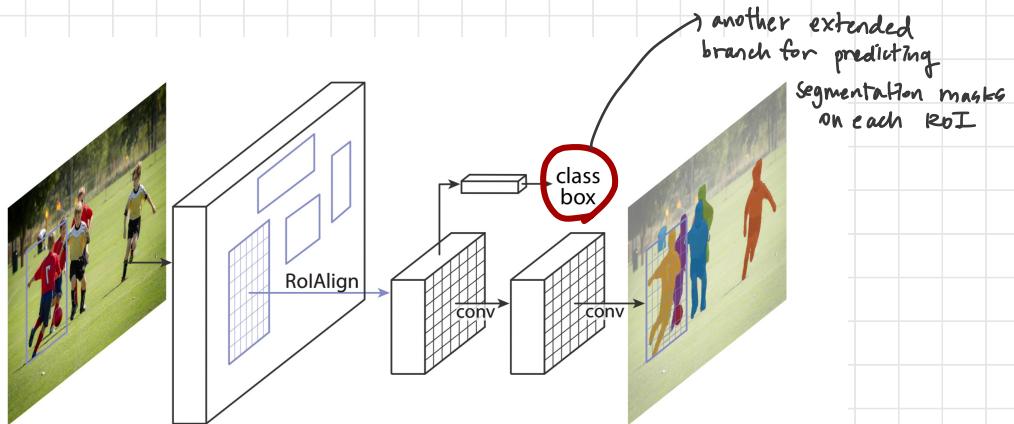
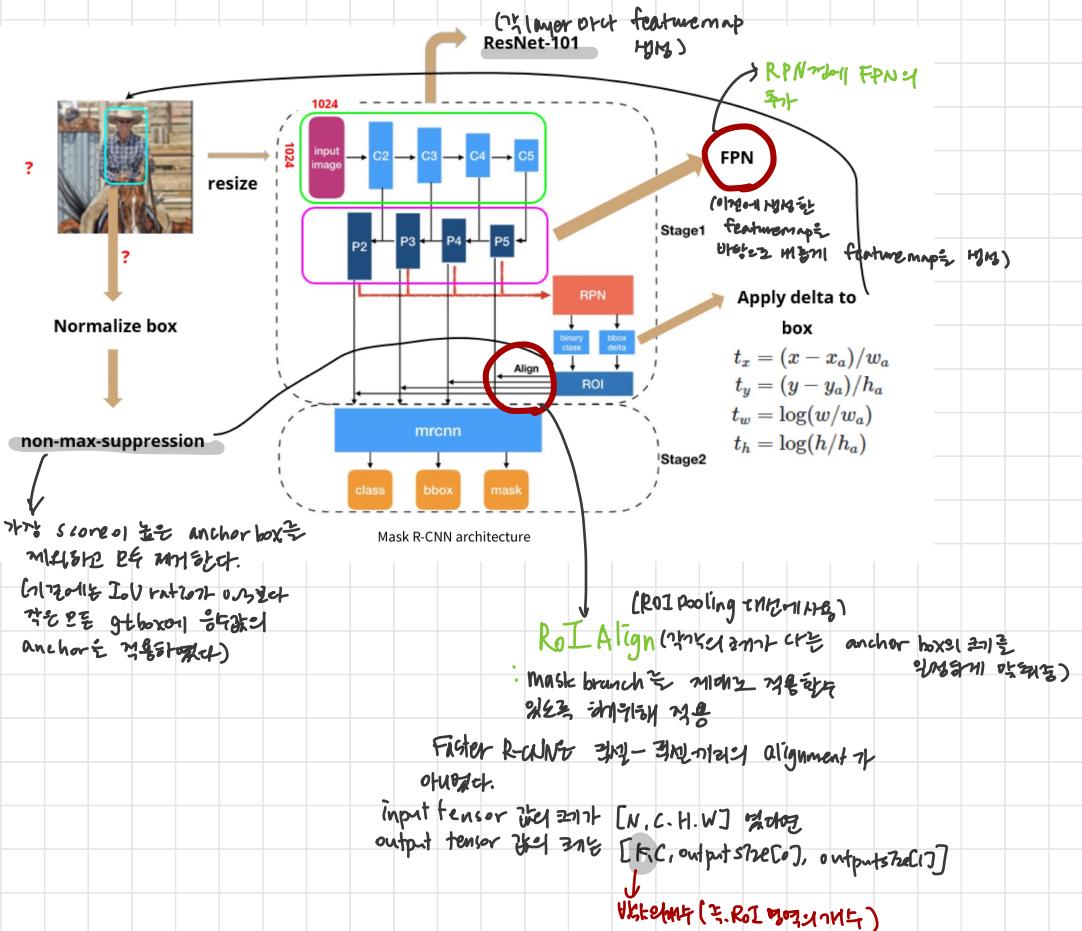




