

YOLACT

Real-time Instance Segmentation

Daniel Bolya

Chong Zhou

Fanyi Xiao

Yong Jae Lee

University of California, Davis

{dbolya, cczhou, fyxiao, yongjaelee}@ucdavis.edu

Abstract

We present a simple, fully-convolutional model for real-time instance segmentation that achieves 29.8 mAP on MS COCO at 33 fps evaluated on a single Titan Xp, which is significantly faster than any previous competitive approach. Moreover, we obtain this result after training on only one GPU. We accomplish this by breaking instance segmentation into two parallel subtasks: (1) generating a set of prototype masks and (2) predicting per-instance mask coefficients. Then we produce instance masks by linearly combining the prototypes with the mask coefficients. We find that because this process doesn't depend on repooling, this approach produces very high-quality masks and exhibits temporal stability for free. Furthermore, we analyze the emergent behavior of our prototypes and show they learn to localize instances on their own in a translation variant manner, despite being fully-convolutional. Finally, we also propose Fast NMS, a drop-in 12 ms faster replacement for standard NMS that only has a marginal performance penalty.

1. Introduction

“Boxes are stupid anyway though, I’m probably a true believer in masks except I can’t get YOLO to learn them.”

– Joseph Redmon, YOLOv3 [34]

What would it take to create a real-time instance segmentation algorithm? Over the past few years, the vision community has made great strides in instance segmentation, in part by drawing on powerful parallels from the well-established domain of object detection. State-of-the-art approaches to instance segmentation like Mask R-CNN [16] and FCIS [22] directly build off of advances in object detection like Faster R-CNN [35] and R-FCN [7]. Yet, these methods focus primarily on performance over speed, leaving the scene devoid of instance segmentation parallels to real-time object detectors like SSD [28] and YOLO [33, 34]. In this work, our goal is to fill that gap with

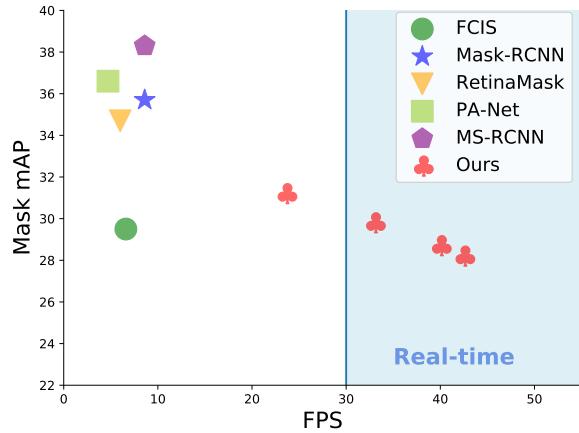


Figure 1: Speed-performance trade-off for various instance segmentation methods on COCO. To our knowledge, ours is the first *real-time* (above 30 FPS) approach with around 30 mask mAP on COCO test-dev.

a fast, one-stage instance segmentation model in the same way that SSD and YOLO fill that gap for object detection.

However, instance segmentation is hard—much harder than object detection. One-stage object detectors like SSD and YOLO are able to speed up existing two-stage detectors like Faster R-CNN by simply removing the second stage and making up for the lost performance in other ways. The same approach is not easily extendable, however, to instance segmentation. State-of-the-art two-stage instance segmentation methods depend heavily on *feature localization* to produce masks. That is, these methods “repool” features in some bounding box region (e.g., via ROI-pool/align), and then feed these now localized features to their mask predictor. This approach is inherently sequential and is therefore difficult to accelerate. One-stage methods that perform these steps in parallel like FCIS do exist, but they require significant amounts of post-processing after localization, and thus are still far from real-time.

To address these issues, we propose YOLACT¹, a real-

¹You Only Look At CoefficienTs

time instance segmentation framework that **forges an explicit localization step**. Instead, YOLACT break ups instance segmentation into two parallel tasks: (1) generating a dictionary of non-local *prototype masks over the entire image*, and (2) predicting a set of *linear combination coefficients per instance*. Then producing a full-image instance segmentation from these two components is simple: for each instance, linearly combine the prototypes using the corresponding predicted coefficients and then crop with a predicted bounding box. We show that by segmenting in this manner, *the network learns how to localize instance masks on its own*, where visually, spatially, and semantically similar instances appear different in the prototypes.

Moreover, since the number of prototype masks is independent of the number of categories (e.g., there can be more categories than prototypes), YOLACT learns a distributed representation in which each instance is segmented with a combination of prototypes that are shared across categories. This distributed representation leads to interesting emergent behavior in the prototype space: some prototypes spatially partition the image, some localize instances, some detect instance contours, some encode position-sensitive directional maps (similar to those obtained by hard-coding a position-sensitive module in FCIS [22]), and most do a combination of these tasks (see Figure 5).

This approach also has several practical advantages. First and foremost, it’s fast: because of its parallel structure and extremely lightweight assembly process, YOLACT adds only a marginal amount of computational overhead to a one-stage backbone detector, making it easy to reach 30 fps even when using ResNet-101 [17]. Second, masks are high-quality: since the masks use the full extent of the image space without any loss of quality from repooling, our masks for large objects are significantly higher quality than those of other methods (see Figure 7). Finally, it’s general: the idea of generating prototypes and mask coefficients could be added to almost any modern object detector.

