

The surprising impact of mask-head architecture on novel class segmentation

Vighnesh Birodkar, Zhichao Lu, Siyang Li, Vivek Rathod, Jonathan Huang
Google

{vighneshb, lzc, siyang, rathodv, jonathanhuang}@google.com

Abstract

Instance segmentation models today are very accurate when trained on large annotated datasets, but collecting mask annotations at scale is prohibitively expensive. We address the partially supervised instance segmentation problem in which one can train on (significantly cheaper) bounding boxes for all categories but use masks only for a subset of categories. In this work, we focus on a popular family of models which apply differentiable cropping to a feature map and predict a mask based on the resulting crop. Within this family, we show that the architecture of the mask-head plays a surprisingly important role in generalization to classes for which we do not observe masks during training. While many architectures perform similarly when trained in fully supervised mode, we show that they often generalize to novel classes in dramatically different ways. We call this phenomenon the strong mask generalization effect, which we exploit by replacing the typical mask-head of 2-4 layers with significantly deeper off-the-shelf architectures (e.g. ResNet, Hourglass models). We also show that the choice of mask-head architecture alone can lead to SOTA results on the partially supervised COCO benchmark without the need of specialty modules or losses proposed by prior literature. Finally, we demonstrate that our effect is general, holding across underlying detection methodologies, (e.g. both anchor-based or anchor free or no detector at all) and across different backbone networks. Code and pre-trained models are available at <https://git.io/deepmac>.

1. Introduction

Large labeled datasets like COCO [31] are crucial for deep neural network based instance segmentation methods [14, 3, 37]. However, collecting groundtruth masks can take $> 10\times$ more time than bounding box annotations. In COCO [31], mask annotations required ≈ 80 seconds on average whereas methods such as Extreme Clicking [36] yield bounding boxes in 7 seconds.

Given that boxes are much cheaper to annotate than masks, we address the “partially supervised” instance seg-



Figure 1: The effect of mask-head architecture on mask predictions for unseen classes. Despite never having seen masks from the ‘parking meter’, ‘pizza’ or ‘mobile phone’ class, the rightmost mask-head architecture can segment these classes correctly. From left to right, we show better mask-head architectures predicting better masks. Moreover, this difference is only apparent when evaluating on unseen classes — if we evaluate on seen classes, all four architectures exhibit similar performance.

mentation problem [18], where all classes have bounding box annotations but only a subset of classes have mask annotations. We will refer to classes with mask annotations as “seen” categories and classes without as “unseen”. Doing well on this task requires the model to generalize in a strong sense, producing correct masks on unseen classes.

We consider a general family of segmentation models where one extracts a feature map over an image, then given a tight bounding box around an instance, performs a differentiable crop (e.g. ROIAlign [14]). The cropped feature map is then fed to a mask-head subnetwork to yield a final

mask prediction. This mask prediction is performed in a class agnostic manner so that a model trained from a subset of classes can be applied unchanged to novel classes.

Previous approaches in this family have used, e.g., offline-trained shape priors [26] or specialty losses [9]. In contrast, we focus on an overlooked part of the problem: the architectural design of the mask-head. While trying alternative mask-head architectures may seem an obvious approach, it is underexplored in the prior work because (1) the choice of mask-head architecture has limited impact in the fully supervised setting and (2) heavier mask-heads adversely impact running time. Thus most prior works in instance segmentation have settled on using shallow (2-4 layer) fully connected or convolution based mask-heads.

Our main finding is that mask-heads that might perform similarly under full supervision can behave differently under partial supervision, generalizing to unseen classes in strikingly different ways. This difference subsequently changes the calculus for deciding whether it’s worth using a heavier mask-head. In our COCO experiments, we find that the difference between worst and best architectures is only 1% (absolute mAP) on seen classes but can be 7% on unseen classes (examples in Figure 1).

We will refer to this effect of certain mask-head architectures on unseen classes as the “*strong mask generalization effect*” and illustrate it with 3 representative model classes: an anchor-free and anchor-based model, and one that discards detection altogether. We show that our effect is general, holding across underlying detection methodologies (or no detector at all) and across different backbone networks. We also identify architectural characteristics (such as depth and encoder-decoder arrangements) that empirically yield strong mask generalization properties.

Our strongest results come from our anchor-free model, which adds a mask-head to CenterNet [54]. Under this family, we show that choosing the right out-of-the-box mask-head architecture alone allows us to surpass the state of the art on the partially supervised COCO benchmark (35.5% mAP) without need of hand designed priors or additional losses. Our best mask-head uses up to 100 layers following an hourglass pattern (we find that deeper mask-heads generalize better despite being counter-intuitively more over-parameterized than shallower ones). We call this model “Deep Mask-heads Above CenterNet” (*Deep-MAC*) and focus our most exhaustive experiments on Deep-MAC.

We also study Mask R-CNN [14], which continues to be one of the most popular and enduring anchor-based segmentation models. Prior works have shown canonical implementations of Mask R-CNN (with a class agnostic mask-head) to perform poorly on the partially supervised segmentation task. We show that a major culprit for this poor performance is the way that Mask R-CNN’s mask-head is trained with (typically noisy) proposals instead of

groundtruth boxes. By modifying the training procedure slightly to crop using groundtruth boxes, we dramatically improve the performance of Mask R-CNN on unseen categories with the standard mask-head (+7% mask mAP). With this change, we evaluate alternative mask-head architectures and show that Mask R-CNN also exhibits the strong mask generalization effect with the strongest mask-head architectures similarly reaching performance on-par with the state of the art.

Finally, we consider a segmentation only model (which takes boxes as an input rather than requiring boxes to also be detected) and show that the strong mask generalization effect can also be observed in this detection-free setting.

When we isolate segmentation quality from detection quality, we find that with respect to segmentation quality our best models have likely already reached a saturation point close to human performance. The implication of this finding is that any future improvements on this task are likely to come from better detection (which does not require strong generalization since we assume that all classes have box annotations). To illustrate, we use a two-stage training procedure, employing Deep-MAC to label masks for unseen categories and training a stronger detector via a second phase in fully supervised mode on these pseudo-labels. This achieves 40.4% non-VOC mask mAP, outperforming the previous best result [9] by 6.4%.

Going in the opposite direction, another benefit of this two stage procedure is that it allows us to train a lightweight mask-head based on pseudo-masks generated from a heavier mask-head. While the heavier mask-head has better mask generalization, we can exploit the fact that both heads perform equally well when fully supervised — thus in the second stage re-training phase, we obtain a model that performs well on the unseen categories while retaining the computational benefits of a lightweight mask-head.

We summarize our main contributions as follows:

- We identify the *strong mask generalization effect* in partially supervised instance segmentation architectures and show that it is general, holding across underlying detectors like CenterNet [54] (Section 4), Mask R-CNN [14] or with no detector at all and across different backbones (Section 5).
- We identify a simple but critical tweak to training for Mask R-CNN which dramatically improves performance on unseen classes (Table 4) and unlocks strong mask generalization.
- We identify characteristics of mask-head architectures (Section 6) that lead to strong mask generalization. Among other things, we find that Hourglass [35] architectures offer excellent performance. We use these findings to achieve state-of-the-art results on the COCO partial supervision task (Table 2, Table 4) with models based on CenterNet and Mask R-CNN.

- Finally, we argue that with respect to mask quality, we are close to saturating performance on COCO and show via experiments that it is much easier to achieve gains on this task simply by using a stronger detector (that need not exhibit strong mask generalization) (Section 7).

2. Related work

Object detection and instance segmentation. There has been a significant progress over the last decade in detection with successful convolutional models like OverFeat [43], YOLO [38, 39, 40, 2], Multibox [44, 8], SSD [34], RetinaNet [30], R-CNN and Fast/Faster versions [11, 10, 42, 15], EfficientDet [45], etc. While many of these works initially focused on box detection, more recently, many benchmarks have focused on the more detailed problem of instance segmentation (COCO [31], OID v5 [27, 1], LVIS [13]) and panoptic segmentation (COCO-Panoptic, [24]) which are arguably more useful tasks in certain applications. A major milestone in this literature was Mask R-CNN [14] which influenced many SOTA approaches today (e.g., [37, 33]) and by itself continues to serve as a strong baseline.

Anchor-free methods. State-of-the-art methods today are predominantly built on anchor-based approaches which predict classification/box offsets relative to a collection of fixed boxes arranged in sliding window fashion (called “anchors”). While effective, the performance of anchor-based methods often depend on manually-specified design decisions, e.g. anchor layouts and target assignment heuristics, a complex space to navigate for practitioners.

In recent years, however, this monopoly has been broken with the introduction of competitive “anchor-free” approaches [28, 7, 46, 55, 54, 25, 56, 4]. These newer anchor-free methods are simpler, more amenable to extension, offer competitive performance and consequently are beginning to be popular. Our anchor-free model (Section 3) in particular builds on the “CenterNet” architecture [54].

Due to the recency of competitive anchor-free methods there are fewer anchor-free instance segmentation approaches in literature. [29, 50, 51, 9] all add mask prediction capabilities on top of the (anchor-free) FCOS [46] framework. While the primary focus of our work is partial supervision, the fully supervised version of our model adds to this growing body of work, offering strong performance among anchor-free instance segmentation approaches.

