

### 2018 Panoptic Challenge

European Conference on Computer Vision

# Joint COCO and Mapillary Recognition Challenge Workshop

Sunday, September 9th, ECCV 2018

Alexander Kirillov, Facebook AI Research



# COCO Panoptic Dataset







# 2018 Panoptic Segmentation Dataset



- $\succ$  For each pixel *i* predict semantic label *l* and instance id *z*
- ➤ no overlaps between segments by design



# 2018 Panoptic Segmentation Dataset



- ➤ COCO annotations have overlaps
- $\succ$  Most overlaps can be resolved automatically
- $\succ \sim \! 25 \mathrm{k}$  overlaps require manual resolution



# 2018 Panoptic Segmentation Dataset

#### Instructions:

In each row click on the image with better objects layout.



#### Task:



# Ø

## 2018 Panoptic Segmentation Dataset



train: 118k, val: 5k, test-dev: 20k, test-challenge: 20k
80 things categories, 53 stuff categories



# Panoptic Quality Measure



Ground Truth



Prediction

PQ Computation:

- Step 1: Matching
- Step 2: Calculation

# Panoptic Quality (PQ): Matching



Ground Truth



Prediction

**Theorem:** For panoptic segmentation problem each ground truth segment can have at most one corresponding predicted segment with IoU greater than 0.5 **Proof sketch:** 



then there is no other non overlapping object that has IoU > 0.5.



# Panoptic Quality (PQ): Matching



Ground Truth



Prediction

$$TP = \{( \bigcirc, \bigcirc), ( \bigcirc, \bigcirc) \}$$
$$FP = \{ \bigcirc \}$$
$$FN = \{ \bigcirc \}$$

# Panoptic Quality (PQ):Calculation



Ground Truth



Prediction

$$\mathrm{PQ} = \frac{\sum_{(p,g)\in TP} \mathrm{IoU}(p,g)}{|TP| + \frac{1}{2}|FP| + \frac{1}{2}|FN|}$$

# Panoptic Quality (PQ):Calculation



# COCO Panoptic Metrics

Average Panoptic Metrics:	
PQ % Panoptic Quality (primary challenge metric)	
SQ % Segmentation Quality component of PQ	
RQ % Recognition Quality component of PQ	
Panoptic Metrics for Things Categories:	
PQ <sup>Th</sup> % PQ for things categories only	
SQ <sup>Th</sup> % SQ for things categories only	
RQ <sup>Th</sup> % RQ for over things categories only	
Panoptic Metrics for Stuff Categories:	
PQ <sup>St</sup> % PQ for stuff categories only	
SQ <sup>St</sup> % SQ for stuff categories only	
RQ <sup>St</sup> % RQ for stuff categories only	



COCO images were annotated independently twice



COCO images were annotated independently twice

	PQ	SQ	RQ
All	53.5	82.6	63.9
Things	57.8	81.4	69.7
Stuff	47.1	84.3	55.2

➤ Crowd sourced annotations are very noisy

#### 5000 COCO images were annotated independently twice

	PQ	SQ	RQ		PQ	SQ	RQ
All	53.5	82.6	63.9	Small	25.2	62.1	32.8
Things	57.8	81.4	69.7	Medium	53.5	81.7	64.6
Stuff	47.1	84.3	55.2	Large	69.6	87.5	78.3

➤ Crowd sourced annotations are very noisy

 $\succ$  Annotations are highly inconsistent for small objects



### Annotations Consistency < Human Performance

### Annotations Consistency < Human Performance

real GT = [0, 0, 0, 0, 0, 1, 1, 1, 1, 1]noisy annotator 1 = [0, 0, 0, 0, 0, 1, 1, 1, 1, 1]noisy annotator 2 = [0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1]

### Annotations Consistency < Human Performance

real GT = [0, 0, 0, 0, 0, 1, 1, 1, 1, 1]noisy annotator 1 = [0, 0, 0, 0, 0, 1, 1, 1, 1, 1]noisy annotator 2 = [0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1]

### Annotations Consistency < Human Performance

- real GT = [0, 0, 0, 0, 0, 1, 1, 1, 1, 1]noisy annotator 1 = [0, 0, 0, 0, 0, 1, 1, 1, 1, 1]noisy annotator 2 = [0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1]

Accuracy(real GT, ideal annotator) = ?





 $\succ$  11 teams joined the competition





- $\succ$  11 teams joined the competition
- ▶ 4 teams achieved better performance than the baseline (RN50 Mask R-CNN + RN50 FPN-FCN)





- $\succ$  11 teams joined the competition
- ▶ 4 teams achieved better performance than the baseline (RN50 Mask R-CNN + RN50 FPN-FCN)



# Summary of Findings

### 2018 COCO Panoptic Challenge Take-aways

➢ All submission above the baseline combined the outputs of two separate networks for stuff and things



# Summary of Findings

#### 2018 COCO Panoptic Challenge Take-aways

- ➢ All submission above the baseline combined the outputs of two separate networks for stuff and things
- Best submission showed better PQ for things categories than the human consistency experiment

	PQ	SQ	RQ	$PQ^{\text{Th}}$	$SQ^{\text{Th}}$	$RQ^{\text{Th}}$	$PQ^{St}$	$SQ^{St} \\$	RQ <sup>St</sup>
Human Consistency	53.5	82.6	63.9	57.8	81.4	69.7	47.1	84.3	55.2
Megvii (Face++)	53.8	83.4	63.6	62.8	85.7	73.1	40.2	80	49.2



# Summary of Findings

### 2018 COCO Panoptic Challenge Take-aways

- ➢ All submission above the baseline combined the outputs of two separate networks for stuff and things
- Best submission showed better PQ for things categories than the human consistency experiment

Things	$PQ^{\text{Th}}$	$SQ^{\text{Th}}$	$RQ^{\text{Th}}$	all  TP	all  FP	all  FN	Precision	Recall
Human Consistency	57.8	81.4	69.7	24890	8860	9628	72.6%	69.5%
Megvii (Face++)	62.8	85.7	73.1	24205	4929	10313	81.1%	67.1%

- ➢ The result suggests ability to learn models with low noise level from large-scale noisy data
- $\succ$  Accuracy of the test set ground truth needs to be improved in the future





Image

Prediction Megvii (Face++)









Image

Prediction Megvii (Face++)









Image

Prediction Megvii (Face++)









Image

Prediction Megvii (Face++)

> dog Vice Andrew Constant of the second secon





#### Image

Prediction Megvii (Face++)





30





- $\succ$  11 teams joined the competition
- ▶ 4 teams achieved better performance than the baseline (RN50 Mask R-CNN + RN50 FPN-FCN)





- $\succ$  11 teams joined the competition
- ▶ 4 teams achieved better performance than the baseline (RN50 Mask R-CNN + RN50 FPN-FCN)

\*External segmentation datasets were used







Team	Position
Megvii (Face++)	$1^{\mathrm{st}}$
Caribbean	$2^{ m nd}$
PKU_360	$3^{ m rd}$