

000  
001       **VX-Mask R-CNN: Multimodal Semantic Object**  
002       **Segmentation for Human Upper Extremity**

004                  Anonymous ACCV 2018 submission

006                  Paper ID 1193

009       **Abstract.** In this paper we propose a modality invariant method to  
010       obtain high quality semantic object segmentation of human body parts,  
011       where the considered modalities are X-ray images and visible images. Due  
012       to the intrinsic difference between images from two modalities, state-of-  
013       the-art approaches such models as Mask R-CNN do not perform satisfac-  
014       torily. We analyse performance of intermediate layers within the Mask  
015       R-CNN on both X-ray and visible modalities. The insights from this  
016       analysis has led us to propose a new and efficient network architecture  
017       which yields highly accurate semantic segmentation results across both  
018       X-ray and visible domains. In the proposed network, dubbed as VX-  
019       Mask R-CNN, multi-task losses are combined to train the network. By  
020       conducting multiple experiments across visible and X-ray images of the  
021       human upper extremity, we validate the proposed approach, which out-  
022       performs the traditional Mask R-CNN method through better exploiting  
023       the output features of CNNs. Our method can be applied to other modal-  
024       ties and can be effectively utilized for medical image analysis tasks such  
as image registration and 3D reconstruction across modalities.

026       **1 Introduction**

028       Leveraging the performance of medical imaging systems with computer vision  
029       techniques has progressed over many years with the aim to improve both diag-  
030       nosis and therapy outcomes [1–3]. X-ray images are widely used in the medical  
031       domain to observe the structures of human anatomy including bones, tissues  
032       and muscles [4, 5]. Some computer vision tasks such as 3D reconstruction, may  
033       require more than one X-ray image of the patient. However, exposing patients  
034       to many X-ray images has adverse health affects and can result in a higher can-  
035       cer risk [6]. In such cases the visible images (RGB images) of the same body  
036       part can be used in combination with X-ray images to mitigate this problem.  
037       The advantages of using visible images in this task are that they can be cap-  
038       tured conveniently, with very low cost while incurring minimal harm to patients.  
039       However, this requires innovative cross-modal processing in low level image pro-  
040       cessing tasks such as image registration, object identification, segmentation and  
041       correspondence.

042       Compared with single-modality image processing, cross-modal image process-  
043       ing exhibits some inherent challenges, arising from potential visual dissimilarity  
044       where the regions in two images can be differently textured or one image may be

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

2 ACCV-18 submission ID 1193

045 textured while the other is homogeneous [7]. In addition, the level of information  
046 presented in each image modality is different. For instance, considering images of  
047 a hand from visible and X-ray modalities, visible images contain details of skin,  
048 nails, clothes and accessories whereas X-ray images contain details of bones and  
049 muscles (Figure 1).

050

051

052

053

054

055

056

057

058

059

060

061

062

063

064

065

066

067

068

069

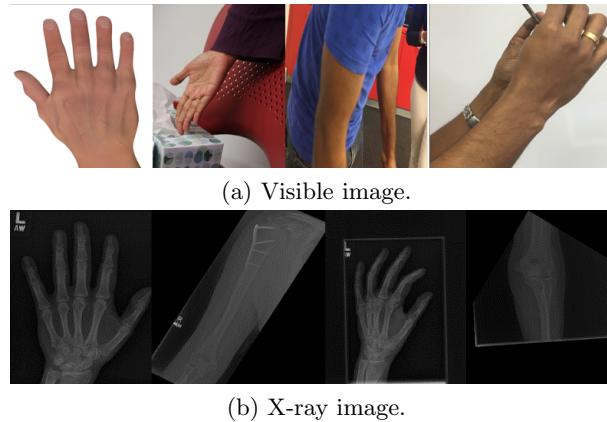


Fig. 1: Sample pairs of images from visible (first row) and X-ray (second row) images, illustrate the difference in visual texture and level of information.

In this paper we present a method to establish object boundary and pixel wise semantic object segmentation of human body parts regardless of the modality, whether it is an X-ray or a visible image. Object boundary identification and semantic segmentation underpin many applications including human pose estimation [8], semantic understanding [9], 3D object reconstruction [10], human tracking [11] and action recognition [12]. Moreover, recent advancements in neural networks have motivated their applicability for semantic segmentation tasks, resulting in robust and impressive baseline frameworks including Fast R-CNN [13], Faster R-CNN [14], and Mask R-CNN [15]. Inspired by this progression, we develop a new method coined as “VX Mask R-CNN” where we exploit the layers which generate the most informative features the capabilities of Mask-RCNN to achieve effective multi-modal and robust object segmentation. For the X-ray images, we use the **musculoskeletal radiograph** (MURA) dataset [16]. For visible images, we have collected our own dataset, which contains human body parts (i.e. elbow, forearm, hand, humerus, shoulder and wrist) images under challenging environments such as different lighting conditions, different resolutions and different backgrounds which we will make publicly available to enable researchers to replicate our results in this paper. In addition to the robust object segmentation method introduced in this paper, we provide the groundtruth region marking by manual annotation for the MURA X-ray image dataset for the upper extremity of the human body which will also be available for researchers

045

046

047

048

049

050

051

052

053

054

055

056

057

058

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

079

080

081

082

083

084

085

086

087

088

089

090 to follow-up on our work. This work has targeted on human upper extremity  
091 due to the unavailability of large X-ray image datasets for other body parts; the  
092 method we propose however can easily be applied to the other body parts as  
093 well.