As a bonus, breaking up instance segmentation in this way is also loosely related to the ventral (“what”) and dorsal (“where”) streams hypothesized to play a prominent role in human vision [14]. The linear coefficients and corresponding detection branch can be thought of as recognizing individual instances (“what”), while the prototype masks can be seen as localizing instances in space (“where”). This is closer to, albeit still far away from, human vision than the two-stage “localize-then-segment” type approaches.

Our main contribution is the first real-time (> 30 fps) instance segmentation algorithm with competitive results on the challenging MS COCO dataset [26] (see Figure 1). In addition, we analyze the emergent behavior of YOLACT’s prototypes and provide experiments to study the speed vs. performance trade-offs obtained with different backbone architectures, numbers of prototypes, and image reso-

lutions. We also provide a novel Fast NMS approach that is 12ms faster than traditional NMS with a negligible performance penalty. *We will make our code publicly available.*

2. Related Work

Instance Segmentation Given its importance, a lot of research effort has been made to push instance segmentation *accuracy*. Mask-RCNN [16] is a representative two-stage instance segmentation approach that first generates candidate region-of-interests (ROIs) and then classifies and segments those ROIs in the second stage. Follow-up works try to **improve its accuracy** by e.g., enriching the FPN features [27] or addressing the incompatibility between a mask’s **confidence score** and its **localization accuracy** [18]. These two-stage methods require re-pooling features for each ROI and processing them with subsequent computations, which make them unable to **obtain real-time speeds** (30 fps) even when decreasing image size (see Table 2c).

One-stage instance segmentation methods generate position sensitive maps that are assembled into final masks with position-sensitive pooling [6, 22] or combine semantic segmentation logits and direction prediction logits [4]. Though conceptually faster than two-stage methods, they still require repooling or other non-trivial computations (e.g., mask voting). This severely limits their speed, placing them far from real-time. In contrast, our assembly step is much more lightweight (only a linear combination) and can be implemented as one GPU-accelerated matrix-matrix multiplication, making our approach very fast.

Finally, some methods first perform semantic segmentation followed by boundary detection [20], pixel clustering [3, 23], or learn an embedding to form instance masks [30, 15, 8, 11]. Again, these methods have multiple stages and/or involve expensive clustering procedures, which limits their viability for real-time applications.

Real-time Instance Segmentation While real-time object detection [28, 32, 33, 34], and semantic segmentation [2, 39, 31, 10, 45] methods exist, few works have focused on real-time instance segmentation. Straight to Shapes [19] can perform instance segmentation with learned encodings of shapes at 30 fps, but its accuracy is far from that of modern baselines. Box2Pix [40] relies on an extremely light-weight backbone detector (GoogLeNet v1 [38] and SSD [28]) combined with a hand-engineered algorithm to obtain 10.9 fps on Cityscapes [5] and 35 fps on KITTI [13]. However, they report no results on the more challenging and semantically-rich COCO dataset, which has 80 classes compared to the 8 of KITTI and Cityscapes. Furthermore, they observe a large drop in relative performance going from a semantically simple dataset (KITTI) to a more complex one (Cityscapes), so an even more difficult dataset (COCO) would pose a challenge. In fact, Mask R-CNN [16] remains one of the fastest

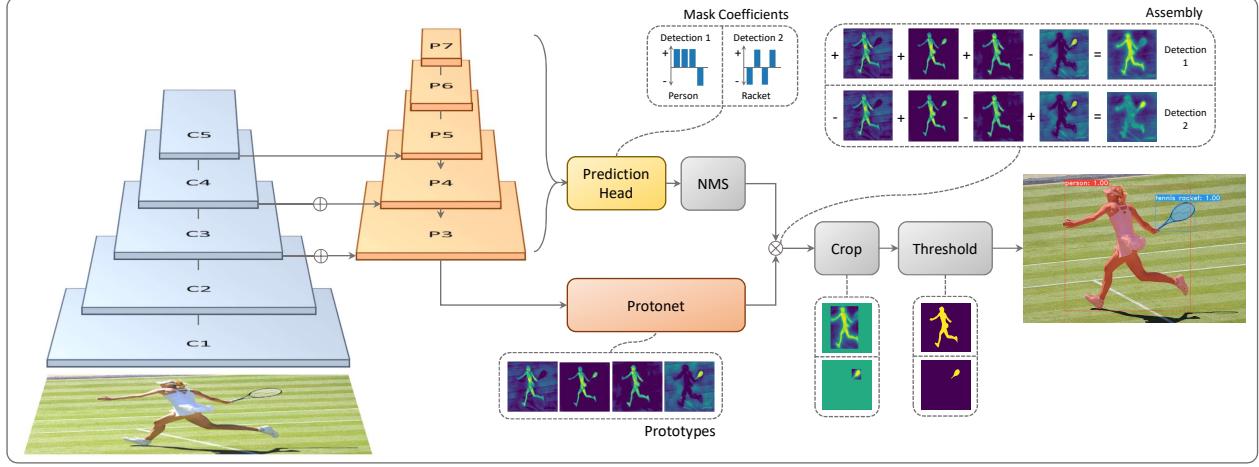


Figure 2: **YOLACT Architecture** Blue/yellow indicates low/high values in the prototypes, gray nodes indicate functions that are not trained, and $k = 4$ in this example. We base this architecture off of RetinaNet [25] using ResNet-101 + FPN.

instance segmentation methods on semantically challenging datasets (13.5 fps on 550² px images; see Table 2c).

