کانولوشنی Convolutional Neural Network با ۲۴ لایه کانولوشنی برای استخراج ویژگی و همچنین ۲ لایه فولی کانکتد Fully Connected برای پیش‌بینی احتمال و مختصات اشیا است.

همچنین، یک نسخه سریع از YOLO برای جابجایی مرزهای تشخیص اشیای سریع طراحی شده است. YOLO سریع، یک شبکه عصبی با تعداد لایه‌های کانولوشنی کمتر است که در آن از ۹ لایه کانولوشنی بجای ۲۴ لایه کانولوشنی (YOLO اصلی) استفاده شده و البته تعداد فیلترهای هر لایه در YOLO سریع نسبت به YOLO اصلی کمتر است. اندازه ورودی هر دو شبکه ۴۴۸×۴۴۸×۳ و خروجی شبکه نیز یک تانسور ۷×۷×۳۰ از پیش‌بینی‌ها است. در تمامی لایه‌ها از Leaky ReLU استفاده شده است.



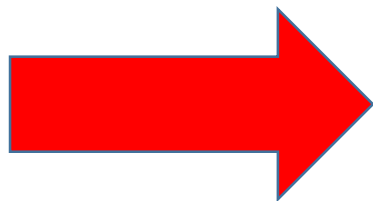
بسیار سریع است. در اینجا، تنها یک شبکه وجود دارد که خیلی ساده به آن ورودی تصویر داده می‌شود (برخلاف RCNN و RFCN) تا شبکه پیش‌بینی‌های تشخیص اشیاء را به ما نشان دهد. دو نسخه شبکه YOLO شامل YOLO اصلی و YOLO سریع طراحی شده است.

بنابراین با توجه به نکات گفته شده الگوریتم YOLO برای اپلیکیشن‌های Real-time بسیار مناسب‌تر از الگوریتم‌های دیگر است. نمونه‌ای مقایسه‌ای سرعت YOLO با دیگر الگوریتم‌ها در شکل زیر وجود دارد.





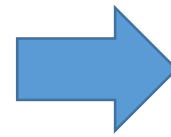
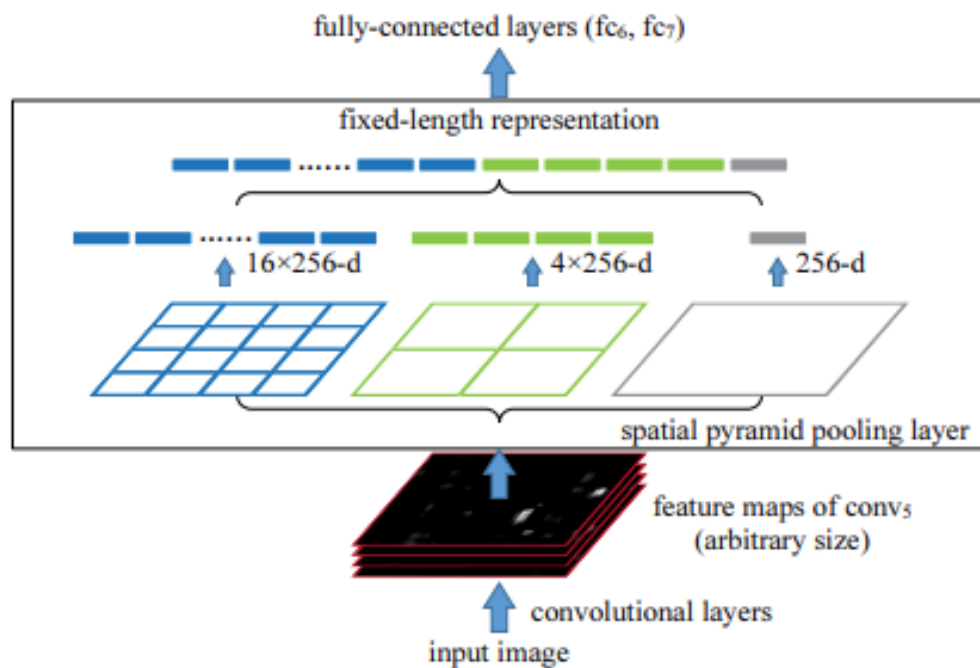
یولو اصلی با کارت گرافیک Titan X با سرعت ۴۵ فریم بر ثانیه اجرا می‌شود. نسخه سریع YOLO هم سرعتی بیش از ۱۵۰ فریم بر ثانیه دارد. یعنی YOLO می‌تواند در یک ویدئوی ۴۰ فریم بر ثانیه در حالت بلادرنگ به تشخیص اشیاء بپردازد. YOLO نسبت به دیگر سیستم‌های تشخیص اشیاء بلادرنگ، به mAP یا همان mean Average Precision دو برابر دست یافته است. عملکرد بهتر نسبت به سایر سیستم‌های بلادرنگ و نه سیستم‌های تشخیص اشیاء قدرتمند مانند Faster R-CNN که بلادرنگ نیستند.