094 The remainder of the paper is organized as follows. In Section 2 we analyse  
095 the recent literature on object detection and segmentation. Section 3 describes  
096 the proposed system and its subsections elaborate each component of the system.  
097 In Section 4, we present our experimental results and in Section 5, we conclude  
098 the paper with a discussion of the results on the presented method.

## 100 2 Related Work 101

102 Object detection and segmentation literature has progressed through decades.  
103 Early stages of image segmentation include approaches that uses the low level  
104 image features such as brightness and texture gradient [17]. However, when con-  
105 sidering the object segmentation in contrast to the image segmentation, it re-  
106 quires more details on the structure and the shape of the object category. Many  
107 of the recent approaches for object segmentation are preceded by object de-  
108 tection, where the object boundary identification can operate as a prior for the  
109 detection task [18–20]. Therefore efforts on improving the detection process have  
110 paramount effect on the improvement of segmentation process.

111 It is notable that neural network based approaches have granted significant  
112 enhancement for the object detection and segmentation task [21–23]. However, in  
113 the early approaches based on CNNs they dealt with specific object categories in  
114 constrained applications such as faces [21, 23] and pedestrians [22]. In the most  
115 recent literature, in rich feature hierarchies for accurate object detection and  
116 semantic segmentation [24], has introduced the concept of Regions with CNN  
117 features coined as “R-CNN” with the target of accomplishing higher accuracies  
118 for object detection. They have first extracted the region proposals for the whole  
119 image, where then the CNN based features are extracted on each of these regions  
120 which are then subjected to SVM based classification.

121 One major drawback with R-CNN is that, to extract features of each candi-  
122 date region proposal it requires a forward pass of the CNN. To overcome this ex-  
123 tensive number of forward passes “Fast R-CNN” [13] has been suggested, where  
124 a Region of Interest (RoI) pooling layer is used to extract feature vectors for  
125 object proposals. Feature vectors use the feature map generated by the froward  
126 pass of the input image over the CNN. To alleviate the bottleneck caused by  
127 having to obtain the candidate region proposals using a separate selective search  
128 process in Fast R-CNN, “Faster R-CNN” [14] has been introduced. In Faster  
129 R-CNN, a single CNN is used to estimate the candidate proposals as well as to  
130 extract features for the classification.

131 The main objectives of R-CNN, Fast R-CNN and Faster R-CNN are object  
132 detection and classification. Mask R-CNN [15] has extended the Faster R-CNN  
133 by adding another branch for object mask detection resulting in state-of-art ac-  
134 curacies for object segmentation. However, the approaches suggested in the liter-

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

4 ACCV-18 submission ID 1193

135 ature for detection and segmentation have focused on visible images. In contrast  
136 to the previous approaches we develop an architecture to detect the objects and  
137 carryout the segmentation while being invariant to the image modality, whether  
138 they are X-ray or visible. To the best of our knowledge, this is the first paper to  
139 consider model invariant semantic segmentation for visible and X-ray images.  
140

### 141 3 VX-Mask R-CNN

142 In this section we provide an overview of Mask R-CNN and then elaborate our  
143 method on deriving VX-Mask R-CNN.  
144

#### 145 3.1 Mask R-CNN

146 As described in Section 2, state-of-the-art object detection and segmentation are  
147 driven by CNNs and Mask R-CNN has depicted robust results for both these  
148 tasks. In this section we will elaborate strengths and weaknesses of Mask R-CNN  
149 and how we can exploit the Mask R-CNN’s strengths to make it invariant for  
150 object segmentation for visible and X-ray images.  
151

152 The foundation of our approach is based on the observations that we en-  
153 countered while attempting to use Mask R-CNN for object segmentation. When  
154 the original Mask R-CNN was trained on a dataset which consisted of human  
155 body parts in the visible image modality it was identified that, though it can  
156 generalize well for visible images, it does not generate masks of adequate quality  
157 for X-ray images (Figure 2). In most of the cases, it was not detecting all the  
158 components of the X-ray object. Conversely, when a model trained on X-ray  
159 images was tested on a set of visible images, one major observation that was  
160 made is that it was erroneously detecting other objects in the environment as  
161 human body parts (Figure 3).  
162

163 An intuitive solution to overcome the problems that are mentioned above is  
164 to use a dataset, which contains the images from both the modalities. Besides  
165 the increased quality of the output compared to the models that were trained  
166 on a single modality (Figure 2 and Figure 3), the newly generated masks also  
167 suffered from the problems of not recognizing the complete object or mistakenly  
168 identifying the parts in the image that do not belong to the object (Figure 4).  
169