Mask R-CNN = Segmentation mask + RoI + Faster R-CNN



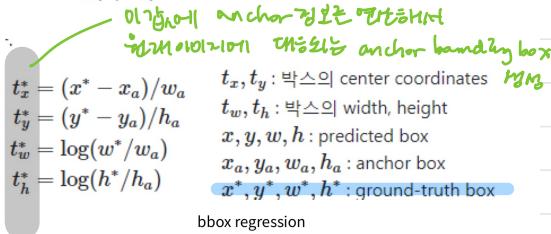
즉각 "Instance segmentation"의 가능성을 증명하는 예



< RPN (Region Proposal Network) >

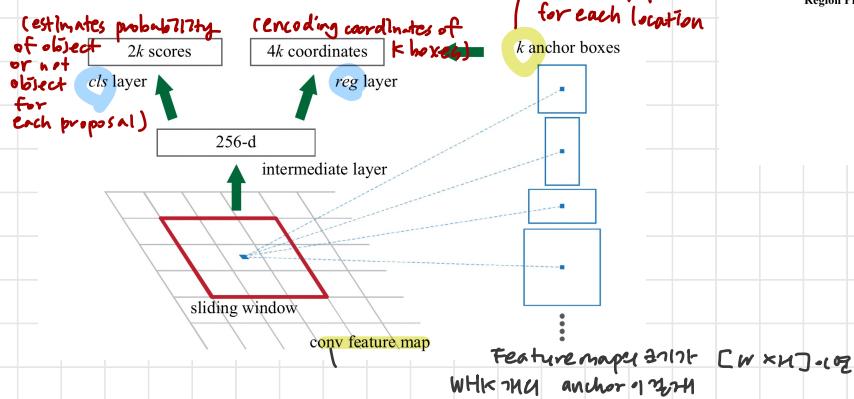
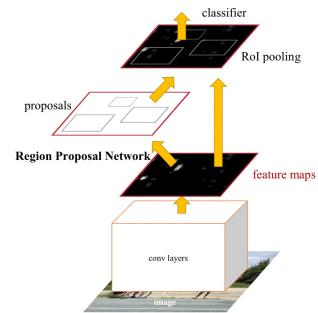
• RPN은 완전히 convolution network로 이루어진다.

$$\begin{aligned} t_x &= (x - x_a)/w_a \\ t_y &= (y - y_a)/h_a \\ t_w &= \log(w/w_a) \\ t_h &= \log(h/h_a) \end{aligned}$$



① Anchors

- 각각의 sliding window location에서 multiple region proposal을 얻는다



② Loss Function

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

log loss over two classes (object vs non-object)

p_i^* :锚이 object인지 여부

$L_{cls}(p_i, p_i^*)$:锚이 object인지 여부

$L_{reg}(t_i, t_i^*)$:锚의 위치 예측 오류

p_i :锚이 object인 확률

index of the anchor in a mini-batch

t_i^* :锚의 실제 위치

\rightarrow Smooth L1 defined as robust loss function

p_i^* 이 0.5가 18-23%의 비율로
 锚이 object인 경우 IoU score > 0.3 일 때는 anchor
 를 제거하는 경우가 많다 (0.3은锚이 object인 경우)
 10% 경우에
 73%의锚은 reg layer에 속하지 않거나 21%의锚은

Mask R-CNN: Mask R-CNN adopts the same two-stage procedure, with an identical first stage (which is RPN). In the second stage, in parallel to predicting the class and box offset, Mask R-CNN also outputs a binary mask for each RoI. This is in contrast to most recent systems, where classification *depends* on mask predictions (*e.g.* [33, 10, 26]). Our approach follows the spirit of Fast R-CNN [12] that applies bounding-box classification and regression in parallel (which turned out to largely simplify the multi-stage pipeline of original R-CNN [13]).

