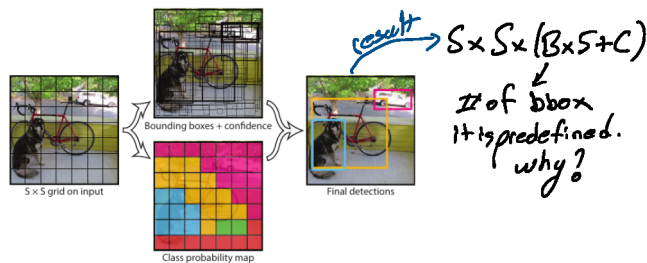


* Reason globally on the entire image & all objects.

Confidence:
↳ Confidence = $\text{Pr}(\text{Obj}) \times 10^{\frac{\text{pred}}{\text{truth}}}$
↳ BBox: x, y, w, h, conf
 center relative to whole img.

↳ Also, predict conditional class probs: $P(c_i | \text{obj})$
↳ per grid cell

$$\rightarrow P(c|s; \text{obj}) \times P(\text{obj}) \times \text{vol}_p^t = P(c|s; \text{obj}) \times \text{vol}_p^t$$



→ Reduction layers (1×1) after conv layers. (3×3)

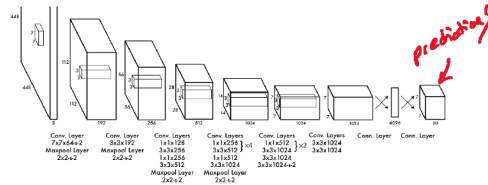


Figure 3: The Architecture. Our detection network has 24 convolutional layers followed by 2 fully connected layers. Alternating 1×1 convolutional layers reduce the features space from preceding layers. We pretrain the convolutional layers on the ImageNet classification task at half the resolution (224×224 input image) and then double the resolution for detection.

- Pretrained on ImageNet.

- Loss: sum-squared error

- ↳ penalizes localization errors and cls. err equally
- ↳ to remedy loss diff. between small and large bboxes, predict $\sqrt{w, h}$.

$$\begin{aligned} \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{I}_{ij}^{\text{obj}} \left[(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\ + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{I}_{ij}^{\text{obj}} \left[\left(\sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left(\sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \\ + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{I}_{ij}^{\text{obj}} (C_i - \hat{C}_i)^2 \\ + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{I}_{ij}^{\text{obj}} (C_i - \hat{C}_i)^2 \\ + \sum_{i=0}^{S^2} \mathbb{I}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \end{aligned}$$

Limitations

↳ struggles w/ small objects

- ↳ group of obj. close to each other like flocks of birds

1) main source of errors is incorrect localization

? R-CNN, region proposals, Selective Search

2 Grasp Detection