Box-only supervision for instance segmentation. The above methods rely on access to massive labeled datasets which are costly to develop, with mask annotations especially so compared to box annotations. Researchers have thus begun to develop methods that are less reliant on mask annotations. In one formulation of this problem (which we

might call *strictly box-supervised*) we ask to learn an instance segmentation model given only box annotations and no masks [22, 41, 17, 21, 47]. However this is intuitively a difficult approach and the performance of all of these methods is still a far cry from fully supervised performance of a strong baseline particularly at high IOU thresholds for mAP.

Partial supervision for instance segmentation Instead of going to the extreme end of discarding all mask annotations, Hu et al. [18] introduced the *partial supervision* formulation which allows for mask annotations from a small subset of classes to be used along with all box annotations. [18] observed that the “obvious” baseline of using a class-agnostic version of Mask R-CNN yielded poor results and proposed a method (Mask^X) for learning to predict mask-head weights given box-head weights hoping that this learned function will generalize to classes whose masks are not observed at training time.

Later papers [26, 9] however revisited the approach of attaching a class-agnostic mask-head on top of a detector, in both cases introducing novel architectures and additional losses to significantly improve generalization to novel classes. ShapeMask [26] builds on RetinaNet, learning a low dimensional shape space from observed masks and uses projections to this space to guide mask estimation; they also introduce a simple method to “condition” features cropped from the backbone on the instance that is being segmented. CP-Net [9], which is the current state of the art on this problem builds on FCOS [46], adding boundary prediction and attention-based aggregation in the mask branch.

We take a similar approach of using a class-agnostic mask-head, but while the ideas explored in these prior works are clearly beneficial, our objective is to demonstrate that mask-head architecture itself plays an underappreciated but significant role in generalization. Notably, by exploiting out-of-the-box architectures with strong mask generalization properties, we show that with only minor tweaks to the training procedure (Section 5.2), even Mask R-CNN has state of the art performance in the partial supervision task.

3. Deep-MAC : an anchor-free model

As our first illustration of the strong mask generalization effect, we introduce the Deep-MAC architecture, which builds instance segmentation capabilities on top of CenterNet [54],¹ a popular anchor-free detection model.

The CenterNet architecture. CenterNet models objects relative to their centers. For predicting bounding boxes it outputs 3 tensors: (1) a class-specific heatmap which indicates the probability of the center of a bounding box being present at each location, (2) a class-agnostic 2-channel tensor indicating the height and width of the bounding box at

¹Not to be confused for the CenterNet from Duan et al. [7].

Model	Backbone	AP	AP_S	AP_M	AP_L
ShapeMask [26]	RF101	37.4	16.1	40.1	53.8
ShapeMask [26]	RNF101	40.0	18.3	43.0	57.1
CPMask [9]	RF101	39.2	22.2	41.8	50.1
Deep-MAC	HG104	39.4	20.5	41.9	54.0

Table 1: Fully supervised instance segmentation performance on COCO test-dev2017. Backbones include RF=ResNet-FPN, RNF=ResNet-NAS-FPN, HG=Hourglass. Deep-MAC (our model) is trained at 1024 \times 1024 resolution with an HG-100 mask-head.

each center pixel, and (3) since the output feature map is typically smaller than the image (stride 4 or 8), CenterNet also predicts an x and y direction offset to recover this discretization error at each center pixel.

Predicting instance masks with CenterNet (Deep-MAC).

In parallel to the box-related prediction heads, we add a fourth *pixel embedding* branch P . For each bounding box b , we crop a region P_b from P corresponding to b via ROIAlign [14] which results in a 32×32 tensor. As in Mask R-CNN [14], we feed each P_b to a second stage mask-head network whose architecture is discussed in the next section. Our final prediction at the end is a class-agnostic 32×32 tensor which we postprocess into a binary mask at test time by applying a sigmoid and thresholding at 0.5. We train this mask-head via a per-pixel sigmoid cross-entropy loss averaged over all pixels and instances.

In addition to this 32×32 cropped feature map, we add two inputs for improved stability of some mask-heads (but note that our main findings *do not depend* on having these additional inputs; see Appendix for ablations): **(1) Instance embedding:** We add an additional head to the backbone that predicts a per-pixel embedding. For each bounding box b we extract its embedding from the center pixel. This embedding is tiled to a size of 32×32 and concatenated to the pixel embedding crop. This helps condition the mask-head on a particular instance and disambiguate it from others. **(2) Coordinate Embedding:** Inspired by CoordConv[32], we add a $32 \times 32 \times 2$ tensor holding normalized (x, y) coordinates relative to the bounding box b .

During training, we use groundtruth bounding boxes and their masks to train the second stage network (this is an important detail which we revisit in Section 5.2 when discussing Mask R-CNN). At test time, we use the boxes predicted by the detection branch. We use an input resolution of 512×512 for our ablations and 1024×1024 for our best results. Other hyperparameters are unchanged across all Deep-MAC experiments (full details in Appendix). Table 1 compares the fully supervised performance of Deep-MAC with other partial supervision capable models.

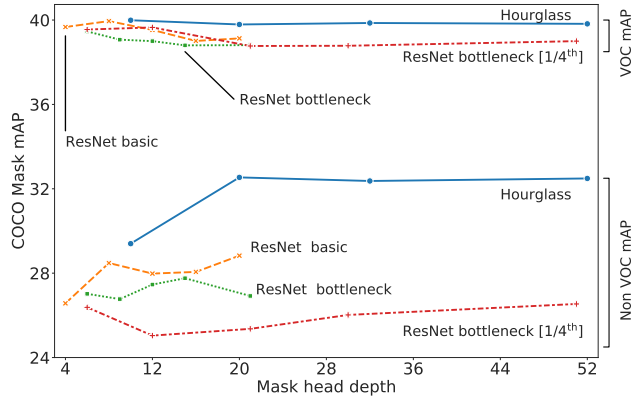


Figure 2: Effect of mask-head architecture and depth on instance segmentation performance over seen (VOC) and unseen (Non-VOC) classes. Although the performance on seen classes does not vary much across different architectures, there is significant variation in the performance on unseen classes. All models are trained with masks only from VOC classes with an input image resolution of 512×512 . Performance is evaluated on the coco-val2017 set.

4. Strong mask generalization in Deep-MAC

We now focus on our main topic, partially supervised instance segmentation where mask annotations are available for a subset of *seen* categories and unavailable for the remaining *unseen* categories. Our main finding in this section is that *mask-head architectures for Deep-MAC affect generalization behavior on unseen classes to a surprising extent*.

For all experiments in this paper (unless otherwise mentioned), we follow the typical partially supervised experimental setup for the COCO dataset with 20 (VOC) categories assumed to be seen and the remaining 60 unseen. In Figure 2 we show the effect of using different mask-heads on the mask mAP (evaluated on coco-val2017 set) of seen (VOC) and unseen (Non-VOC) classes. We experiment with 4 mask-head architectures: (a) Hourglass, a stacked hourglass [35] mask-head, (b) ResNet basic, a ResNet [16] v1 mask-head composed of basic ResNet building blocks, (c) ResNet bottleneck, a ResNet v1 mask-head composed of bottleneck building blocks, and (d) ResNet bottleneck [1/4th], a variant of the ResNet bottleneck mask-head with $4 \times$ fewer channels.

We first observe that while mAP on the seen classes depends a little bit on the specific mask-head architecture, the effect is small ($38.8 \rightarrow 40.0$). However, for the same settings, the mAP on unseen classes varies significantly ($25.0 \rightarrow 32.5$). This indicates that mask-head architectures play a critical role in generalization to unseen classes and not just by virtue of fitting the training data better. For one thing, depth plays a role: empirically, it is important to go significantly beyond 4 layers to achieve the best performance. From a classical perspective, this is counterintuitive given the over-parameterization of very deep mask-heads, but perhaps is not so surprising in light of recent ways of

Model	b-box.	VOC \rightarrow Non-VOC (mask)						Non-VOC \rightarrow VOC (mask)					
	AP	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
Mask R-CNN [14]	38.2	18.5	24.8	18.1	11.3	23.4	21.7	24.7	43.5	24.9	11.4	25.7	35.1
Mask GrabCut [18]	38.2	19.7	39.7	17.0	6.4	21.2	35.8	19.6	46.1	14.3	5.1	16.0	32.4
Mask ^X R-CNN	38.2	23.8	42.9	23.5	12.7	28.1	33.5	29.5	52.4	29.7	13.4	30.2	41.0
ShapeMask [26]	45.4	33.2	53.1	35.0	18.3	40.2	43.3	35.7	60.3	36.6	18.3	40.5	47.3
CPMask [9]	41.5	34.0	53.7	36.5	18.5	38.9	47.4	36.8	60.5	38.6	17.6	37.1	51.5
Deep-MAC (ours)	44.5	35.5	54.6	38.2	19.4	40.3	50.6	39.1	62.6	41.9	17.6	38.7	54.0
Oracle Mask-RCNN [14]	38.2	34.4	55.2	36.3	15.5	39.0	52.6	39.1	64.5	41.4	16.3	38.1	55.1
Oracle ShapeMask [26]	45.4	37.6	57.7	40.2	20.1	44.4	51.1	43.1	67.9	45.8	20.1	44.3	57.8
Oracle CPMask [9]	41.5	37.6	58.2	40.2	19.9	42.6	54.2	42.9	67.6	46.6	21.6	42.1	58.9
Oracle Deep-MAC	44.5	38.3	57.5	41.8	20.7	42.1	52.7	42.6	66.6	46.1	19.8	41.2	58.0

Table 2: Partially supervised performance of Deep-MAC compared to other models. We measure mask mAP on the `coco-val2017` set. The top row with label A \rightarrow B indicates that we train on masks from set A and evaluate our masks on set B. The Oracle models indicate the performance when trained with masks from all categories. We also report the detection (b-box) AP of the base detector used in each of the models on `coco-testdev2017`. Deep-MAC outperforms previous partially supervised approaches and even exceeds the Oracle performance of Mask R-CNN.

rethinking generalization for deep learning [53, 52]. However, depth is not the only factor that drives generalization.

