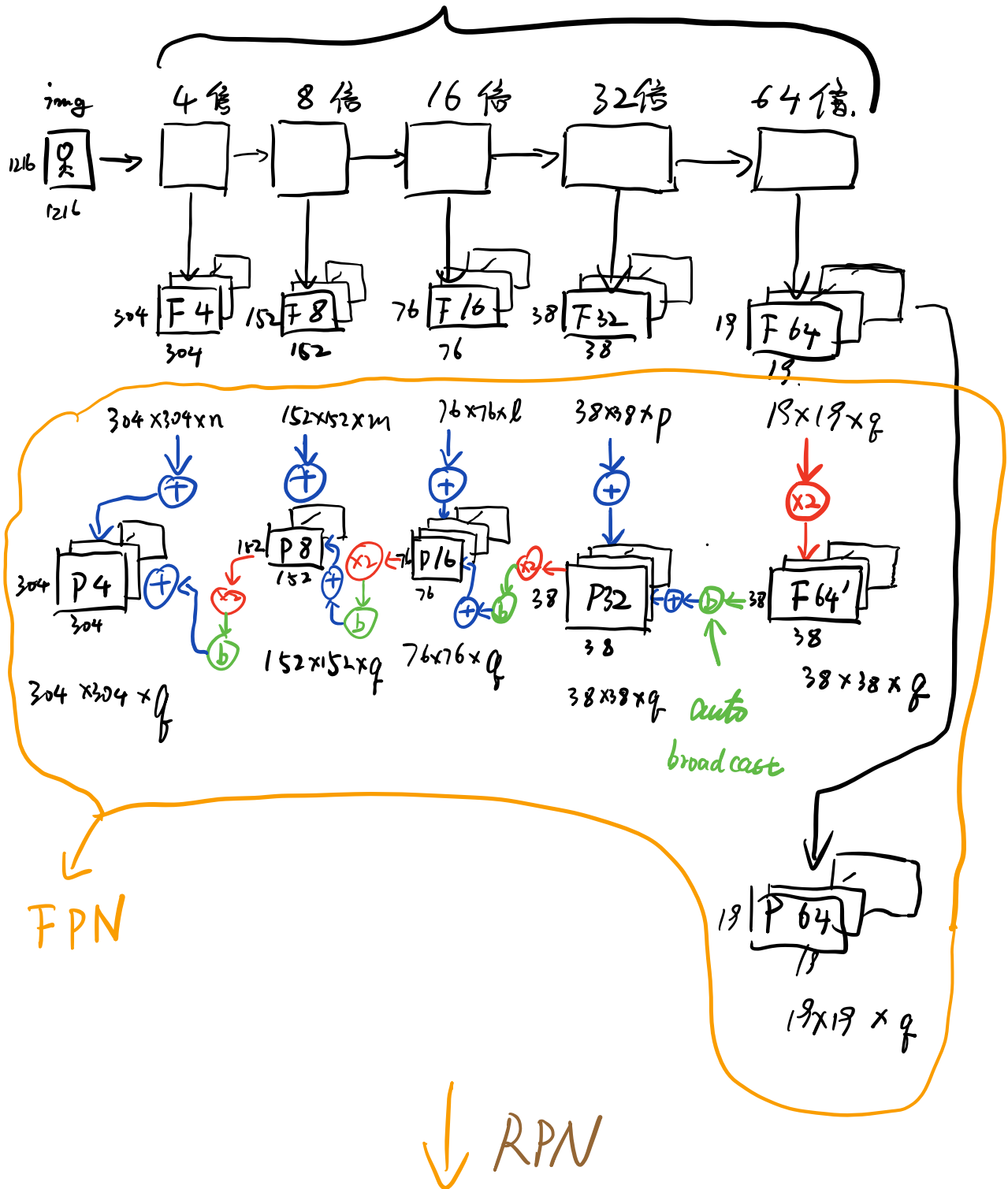
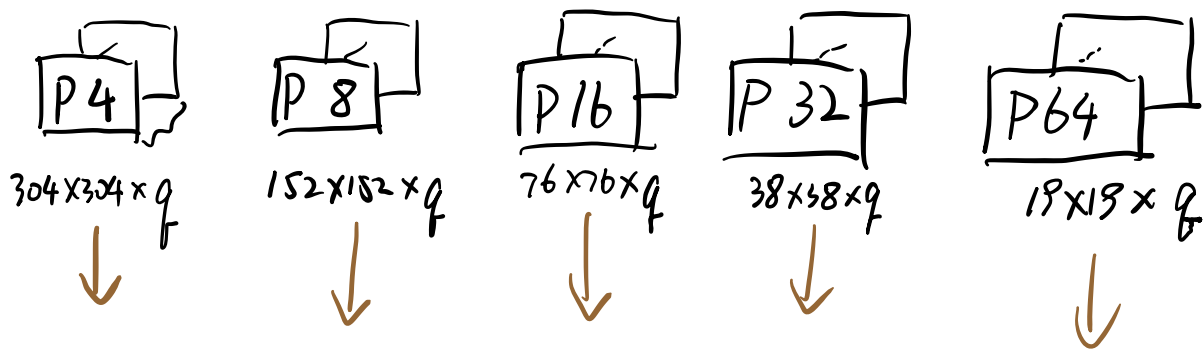