170 With the aim of gaining an insight for Mask R-CNN, when the output of  
171 the layers were inspected, a simple yet crucial observation was made. That is,  
172 only some of the layers in Mask R-CNN generates discerning features for both  
173 the image modalities. A subset of feature maps, extracted from some random  
174 layers of the Mask R-CNN is depicted in Figure 5. The first column of Figure 5  
175 shows the input image, where the first row of it is the input visible image and  
176 the second row shows the input X-ray image. The rest of the columns depict  
177 the feature output from some random layers. It can be seen that some layers  
178 have generated features for both the input images, whereas some layers only  
179 have generated sensible features for either of input images. The layers that were  
generating sensible features for both image modalities were consistent for all the  
body parts under our consideration.  
180

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

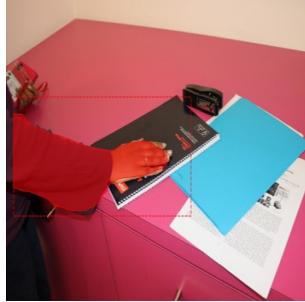
220

221

222

223

224



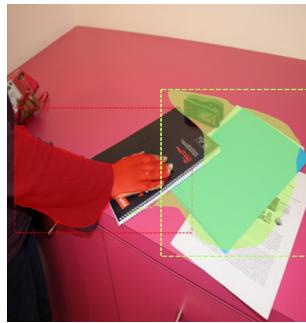
(a) Tested on a visible image.



(b) Tested on an X-ray image.



(a) Tested on an X-ray image.



(b) Tested on a visible image.

Fig. 2: Results that were obtained from the model trained on visible images illustrating its poor performance in mask identification in X-ray images.



Fig. 3: Results that were obtained from the model trained on X-ray images illustrating its poor performance in mask identification in visible images.

212

213

214

215

216

217

218

219

220

221

222

223

224

Fig. 4: Results on images which were tested on the model trained on multimodal image dataset.

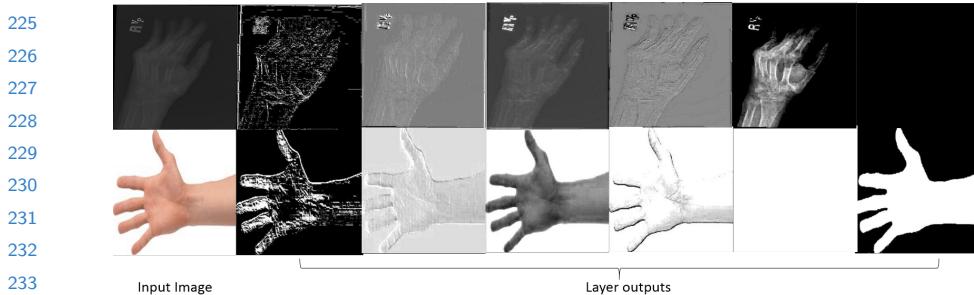


Fig. 5: Visualization of layer outputs of Mask R-CNN for a visible image and for an X-ray image, the first column depicts the input images. The 2<sup>nd</sup>, 3<sup>rd</sup>, 4<sup>th</sup> and 5<sup>th</sup> columns show the layer outputs where informative features have been generated for both X-ray and visible images, 6<sup>th</sup> column shows an example for a layer output where for the X-ray image the generated feature is informative but for the visible image it does not, and the 7<sup>th</sup> column shows an example for the opposite scenario.

### 3.2 VX-Mask R-CNN

Motivated by the observations described in Section 3.1, in this paper we introduce the architecture of VX-Mask R-CNN, which generates robust results for visible images as well as X-ray images.

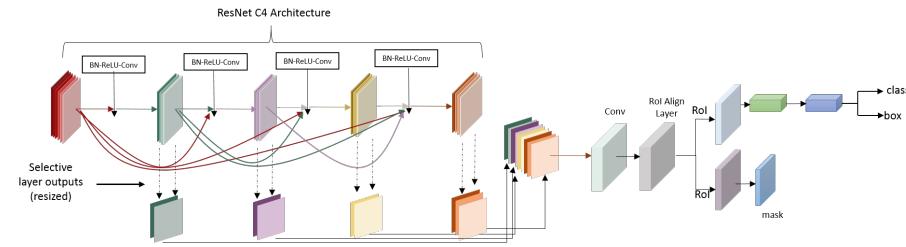


Fig. 6: VX-Mask architecture, where the selected layer outputs are resized and fused to generate the new layer, which is then subjected for ROI extraction, mask extraction and class label extraction.

VX-Mask introduces a new feature concatenation layer where the feature output of the layers that generate features for both the modalities from Mask R-CNN are concatenated. Depending on the backbone architecture used in the Mask R-CNN the dimension of feature concatenation layer vary. The architecture that is depicted in Figure 6 corresponds to the ResNet101-C4 [25] architecture. Following the previous work on Faster R-CNN we employed ResNet architecture and ResNeXt architecture [26] with depth of 50 or 101 layers. Following the

270 commonly adopted notation [27, 25, 28, 29] we denote the ResNet architecture  
271 with 50 layers where the features are extracted from 4th stage convolution layer  
272 as “ResNet-50-C4”. In addition, we analyse Feature Pyramid Network (FPN)  
273 [30] in combination with ResNet architecture.