Among the above alternatives, the Hourglass [35] architecture for the mask-head provides the best generalization performance to unseen classes. We note that this is fortunate since, compared to the other mask-heads, it is also the most memory-efficient due to its successive downsampling layers, which make the feature maps smaller as depth increases. As we next show, this out-of-the-box architecture allows us, with no extra losses or additional priors, to surpass current state-of-the-art results by a significant margin.

Comparison with previous approaches In the previous experiments, we trained Deep-MAC on 512x512 input images. We now train models on 1024x1024 inputs and compare against previous state-of-the-art approaches in Table 2. For our best result we used an Hourglass-100 mask-head, which exceeds the previous state of the art [9] on VOC to Non-VOC transfer by 1.5% and Non-VOC to VOC transfer by 2.3%. Compared to prior approaches, our method is end-to-end trainable and does not require auxiliary losses or specialty modules. Moreover, Deep-MAC is the first partially supervised model to exceed the Oracle (fully supervised) Mask R-CNN performance in the VOC \rightarrow Non-VOC setting and we match its overall mAP in the Non-VOC \rightarrow VOC setting. Although Deep-MAC surpasses the state of the art by itself, we show in Section 7 that we can do significantly better by using stronger detection architectures.

5. How general is the strong mask generalization effect?

Having established that mask-head architecture significantly affects strong mask generalization in Deep-MAC, we now go beyond the particulars of Deep-MAC and show that

Mask-Head	ResNet-101-FPN		Hourglass-104	
	Box	Mask	Box	Mask
ResNet-4	32.6	22.6	39.7	26.6
Hourglass-10	32.2	24.8	39.9	29.4
Hourglass-20	32.5	26.7	39.7	32.5

Table 3: Effect of backbones on the relative performance of mask-heads. We report the mAP on non-VOC classes. Training and evaluation setup is the same as Figure 2. Note that the box mAP is relatively unchanged for each backbone because we train using all boxes.

this is in fact a highly general phenomenon. We first demonstrate that the effect occurs independently of the backbone network used by Deep-MAC. Second, we show that it is not critical that the model be based on the CenterNet detector — Mask R-CNN works as well. Finally we show that there need not even be an underlying detection model.

5.1. Impact of different backbones

Our first question is whether the strong performance of the Hourglass based mask-heads from the previous section are somehow tied to our choice of Hourglass backbone in the CenterNet model. To address this, we compare (Table 3) the performance of various mask-heads when using a ResNet-FPN [16] and an Hourglass [35] backbone. Although we observe that absolute mask mAP values on Non-VOC classes are lower using the ResNet backbone, we find that the same trends hold; specifically, mask-head architecture impacts generalization to unseen classes, and the Hourglass mask-head is still the best relative to other alternatives. Thus, our strong mask generalization effect is not tied to the specific choice of Hourglass backbone.

5.2. Strong mask generalization in Mask R-CNN

Next we show strong mask generalization is *not* tied to our choice of the CenterNet detector. In the following, we

Mask-Head	ResNet-50-FPN		ResNet-101-FPN	
	Prop.	GT.	Prop.	GT.
4x Conv	18.2	25.9	19.6	29.3
ResNet-4	17.4	24.6	21.0	27.4
HG-20	22.1	32.8	20.6	33.8
HG-52	20.2	33.1	20.6	34.4

Table 4: Performance of Mask RCNN with its mask-head trained in a class agnostic manner with different mask-heads. All models were trained only on VOC classes and we report mask mAP on non-VOC classes on `coco-val2017` set. We compare the performance when training the proposed boxes (Prop.) and groundtruth (GT.) boxes and note that performance when training with groundtruth boxes is markedly better. A ResNet-101-FPN backbone combined with an HG-52 mask-head also manages to beat the previous SOTA [9].

modify Mask R-CNN [14] to use a class-agnostic mask-head and show that here too (with a simple modification to the typical training procedure), mask-head architectures significantly impact strong mask generalization.

Whereas canonical implementations of Mask R-CNN use (typically noisy) RPN-generated proposals as input to the mask-head during training, we find that it is necessary to train the mask-head with *only* groundtruth boxes. Table 4 summarizes our results and shows in particular the dramatic increase in performance when training the mask-head with groundtruth boxes instead of proposals. Consistent with prior works (e.g., [18]), Mask R-CNN trained with proposals generalizes poorly to unseen categories, and we find that better mask-heads do not improve the situation. Implementing our simple modification of training with groundtruth boxes on the other hand leads to +7.7 and +9.7 mAP gains using the standard convolutional mask-head and ResNet-50-FPN and ResNet-101-FPN backbones respectively. Interestingly, we note that training using groundtruth boxes instead of proposals does not significantly impact performance when fully supervised (see Appendix) — why this makes such a big difference for unseen classes is thus an interesting question further discussed in Section 6.

In the last column of Table 4, we see that better mask-heads also dramatically improve generalization on unseen classes following a similar trend compared to Deep-MAC, thus demonstrating strong mask generalization in the Mask R-CNN setting. In fact, we find that our ResNet-101-FPN based Mask-RCNN model (with HG-52 mask-head) also outperforms the current state of the art, CPMask [9]. For detailed experiment setup, see the Appendix.

5.3. Strong mask generalization without a detector

To further make the point that detection does not play a critical role in our story, we consider a “detection-free” incarnation of our model family, in which we do not even require the model to produce detections. In this most basic of settings, we assume that a groundtruth box for each instance is provided as input and the task is to simply pro-

Mask-Head	mIOU		
	Overall	VOC	Non-VOC
ResNet-4	67.0	78.6	62.1
Hourglass-20	78.6	81.0	77.8
Hourglass-52	78.9	81.1	79.2

Table 5: mIOU of Deep-MAC trained without any detection losses. Because we cannot compute Mask mAP without a detector, we report mIOU computed over the full dataset and over VOC/non-VOC class splits. The models were trained only on VOC masks with the same setup as Figure 2. Hourglass mask-heads continue to show strong mask generalization on non-VOC classes, even when they are not coupled with a detector.

duce the correct segmentation mask. For this setting, we use the same architecture family as above, using an Hourglass backbone, cropping to each groundtruth box and passing the result to the mask-head. Since detection is no longer a task of interest, we drop all detection related losses and train only with sigmoid cross entropy loss for the masks. We also evaluate using the mean IOU metric instead of mask mAP.

Table 5 shows the results of this experiment (evaluating ResNet-4, and Hourglass 52, 100) mask-head architectures. We observe that all of these architectures have similar performance on the seen categories ($\sim 2.5\%$ spread) whereas on unseen categories, the best mask-head (HG-100) outperforms the worst (ResNet-4) by $>16\%$. This confirms the strong mask generalization effect occurs in the detection-free setting and together with our results for Deep-MAC and Mask R-CNN provide strong evidence that we would find similar effects using other detection architectures.

6. Mask-head architecture ablations

In Section 4 we illustrated the role of depth in mask-head architectures. Next, we consider mask-head architectural characteristics in some more detail via a series of ablations. All experiments in this section are performed using Deep-MAC at 512x512 resolution with Hourglass backbones.

What makes Hourglass mask-heads so good? We first address the question of which architectural elements are most responsible for the superior generalization of Hourglass networks. To investigate, we focus on 20-layer Hourglass and ResNet basic mask-heads. The Hourglass architecture differs in two main ways from ResNet, having a) an encoder-decoder structure in which the encoder down-samples the input and the decoder upsamples the result of the encoder, and b) long range skip connections connecting feature maps of the same size in the encoder and decoder.

To understand each difference in isolation, we explore the effect of (a) replacing the downsampling/upsampling layers with layers that do not change the feature map resolution, and (b) severing long range skip connections.

Table 7 shows the corresponding ablation results. We see that removing the long range skip connections (No

Backbone	Mask-Head	mIOU		
		Overall	VOC	Non-VOC
HG-52	HG-52	78.43	80.44	77.76
HG-104	ResNet-4	71.43	79.19	68.84

Table 6: Can we reproduce strong mask generalization by adding an hourglass network to the shared backbone instead of using it in the per-proposal mask-head? We compare two networks of similar depth where the first network has a deeper mask-head. For fair comparison, we use groundtruth boxes as input and report mIOU.