Formally, during training, we define a multi-task loss on each sampled RoI as $L = L_{cls} + L_{box} + L_{mask}$. The classification loss L_{cls} and bounding-box loss L_{box} are identical as those defined in [12]. The mask branch has a Km^2 -dimensional output for each RoI, which encodes K binary masks of resolution $m \times m$, one for each of the K classes. To this we apply a per-pixel sigmoid, and define L_{mask} as the average binary cross-entropy loss. For an RoI associated with ground-truth class k , L_{mask} is only defined on the k -th mask (other mask outputs do not contribute to the loss).

Our definition of L_{mask} allows the network to generate masks for every class without competition among classes; we rely on the dedicated classification branch to predict the

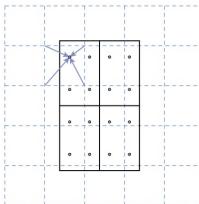


Figure 3. **RoIAlign:** The dashed grid represents a feature map, the solid lines an RoI (with 2×2 bins in this example), and the dots the 4 sampling points in each bin. RoIAlign computes the value of each sampling point by bilinear interpolation from the nearby grid points on the feature map. No quantization is performed on any coordinates involved in the RoI, its bins, or the sampling points.

class label used to select the output mask. This *decouples* mask and class prediction. This is different from common practice when applying FCNs [30] to semantic segmentation, which typically uses a per-pixel *softmax* and a *multinomial* cross-entropy loss. In that case, masks across classes compete; in our case, with a per-pixel *sigmoid* and a *binary* loss, they do not. We show by experiments that this formulation is key for good instance segmentation results.

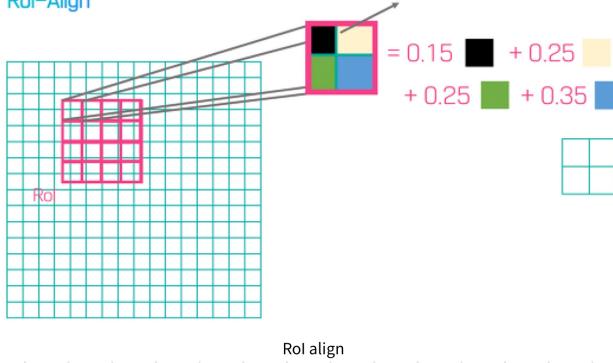
기존의 R-CNN은 레이어별로 ROI를 처리하는
계산은 $L_{cls} + L_{box}$ 외에 각각의 Mask의 Loss를 처리하는 계산.
K개의 class마다 K개의 binary mask
resolution은 각각의 ROI에 Km^2 의 크기로 나누어
loss를 처리하는 계산.

→ layer별로 처리
classification의 mask prediction 계산
각각의 Mask R-CNN 계산하는
classification layer 계산하는
output mask 계산하는 계산 계산.

→ class별로 처리하는 계산
[per pixel sigmoid]
[binary loss]
→ 각각의 Mask를 계산하는 계산입니다.