274 When the features generated by Resnet-50-C4, Resnet-101-C4, ResNet-50-  
275 FPN, ResNet-101-FPN, ResNeXt-50-C4, ResNeXt-101-C4, ResNeXt-50-FPN and  
276 ResNeXt-101-FPN were investigated it was identified that there are 2100, 3435,  
277 2421, 3606, 2241, 3527, 2627, and 3677 feature outputs that create features for  
278 both visible and X-ray images. The concatenation layer is followed by a convolutional  
279 layer, which is then subjected to RoI align layer, which then get segregated  
280 into mask prediction branch and to class and bounding box prediction branch.

281 We use a multi-task loss function to train the network,

$$282 \quad L = L_{class} + \lambda L_{box} + \mu L_{mask}, \quad (1)$$

283

284 where  $L_{class}$  denotes the classification loss,  $L_{box}$  denotes the loss defined  
285 on bounding box identification and  $L_{mask}$  denotes the loss defined on mask  
286 identification. The combined loss is calculated as in Equation 1, where  $\lambda$  and  $\mu$   
287 are the hyper-parameters that are used for maintaining the balance between the  
288 three losses that are aggregating for the final loss value.

289 The  $L_{class}$  is computed based on the softmax over 7 outputs of the fully  
290 connected layer in the classification branch. The last layer of the classification  
291 branch contains 7 outputs as the dataset we have used 6 classes and the default  
292 background layer is also taken into the consideration. For a given image, for a  
293 given Region of Interest (RoI) with groundtruth class label  $u$ , the  $L_{class}$  is the  
294 negative log softmax value for the groundtruth class label (Equation 2),  
295

$$296 \quad L_{class} = -\log(p_u). \quad (2)$$

297

298 A bounding box is defined using four parameters  $t_x$ ,  $t_y$ ,  $t_w$  and  $t_h$  which  
299 denote the  $x$  coordinate of the right corner of the bounding box,  $y$  coordinate  
300 of the right corner of the bounding box, width of the bounding box and height  
301 of the bounding box respectively. Let  $t^u = (t_x^u, t_y^u, t_w^u, t_h^u)$  be the bounding box  
302 parameter tuple of instance  $u$  of groundtruth annotation, and  $t^v$  be the predicted  
303 tuple the  $L_{box}$  is defined as,

$$304 \quad L_{box} = \sum_{i \in x, y, w, h} smooth_{L1}(t_i^v - t_i^u). \quad (3)$$

305

306 In the equation 3, the term  $smooth_{L1}$  is defined as,

$$309 \quad smooth_{L1}(r) = \begin{cases} 0.5q^2 & if |q| < 1 \\ |q| - 0.5 & otherwise. \end{cases} \quad (4)$$

310

311 It should be noted that the  $L_{box}$  is not calculated for the *background* class.

312 In this work, the mask is defined as a binary array where the foreground  
313 pixels of the objects are marked. The mask branch output consists of  $K$  masks  
314 in  $m \times m$  resolutions, where for each of the entries in  $m \times m$  array a per-pixel

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

8 ACCV-18 submission ID 1193

315 sigmoid is applied and when an object belongs to the class  $k$ , the  $L_{mask}$  is defined  
316 by the average binary cross entropy,

317  
318 
$$\text{Binary Cross Entropy} = -(y \log(p)) + (1 - y) \log(1 - p), \quad (5)$$

319 of the  $k^{\text{th}}$  mask. In this process K is the number of classes under consideration.  
320 In Equation 5,  $y$  is the label for the mask that can either be 1 or 0, depending on  
321 whether that point correspondence to a foreground pixel or a background pixel  
322 and  $p$  stands for the output of the neuron.

## 324 4 Experiments

### 325 4.1 Dataset

326 For the X-ray image dataset we used the MURA dataset[16] which contains X-ray  
327 images of the human body parts. The MURA dataset contains X-ray images of  
328 elbow, forearm, hand, humerus, shoulder and wrist. For this dataset we generated  
329 manually annotated mask groundtruth. For the visible images, we generated a  
330 dataset which contains visible images of the same body parts that are included  
331 in MURA dataset and the corresponding groundtruth masks. The visible image  
332 dataset contains the images captured under challenging conditions including  
333 different backgrounds, resolutions, occlusions and illumination conditions.

### 338 4.2 Experiment Setting

339 For the VX-Mask training we use the same learning rates and the weight decay  
340 values that have been used in the preceding approaches including Fast R-CNN  
341 [13, 14]. We use the pretrained backbone architectures where the training has  
342 been performed on COCO [31] dataset and the convolution layer was initiated  
343 with random values. The  $\lambda$  and  $\mu$  were set to 1, and the learning rate, weight  
344 decay and the momentum was set to 0.02, 0.00001 and 0.9 accordingly.