LRS) has a small negative impact on the performance. More importantly, we find that the majority of the gap between ResNet and Hourglass is closed by getting rid of the encoder-decoder structure in the Hourglass mask-head (No ED). Given these results, we conclude that this style of downsampling followed by upsampling likely captures particularly appropriate inductive biases for segmentation.

What’s so special about the mask-head? Next, given that an hourglass mask-head offers generalization advantages, we ask: could we reproduce these advantages by adding an hourglass network to the shared backbone instead of using it in the per-proposal mask-head? In other words, what is so special about the mask-head? Here we show that the answer is negative and that it is indeed the mask-head’s architecture that impacts strong mask generalization.

Consider an HG-104 network which is a stack of two hourglass modules each with 52 layers. We compare (a) a model where all 104 layers are situated in the backbone and we use a simple ResNet-4 mask-head versus (b) a model with an HG-52 backbone and HG-52 mask-head. In both cases, inputs undergo roughly 100 layers contained within two hourglass modules but in the second case, the 52 layer mask-head is applied on a per-proposal basis.

Since using 52 layers in the backbone in general yields inferior detection quality compared to the 104 layer backbone, we use groundtruth boxes as input so that both models are on equal footing and we evaluate mIOU.

Our finding (Table 6) is that despite having slightly fewer total layers, our model with the 52 layer mask-head outperforms the model with the 4 layer mask-head by 9% mIOU on unseen classes (both models have similar performance on seen classes). More generally, this supports our hypothesis that within the entire architecture the mask-head plays a disproportionately significant role with respect to generalization to unseen classes.

Is it sufficient to have a large receptive field? Finally, given that depth and encoder/decoder structures do so well, it seems natural to conjecture that increased receptive field in these architectures may play a significant role.

To evaluate this hypothesis, we explore two additional families of mask-heads: (a) we replace the vanilla convolu-

Mask-Head	Variant	Mask mAP		
		Overall	VOC	Non-VOC
ResNet-20	Default	31.4	39.1	28.8
Hourglass-20	Default	34.1	39.8	32.2
	No LRS	33.6	39.2	31.7
	No ED	31.7	39.1	29.2

Table 7: Isolating what makes Hourglass architectures achieve strong mask generalization. No LRS = No long range skip connections. No ED = No encoder-decoder structure, i.e., no downsampling or upsampling layers. Training and evaluation setup is the same as Figure 2.

# Dilated conv layers	Mask mAP		
	Overall	VOC	non-VOC
0	32.2	39.4	29.9
10	32.7	39.1	30.6
20	32.8	39.3	30.7

Table 8: Replacing different numbers of regular convolutional layers with dilated convolutions (rate=2) to isolate the effect of receptive field. For this experiment we use a ResNet-20 mask-head with HG-104 backbone.

tions in a ResNet mask-head with dilated convolutions (w/ rate 2), which has the effect of expanding receptive field without changing the depth or number of parameters, and (b) we use fully connected (MLP) mask-heads which have full receptive field. See Table 8 and Appendix for dilated and FC results, respectively.

Our experiments using both families of models show, first, that none of these models are able to reach the performance of Hourglass mask-heads, so there must be further factors at play beyond receptive field. On the other hand, growing the receptive field early seems to have positive benefits for generalization to some extent (e.g., a shallow FC mask-head seems to outperform a shallow convolution based mask-head).

And this raises an interesting question which we leave for further study: what about receptive field would help unseen classes without simultaneously helping seen classes? Here we offer one conjecture based on our Mask R-CNN finding (Section 5.2), that it is important to train using groundtruth boxes instead of proposals. A groundtruth box, when tight on an instance, acts as a cue, indicating the object that is meant to be segmented. When trained on noisy proposals, we conjecture that Mask R-CNN tries to memorize the types of foreground classes seen at training time and thus struggles to generalize to unseen classes. With a precise cue, however, perhaps the model learns to compare interior pixels to boundary pixels to make this decision, a strategy that is more generalizable across categories and requires a large enough receptive field so that boundary pixels can interact with interior pixels.

Model	B.B.	Mask mAP		
		Overall	VOC	non-VOC
Deep-MAC [R4]	HG104	37.8	42.2	36.3
Mask R-CNN	RF50	36.1	40.2	34.7
Mask R-CNN	SN143	41.9	46.4	40.4

Table 9: Using Deep-MAC generated pseudo labels to train other models. Deep-MAC is trained as described in Table 1 on pseudo labels and evaluated on the `coco-val2017` set. Other models are trained with their default settings. Backbones(B.B.) include HG=Hourglass, RF=ResNet-FPN, SN=SpineNet. R4=ResNet-4 mask-heads. For reference, the “teacher” Deep-MAC model achieves a non-VOC mAP of 35.5% (c.f. Table 2).

7. Using Deep-MAC just for its masks

If, as with the detector-free model from Section 5.3, we run the same mIOU evaluation on the Deep-MAC model using groundtruth boxes, we obtain 81.4% which is slightly better than the detection-free model. To put this number in perspective, [13] showed that COCO groundtruth masks achieve 83%-87% mIOU when compared to expert labels. Thus our finding suggests that remaining headroom on improving segmentation quality is quite limited (we are likely close to a saturation point) and future improvements on the partially supervised task on COCO as measured by mean AP will be much easier to come by via improvements to detection quality as opposed to segmentation quality.

To illustrate, we use Deep-MAC just for its masks (and not its boxes), first segmenting unseen categories and then training a stronger detection model (Mask R-CNN with SpineNet [6], which reaches 48.9% box AP compared to Deep-MAC which reaches 44.1% box AP) in fully supervised mode on these pseudo labels. Table 9 (last row) shows the result of this experiment — specifically, Mask R-CNN with SpineNet is able to leverage the pseudo labels to get to a 40.4% non-VOC mask mAP, which is significantly higher than the original model that generated the pseudo labels. Thus improving box detection quality leads to a significantly increased final non-VOC mAP which is not upper bounded by the non-VOC mAP of Deep-MAC itself. This is also the highest performance ever reported on the partially supervised task by a margin of 6.4% (but only by virtue of better detection and without improving generalization to novel classes).

Our recommendation, consequently, is that the community should focus on harder tasks either by training with even fewer mask annotations, or evaluate partially supervised performance on LVIS [13] which has more classes and higher quality masks. As an initial step, we train Deep-MAC on COCO masks from all 80 categories and evaluate mIOU on LVIS masks (from the `val` set) cropping to LVIS groundtruth boxes. Here our models using ResNet-4 and HG-100 mask-heads achieve 70.3% and 79.9% mIOU respectively, showing that architecture continues to matter

for strong mask generalization even on LVIS. Comparing to [13] who report 90-92% mIOU dataset-to-expert agreement, we also see that there is still a gap between Deep-MAC and human performance (but this is likely at least in part due to COCO’s lower quality masks).

Another application of two stage training is to train a cheaper instance segmentation model on masks produced by Deep-MAC. The first two rows of Table 9 demonstrate results using a cheaper Mask R-CNN model or Deep-MAC model with a shallower (4 layer) mask-head. This experiment is particularly interesting in the case of the student Deep-MAC model with the shallow head since in this two stage setting, the student trains as if it were being fully supervised. According to our findings in Section 4, we should therefore expect it to achieve the same performance as Deep-MAC with the heavier mask-head (which it does, and even exceeds). Thus for COCO categories we are able to leverage the strong mask generalization properties of the heavier mask-head while retaining the computational benefits of the cheaper mask-head. When running at 1024×1024 resolution on a V100 GPU, Deep-MAC with an HG-100 mask-head takes 232 ms per image, whereas the cheaper student model with a ResNet-4 mask-head is faster (201 ms per image). Notably, this cheaper student model is on par with ShapeMask [26] in terms of speed (204 ms) while achieving a 2.1 % improvement on non-VOC mAP.²

8. Conclusions

In this work, we have identified and studied the surprising extent to which the mask-head architecture impacts generalization to unseen categories. Through extensive experiments, we demonstrated the generality of this effect across detection methodologies and backbone networks. And by exploiting this strong mask generalization effect, we established a new state of the art on this problem by a significant margin using a conceptually simple model and we argue that performance on this problem (at least with respect to segmentation quality on COCO) is likely near saturation.

While we have taken initial steps in understanding strong mask generalization, how to better understand the inductive biases encoded within mask-head architectures and how to explain our results theoretically remain important directions. Along these lines, we leave readers with pointers to two papers which have noted similar empirical phenomena where certain architectures generalize effectively to data outside of the training distribution. The Deep Image Priors work [48] similarly observed that Hourglass-style networks seem to automatically capture image level statistics in a natural way without being trained on data. [52] showed that sufficiently deep networks unlock a certain strong generalization behavior. We conjecture that there may be a com-

²Inference speed is not reported in CPMask [9]

mon denominator at play and that exploring these synergies further would be a fruitful area of further research potentially yielding insights useful beyond segmentation.

Acknowledgements

We would like to thank David Ross for thoughtful comments on our draft and Pengchong Jin, Abdullah Rashwan and Xianzhi Du for help with setting up our Mask R-CNN experiments.