Prototypes Learning prototypes (aka vocabulary or code-book) has been extensively explored in computer vision. Classical representations include textons [21] and visual words [37], with advances made via sparsity and locality priors [42, 41, 44]. Others have designed prototypes for object detection [1, 43, 36]. Though related, these works use prototypes to represent features, whereas we use them to assemble masks for instance segmentation. Moreover, we learn prototypes that are specific to each image, rather than global prototypes shared across the entire dataset.

3. YOLACT

Our goal is to add a mask branch to an existing one-stage object detection model in the same vein as Mask R-CNN [16] does to Faster R-CNN [35], but without an explicit localization step (e.g., feature repooling). To do this, we break up the complex task of instance segmentation into two simpler, parallel tasks that can be assembled to form the final masks. The first branch uses an FCN [29] to produce a set of image-sized “prototype masks” that do not depend on any one instance. The second adds an extra head to the object detection branch to predict a vector of “mask coefficients” for each anchor that encode an instance’s representation in the prototype space. Finally, for each instance that survives NMS we construct a mask for that instance by linearly combining the work of these two branches.

Rationale We perform instance segmentation in this way primarily because masks are spatially coherent; i.e., pixels close to each other are likely to be part of the same instance. While a convolutional (*conv*) layer naturally takes advantage of this coherence, a fully-connected (*fc*) layer does not.

That poses a problem, since one-stage object detectors produce class and box coefficients for each anchor as an output of an *fc* layer.² Two stage approaches like Mask R-CNN get around this problem by using a localization step (e.g. ROI-Align), which preserves the spatial coherence of the features while also allowing the mask to be a *conv* layer output. However, doing so requires a significant portion of the model to wait for a first-stage RPN to propose localization candidates, inducing a significant speed penalty.

Thus, we break the problem into two parallel parts, making use of *fc* layers, which are good at producing semantic vectors, and *conv* layers, which are good at producing spatially coherent masks, to produce the “mask coefficients” and “prototype masks”, respectively. Then, because prototypes and mask coefficients can be computed independently, the computational overhead over that of the backbone detector comes mostly from the assembly step, which can be implemented as a single matrix multiplication. In this way, we can maintain spatial coherence in the feature space while still being one-stage and *fast*.

3.1. Prototype Generation

The prototype generation branch (protonet) predicts a set of k prototype masks for the entire image. We implement protonet as an FCN whose last layer has k channels (one for each prototype) and attach it to a backbone feature layer (see Figure 3 for an illustration). While this formulation is similar to standard semantic segmentation, it differs in that we exhibit no explicit loss on the prototypes. Instead, all supervision comes from the final mask loss after assembly.

²To show that this is an issue, we develop an “*fc*-mask” model that produces masks for each anchor as the reshaped output of an *fc* layer. As our experiments in Table 2c show, simply adding masks to a one-stage model as *fc* outputs only obtains 20.7 mAP and is thus very much insufficient.

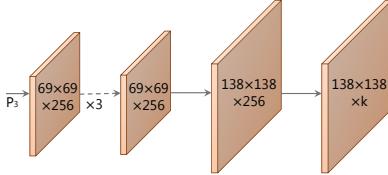


Figure 3: Protonet Architecture The labels denote feature size and channels for an image size of 550×550 . Arrows indicate 3×3 *conv* layers, except for the final *conv* which is 1×1 . The increase in size is an upsample followed by a *conv*. Inspired by the mask branch in [16].

We note two important design choices: taking protonet from deeper backbone features produces more robust masks, and higher resolution prototypes result in both higher quality masks and better performance on smaller objects. Thus, we use FPN [24] because its largest feature layers (P_3 in our case; see Figure 2) are the deepest. Then, we upsample it to one fourth the dimensions of the input image to increase performance on small objects.

Finally, we find it important that the protonet’s output being unbounded, as this allows the network to produce large, overpowering activations on prototypes it is very confident about (e.g., obvious background). Thus, we have the option of following protonet with either a ReLU or no nonlinearity. We choose ReLU for more interpretable prototypes.

3.2. Mask Coefficients

Typical anchor-based object detectors have two branches in their prediction heads: one branch to predict c class confidences, and the other to predict 4 bounding box regressors. For mask coefficient prediction, we simply add a third branch in parallel that predicts k mask coefficients, one corresponding to each prototype. Thus, instead of producing $4 + c$ coefficients per anchor, we produce $4 + c + k$.

Then for nonlinearity, we find it important to be able to subtract out prototypes from the final mask. Thus, we apply \tanh to the k mask coefficients, which produces more stable outputs over no nonlinearity. The relevance of this design choice is apparent in Figure 2, as neither mask would be constructable without allowing for subtraction.