345 The training accuracies and the validation for the combined dataset are de-  
346 picted in Figure 7.

### 349 4.3 Evaluation Matrices

350 The evaluation matrix we used the average precision (AP), which is been calcu-  
351 lated using Intersection over Union (IoU),

352  
353  
354 
$$\text{IoU} = \frac{\text{Area}(b_p \cap b_t)}{\text{Area}(b_p \cup b_t)}, \quad (6)$$

355 based on a defined threshold. In Equation 6,  $b_p$  is the prediction, which is the  
356 bounding box prediction or the mask prediction, and  $b_t$  is the corresponding  
357 groundtruth value. In this evaluation we used the thresholds of 0.5 and 0.75.

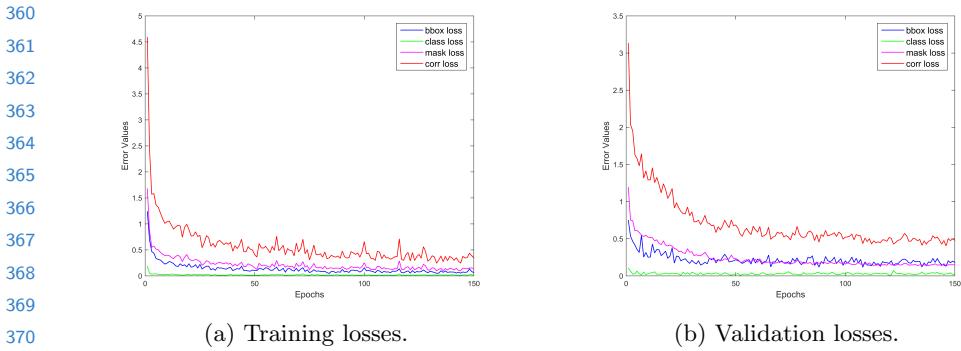


Fig. 7: Training and validation losses.

#### 4.4 Results

To compare our method with original Mask R-CNN we conducted a set of experiments where previously unseen images were fed into our architecture as well as to the Mask R-CNN. When the original Mask R-CNN and the VX-Mask were tested on the same datasets for mask identification, the obtained average precision values are indicated in Table 1. By analysing the table it can be understood that the VX-Mask R-CNN outperforms the Mask R-CNN for both the X-ray and visible images. In addition, it can be identified that the features extracted from deeper and advanced networks can perform better comparatively. From Table 1 it can be understood that ResNeXt-101-FPN has yielded the highest average precision for both the modalities and for both Mask R-CNN and VX-Mask R-CNN, except for the  $AP_{50}$  in visible modality on which the Mask R-CNN with backbone ResNeXt-101-C4 has performed at the best. However, when comparing the results of Mask R-CNN and VX-Mask R-CNN it was identified that the later with all the backbone architectures have performed better than the best performing setting of Mask R-CNN.

The average precision for bounding box identification using the original Mask R-CNN and VX-Mask R-CNN, which have been trained on the same datasets is recorded in Table 2. It can be identified that the same observations that were made for the mask prediction are presented in bounding box identification as well. Some qualitative results that we obtained for both visible and X-ray images are depicted in Figure 8. From the qualitative results it can be identified that the when the results gained by VX-Mask R-CNN are compared with the original Mask R-CNN results, that the VX-Mask R-CNN is capable of generating more precise results, while identifying more components and regions of upper extremity of human body. However, it should be noted that even when using VX-Mask R-CNN, there are circumstances that where all the components that are presented in the image are not properly identified. For an instance, in the 3<sup>rd</sup> row of Figure 8 the elbow of the person has not been identified by both the

405 approaches, instead the combination of forearm and the elbow has been detected  
406 as the forearm.  
407  
408  
409