References

- [1] Rodrigo Benenson, Stefan Popov, and Vittorio Ferrari. Large-scale interactive object segmentation with human annotators. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 11700–11709, 2019. 3
- [2] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934*, 2020. 3
- [3] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: high quality object detection and instance segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019. 1
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. *arXiv preprint arXiv:2005.12872*, 2020. 3
- [5] Chen Chen, Xianzhi Du, Le Hou, Jaeyoun Kim, Pengchong Jin, Jing Li, Yeqing Li, Abdullah Rashwan, and Hongkun Yu. Tensorflow official model garden, 2020. 12
- [6] Xianzhi Du, Tsung-Yi Lin, Pengchong Jin, Golnaz Ghiasi, Mingxing Tan, Yin Cui, Quoc V Le, and Xiaodan Song. Spinenet: Learning scale-permuted backbone for recognition and localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11592–11601, 2020. 8
- [7] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. Centernet: Keypoint triplets for object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6569–6578, 2019. 3
- [8] Dumitru Erhan, Christian Szegedy, Alexander Toshev, and Dragomir Anguelov. Scalable object detection using deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2147–2154, 2014. 3
- [9] Qi Fan, Lei Ke, Wenjie Pei, Chi-Keung Tang, and Yu-Wing Tai. Commonality-parsing network across shape and appearance for partially supervised instance segmentation. *arXiv preprint arXiv:2007.12387*, 2020. 2, 3, 4, 5, 6, 8
- [10] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015. 3
- [11] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Region-based convolutional networks for accurate object detection and segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 38(1):142–158, 2015. 3
- [12] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, 2010. 12
- [13] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5356–5364, 2019. 3, 8
- [14] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 1, 2, 3, 4, 5, 6, 12
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 3
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 4, 5, 12
- [17] Cheng-Chun Hsu, Kuang-Jui Hsu, Chung-Chi Tsai, Yen-Yu Lin, and Yung-Yu Chuang. Weakly supervised instance segmentation using the bounding box tightness prior. In *Advances in Neural Information Processing Systems*, pages 6586–6597, 2019. 3
- [18] Ronghang Hu, Piotr Dollár, Kaiming He, Trevor Darrell, and Ross Girshick. Learning to segment every thing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4233–4241, 2018. 1, 3, 5, 6
- [19] Jonathan Huang, Vivek Rathod, Chen Sun, Menglong Zhu, Anoop Korattikara, Alireza Fathi, Ian Fischer, Zbigniew Wojna, Yang Song, Sergio Guadarrama, et al. Speed/accuracy trade-offs for modern convolutional object detectors. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7310–7311, 2017. 12
- [20] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. 12
- [21] Hoel Kervadec, Jose Dolz, Shanshan Wang, Eric Granger, and Ismail Ben Ayed. Bounding boxes for weakly supervised segmentation: Global constraints get close to full supervision. *arXiv preprint arXiv:2004.06816*, 2020. 3
- [22] Anna Khoreva, Rodrigo Benenson, Jan Hosang, Matthias Hein, and Bernt Schiele. Simple does it: Weakly supervised instance and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 876–885, 2017. 3
- [23] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 12
- [24] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9404–9413, 2019. 3

- [25] Tao Kong, Fuchun Sun, Huaping Liu, Yuning Jiang, Lei Li, and Jianbo Shi. Foveabox: Beyond anchor-based object detection. *IEEE Transactions on Image Processing*, 29:7389–7398, 2020. 3
- [26] Weicheng Kuo, Anelia Angelova, Jitendra Malik, and Tsung-Yi Lin. Shapemask: Learning to segment novel objects by refining shape priors. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9207–9216, 2019. 2, 3, 4, 5, 8, 13
- [27] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Tom Duerig, et al. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *arXiv preprint arXiv:1811.00982*, 2018. 3
- [28] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 734–750, 2018. 3
- [29] Youngwan Lee and Jongyoul Park. Centermask: Real-time anchor-free instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13906–13915, 2020. 3
- [30] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 3
- [31] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 1, 3
- [32] Rosanne Liu, Joel Lehman, Piero Molino, Felipe Petroski Such, Eric Frank, Alex Sergeev, and Jason Yosinski. An intriguing failing of convolutional neural networks and the coordconv solution. In *Advances in Neural Information Processing Systems*, pages 9605–9616, 2018. 4
- [33] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. Path aggregation network for instance segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8759–8768, 2018. 3
- [34] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016. 3
- [35] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *European conference on computer vision*, pages 483–499. Springer, 2016. 2, 4, 5, 12
- [36] Dim P Papadopoulos, Jasper RR Uijlings, Frank Keller, and Vittorio Ferrari. Extreme clicking for efficient object annotation. In *Proceedings of the IEEE international conference on computer vision*, pages 4930–4939, 2017. 1
- [37] Siyuan Qiao, Liang-Chieh Chen, and Alan Yuille. Detectors: Detecting objects with recursive feature pyramid and switchable atrous convolution. *arXiv preprint arXiv:2006.02334*, 2020. 1, 3
- [38] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016. 3
- [39] Joseph Redmon and Ali Farhadi. Yolo9000: better, faster, stronger. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7263–7271, 2017. 3
- [40] Joseph Redmon and Ali Farhadi. Yolo3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018. 3
- [41] Tal Remez, Jonathan Huang, and Matthew Brown. Learning to segment via cut-and-paste. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 37–52, 2018. 3
- [42] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. 3
- [43] Pierre Sermanet, David Eigen, Xiang Zhang, Michaël Mathieu, Rob Fergus, and Yann LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *arXiv preprint arXiv:1312.6229*, 2013. 3
- [44] Christian Szegedy, Alexander Toshev, and Dumitru Erhan. Deep neural networks for object detection. In *Advances in neural information processing systems*, pages 2553–2561, 2013. 3
- [45] Mingxing Tan, Ruoming Pang, and Quoc V Le. Efficientdet: Scalable and efficient object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10781–10790, 2020. 3
- [46] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 9627–9636, 2019. 3
- [47] Zhi Tian, Chunhua Shen, Xinlong Wang, and Hao Chen. Boxinst: High-performance instance segmentation with box annotations. *arXiv preprint arXiv:2012.02310*, 2020. 3
- [48] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Deep image prior. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9446–9454, 2018. 8
- [49] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2. <https://github.com/facebookresearch/detectron2>, 2019. 12
- [50] Enze Xie, Peize Sun, Xiaoge Song, Wenhai Wang, Xuebo Liu, Ding Liang, Chunhua Shen, and Ping Luo. Polarmask: Single shot instance segmentation with polar representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12193–12202, 2020. 3
- [51] Hui Ying, Zhaojin Huang, Shu Liu, Tianjia Shao, and Kun Zhou. Embedmask: Embedding coupling for one-stage instance segmentation. *arXiv preprint arXiv:1912.01954*, 2019. 3
- [52] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Michael C Mozer, and Yoram Singer. Identity crisis: Memorization and generalization under extreme overparameterization. *arXiv preprint arXiv:1902.04698*, 2019. 5, 8

- [53] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016. [5](#)
- [54] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019. [2](#), [3](#), [12](#), [13](#)
- [55] Xingyi Zhou, Jiacheng Zhuo, and Philipp Krahenbuhl. Bottom-up object detection by grouping extreme and center points. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 850–859, 2019. [3](#), [12](#)
- [56] Chenchen Zhu, Yihui He, and Marios Savvides. Feature selective anchor-free module for single-shot object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 840–849, 2019. [3](#)

Appendix

A. Experimental details

A.1. Deep-MAC

Architecture settings. We use a pixel embedding layer with 16 channels and an instance embedding layer with 32 channels. For mask head architectures, we experiment with Hourglass networks [35] and Residual Networks [16] (using both basic and bottleneck variants). The number of channels in the first layer of all mask heads is set to 64 and increased gradually in the deeper layers (see Section F). The number of convolution layers of each dimensionality is kept roughly similar between mask heads of similar depth.³

Training details and hyperparameters. We train using Adam[23] with synchronized batch normalization [20] with batch size 128 for 50K steps and apply a loss weight of 5 to the mask loss. We use a learning rate of 10^{-3} which is cosine decayed to 0 after linear warm-up from 2.5×10^{-4} for 5K steps. All detection related hyperparameters are unchanged from CenterNet defaults as described in [54] including data augmentation. We do not use test time augmentation. Finally, models trained at 512×512 resolution are initialized from an ExtremeNet [55] checkpoint, whereas models trained at 1024×1024 resolution are initialized from a COCO (detection) trained checkpoint.

Implementation details Deep-MAC is built on top of the publicly available CenterNet implementation in the Tensorflow Object Detection API [19]. For our best results as well as most ablations, we use an Hourglass-104 [35] backbone, though we show that our findings hold when using ResNet-FPN backbones as well. Since the Stacked Hourglass model produces predictions at the end of multiple Hourglass modules, we follow the common approach of applying prediction heads and losses at the end of each such module, using only predictions from the final stage at test time.

A.2. Mask R-CNN

For all Mask R-CNN experiments, we train at an input resolution of 1024×1024 and follow the 3x learning rate schedule as outlined in [49]. We do not use any test-time augmentations. For the Hourglass-20 and Hourglass-52 mask-heads, we use RoI crops of sizes 16×16 and 32×32 respectively, to accommodate the successive downsampling and upsampling operations in these architectures. Otherwise, we use the default Mask R-CNN [14] RoI crop size of 14. The output resolution of the mask is always fixed to be $2 \times$ the resolution of the input RoI crop. We verified that

³Sometimes depths cannot be exactly matched between different mask heads.