<RoI-Align>

RoI-Align

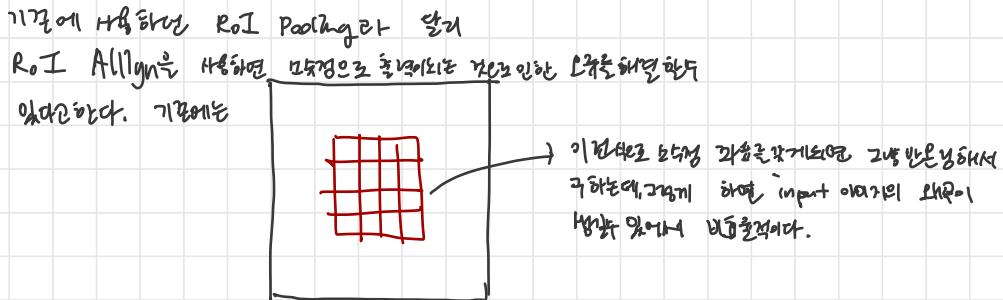


RoIPool: RoIPool [12] is a standard operation for extracting a small feature map (e.g., 7×7) from each ROI. RoIPool first quantizes a floating-number ROI to the discrete granularity of the feature map, this quantized ROI is then subdivided into spatial bins which are themselves quantized, and finally feature values covered by each bin are aggregated (usually by max pooling). Quantization is performed, e.g., on a continuous coordinate x by computing $[x/16]$, where 16 is a feature map stride and $[.]$ is rounding; likewise, quantization is performed when dividing into bins (e.g., 7×7). These quantizations introduce misalignments between the ROI and the extracted features. While this may not impact classification, which is robust to small translations, it has a large negative effect on predicting pixel-accurate masks.

To address this, we propose an *RoIAlign* layer that removes the harsh quantization of RoIPool, properly aligning the extracted features with the input. Our proposed change is simple: we avoid any quantization of the ROI boundaries

or bins (i.e., we use $x/16$ instead of $[x/16]$). We use bilinear interpolation [22] to compute the exact values of the input features at four regularly sampled locations in each ROI bin, and aggregate the result (using max or average), see Figure 3 for details. We note that the results are not sensitive to the exact sampling locations, or how many points are sampled, as long as no quantization is performed.

RoIAlign leads to large improvements as we show in §4.2. We also compare to the RoIWarp operation proposed in [10]. Unlike RoIAlign, RoIWarp overlooked the alignment issue and was implemented in [10] as quantizing ROI just like RoIPool. So even though RoIWarp also adopts bilinear resampling motivated by [22], it performs on par with RoIPool as shown by experiments (more details in Table 2c), demonstrating the crucial role of alignment.

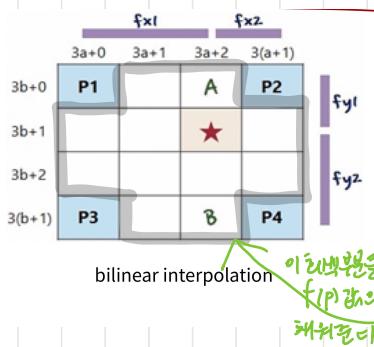


<Non-Max Suppression>

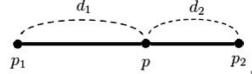
~ object detection 과정에서 사용되는 NMS의 anchor을 살펴보자

기본 맵의 각 위치에 정복가능한 anchor를 놓고 나머지는
겹치지 않는 방향을 사용된다.

<Resize Input Image>



Mask R-CNN에서는 backbone을 ResNet 101을 사용하는데
ResNet은 input 이미지의 크기가 800 - 1024일 때 성능이 좋다.
따라서 이미지를 해상도 size로 1024x1024로 고정
필요하다.

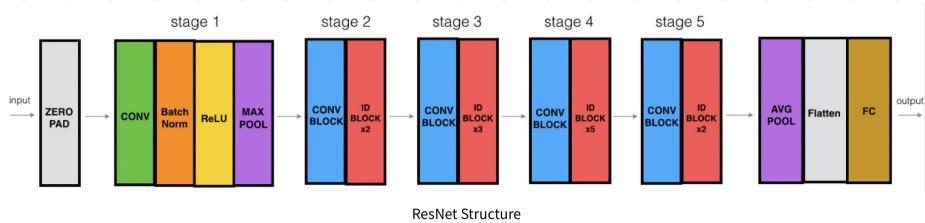


일반적으로 두 지점 p_1, p_2 에서의 데이터 값 각각 $f(p_1), f(p_2)$ 일 때, p_1, p_2 사이의 임의의
지점 p 에서의 데이터 값 $f(p)$ 은 다음과 같이 계산할 수 있다.

$$f(p) = \frac{d_2}{d_1 + d_2} f(p_1) + \frac{d_1}{d_1 + d_2} f(p_2)$$

대부분 1024x1024짜리 size
에서는 이미지 크기를 zero padding을
해당 사이즈로 맞춘다.

