

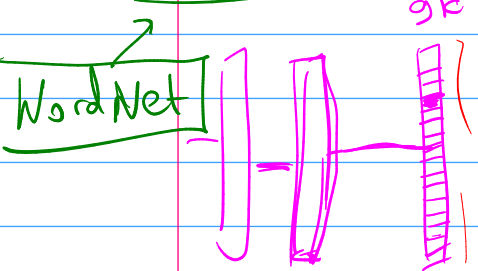
$$L = \lambda_1 \sum_s \sum_b [bbox] \cdot \left( (\sqrt{w} - \sqrt{\hat{w}})^2 + (\sqrt{h} - \sqrt{\hat{h}})^2 + (x - \hat{x})^2 + (y - \hat{y})^2 \right)$$

$$+ \lambda_2 \sum_s [bbox] \cdot L_{class}(c, \hat{c})$$

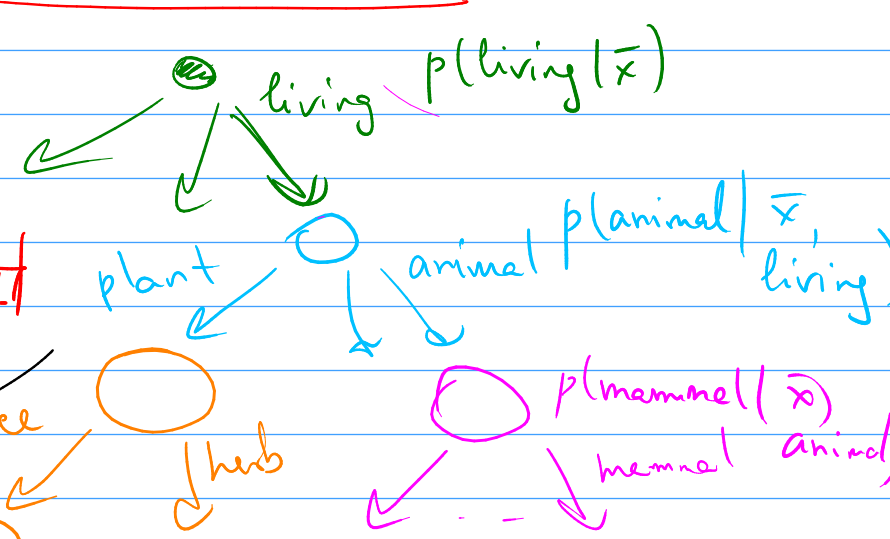
$$+ \lambda_3 \sum_s \sum_b [bbox] \cdot (1-p)^2 + \lambda_4 \sum_s \sum_b [bbox] \cdot p^2$$

ImageNet - 120 megapixel

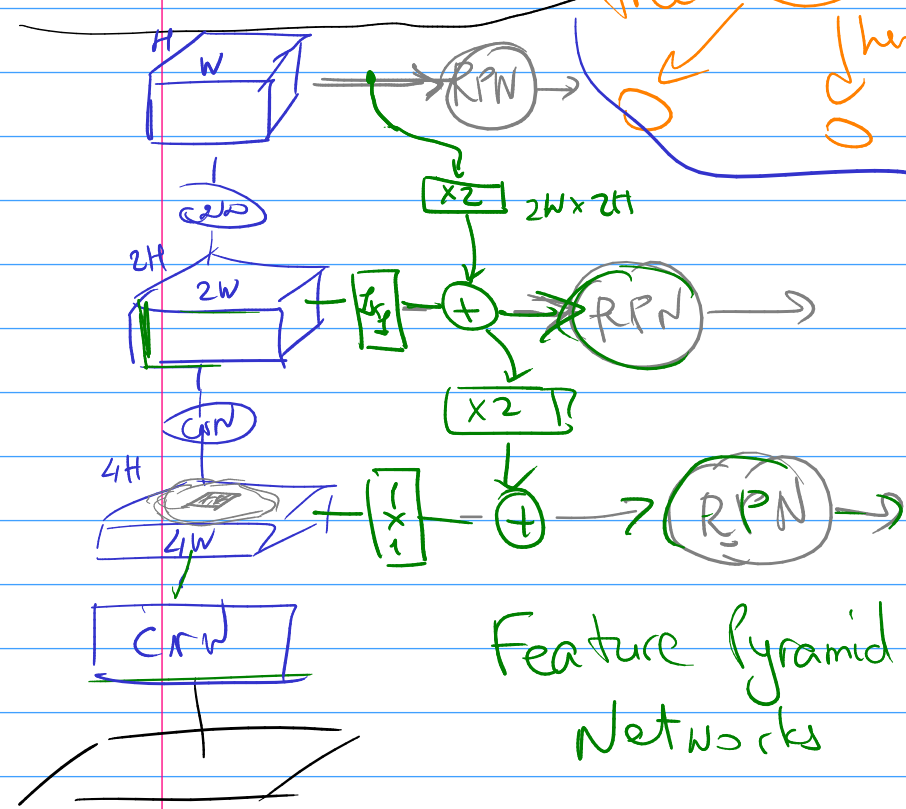
9000



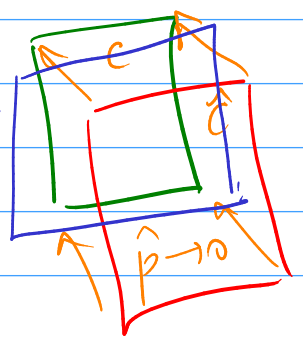
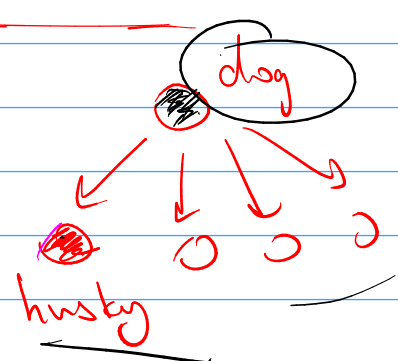
hierarchical softmax



$$L = l_1 + l_2 + l_4 + \dots + l_n$$



Feature Pyramid Networks



$\rightarrow p \rightarrow 1$   
 $\rightarrow 0$

$$\text{BCE} = -t \log p - (1-t) \log(1-p)$$

5000

$$p = 0.2 \quad \text{BCE} = -\log 4/5 = 0.32$$

$(1-p)^{\delta}$

5

$$p = 0.8 \quad \text{BCE} = -\log 1/5 = 2.32 \dots$$

$$\text{Focal loss} = -t (1-p)^{\delta} \log - (1-t) p^{\delta} \log(1-p)$$