Mask Head	C	I	Mask mAP		
			Overall	VOC	Non-VOC
ResNet-20			–	–	–
		✓	–	–	–
	✓	✓	30.9	39.1	28.2
Hourglass-20	✓	✓	31.4	39.1	28.8
			34.1	39.8	32.2
		✓	34.5	39.9	32.7
	✓		33.6	39.9	31.5
	✓	✓	34.3	39.8	32.5

Table 10: Effect of Coordinate Embedding (C) and Instance Embedding (I) on the generalization ability of Deep-MAC on unseen classes. A ‘–’ indicates that the model failed to converge. All models are trained with masks only from VOC classes at an input image resolution of 512×512 . Performance is reported on the coco-val2017 set.

training with only groundtruth boxes does not impact fully supervised performance (See Table 12). We implement our Mask R-CNN variant in the TF-vision code base [5].

B. Additional Ablations

B.1. Deep-MAC

B.1.1 Effect of instance and coordinate embedding

Table 10 shows the effects of coordinate embedding and instance embedding on ResNet and Hourglass mask heads. We notice that the additional embeddings do not make a significant difference to the Hourglass model, but coordinate embedding is required for the ResNet based mask heads to converge. For uniformity, we have thus used both components in all Deep-MAC variants.

B.1.2 Effect of using fully connected layers

See Table 11 for experiments with fully connected layers. We used Glorot normal initialization [12] the mask-head weights. Based on these results, we see that the fully connected mask-head models, which have full receptive field with respect to the input tensor, do not offer competitive performance compared to the HG-based mask-heads. However, early large receptive fields may still be beneficial to some extent as these fully connected mask-heads do outperform our shallowest convolution-only mask-heads (e.g. Resnet-4).

B.2. Mask R-CNN

Table 12 shows the impact of using groundtruth boxes (instead using proposals, which is the standard approach) for training the mask-head of a fully supervised Mask R-CNN model on COCO. First we see that using a class-agnostic mask head results in a slightly lower mask mAP

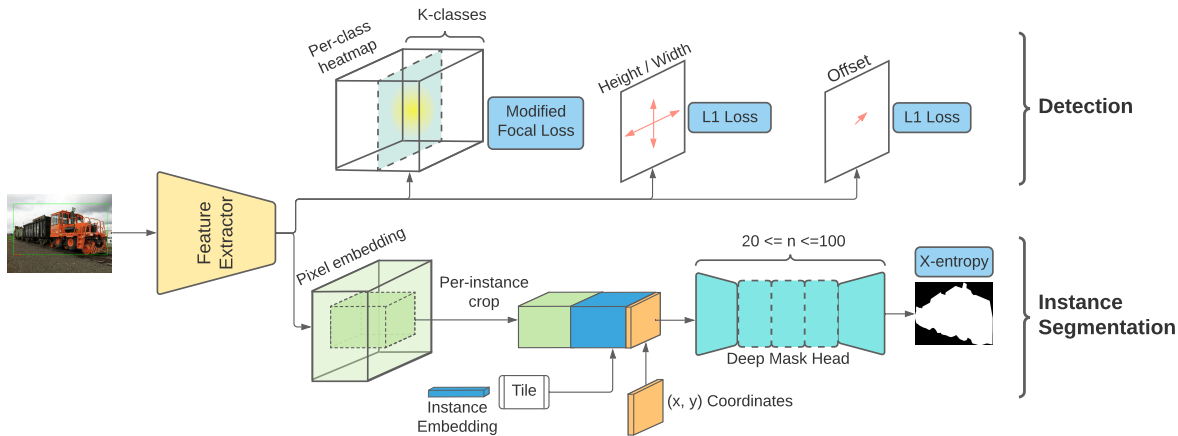


Figure 3: Schematic of the Deep-MAC architecture. The top-half is kept identical to CenterNet [54] and the bottom-half uses an RoI crop followed by a deep mask head. In our experiments, it was crucial to train the mask head with only groundtruth boxes.

FCN layers	Mask mAP		
	Overall	VOC	Non-VOC
2	29.1	38.4	26.0
4	30.5	37.5	28.2

Table 11: Effect of using fully connected layers as mask heads. All models are trained with masks only from VOC classes at an input image resolution of 512×512 . Performance is reported on the `coco-val2017` set. For easy reference, the VOC/non-VOC mask mAP values for Resnet-4 and HG-52 mask-heads are 39.7/26.6 and 39.8/32.5 respectively.

Variant	Mask mAP
Class-specific (standard)	37.2
Class-agnostic	36.7
Class-agnostic + GT boxes	36.4

Table 12: Fully supervised mask mAP of Mask-RCNN variants with a ResNet-50-FPN backbone.

compared to the standard class-specific mask-head. Training with groundtruth boxes instead of proposals does not further impact the performance of the class agnostic mask head significantly.

C. Generalizing from a single class

In the majority of our experiments, we assumed the standard setup of “train-on-VOC, test-on-non-VOC”. In this section, we restrict further, training on a single “source” class at a time, in order to better understand when Deep-MAC can be expected to strongly generalize to a novel class. In Figure 4a we plot results from this experiment, training on each of the VOC categories with 512×512 resolution inputs and an Hourglass-52 mask head. We observe that while some classes lead to strong performance, there is high variance depending on the source category (ranging

from 12.5% mAP to 27.8% mAP). Notably, a single class can achieve strong results — as one datapoint, training only on the chair category with higher resolution 1024×1024 inputs yields a non-VOC mask mAP of 31.5, which is competitive with previous high-performance methods (e.g., ShapeMask [26]) when trained on all VOC categories.

In some cases it is easy to guess why a class might be a poor source — on the worst classes, we see that the quality of groundtruth masks is uneven in COCO. For example, labels were not consistent about excluding objects that were on but not part of a dining table (see Section E).

For more detailed insight, we ask how training on a specific source class might generalize to a specific new target class. For source-target pairs (i, j) , Figure 4b visualizes this relationship via the ratio between mAP on target class j if we were to train on just the source class i to mAP on target class j if we were to train on all classes. Here we cluster the rows and columns by similarity and truncate ratios to be at most 1.0 for visualization purposes.⁴

Figure 4b illustrates that some classes (e.g., apple [52], umbrella [45], stop sign [42]) are universally easy transfer targets likely due to being visually salient, having consistent appearance and not typically co-occurring with other examples of their own class. We also see that co-occurrence of source and target classes does not always lead to improved ratios (i.e. close to 1). For example, training on car does not yield strong performance on stop signs [42] or parking meters [43] and training on person does not yield strong performance on bench [01] or baseball bat [09]. On the other hand, categories that are similar semantically seem to function similarly as source categories, and with a few exceptions, the source categories cluster naturally into two broad groups: man-made and natural objects.

⁴It is worth noting in several cases (most notably, hair drier [34]) that it is better to train on other source classes than it is to include the target class annotations during training.