<Backbone Network – ResNet 101>



Feature Pyramid Network

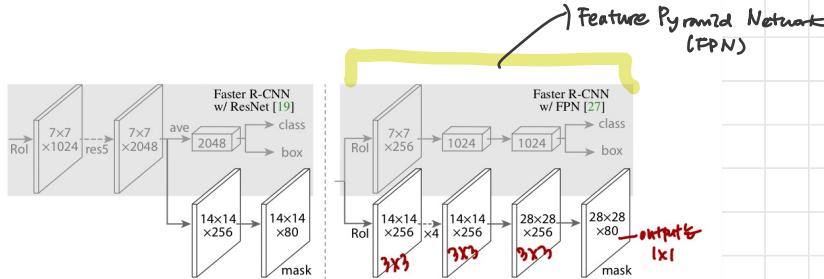
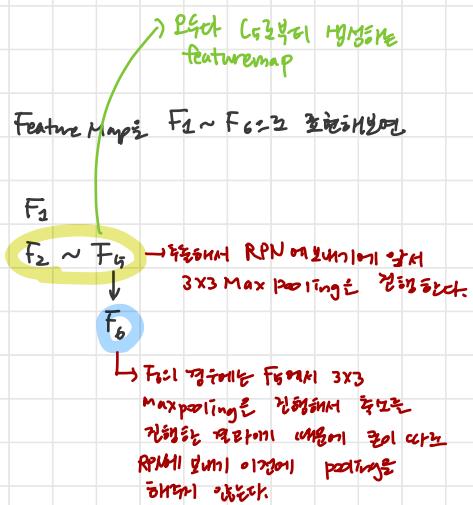
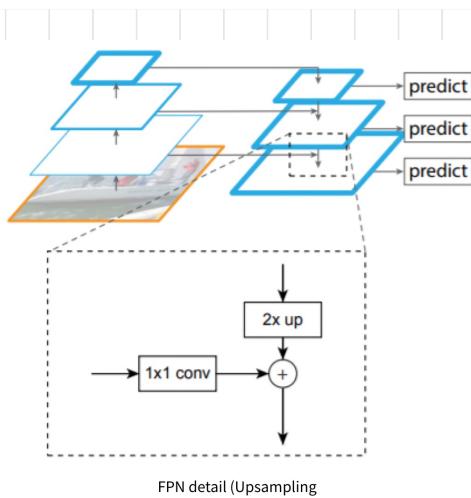


Figure 4. Head Architecture: We extend two existing Faster R-CNN heads [19, 27]. Left/Right panels show the heads for the ResNet C4 and FPN backbones, from [19] and [27], respectively, to which a mask branch is added. Numbers denote spatial resolution and channels. Arrows denote either conv, deconv, or *fc* layers as can be inferred from context (conv preserves spatial dimension while deconv increases it). All convs are 3×3 , except the output conv which is 1×1 , deconv are 2×2 with stride 2, and we use ReLU [31] in hidden layers. *Left*: ‘res5’ denotes ResNet’s fifth stage, which for simplicity we altered so that the first conv operates on a 7×7 RoI with stride 1 (instead of 14×14 / stride 2 as in [19]). *Right*: ‘ $\times 4$ ’ denotes a stack of four consecutive convs.



We denote the **backbone** architecture using the nomenclature *network-depth-features*. We evaluate ResNet [19] and ResNeXt [45] networks of depth 50 or 101 layers. The original implementation of Faster R-CNN with ResNets [19] extracted features from the final convolutional layer of the 4-th stage, which we call C4. This backbone with ResNet-50, for example, is denoted by ResNet-50-C4. This is a common choice used in [19, 10, 21, 39].

↑ 모델이 나온 후 모델 성능의 경우에 는
mobilenet_v2를 (192x152), 그것의
feature 맵은 잘 생겼다.

↑ backbone network는 ResNet을 사용하는
(276-11M) Resnet + FPN을 사용하는 것이 성능 좋다)