3.3. Mask Assembly

To produce instance masks, we combine the work of the prototype branch and mask coefficient branch, using a linear combination of the former with the latter as coefficients. We then follow this by a sigmoid nonlinearity to produce the final masks. These operations can be implemented efficiently using a single matrix multiplication and sigmoid:

$$M = \sigma(PC^T) \quad (1)$$

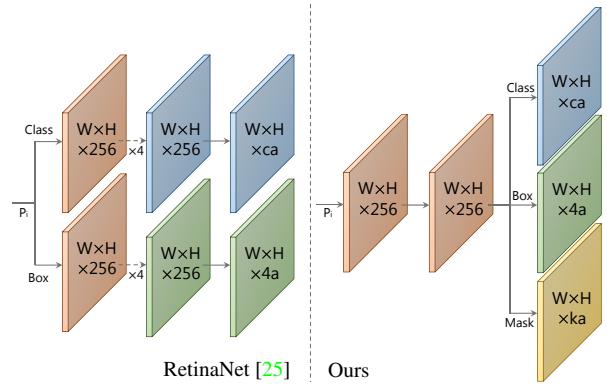


Figure 4: Head Architecture We use a shallower prediction head than RetinaNet [25] and add a mask coefficient branch. This is for c classes, a anchors for feature layer P_i , and k prototypes. See Figure 3 for a key.

where P is an $h \times w \times k$ matrix of prototype masks and C is a $n \times k$ matrix of mask coefficients for n instances surviving NMS and score thresholding. Other, more complicated combination steps are possible; however, we keep it simple (and fast) with a basic linear combination.

Losses We use three losses to train our model: classification loss L_{cls} , box regression loss L_{box} and mask loss L_{mask} . Both L_{cls} and L_{box} are defined in the same way as in [28]. Then to compute mask loss, we simply take the pixel-wise binary cross entropy between assembled masks M and the ground truth masks M_{gt} : $L_{mask} = \text{BCE}(M, M_{gt})$.

Cropping Masks To preserve small objects in the prototypes, we crop the final masks with the predicted bounding box during evaluation. During training, we instead crop with the ground truth bounding box and divide L_{mask} by the ground truth bounding box area.

3.4. Emergent Behavior

Our approach might seem surprising, as the general consensus around instance segmentation is that because FCNs are translation invariant, the task needs translation variance added back in [22]. Thus methods like FCIS [22] and Mask R-CNN [16] try to explicitly add translation variance, whether it be by directional maps and position-sensitive re-pooling, or by putting the mask branch in the second stage so it does not have to deal with localizing instances. In our method, the only translation variance we add is to crop the final mask with the predicted bounding box. However, we find that our method also works without cropping for medium and large objects, so this is not a result of cropping. Instead, YOLACT *learns how to localize instances on its own* via different activations in its prototypes.

To see how this is possible, first note that the prototype activations for the solid red image (image a) in Figure 5 are

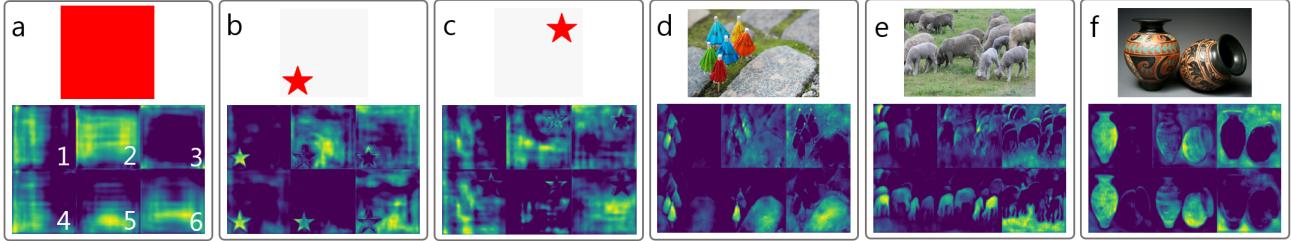


Figure 5: Prototype Behavior The activations of the same six prototypes across different images. Prototypes 1, 4, and 5 are partition maps with boundaries clearly defined in image a, prototype 2 is a bottom-left directional map, prototype 3 segments out the background and provides instance contours, and prototype 6 segments out the ground.

actually not possible in an FCN without padding. Because a convolution outputs to a single pixel, if its input everywhere in the image is the same, the result everywhere in the *conv* output will be the same. On the other hand, the consistent rim of padding in modern FCNs like ResNet gives the network the ability to tell how far away from the image’s edge a pixel is.) Conceptually, one way it could accomplish this is to have multiple layers in sequence spread the padded 0’s out from the edge toward the center (e.g., with a kernel like [1, 0]). This means ResNet, for instance, is *inherently translation variant*, and our method makes heavy use of that property (images b and c exhibit clear translation variance).

We observe many prototypes to activate on certain “partitions” of the image. That is, they only activate on objects that are on one side of an implicitly learned boundary. In Figure 5, prototypes 1, 4, 5, and 6 are such examples (with 6 partitioning the background and not the foreground). By combining these partition maps, the network can distinguish between different (even overlapping) instances of the same semantic class. For instance, in image d, the green umbrella can be separated from the red one by subtracting prototype 5 from prototype 4.

Furthermore, being learned objects, prototypes are compressible. That is, if protonet combines the functionality of multiple prototypes into one, the mask coefficient branch can learn which situations call for which functionality. For instance, in Figure 5, prototype 2 encodes the bottom-left side of objects, but also fires more strongly on instances in a vertical strip down the middle of the image. Prototype 4 is a partitioning prototype but also fires most strongly on instances in the bottom-left corner. Prototype 5 is similar but for instances on the right. This explains why in practice, the model does not degrade in performance even with as low as 32 prototypes (see Table 2b).