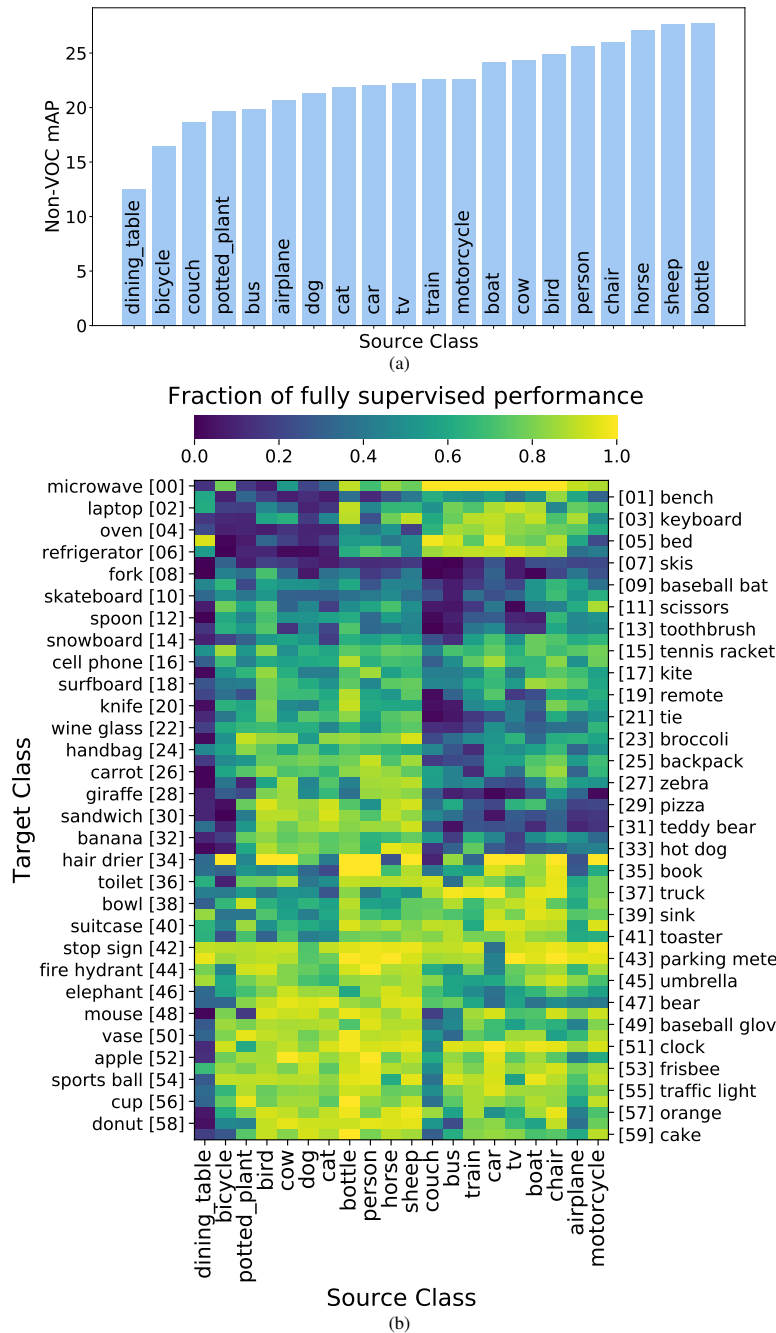


Figure 4: (a) Mask mAP on Non-VOC classes when training with masks from only a single source class from the VOC set; (b) Performance on specific (Non-VOC) target classes when training with masks from only a single class. We visualize performance relative to full supervision.

It remains an open question why a class might excel as a source class in general. Intuitively one might think that person, car or chair categories might be the best because they have the most annotations and are visually diverse, but perhaps surprisingly, using the bottle category is the best. This may be due to the fact that bottles tend to look alike and appear in groups, forcing the model to make non-local decisions about mask boundaries. We leave exploration of

this hypothesis for future work.

D. Example images on unknown classes

See Figure 5 for example outputs of Deep-MAC. We look at the output of our model with user-specified boxes around object categories that are not in the COCO dataset. We observe that Deep-MAC generalizes to multiple differ-

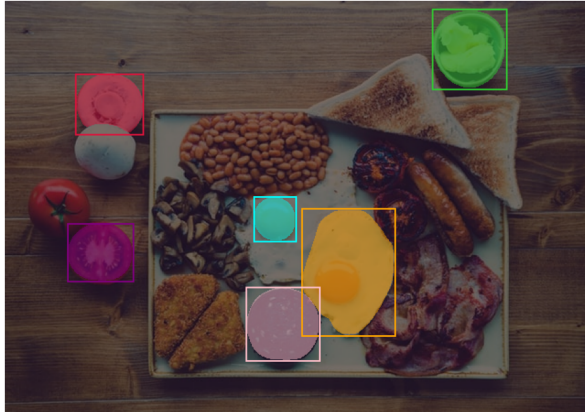
ent domains like biological and camera trap images and does well even in cluttered settings. For this experiment, we used a model trained with all COCO classes in fully supervised mode.

E. Looking at annotation quality

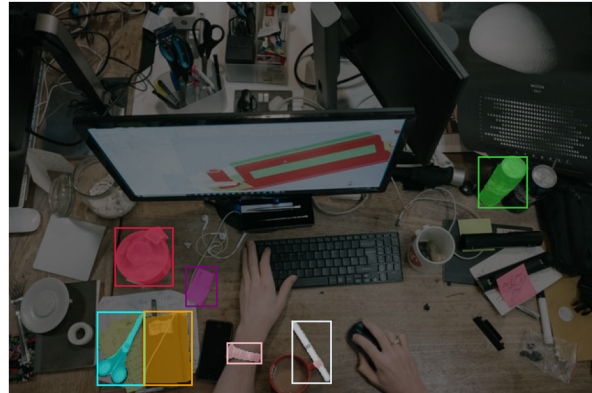
In Figure 6 we show examples of COCO groundtruth annotations from the dining table, bicycle and potted plant categories, the worst three categories to use as source training categories. The examples illustrate the inconsistencies/inaccuracies in mask annotations for these categories — for example, annotators were inconsistent about including or excluding objects on the dining tables.

F. Mask head architecture details

Details of mask head architectures can be found in Table 13, 14, 15 and 16. Figure 7 illustrates the computation graph of an hourglass mask head.



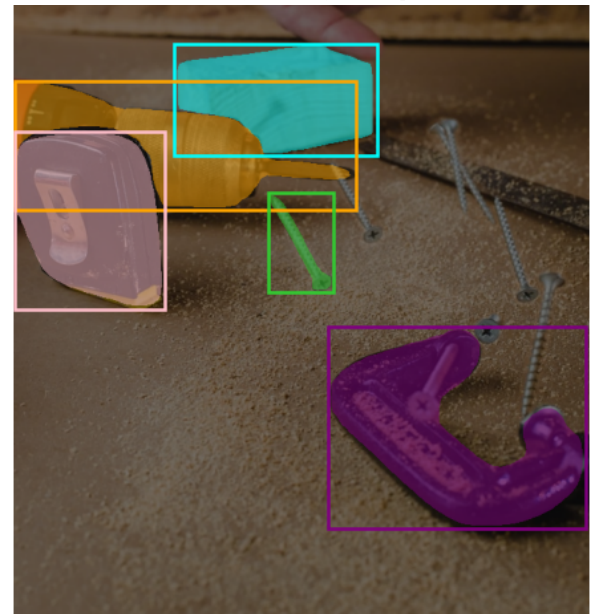
(a) Photo by Jonathan Farber on Unsplash.



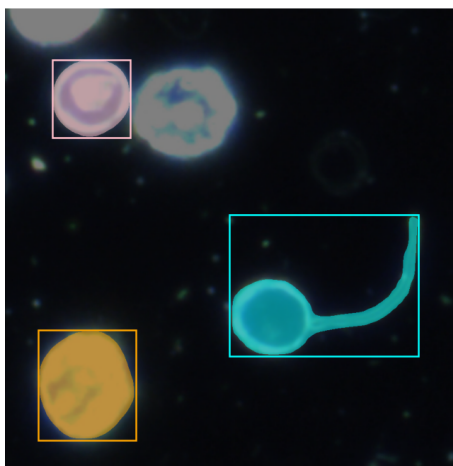
(b) Photo by Robert Bye on Unsplash.



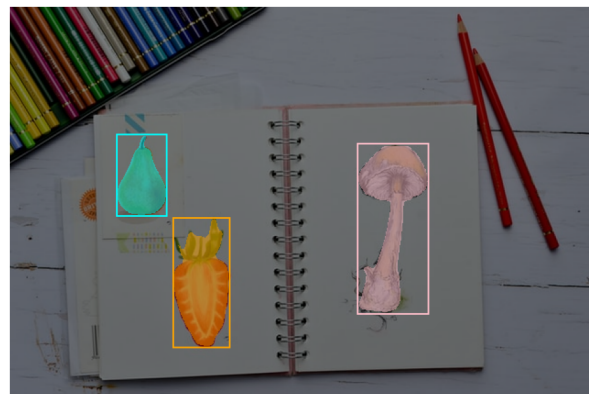
(c) Sample from the Snapshot Serengeti dataset.



(d) Photo by Chris Briggs on Unsplash.

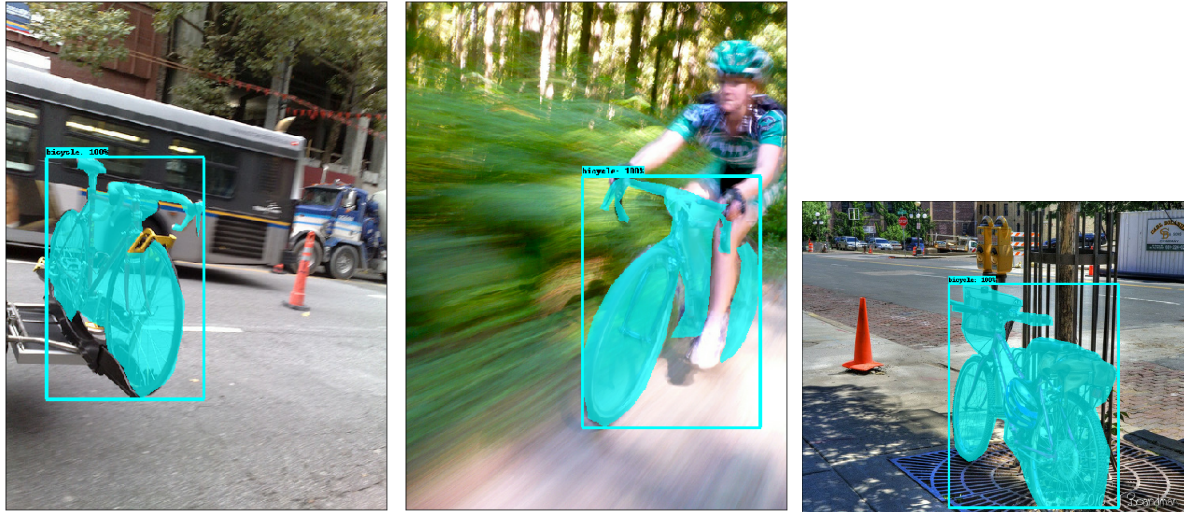


(e) Sample from the Bacteria detection with darkfield microscopy competition on Kaggle.