MaskRCNN(
 (transform): GeneralizedRCNNTransform(
 Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
 Resize(min_size=(800,), max_size=1333, mode='bilinear')
)
 (backbone): Sequential(
 (0): ConvBNActivation(
 (0): Conv2d(3, 32, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1),
 bias=False)
 (1): BatchNorm2d(32, eps=1e-05, momentum=0.1, affine=True,
 track_running_stats=True)
 (2): ReLU6(inplace=True)
)
 (1): InvertedResidual(
 (conv): Sequential(
 (0): ConvBNActivation(
 (0): Conv2d(32, 32, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1),
 groups=32, bias=False)
 (1): BatchNorm2d(32, eps=1e-05, momentum=0.1, affine=True,
 track_running_stats=True)
 (2): ReLU6(inplace=True)
)
 (1): Conv2d(32, 16, kernel_size=(1, 1), stride=(1, 1), bias=False)
 (2): BatchNorm2d(16, eps=1e-05, momentum=0.1, affine=True,
 track_running_stats=True)
)
)
 (2): InvertedResidual(
 (conv): Sequential(
 (0): ConvBNActivation(
 (0): Conv2d(16, 96, kernel_size=(1, 1), stride=(1, 1), bias=False)
 (1): BatchNorm2d(96, eps=1e-05, momentum=0.1, affine=True,
 track_running_stats=True)
 (2): ReLU6(inplace=True)
)
 (1): ConvBNActivation(
 (0): Conv2d(96, 96, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1),
 groups=96, bias=False)
 (1): BatchNorm2d(96, eps=1e-05, momentum=0.1, affine=True,
 track_running_stats=True)
 (2): ReLU6(inplace=True)
)
 (2): Conv2d(96, 24, kernel_size=(1, 1), stride=(1, 1), bias=False)
 (3): BatchNorm2d(24, eps=1e-05, momentum=0.1, affine=True,
 track_running_stats=True)
)
)
 (3): InvertedResidual(
 (conv): Sequential(
 (0): ConvBNActivation(
 (0): Conv2d(24, 144, kernel_size=(1, 1), stride=(1, 1), bias=False)
 (1): BatchNorm2d(144, eps=1e-05, momentum=0.1, affine=True,

ResNet 9x
FPN은 192x152
생겼다

For the network *head* we closely follow architectures presented in previous work to which we add a fully convolutional mask prediction branch. Specifically, we extend the Faster R-CNN box heads from the ResNet [19] and FPN [27] papers. Details are shown in Figure 4. The head on the ResNet-C4 backbone includes the 5-th stage of ResNet (namely, the 9-layer ‘res5’ [19]), which is compute-intensive. For FPN, the backbone already includes res5 and thus allows for a more efficient head that uses fewer filters.

We note that our mask branches have a straightforward structure. More complex designs have the potential to improve performance but are not the focus of this work.

```
(rpn): RegionProposalNetwork(  
    (anchor_generator): AnchorGenerator()  
    (head): RPNAutoDeconvHead(  
        (conv): Conv2d(1280, 1280, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (cls_logits): Conv2d(1280, 15, kernel_size=(1, 1), stride=(1, 1))  
        (bbox_pred): Conv2d(1280, 60, kernel_size=(1, 1), stride=(1, 1))  
    )  
)  
  
(roi_heads): ROIHeads(  
    (box_roi_pool): MultiScaleRoIAlign(featmap_names=['0'], output_size=(7, 7),  
    sampling_ratio=2)  
    (box_head): TwoMLPHead(  
        (fc6): Linear(in_features=62720, out_features=1024, bias=True)  
        (fc7): Linear(in_features=1024, out_features=1024, bias=True)  
    )  
    (box_predictor): FastRCNNPredictor(  
        (cls_score): Linear(in_features=1024, out_features=2, bias=True)  
        (bbox_pred): Linear(in_features=1024, out_features=8, bias=True)  
    )  
    (mask_roi_pool): MultiScaleRoIAlign(featmap_names=['0'], output_size=(14, 14),  
    sampling_ratio=2)  
    (mask_head): MaskRCNNHeads(  
        (mask_fcn1): Conv2d(1280, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (relu1): ReLU(inplace=True)  
        (mask_fcn2): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (relu2): ReLU(inplace=True)  
        (mask_fcn3): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (relu3): ReLU(inplace=True)  
        (mask_fcn4): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (relu4): ReLU(inplace=True)  
    )  
    (mask_predictor): MaskRCNNPredictor(  
        (conv5_mask): ConvTranspose2d(256, 256, kernel_size=(2, 2), stride=(2, 2))  
        (relu): ReLU(inplace=True)  
        (mask_fcn_logits): Conv2d(256, 2, kernel_size=(1, 1), stride=(1, 1))  
    )  
)
```