4. Backbone Detector

For our backbone detector we prioritize speed as well as feature richness, since predicting these prototypes and coefficients is a difficult task that requires good features to do well. Thus, the design of our backbone detector closely

follows RetinaNet [25] with an emphasis on speed.

YOLOACT Detector We use ResNet-101 [17] with FPN [24] as our default feature backbone and a base image size of 550×550 . We do not preserve aspect ratio in order to get consistent evaluation times per image. Like RetinaNet, we modify FPN by not producing P_2 and producing P_6 and P_7 as successive 3×3 stride 2 *conv* layers starting from P_5 and place 3 anchors with aspect ratios [1, 1/2, 2] on each. The anchors of P_3 have areas of 24 pixels squared, and every subsequent layer has double the scale of the previous (resulting in the scales [24, 48, 96, 192, 384]). For the prediction head attached to each P_i , we have one 3×3 *conv* shared by all three branches, and then each branch gets its own 3×3 *conv* in parallel. Compared to RetinaNet, our prediction head design (see Figure 4) is more lightweight and much faster. We apply smooth- L_1 loss to train box regressors and encode box regression coordinates in the same way as SSD [28]. To train class prediction, we use softmax cross entropy with c positive labels and 1 background label, selecting training examples using OHEM with 3:1 neg:pos ratio. Thus, unlike RetinaNet we do not use focal loss, which we found not to be viable in our situation.

With these design choices, we find that this backbone performs better and faster than SSD [28] modified to use ResNet-101 [17], with the same image size.

5. Other Improvements

We also discuss other improvements that either increase speed with little effect on performance or increase performance with no speed penalty.

Fast NMS After producing bounding box regression coefficients and class confidences for each anchor, like most object detectors we perform NMS to suppress duplicate detections. In many previous works [33, 34, 28, 35, 16, 25], NMS is performed sequentially. That is, for each of the c classes in the dataset, sort the detected boxes descending by confidence, and then for each detection remove all those with lower confidence than it that have an IoU overlap greater than some threshold. While this sequential approach is fast



Figure 6: **YOLACT** evaluation results on COCO’s test-dev set. This base model achieves 29.8 mAP at 33.0 fps. All images have the confidence threshold set to 0.3.

enough at speeds of around 5 fps, it becomes a large barrier for obtaining 30 fps (for instance, a 10 ms improvement at 5 fps results in a 0.26 fps boost, while a 10 ms improvement at 30 fps results in a 12.9 fps boost).

To fix the sequential nature of traditional NMS, we introduce Fast NMS, a version of NMS where every instance can be decided to be kept or discarded in parallel. To do this, we simply allow already-removed detections to suppress other detections, which is not possible in traditional NMS. This relaxation allows us to implement Fast NMS entirely in standard matrix operations available in most GPU-accelerated libraries.

To perform Fast NMS, we first compute a $c \times n \times n$ pairwise IoU matrix X for the top n detections sorted descending by score for each of c classes. Batched sorting on the GPU is readily available and computing IoU can be easily vectorized. Then, we find which detections to remove by checking if there are any higher-scoring detections with a corresponding IoU greater than some threshold t . We efficiently implement this by first setting the lower triangle and diagonal of X to 0:

$$X_{kij} = 0 \quad \forall k, j, i \geq j \quad (2)$$

which can be performed in one batched `triu` call, and then

taking the column-wise max:

$$K_{kj} = \max_i(X_{kij}) \quad \forall k, j \quad (3)$$

to compute a matrix K of maximum IoU values for each detection. Finally, thresholding this matrix with t ($K < t$) will indicate which detections to keep for each class.

Because of the relaxation, Fast NMS has the effect of removing slightly too many boxes. However, the performance hit caused by this is negligible compared to the stark increase in speed (see Table 2a). In our code base, Fast NMS is 11.8 ms faster than a Cython implementation of traditional NMS while only reducing performance by 0.1 mAP. In the Mask R-CNN benchmark suite [16], Fast NMS is 16.5 ms faster than their CUDA implementation of traditional NMS with a performance loss of only 0.3 mAP.

Semantic Segmentation Loss While Fast NMS trades a small amount of performance for speed, there are ways to increase performance with no speed penalty. One of those ways is to apply extra losses to the model during training using modules not executed at test time. This effectively increases feature richness while at no speed penalty.

Thus, we apply a semantic segmentation loss on our feature space using layers that are only evaluated during training. Note that because we construct the ground truth for this loss from instance annotations, this does not strictly capture

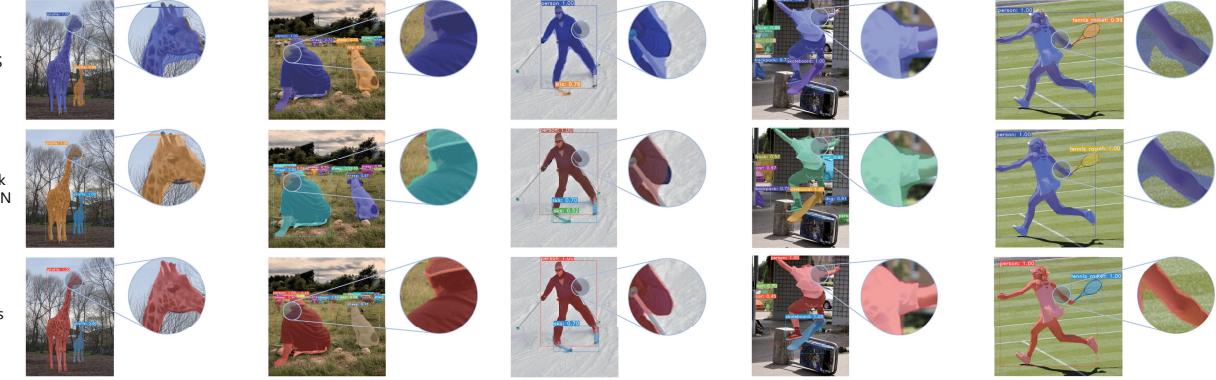


Figure 7: Mask Quality Our masks are typically higher quality than those of Mask R-CNN [16] and FCIS [22] because of the larger mask size and lack of feature repooling. These images were evaluated on the 29.2 mAP version of FCIS and 35.7 mAP version of Mask R-CNN.

semantic segmentation (i.e., we do not enforce the standard one class per pixel). To create predictions during training, we simply attach a 1×1 *conv* layer with c output channels directly to the largest feature map (P_3) in our backbone. Because each pixel can be assigned to more than one class, we use sigmoid and c channels instead of softmax and $c+1$. Training with this loss results in a +0.4 mAP boost.

6. Results

We report results on MS COCO’s instance segmentation task [26] using the standard metrics for the task. We train on `train2017` and evaluate on `val2017` and `test-dev`.

6.1. Instance Segmentation Results

We first compare YOLACT to state-of-the art methods on MS COCO’s `test-dev` set in Table 1. Because our main goal is speed, we compare against other single model results with no test-time augmentations. We report all speeds computed on a single Titan Xp, so some listed speeds may differ from those reported in the original paper.

YOLACT-550 offers competitive instance segmentation performance while at 3.8x the speed of the previous fastest instance segmentation method on COCO. We also note an interesting difference in where the performance of our method lies compared to others, as the performance supports our qualitative findings in Figure 7: observe that the gap between YOLACT-550 and Mask R-CNN at the 50% overlap threshold is 9.5 AP while at the 75% IoU threshold it’s 6.6. This is different from the performance of FCIS, for instance, compared to Mask R-CNN where the gap is consistent (AP values of 7.5 and 7.6 respectively). Furthermore, at the highest (95%) IoU threshold, we outperform Mask R-CNN with 1.6 AP vs. 1.3 AP.

We also report numbers for alternate configurations of our model in Table 1. In addition to our base 550×550

image size model, we train 400×400 (YOLACT-400) and 700×700 (YOLACT-700) models, adjusting the anchor scales for these models accordingly ($s_x = s_{550}/550 * x$). Lowering the image size results in a large decrease in performance, demonstrating that instance segmentation naturally demands larger images. Then, raising the image size decreases speed significantly but also increases performance, as expected.

In addition to our base backbone of ResNet-101 [17], we also test ResNet-50 and DarkNet-53 [34] to obtain even faster results. If higher speeds are preferable we suggest using ResNet-50 or DarkNet-53 instead of lowering the image size, as the performance of these configurations is much better than YOLACT-400, while only being slightly slower.

6.2. Mask Quality

Because we produce a final mask of size 138×138 , and because we create masks directly from the original features (with no repooling step to transform and potentially misalign the features), our masks for large objects are noticeably higher quality than those of Mask R-CNN [16] and FCIS [22]. For instance, in Figure 7, YOLACT produces a mask that cleanly follows the boundary of the arm, whereas both FCIS and Mask R-CNN have more noise. Moreover, despite being 5.9 mAP worse overall, at the 95% IoU threshold, our base model achieves 1.6 AP while Mask R-CNN obtains 1.3. This indicates that repooling does result in a quantifiable decrease in mask quality.

6.3. Temporal Stability

Although we only train our model using static images and do not apply any temporal smoothing, we find that the model produces more temporally stable masks on videos than Mask R-CNN, whose masks jitter across frames even when objects are completely stationary. We believe our masks are more temporally stable in part because they are

Method	Backbone	FPS	Time	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
PA-Net [27]	R-50-FPN	4.7	212.8	36.6	58.0	39.3	16.3	38.1	53.1
RetinaMask [12]	R-101-FPN	6.0	166.7	34.7	55.4	36.9	14.3	36.7	50.5
FCIS [22]	R-101-C5	6.6	151.5	29.5	51.5	30.2	8.0	31.0	49.7
Mask R-CNN [16]	R-101-FPN	8.6	116.3	35.7	58.0	37.8	15.5	38.1	52.4
MS R-CNN [18]	R-101-FPN	8.6	116.3	38.3	58.8	41.5	17.8	40.4	54.4
YOLACT-550	R-101-FPN	33.0	30.3	29.8	48.5	31.2	9.9	31.3	47.7
YOLACT-400	R-101-FPN	44.0	22.7	24.9	42.0	25.4	5.0	25.3	45.0
YOLACT-550	R-50-FPN	42.5	23.5	28.2	46.6	29.2	9.2	29.3	44.8
YOLACT-550	D-53-FPN	40.0	25.0	28.7	46.8	30.0	9.5	29.6	45.5
YOLACT-700	R-101-FPN	23.6	42.4	31.2	50.6	32.8	12.1	33.3	47.1

Table 1: **Mask Performance** We compare our approach to other state-of-the-art methods for mask mAP and speed on COCO test-dev. We denote the backbone architecture with the nomenclature network-depth-features, where R and D refer to ResNet [17] and DarkNet [34], respectively. Our base model, YOLACT-550 with ResNet-101, is 3.8x faster than the previous fastest approach with competitive mask mAP.

Method	NMS	AP	FPS	Time	<i>k</i>	AP			Time
						FPS	Time	AP	
YOLACT	Standard	30.0	23.8	42.1	32	27.7	32.4	30.9	105.3
	Fast	29.9	33.0	30.3	64	27.8	31.7	31.5	
Mask R-CNN	Standard	36.1	8.6	116.0	128	27.6	31.5	31.8	73.9
	Fast	35.8	9.9	101.0	256	27.7	29.8	33.6	

(a) **Fast NMS** Fast NMS performs only slightly worse than standard NMS, while being around 12 ms faster. We also observe a similar trade-off implementing Fast NMS in Mask R-CNN.

(b) **Prototypes** Choices for *k* in our method. YOLACT is robust to varying *k*, so we choose the fastest (*k* = 32).

(c) **Accelerated Baselines** We compare to other baseline methods by tuning their speed-accuracy trade-offs. *fc*-mask is our model but with 16×16 masks produced from an *fc* layer.

Table 2: **Ablations** All models evaluated on COCO val2017 using our servers. Models in Table 2b were trained for 400k iterations instead of 800k. Time in milliseconds reported for convenience.

higher quality (thus there is less room for error between frames), but mostly because our model is one-stage. Masks produced in two-stage methods are highly dependent on their region proposals in the first stage. In contrast for our method, even if the model predicts different boxes across frames, the prototypes are not affected, yielding much more temporally stable masks.

6.4. Implementation Details

We train all models with **batch size 8** on one GPU using ImageNet [9] pretrained weights. We find that this is a sufficient batch size to use batch norm, so we leave the pretrained batch norm unfrozen but do not add any extra *bn* layers. We train with **SGD** for 800k iterations starting at an initial learning rate of 10^{-3} and divide by 10 at iterations 280k, 600k, 700k, and 750k, using a weight decay of 5×10^{-4} and a momentum of 0.9. We also train with all data augmentations used in SSD [28].

7. Discussion

Despite our masks being higher quality and having nice properties like temporal stability, we fall behind state-of-the-art instance segmentation methods in overall perfor-

mance, albeit while being much faster. Most errors are simply caused by mistakes in the detector: misclassification, box misalignment, etc. However, we have identified two typical errors caused by YOLACT’s mask generation algorithm.

Localization Failure If there are too many objects in one spot in a scene, the network can fail to localize each object in its own prototype. In these cases, the network will output something closer to a foreground mask than an instance segmentation for some objects in the group. An example of this can be seen in the first image in Figure 6 (row 1 column 1), where the blue truck under the red airplane is not properly localized.

Leakage Our network leverages the fact that masks are cropped after assembly, and makes no attempt to suppress noise outside of the cropped region. This works fine when the bounding box is accurate, but when it is not, that noise can creep into the instance mask, creating some “leakage” from outside the cropped region. This can also happen when two instances are far away from each other, because the network has learned that it doesn’t need to localize far away instances—the cropping will take care of it. However, if the predicted bounding box is too big, the mask will include



Figure 8: **More YOLACT** evaluation results on COCO’s test-dev set with the same parameters as before. To further support that YOLACT implicitly localizes instances, we select examples with adjacent instances of the same class.

some of the far away instance’s mask as well. For instance, Figure 6 (row 2 column 4) exhibits this leakage because the mask branch deems the three skiers to be far enough away to not have to separate them.

These issues could potentially be mitigated with a mask error down-weighting scheme like in MS R-CNN [18], where masks exhibiting these errors could be ignored. However, we leave this for future works to address.

References

- [1] S. Agarwal and D. Roth. Learning a sparse representation for object detection. In *ECCV*, 2002. 3
- [2] V. Badrinarayanan, A. Kendall, and R. Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *CoRR*, 2015. 2
- [3] M. Bai and R. Urtasun. Deep watershed transform for instance segmentation. In *CVPR*, 2017. 2
- [4] L.-C. Chen, A. Hermans, G. Papandreou, F. Schroff, P. Wang, and H. Adam. Masklab: Instance segmentation by refining object detection with semantic and direction features. In *CVPR*, 2018. 2
- [5] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016. 2
- [6] J. Dai, K. He, Y. Li, S. Ren, and J. Sun. Instance-sensitive fully convolutional networks. In *ECCV*, 2016. 2
- [7] J. Dai, Y. Li, K. He, and J. Sun. R-fcn: Object detection via region-based fully convolutional networks. In *NeurIPS*, 2016. 1
- [8] B. De Brabandere, D. Neven, and L. Van Gool. Semantic instance segmentation with a discriminative loss function. *arXiv preprint arXiv:1708.02551*, 2017. 2
- [9] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A Large-Scale Hierarchical Image Database. In *CVPR*, 2009. 8
- [10] N. Dvornik, K. Shmelkov, J. Mairal, and C. Schmid. Blitznet: A real-time deep network for scene understanding. In *ICCV*, 2017. 2
- [11] A. Fathi, Z. Wojna, V. Rathod, P. Wang, H. O. Song, S. Guadarrama, and K. P. Murphy. Semantic instance segmentation via deep metric learning. *arXiv preprint arXiv:1703.10277*, 2017. 2
- [12] C.-Y. Fu, M. Shvets, and A. C. Berg. Retinamask: Learning to predict masks improves state-of-the-art single-shot detection for free. *arXiv preprint arXiv:1901.03353*, 2019. 8