(f) Photo by Maggie Jaszowska on Unsplash.

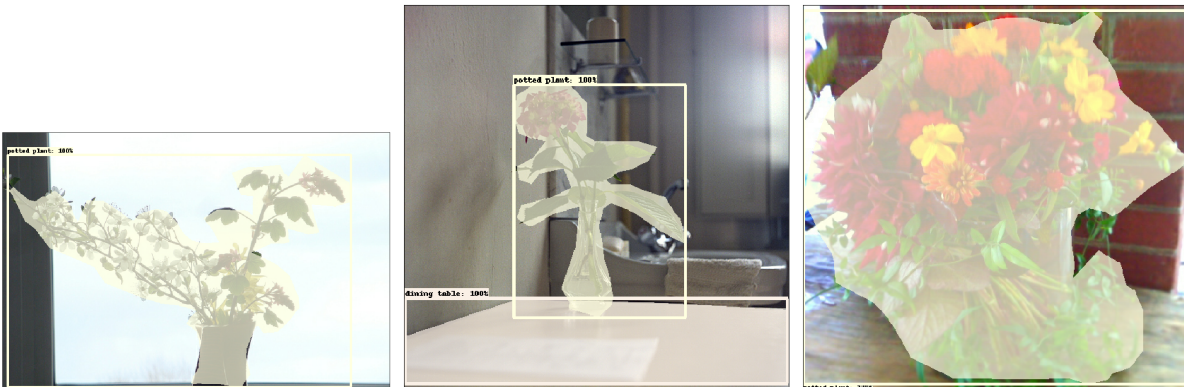
Figure 5: Example outputs of Deep-MAC with hand-drawn boxes on unknown classes.



(a) Bicycle: The annotated masks don't capture the shape correctly, and quite often label parts of the background interspersed with the bicycle frame as foreground.



(b) Dining Table: Inconsistencies in annotated parts of dining tables. Left: Plate and cup with carrots is excluded whereas plate with empty glass is included in the mask. Center: Some glasses on the dining table are included as part of it whereas some classes aren't. Right: Chairs are excluded from the dining table mask in the dining tables near the bottom, whereas they are included in the dining table masks near the top.



(c) Potted plant: Areas of background are included in the foreground masks of potted plants, especially near the leaves.

Figure 6: Example annotations of the 3 worst source classes to train on.

Type	Depth	# of Blocks	Conv Block		
			Size	Channels	
ResNet	4	1	32×32	64	
		2	32×32 32×32	128 128	
	8	1	32×32	64	
		4	32×32 32×32	128 128	
	12	1	32×32	64	
		6	32×32 32×32	128 128	
	16	1	32×32	64	
		8	32×32 32×32	128 128	
	20	1	32×32	64	
		8	32×32 32×32	128 128	
		2	32×32 32×32	128 128	
	ResNet Bottleneck	6	1	32×32	64
			2	32×32 32×32 32×32	128 512 128
			9	1	32×32
3				32×32 32×32 32×32	128 512 128
12		1	32×32	64	
		4	32×32 32×32 32×32	128 512 128	
		15	1	32×32	64
5			32×32 32×32 32×32	128 512 128	
21			1	32×32	64
		6	32×32 32×32 32×32	128 512 128	
		1	32×32	192	
			32×32	384	
			32×32	192	

Table 13: Architecture details of ResNet and ResNet bottleneck mask heads.

Type	Depth	# of Blocks	Conv Block	
			Size	Channels
ResNet Bottleneck [1/4 th]	6	1	32 × 32	16
		2	32 × 32	32
			32 × 32	128
	32 × 32		32	
	12	1	32 × 32	16
		4	32 × 32	32
			32 × 32	128
	32 × 32		32	
	21	1	32 × 32	16
		6	32 × 32	32
			32 × 32	128
	32 × 32		32	
	30	1	32 × 32	16
		5	32 × 32	32
			32 × 32	128
	32 × 32		32	
	30	5	32 × 32	48
			32 × 32	192
			32 × 32	48
	51	1	32 × 32	16
5		32 × 32	32	
		32 × 32	128	
	32 × 32	32		
51	5	32 × 32	48	
		32 × 32	192	
		32 × 32	48	
51	1	32 × 32	16	
	6	32 × 32	32	
		32 × 32	128	
32 × 32		32		
51	8	32 × 32	48	
		32 × 32	192	
		32 × 32	48	
51	3	32 × 32	64	
		32 × 32	256	
		32 × 32	64	

Table 14: Architecture details of ResNet bottleneck [1/4th] mask head.

Type	Depth	# of Blocks	Conv Block	
			Size	Channels
Hourglass	10	1	32 × 32	64
		3	32 × 32	128
			32 × 32	128
		1	16 × 16	128
			16 × 16	128
	1	32 × 32	128	
	20	1	32 × 32	64
		3	32 × 32	128
			32 × 32	128
		4	16 × 16	128
			16 × 16	128
	2	8 × 8	192	
		8 × 8	192	
	1	32 × 32	128	
	32	1	32 × 32	64
		5	32 × 32	128
			32 × 32	128
		4	16 × 16	128
			16 × 16	128
	4	8 × 8	192	
		8 × 8	192	
2	4 × 4	192		
	4 × 4	192		
1	32 × 32	128		
52	1	32 × 32	64	
	5	32 × 32	128	
		32 × 32	128	
	4	16 × 16	128	
		16 × 16	128	
	4	8 × 8	192	
		8 × 8	192	
	4	4 × 4	192	
		4 × 4	192	
	4	2 × 2	192	
2 × 2		192		
4	1 × 1	256		
	1 × 1	256		
1	32 × 32	128		

Table 15: Architecture details of Hourglass mask head (Part 1 of 2).

Type	Depth	# of Blocks	Conv Block	
			Size	Channels
Hourglass	100	1	32×32	64
		9	32×32	128
			32×32	128
		8	16×16	128
			16×16	128
		8	8×8	192
			8×8	192
		8	4×4	192
			4×4	192
		8	2×2	192
			2×2	192
		8	1×1	256
			1×1	256
1	32×32	128		

Table 16: Architecture details of Hourglass mask head (Part 2 of 2).

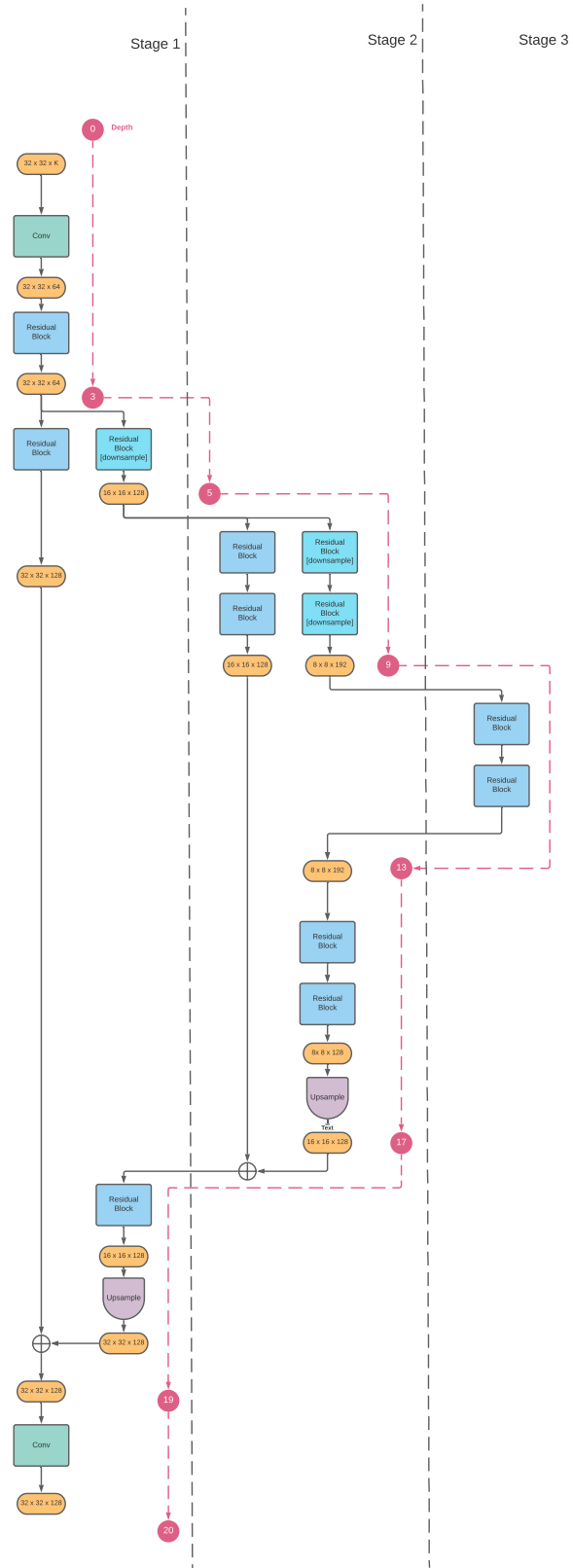


Figure 7: Illustration of the Hourglass 20 mask head computation graph.