Pre-trained ResNet Convolution Block





size = 3

stride = 1

padding = 1

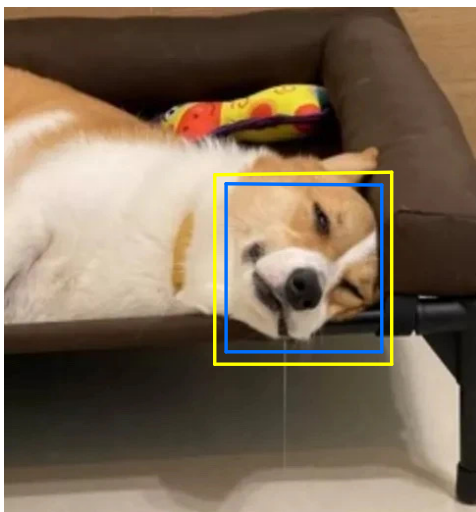
感受野从小到大，

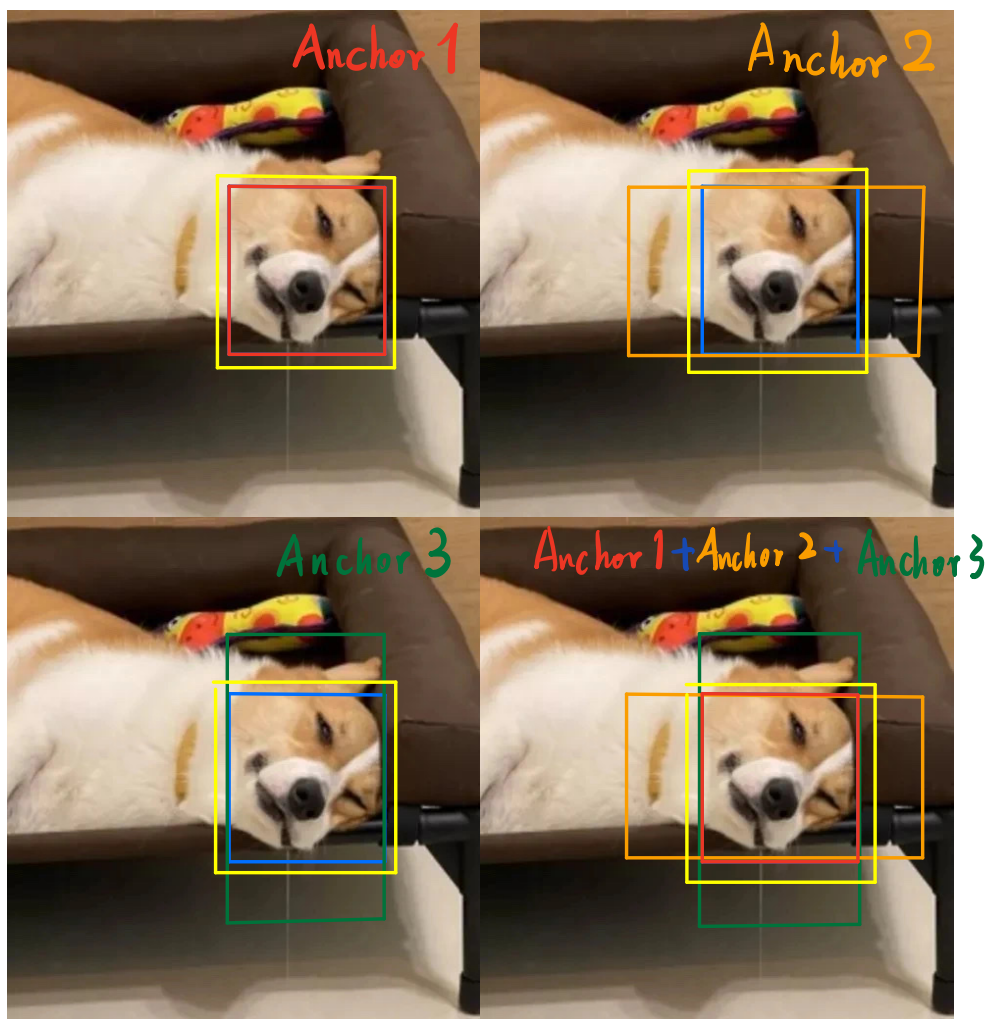
在特征图中每个像素对应的原图中的感受野上

生成 3 个 anchor 比例为 1:1, 1:2, 2:1

设蓝框为感受野

黄框为 Ground Truth



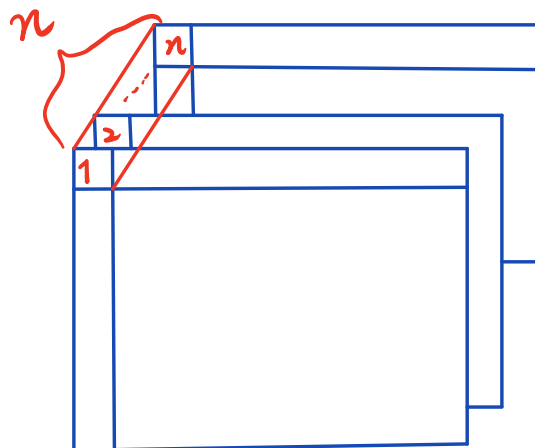


对于 P4 $304 \times 304 \times q$.

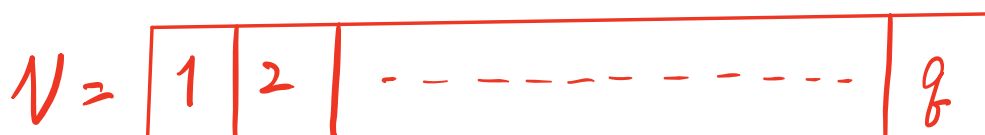
304×304 个像素 \times 3 个 anchor.

$304 \times 304 \times 3 = 277\,248$ 个 anchor.

使用 P_4 中每个像素对应的长度为 n 的向量 v



↓
对应向量 v

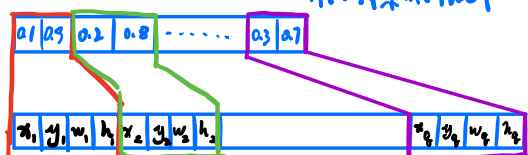


使用 v 计算 { ① 分类 → 输出 2 个值

每个 anchor 中的内容是前景还是背景的概率

② 回归 → 输出 4 个值

分类输出: anchor 中内容属于前、背景的概率



回归输出: anchor 的偏移量和缩放量.

每个 anchor 的中心点偏移量和长宽缩放量.

Proposal 层

- ① 利用回归结果对所有 anchor 进行修正, 得到修正后检测框.
- ② 根据分类结果, 按前景概率^{降序}排序, 取前 6000 个框.
- ③ 限定超出图像边界的框, 防止后续 ROI pooling 超出边界
- ④ 对~~剩余~~框进行非极大值抑制
- ⑤ 输出对应输入图像尺度归一化后的坐标值
 $[x_1, y_1, x_2, y_2]$

ROI pooling.

由于有 FPN 输出了多个 Feature maps

在 ROI pooling 时要计算一个 k

来确定应该使用哪个 Feature map.

$$k = \lfloor k_0 + \log_2(\sqrt{wh} / 224) \rfloor$$

224: 图像输入尺寸

w, h : ROI 区域的长宽 $\neq 7$

k_0 : 基准值 = 4

例

ROI = 112 x 112.

$$k = \lfloor 4 + \log_2\left(\frac{112}{224}\right) \rfloor$$

$$= [4 - 1]$$

$$= 3$$

取 P_3 层.

k 要取整, 且在 2-5 间, 超出就截断