$$\lambda_{coord} \sum_{i=0}^{s^2} \sum_{j=0}^B 1_{ij}^{obj} [(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2] +$$
$$\lambda_{coord} \sum_{i=0}^{s^2} \sum_{j=0}^B 1_{ij}^{obj} [(\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2] +$$
$$\sum_{i=0}^{s^2} \sum_{j=0}^B 1_{ij}^{obj} (C_i - \hat{C}_i)^2 + \lambda_{coord} \sum_{i=0}^{s^2} \sum_{j=0}^B 1_{ij}^{obj} (C_i - \hat{C}_i)^2 +$$
$$\sum_{i=0}^{s^2} 1_{ij}^{obj} \sum_{c \in classes} (p_i(c) - \hat{p}_i(c))^2$$



❖ ژمانی که ما از شبکه های کانوولوشنی برای استخراج ویژگی و از لایه های fully connected برای دسته بندی و پیش بینی استفاده می کنیم به ناچار مجبور می شویم که ابعاد ویژگی های ورودی شبکه ی کانوولوشنی را یکسان در نظر بگیریم و این مسئله به نوعی برای ما محدودیت ایجاد می کند زیرا لایه ی fully connect در غیر اینطورت دچار مشکل می شد. بنابراین گاهی برای اینکه ابعاد ویژگی های ورودی را یکسان کنیم مجبور می شویم داده ها را reshape کنیم که خود این مسئله می تواند ورودی های شبکه و خروجی آن را تحت تاثیر قرار دهد. بنابراین ایده ی استفاده از Spatial Pyramid pooling شکل گرفت. طبق این فرضیه وقتی ما لایه های ورودی را از لایه های کانوولوشنی عبور میدادیم feature map هایی را استخراج میکنیم می توانیم لایه های پولینگ (مثلا ۳ لایه) را در نظر بگیریم و با استفاده از آن ابعاد را یکسان کنیم همانطور که در شکل زیر می بینیم.



همانطور که در شکل مقابل می بینیم پس از استخراج ویژگی ها از لایه های کانوولوشی ما یک سری ویژگی با عمیق  $d$  داریم به ازای هر لایه (هر bin) مثلا ما یک لایه پولینگ با سایز ۱ در نظر میگیریم بنابر این به ازای هر bin و بر اساس متد pooling که در نظر گرفتیم مثلا max pooling یک ویژگی استخراج می کنیم بنابر این پس از اولین pooling ما به اندازه ی  $d$  تا ویژگی داریم. سپس یک لایه pooling دیگر با سایز دو مثلا ایجاد می کنیم که در این صورت به ازای هر bin ما ۴ تا ویژگی استخراج می کنیم که در این صورت کل ویژگی های استخراج شده برابر با  $4*d$  می شود. پس از این لایه یک لایه ی پولینگ دیگر استفاده می کنیم با سایز پنجره ۴ بنابرین در هر bin ما ۱۶ ویژگی استخراج می کنیم که در این صورت کل ویژگی های استخراج شده در این لایه برابر با  $16*d$  می شود. بنابراین کل ویژگی های استخراج شده به ازای لایه ی Spatial Pyramid pooling برابر با  $16d+4d+d$  می شود که یک سایز ثابت است.



در مسئله ی object detection ما نیاز داریم تا بخش های تصویر را crop کنیم یا wrap کنیم و تعداد زیادی تصویر با ابعاد مختلف ایجاد می شود (بر اساس RCNN). بنابراین پیوسته ابعاد تصویر ورودی ما متغیر هستند و روش Spatial Pyramid pooling می تواند بسیار کمک کننده باشد. بر این اساس به ازای هر تصویر با سایز گوناگون لایه ی SSP خروجی با سایز یکسان تولید می کند همانطور که در شکل زیر می بینیم:

