

---

# [Re] Panoptic-DeepLab: A Simple, Strong, and Fast Baseline for Bottom-Up Panoptic Segmentation

---

Anonymous Author(s)

Affiliation

Address

email

## Reproducibility Summary

1

### 2 **Scope of Reproducibility**

3 The original work by Cheng et al. [5] introduces Panoptic-DeepLab - a novel architecture for panoptic segmentation,  
4 claiming to achieve comparable performance to two-stage, top down approaches while yielding fast inference speeds.  
5 At the time of publication, Panoptic-Deeplab claims to have ranked first in all three cityscapes benchmarks (*specifically:*  
6 *mIoU, AP & PQ*).

### 7 **Methodology**

8 As the original paper authors published their source code, our codebase integrates sections of their codebase, while  
9 re-implementing components intrinsic to the main claim we are attempting to evaluate. We also studied the source code,  
10 using information provided from it and the pipelines to augment our understanding from what the paper described.

11 While we initially attempted a code-blind reproduction, it was soon determined to be unfeasible following which a  
12 hybrid approach was instantiated.

### 13 **Results**

14 While we successfully reproduced the given architecture, we have been unable to train it. Therefore: our contributions  
15 currently remaining exclusive to architecture, and certain unit tests within the system itself. We also highlight potential  
16 low-level Tensorflow that were pitfalls to our development, that may be advantageous to investigate.

### 17 **What was easy**

18 The authors of the paper structured their contributions on well-documented and tested frameworks such as ResNet and  
19 DeepLabV3+, while training on popular datasets such as Cityscapes and Mapillary Vistas. Consequently, setting up the  
20 dataset and the environment to reproduce the given research was straightforward.

### 21 **What was difficult**

22 A significant hurdle we came across during our reading of the paper was vagueness within the expected implementation.  
23 This extended from the architecture to the training regime. The descriptions provided, although accurate, were presented  
24 as a high-level overview, with the expectation of a lot of prior domain knowledge. This resulted in a significant time-sink,  
25 following which we looked into the codebase for necessitated context.

26 Despite the well-structured objected oriented implementation through which the code was written, we found certain  
27 sections hard to understand. We observed convoluted re-implementations of high level functions already part of  
28 Tensorflow as part of the codebase. However, this could have been a direct result of the implementation not using the  
29 now-popularised Functional API within Tensorflow, which may have resulted in the required use of custom layers.

30 **Communication with original authors**

31 We communicated with the authors over e-mail, resolving doubts that arose while reading the paper. It is also through  
32 author communication that we were directed to the codebase, as although public - the relevant repository wasn't  
33 mentioned within the paper.

## 34 **1 Introduction**

35 Since its inception in 2018, Panoptic Segmentation[10] has remained a popular task within the domain of computer  
36 vision. It is in effect the unification of two distinct yet related tasks, namely: semantic and instance segmentation.  
37 Semantic segmentation broadly involves the assignment of a class ID to every input pixel, whereas instance segmentation  
38 is the delineation of distinct objects within an input frame. Broadly classified as “stuff” and “things”, the unification of  
39 the two produces the target output known as Panoptic Segmentation.

40 Panoptic-Deeplab[5] aims to establish a strong baseline for a bottom-up approach to the task. Consequently, it places  
41 a focus on simplicity, cleverly incorporating established components within neural architecture to set state-of-art  
42 benchmarks as of the date of publication.

## 43 **2 Scope of reproducibility**

44 We investigate the following claims from the original paper:

- 45 • Panoptic-DeepLab establishes a solid baseline for bottom-up methods that can achieve comparable performance  
46 of two-stage methods while yielding fast inference speed - nearly real-time on the MobileNetV3 backbone.
- 47 • Single Panoptic-DeepLab simultaneously ranks first (at the time of publication) at all three Cityscapes  
48 benchmarks, setting the new state-of-art of 84.2% mIoU, 39.0% AP, and 65.5% PQ on test set.

## 49 **3 Methodology**

50 Initially, we attempted a code-blind reproduction of Panoptic-DeepLab. However, we swiftly determined it to be  
51 unfeasible - primarily as a result of us being unable to fully grasp implementation details from the paper itself. The  
52 paper does incredibly well to provide a high level explanation of how the architecture functions; unfortunately, the lack  
53 of implementation-specific information prevented a blind-paper reproduction without extensive interpolation. Crucially:  
54 we note the importance of a standardized system for presenting architecture diagram. While the current abstract layers  
55 look nicer, we find they lack important information necessary to reproduction.

56 It is important to note here that upon re-reading the paper post implementation - with a prior understanding of the  
57 architecture - we found that just the paper did very well to explain the architecture, enough even, for a code-blind  
58 reproduction. Going through long-expired threads of discussion was an exercise that did well to remind us of implicit  
59 interpolations we made, having already known the architecture.