410 Table 1: Average precision on visible images and X-ray images for image seg-  
411 mentation

	Backbone	On Visible Images		On X-ray Images		Average		
		$AP_{50}$	$AP_{75}$	$AP_{50}$	$AP_{75}$	$AP_{50}$	$AP_{75}$	
412	Mask R-CNN	Resnet-50-C4	39.1	27.4	37.7	21.8	38.4	24.6
413	Mask R-CNN	Resnet-101-C4	43	29.1	44.3	23.3	43.7	26.2
414	Mask R-CNN	ResNet-50-FPN	41.6	28.3	40.3	22.4	41.0	25.4
415	Mask R-CNN	ResNet-101-FPN	43.8	29.5	46.1	25	45.0	27.3
416	Mask R-CNN	ResNeXt-50-C4	41.9	28.7	42.1	22.8	42.0	25.8
417	Mask R-CNN	ResNeXt-101-C4	<b>45.6</b>	29.3	44.8	23.9	45.2	26.6
418	Mask R-CNN	ResNeXt-50-FPN	44.1	29.4	45.7	24.3	44.9	26.9
419	Mask R-CNN	ResNeXt-101-FPN	44.7	<b>29.6</b>	<b>46.2</b>	<b>25.3</b>	<b>45.5</b>	<b>27.5</b>
420	VX-Mask R-CNN	Resnet-50-C4	57.3	31.3	53.4	27.7	55.4	29.5
421	VX-Mask R-CNN	Resnet-101-C4	59.7	33.1	54.4	28.8	57.1	31.0
422	VX-Mask R-CNN	ResNet-50-FPN	58.2	32.4	53.1	28.2	55.7	30.3
423	VX-Mask R-CNN	ResNet-101-FPN	62.6	34.6	54.9	29.3	58.8	32.0
424	VX-Mask R-CNN	ResNeXt-50-C4	62.9	34.9	55.8	30	59.4	32.5
425	VX-Mask R-CNN	ResNeXt-101-C4	64.1	35.3	56.9	31.7	60.5	33.5
426	VX-Mask R-CNN	ResNeXt-50-FPN	65.1	35.7	58.2	32.4	61.7	34.1
427	VX-Mask R-CNN	ResNeXt-101-FPN	<b>66.8</b>	<b>37.1</b>	<b>58.9</b>	<b>33.5</b>	<b>62.9</b>	<b>35.3</b>

## 433 5 Conclusion

434  
435  
436 In this paper we have presented modality invariant Mask R-CNN called "VX  
437 Mask R-CNN", where the most informative layers of CNNs are effectively utilized  
438 as the features for class prediction, region identification and mask identification  
439 across multimodal images. The different textures of the two image modalities  
440 (X-ray and visible) and the features that are generated by the convolutional  
441 neural networks for these different textures are taken into consideration  
442 when devising our architecture. In contrast to the previous approaches where the  
443 successive layers of the network use all the outputs from the previous layers, and  
444 the features of the last convolutional layer is subjected for the mask extraction,  
445 we filter the most informative outputs of all the layers in the network. The current  
446 implementation has outperformed the original Mask R-CNN for the visible and  
447 X-ray images, and this method can be extended to other image modalities by  
448 observing the convolutional neural network's behaviour on the modalities under  
449 consideration.