- [13] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, 2012. 2
- [14] M. A. Goodale and A. Milner. Separate visual pathways for perception and action. *Trends in Neurosciences*, 1992. 2
- [15] A. W. Harley, K. G. Derpanis, and I. Kokkinos. Segmentation-aware convolutional networks using local attention masks. In *ICCV*, 2017. 2
- [16] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *ICCV*, 2017. 1, 2, 3, 4, 5, 6, 7, 8
- [17] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 2, 5, 7, 8
- [18] Z. Huang, L. Huang, Y. Gong, C. Huang, and X. Wang. Mask scoring r-cnn. In *CVPR*, 2019. 2, 8, 9
- [19] S. Jetley, M. Sapienza, S. Golodetz, and P. H. Torr. Straight to shapes: real-time detection of encoded shapes. In *CVPR*, 2017. 2
- [20] A. Kirillov, E. Levinkov, B. Andres, B. Savchynskyy, and C. Rother. Instancecut: from edges to instances with multi-cut. In *CVPR*, 2017. 2
- [21] T. Leung and J. Malik. Representing and recognizing the visual appearance of materials using three-dimensional textures. *IJCV*, 2001. 3
- [22] Y. Li, H. Qi, J. Dai, X. Ji, and Y. Wei. Fully convolutional instance-aware semantic segmentation. In *CVPR*, 2017. 1, 2, 4, 7, 8
- [23] X. Liang, L. Lin, Y. Wei, X. Shen, J. Yang, and S. Yan. Proposal-free network for instance-level object segmentation. *TPAMI*, 2018. 2
- [24] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017. 4, 5
- [25] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *CVPR*, 2017. 3, 4, 5
- [26] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 2, 7
- [27] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia. Path aggregation network for instance segmentation. In *CVPR*, 2018. 2, 8
- [28] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. Ssd: Single shot multibox detector. In *ECCV*, 2016. 1, 2, 4, 5, 8
- [29] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 3
- [30] A. Newell, Z. Huang, and J. Deng. Associative embedding: End-to-end learning for joint detection and grouping. In *NeurIPS*, 2017. 2
- [31] A. Paszke, A. Chaurasia, S. Kim, and E. Culurciello. Enet: A deep neural network architecture for real-time semantic segmentation. *CoRR*, 2016. 2
- [32] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, 2016. 2
- [33] J. Redmon and A. Farhadi. Yolo9000: Better, faster, stronger. In *CVPR*, 2017. 1, 2, 5
- [34] J. Redmon and A. Farhadi. Yolov3: An incremental improvement. *arXiv:1804.02767*, 2018. 1, 2, 5, 7, 8
- [35] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NeurIPS*, 2015. 1, 3, 5
- [36] X. Ren and D. Ramanan. Histograms of sparse codes for object detection. In *CVPR*, 2013. 3
- [37] J. Sivic and A. Zisserman. Video google: A text retrieval approach to object matching in videos. In *ICCV*, 2003. 3
- [38] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, 2015. 2
- [39] M. Treml, J. Arjona-Medina, T. Unterthiner, R. Durgesh, F. Friedmann, P. Schubert, A. Mayr, M. Heusel, M. Hofmarcher, M. Widrich, et al. Speeding up semantic segmentation for autonomous driving. In *NeurIPS Workshops*, 2016. 2
- [40] J. Uhrig, E. Rehder, B. Fröhlich, U. Franke, and T. Brox. Box2pix: Single-shot instance segmentation by assigning pixels to object boxes. In *IEEE Intelligent Vehicles Symposium*, 2018. 2
- [41] J. Wang, J. Yang, K. Yu, F. Lv, T. Huang, and Y. Gong. Locality-constrained linear coding for image classification. In *CVPR*, 2010. 3
- [42] J. Yang, J. Wright, T. S. Huang, and Y. Ma. Image super-resolution via sparse representation. *IEEE Transactions on Image Processing*, 2010. 3
- [43] X. Yu, L. Yi, C. Fermüller, and D. S. Doermann. Object detection using shape codebook. In *BMVC*, 2007. 3
- [44] T. Zhang, B. Ghanem, S. Liu, C. Xu, and N. Ahuja. Low-rank sparse coding for image classification. In *ICCV*, 2013. 3
- [45] H. Zhao, X. Qi, X. Shen, J. Shi, and J. Jia. Icnet for real-time semantic segmentation on high-resolution images. In *ECCV*, 2018. 2