### 60 **3.1 Model description**

61 Panoptic-DeepLab[5] incorporates an encoder-decoder architecture to generate target inference, with our implementation  
62 encapsulating 6,547,894 total parameters, of which 6,534,998 are trainable, while the remaining 12,896 are non-  
63 trainable. Broadly - it sequentially incorporates the modules discussed in the following subsections.

#### 64 **3.1.1 Image DataGenerator**

65 To the extent of our understanding, Panoptic-Deeplab[5] did not discuss the implementation of its dataset loaders. As  
66 a result, we entirely used a custom implementation of Tensorflow’s ImageDataGenerator[1] class, to function as an  
67 iterator for the training regime. Since we did not find it highlighted within the paper to generate ground-truth center  
68 heatmaps and centerpoint predictions, we discuss this in the following paragraph.

69 **Center Heatmaps & Prediction[13]** The center heatmaps & prediction maps are representations of the ground truth  
70 instance ID images. These images are effective data representations of instances within the frame. Each ‘thing’ has an  
71 encoded value, for instance: each pixel representing car#1 may be labeled 10001, while car#2 is labelled 10002. The  
72 first two digits encode one of the 19 different objects tracked by Cityscapes - in this case, the car - while the final three  
73 digits refer to the instance of the given object. The representation in specific are the computed averages of each of the  
74 instances - producing the center prediction. The center heatmaps are a gaussian distribution applied over the centerpoint  
75 predictions with *standard deviation* =  $8px$ .

### 76 3.1.2 Encoder

77 Panoptic-DeepLab is trained on three popular encoder ImageNet pre-trained backbones, namely: Xception-71[6],  
78 ResNet-50[8] & MobileNetV3[9]. The backbone works to generate feature maps from input images. For the purpose  
79 of this reproduction, we use Xception-71 as our encoder backbone, as this is the primary implementation used by the  
80 original authors. We integrate our own implementation of the Xception-71 module as part of the paper reproduction.

### 81 3.1.3 Atrous Spatial Pyramid Pooling

82 From the encoder, the feature maps are split into dual modules. The first layer to run the decoupled modules is Atrous  
83 Spatial Pyramid Pooling[4], abbreviated - ASPP, is a module that concurrently resamples encoded feature layers at  
84 different rates, finally pooled together to capture objects and relevant context at multiple scales.

85 We derived the ASPP block directly from the tensorflow implementation maintained by the paper authors, with no  
86 modifications made to the architecture.

### 87 3.1.4 Decoder

88 Panoptic-DeepLab is a fork of the DeepLabV3+[4] decoder architecture. It incorporates two fundamental contributions,  
89 specifically: an additional skip connection in the upsampling stage, and an additional upsampling layer with output  
90 `stride = 8`. We developed a custom implementation of this utilizing the modern Keras Functional[2] API. Through  
91 our development of the decoder, we ran into a prominent problem, that delayed significantly our progress within model  
92 architecture. This is in direct correlation with how Tensorflow handles internal API calls, type conversion.

93 **tf.Tensor v KerasTensor** `KerasTensor` is an internal class within the Keras API. It is generated during layer  
94 definition, during the construction of a neural architecture. When latent features are passed during the function calls, the  
95 `KerasTensor` object is converted implicitly to the `tf.Tensor` format - covering up significant type discrepancies. As  
96 part of testing the original Panoptic-Deeplab code, we evaluated that as part of the model conversion to the Functional  
97 API, it was unable to retrace inputs to the decoder. This resulted in a Graph Disconnected error. In an attempt to  
98 allow traceback to work, we devised an approach wherein skip connections were made instance variables within the  
99 Decoder class, and passed separately to the functional call. It is here that we discovered that the lack of the implicit type  
100 conversion, while transferring precisely the same set of data resulted in a `TypeError`. We were unable to manually make  
101 the necessary conversion, highlighting a lack of documentation as `KerasTensor` is a backend class. Consequently, we  
102 were unable to patch the approach and proceeded to a full rewrite.

103 **Graph Disconnected** An error we struggled to get past - the Graph Disconnected error is thrown when the traceback  
104 method within the functional API is unable to generate the necessary I/O graph to create a valid architecture. While in  
105 retrospect: the information provided was enough to debug effectively the point of failure, we would like to highlight  
106 that we believe a more visual or verbose representation - for instance, a plot describing the graph upto the point of  
107 failure - may allow the quicker & clearer identification of the issue.

### 108 3.1.5 Prediction Heads

109 The decoupled decoder modules further split into three separate prediction heads. These generate the final deep-learning  
110 based output within our implementation. They are a final set of convolutional followed by fully connected layers  
111 generating the final result.

112 Similar to ASPP[4], we derived prediction heads directly from the tensorflow implementation maintained by the paper  
113 authors, with no modifications made to the architecture.

### 114 3.1.6 Loss Function

Panoptic-DeepLab employs a collective loss function intended to train resultant outputs.

$$L = \lambda_{sem}L_{sem} + \lambda_{heatmap}L_{heatmap} + \lambda_{offset}L_{offset}$$

115 This was a straightforward function, the implementation of which was just as straightforward, and did not require any  
116 effort above the requisite minimum.

### 117 3.1.7 Post Processing

118 Post processing of the outputs heads in effect involves stitching the instance and semantic segmentation outputs via a  
119 majority vote, generating the final panoptic segmentation. Since output post processing involves a traditional script with  
120 no trainable parameters, we have used post-processing code directly from the original tensorflow implementation, as  
121 put forward by the authors of the paper.

## 122 3.2 Datasets

123 Panoptic-DeepLab used Cityscapes[7], Mapillary Vistas[12] & COCO[11] datasets over the proposed architecture. For  
124 the purpose of our implementation, we train our model on the Cityscapes dataset, as examples are referenced from it  
125 through the evaluation stages of the model. Each image is of size (1025, 2049), and utilizes an odd crop size to allow  
126 centering, aligning features across spatial resolutions.

## 127 3.3 Hyperparameters

128 Panoptic-DeepLab uses a training protocol similar to that of the original DeepLab, specifically: the ‘poly’ learning  
129 rate policy. It uses the Adam optimizer with a learning rate of 0.001 without weight decay, with fine-tuned batch  
130 normalization parameters and random scale data augmentation. While we prepared our re-implementation with the  
131 same set of hyperparameters, we were unable to validate our approach, further discussed in Section 3.5.

## 132 3.4 Experimental setup and code

133 Alongside git for code tracking, we also employ data science specific tools such as DVC (Data Version Control) and  
134 MLFlow[3] with DAGsHub as the platform operating the relevant stack of services. DVC requires S3 buckets, that  
135 maintain the dataset, models, visualization and high storage binaries utilized during training. MLFlow - specifically,  
136 MLFlow tracking was the service we utilized as part of documenting the training lifecycle, including experimentation,  
137 and the relevant comparison between training cycles.

## 138 3.5 Computational requirements

139 By an astronomical margin, the computational requirements necessary for training Panoptic-DeepLab was the factor  
140 that prevented us from successfully testing our target reproduction. Originally, the architecture was trained on a cluster  
141 of 32 TPUs. In a technical report that detailed a PyTorch re-implementation of Panoptic-DeepLab, they coupled runtime  
142 optimization techniques alongside smaller batch size to reduce the training size to 4-8 GPUs. While a significant  
143 improvement, we find that stating it enables ‘everyone be able to reproduce state-of-the-art results with limited resources’  
144 a vast extension.

145 The computational stack under active access to our team includes a single GPU on a docker container, personal  
146 workstations as well as any GPUs provisioned by cloud notebook service *Google Colaboratory*. Even considering the  
147 use of cloud compute services such as AWS - that are estimated to cost upwards of 2,000 USD - for the acquisition of  
148 necessary compute, it is not possible to acquire access to the high performance GPU-enabled G3 instances without  
149 explicit approval from AWS customer support. Through a back-and-forth that extended across weeks, we have been  
150 unable to acquire the approval necessary to create stated instances.

151 We therefore attempted the utilization of CPU resources to train the model to the best of our ability. We theorized the  
152 use of high learning rates in an attempt to overfit the model in a single epoch as a sanity check; to ensure the pipeline  
153 for our re-implementation worked as intended. Predictably, the training failed, and python was killed as the memory  
154 usage exceeded the cap permitted by the system, causing it to crash.

## 155 4 Results

156 As a result of the scenario detailed in the previous section: while we did manage to reproduce the architecture, we  
157 have been - as of now - unable to train it. Therefore, to this degree, our reproduction has not been a success, with  
158 our contributions currently remaining exclusive to architecture and the challenges encountered by us through our  
159 reproduction of the paper.

## 160 5 Discussion

161 Through the constant cycle of updates across which the languages on which neural architectures are written, the  
162 Reproducibility Challenge presents the fantastic opportunity to (1) take a step back, and (2) re-approach a pre-existing  
163 codebases with an entirely different perspective. It allows us the opportunity to fine-tune both past research and research  
164 in the near future. The insights our team has generated from our work on Panoptic-DeepLab itself, has done immensely  
165 to broaden our own perspective on the state of our field at the moment.

### 166 5.1 What was easy

167 The authors of the paper structured their contributions on well-documented frameworks such as ResNet and  
168 DeepLabV3+, while training on popular datasets such as Cityscapes and Mapillary Vistas. Consequently, setting up the  
169 dataset and the environment to reproduce the given research was straightforward.