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

ACCV-18 submission ID 1193

11

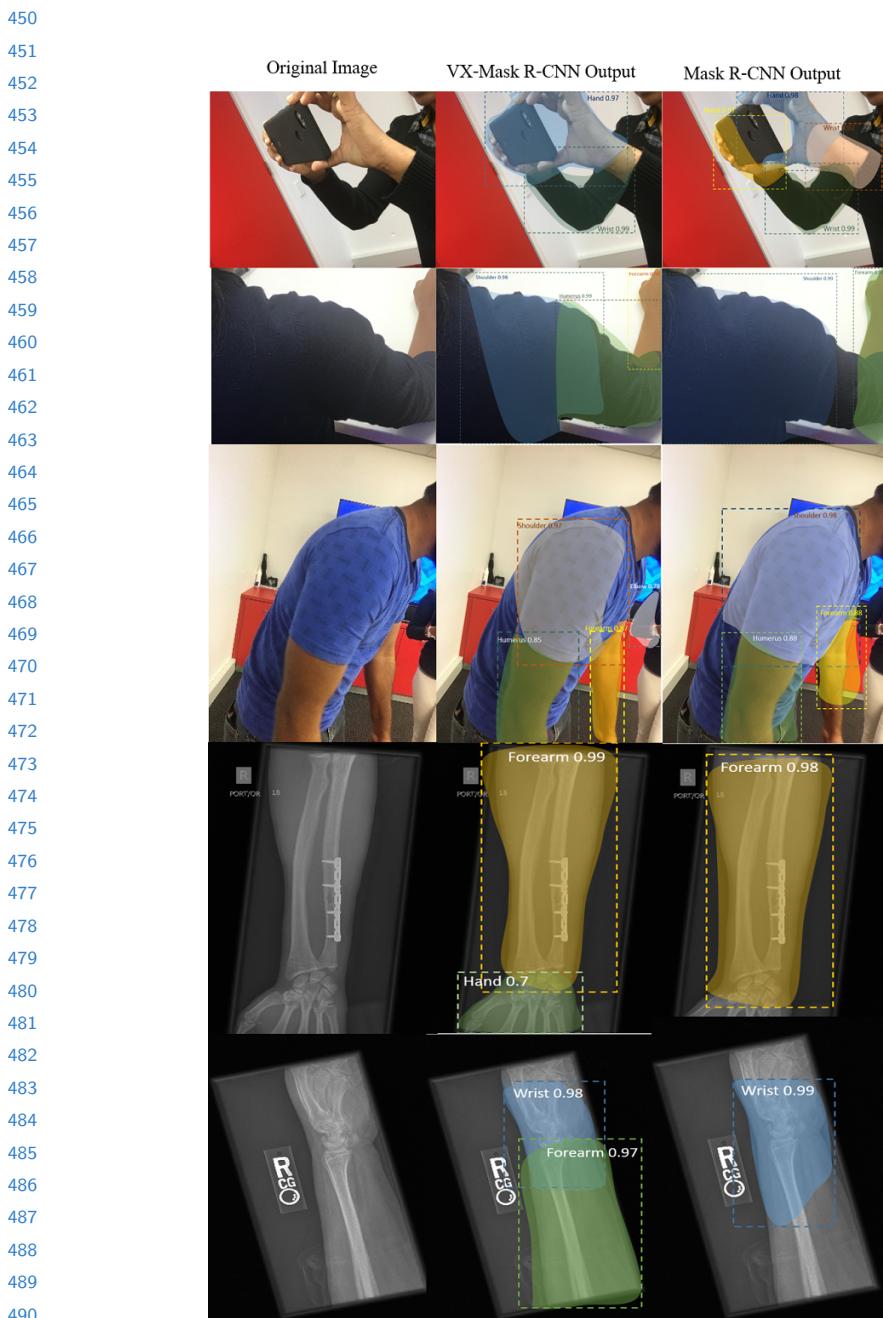


Fig. 8: Qualitative results on visible images and X-ray images on object identification, mask identification and bounding box identification.

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

12 ACCV-18 submission ID 1193

495  
 496 Table 2: Average precision on visible images and X-ray images for bounding box  
 497 identification

	Backbone	On Visible Images		On X-ray Images		Average	
		$AP_{50}$	$AP_{75}$	$AP_{50}$	$AP_{75}$	$AP_{50}$	$AP_{75}$
Mask R-CNN	Resnet-50-C4	42.10	29.05	42.05	25.46	42.07	27.26
Mask R-CNN	Resnet-101-C4	46.92	34.98	49.34	28.36	48.13	31.67
Mask R-CNN	ResNet-50-FPN	42.61	29.94	44.06	27.54	43.34	28.74
Mask R-CNN	ResNet-101-FPN	48.90	36.62	51.80	33.03	50.35	34.83
Mask R-CNN	ResNeXt-50-C4	44.62	34.50	48.63	27.83	46.63	31.17
Mask R-CNN	ResNeXt-101-C4	<b>50.48</b>	35.48	50.34	31.13	50.41	33.30
Mask R-CNN	ResNeXt-50-FPN	49.25	36.07	50.40	33.02	49.83	34.54
Mask R-CNN	ResNeXt-101-FPN	49.54	<b>39.33</b>	<b>52.95</b>	<b>33.62</b>	<b>51.24</b>	<b>36.48</b>
VX-Mask R-CNN	Resnet-50-C4	63.29	39.75	57.28	30.09	60.28	34.92
VX-Mask R-CNN	Resnet-101-C4	65.56	38.99	61.70	33.01	63.63	36.00
VX-Mask R-CNN	ResNet-50-FPN	63.34	33.13	56.44	30.81	59.89	31.97
VX-Mask R-CNN	ResNet-101-FPN	65.81	38.70	61.75	33.86	63.78	36.28
VX-Mask R-CNN	ResNeXt-50-C4	68.40	44.18	62.49	34.26	65.44	39.22
VX-Mask R-CNN	ResNeXt-101-C4	69.13	42.84	63.58	36.39	66.35	39.62
VX-Mask R-CNN	ResNeXt-50-FPN	69.27	41.36	63.95	36.91	66.61	39.13
VX-Mask R-CNN	ResNeXt-101-FPN	<b>69.80</b>	<b>46.31</b>	<b>64.16</b>	<b>36.94</b>	<b>66.98</b>	<b>41.63</b>

## 514 515 References 516 517

- Chen, C.h.: Computer vision in medical imaging. Volume 2. World scientific (2014)
- Cootes, T.F., Taylor, C.J.: Statistical models of appearance for medical image analysis and computer vision. In: Medical Imaging 2001: Image Processing. Volume 4322., International Society for Optics and Photonics (2001) 236–249
- Florack, L.: Image structure. computational imaging and vision (1997)
- Farber, L., Tardos, G., Michaels, J.N.: Use of x-ray tomography to study the porosity and morphology of granules. Powder Technology **132** (2003) 57–63
- Hu, S., Hoffman, E.A., Reinhardt, J.M.: Automatic lung segmentation for accurate quantitation of volumetric x-ray ct images. IEEE transactions on medical imaging **20** (2001) 490–498
- Rao, V.M., Levin, D.C.: The overuse of diagnostic imaging and the choosing wisely initiative. Annals of internal medicine **157** (2012) 574–576
- Torabi, A., Bilodeau, G.A.: A lss-based registration of stereo thermal-visible videos of multiple people using belief propagation. Computer Vision and Image Understanding **117** (2013) 1736–1747
- Xiao, B., Wu, H., Wei, Y.: Simple baselines for human pose estimation and tracking. arXiv preprint arXiv:1804.06208 (2018)
- Neuhold, G., Ollmann, T., Bulo, S.R., Kontschieder, P.: The mapillary vistas dataset for semantic understanding of street scenes. In: Proceedings of the International Conference on Computer Vision (ICCV), Venice, Italy. (2017) 22–29
- Yang, B., Rosa, S., Markham, A., Trigoni, N., Wen, H.: 3d object dense reconstruction from a single depth view. arXiv preprint arXiv:1802.00411 (2018)
- Li, J., Zhao, J., Wei, Y., Lang, C., Li, Y., Sim, T., Yan, S., Feng, J.: Multi-human parsing in the wild. arXiv preprint arXiv:1705.07206 (2017)

- 540 12. Liu, Y., Wang, Q.: A new simple human abnormal action detection based on static  
541 images. In: Computer Science and Automation Engineering (CSAE), 2011 IEEE  
542 International Conference on. Volume 1., IEEE (2011) 578–581  
543 13. Girshick, R.: Fast r-cnn. arXiv preprint arXiv:1504.08083 (2015)  
544 14. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detec-  
545 tion with region proposal networks. In: Advances in neural information processing  
systems. (2015) 91–99  
546 15. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: Computer Vision  
547 (ICCV), 2017 IEEE International Conference on, IEEE (2017) 2980–2988  
548 16. Rajpurkar, P., Irvin, J., Bagul, A., Ding, D., Duan, T., Mehta, H., Yang, B., Zhu,  
549 K., Laird, D., Ball, R.L., et al.: Mura dataset: Towards radiologist-level abnor-  
550 mality detection in musculoskeletal radiographs. arXiv preprint arXiv:1712.06957  
551 (2017)  
552 17. Martin, D.R., Fowlkes, C.C., Malik, J.: Learning to detect natural image bound-  
553 aries using brightness and texture. In: Advances in Neural Information Processing  
Systems. (2003) 1279–1286  
554 18. Elgammal, A., Harwood, D., Davis, L.: Non-parametric model for background  
555 subtraction. In: European conference on computer vision, Springer (2000) 751–767  
556 19. Mittal, A., Paragios, N.: Motion-based background subtraction using adaptive ker-  
557 nel density estimation. In: Computer Vision and Pattern Recognition, 2004. CVPR  
558 2004. Proceedings of the 2004 IEEE Computer Society Conference on. Volume 2.,  
Ieee (2004) II–II  
559 20. Gonfaus, J.M., Boix, X., Van de Weijer, J., Bagdanov, A.D., Serrat, J., Gonzalez,  
560 J.: Harmony potentials for joint classification and segmentation. In: Computer  
561 Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on, IEEE (2010)  
562 3280–3287  
563 21. Rowley, H.A., Baluja, S., Kanade, T.: Neural network-based face detection. IEEE  
564 Transactions on pattern analysis and machine intelligence **20** (1998) 23–38  
565 22. Sermanet, P., Kavukcuoglu, K., Chintala, S., LeCun, Y.: Pedestrian detection  
566 with unsupervised multi-stage feature learning. In: Computer Vision and Pattern  
Recognition (CVPR), 2013 IEEE Conference on, IEEE (2013) 3626–3633  
567 23. Vaillant, R., Monrocq, C., Le Cun, Y.: Original approach for the localisation  
568 of objects in images. IEE Proceedings-Vision, Image and Signal Processing **141**  
569 (1994) 245–250  
570 24. Girshick, R., Donahue, J., Darrell, T., Malik, J.: Rich feature hierarchies for ac-  
571 curate object detection and semantic segmentation. In: Proceedings of the IEEE  
572 conference on computer vision and pattern recognition. (2014) 580–587  
573 25. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recogni-  
574 tion. In: Proceedings of the IEEE conference on computer vision and pattern recogni-  
575 tion. (2016) 770–778  
576 26. Xie, S., Girshick, R., Dollár, P., Tu, Z., He, K.: Aggregated residual transfor-  
577 mations for deep neural networks. In: Computer Vision and Pattern Recognition (CVPR),  
578 2017 IEEE Conference on, IEEE (2017) 5987–5995  
579 27. Dai, J., He, K., Sun, J.: Instance-aware semantic segmentation via multi-task  
580 network cascades. In: Proceedings of the IEEE Conference on Computer Vision  
and Pattern Recognition. (2016) 3150–3158  
581 28. Huang, J., Rathod, V., Sun, C., Zhu, M., Korattikara, A., Fathi, A., Fischer, I.,  
582 Wojna, Z., Song, Y., Guadarrama, S., et al.: Speed/accuracy trade-offs for modern  
583 convolutional object detectors. In: IEEE CVPR. (2017)  
584 29. Shrivastava, A., Sukthankar, R., Malik, J., Gupta, A.: Beyond skip connections:  
Top-down modulation for object detection. arXiv preprint arXiv:1612.06851 (2016)

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

14 ACCV-18 submission ID 1193

- 585 30. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature 585  
586 pyramid networks for object detection. In: CVPR. Volume 1. (2017) 4 586  
587 31. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., 587  
588 Zitnick, C.L.: Microsoft coco: Common objects in context. In: European conference 588  
589 on computer vision, Springer (2014) 740–755 589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629