170 Additionally, various modules within the architecture were concisely and concretely defined - which enabled us to  
171 re-implement them without additional effort, above the minimum requisite. We found several sections of the paper were  
172 written with meticulous detail, and we especially appreciated the exhaustive, vast array of experiments and benchmarks  
173 provided as part of the research, which led our primary motivations towards attempting the reproduction.

### 174 5.2 What was difficult

175 A significant hurdle we came across during our reading of the paper was vagueness within the expected implementation.  
176 This extended from the architecture to the training regime. The descriptions provided, although accurate, were presented  
177 as a high-level overview, with the expectation of a lot of prior domain knowledge. This resulted in a significant time-sink,  
178 following which we looked into the codebase for necessitated context.

179 Despite the well-structured objected oriented implementation through which the code was written, we found certain  
180 sections hard to understand. We observed convoluted re-implementations of high level functions already part of  
181 Tensorflow as part of the codebase. However, this could have been a direct result of the implementation not using the  
182 now-popularised Functional API within Tensorflow, which may have resulted in the required use of custom layers.

183 Additionally, we would also like to highlight the importance of excessive computational requirements within the  
184 machine learning space, and it's relation to the reproducibility of a paper. With the exploding costs of GPUs owing  
185 to extensive crypto-mining farms[14], and the ever increasing complexity of models being trained over time, it is  
186 imperative to consider designing systems that adhere to development policies ranging beyond the best-funded labs, and  
187 represents an important milestone within the democratization of research within high-throughput deep learning.

### 188 5.3 Communication with original authors

189 We enjoyed minimal yet significant communication with the original authors of the research. We communicated over  
190 e-mail, resolving doubts we came across as we read the paper. We found valuable insight through this communication,  
191 which has consequently been imperative to the success of our project. It has enabled discovering an additional suite of  
192 supplementary literature written with respect to the target architecture, which we may have potentially been unable to  
193 find without significant delay.

## 194 References

- 195 [1] Shervine Amidi Afshine Amidi. *A detailed example of how to use data generators with Keras*. <https://github.com/afshinea/keras-data-generator/>. 2018.
- 196 [2] Ekaba Bisong. "Tensorflow 2.0 and keras". In: *Building Machine Learning and Deep Learning Models on Google Cloud Platform*. Springer, 2019, pp. 347–399.
- 197 [3] Andrew Chen et al. "Developments in MLflow: A System to Accelerate the Machine Learning Lifecycle". In: *Proceedings of the Fourth International Workshop on Data Management for End-to-End Machine Learning (2020)*.
- 198 [4] Liang-Chieh Chen et al. "DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 40 (2018), pp. 834–848.
- 199
- 200
- 201
- 202
- 203
- 204

- 205 [5] Bowen Cheng et al. “Panoptic-DeepLab: A Simple, Strong, and Fast Baseline for Bottom-Up Panoptic Seg-  
206 mentation”. In: *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2020)*,  
207 pp. 12472–12482.
- 208 [6] François Chollet. “Xception: Deep Learning with Depthwise Separable Convolutions”. In: *2017 IEEE Conference*  
209 *on Computer Vision and Pattern Recognition (CVPR) (2017)*, pp. 1800–1807.
- 210 [7] Marius Cordts et al. “The Cityscapes Dataset for Semantic Urban Scene Understanding”. In: *2016 IEEE*  
211 *Conference on Computer Vision and Pattern Recognition (CVPR) (2016)*, pp. 3213–3223.
- 212 [8] Kaiming He et al. “Deep Residual Learning for Image Recognition”. In: *2016 IEEE Conference on Computer*  
213 *Vision and Pattern Recognition (CVPR) (2016)*, pp. 770–778.
- 214 [9] Andrew G. Howard et al. “MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications”.  
215 In: *ArXiv abs/1704.04861 (2017)*.
- 216 [10] Alexander Kirillov et al. “Panoptic Segmentation”. In: *2019 IEEE/CVF Conference on Computer Vision and*  
217 *Pattern Recognition (CVPR) (2019)*, pp. 9396–9405.
- 218 [11] Tsung-Yi Lin et al. “Microsoft COCO: Common Objects in Context”. In: *ECCV*. 2014.
- 219 [12] Gerhard Neuhold et al. “The Mapillary Vistas Dataset for Semantic Understanding of Street Scenes”. In: *2017*  
220 *IEEE International Conference on Computer Vision (ICCV) (2017)*, pp. 5000–5009.
- 221 [13] Jonathan Tompson et al. “Joint Training of a Convolutional Network and a Graphical Model for Human Pose  
222 Estimation”. In: *NIPS*. 2014.
- 223 [14] Linus Wilson. “GPU Prices and Cryptocurrency Returns”. In: *ERN: Technology (Topic) (2021)*.