# **Scaling Object Detection up** to More Categories

Y. Lu, B. Xiao, Z. Gong, X. Zhang, H. Guo, L. Wen **Bytedance Al Lab** 



- More object categories: 80 -> 365
- More training images: 11W -> 60W
- More data  $\rightarrow$  more gains
- But...

#### Object365 dataset has a longer tail



Class ID (sorted by # instances)

### Class instance distribution of Object365 2250000 1687500 number of instances 1125000 long tail 562500 0 **In ByteDance** class ID

Class imbalance problem is more severe on Object365 

	COCO	Object365
Max #Instance	262465	2120895
Min #Instance	198	28
Max / Min	1326	75746

- More object classes: 80 -> 365
- More training images: 11W -> 60W
- But longer tail and more imbalance data
- What if we simply apply COCO models onto 365 classes?



- mAP of 44.7 on COCO  $\bigcirc$
- Achieve only mAP of 29.5 on the validation set of Object365



[1] Cai Z, Vasconcelos N. Cascade r-cnn: Delving into high quality object detection. CVPR 2018. [2] Xie S, Girshick R, Dollár P, et al. Aggregated residual transformations for deep neural networks. CVPR 2017.

### Start from Cascade R-CNN [1] with ResNext101 64x4d [2] backbone

**by ByteDance** 



## Class AP distribution on Object365

The AP is worse for the classes with less instances 



Class AP Distribution on Object365

## A detailed look on class 301-365

• 39 out of 65 classes has 0 AP !



AP of Class 301-365

Class ID

# A detailed look on class 301-365

### • Zero AP classes: okra, scallop, pitaya





#### Most small things with heavy clustering







#### hal ByteDance

### A detailed look on class 301-365 High AP classes: donkey, polar bear, seal



Most animals, with large scales and simple appearance





#### **h** ByteDance

## **Possible solutions**

- Expert models
- Data distribution resampling

## Expert models

- Fine-tuning the full classes model on class 301-365
- mAP on Class 301-365: 18.4 → **29.5\*** 
  - APs of **46** classes increase  $\bigcirc$



Class ID



## Expert models

#### Introducing expert models improves overall mAP by 1.1

- Expert 1: 301-365 classes  $\bigcirc$
- Expert 2: 151-300 classes  $\bigcirc$



	mAP
	29.6
	29.9
t 2	30.7



## Data distribution resampling

#### Down-sample classes with huge number of instances

Number of Instances of Class 151-365





# 

Class ID



# Data distribution resampling

- Down-sample classes with huge number of instances
  - mAP of Class 301-365: 18.4 -> 23.3\*  $\bigcirc$
  - overall mAP: 31.3 -> 31.0  $\bigcirc$
- No gain on overall mAP





mAP on validation set

A better pretrained backbone improves mAP by 0.6 



mAP on validation set

31.3 +0.6



Multi-scale training improves mAP by 0.9 



mAP on validation set

#### hul ByteDance

Multi-scale testing and soft NMS improve mAP by 1.4 



Model ensemble improves mAP by 0.9 



## Tiny track experiments

- Pretraining on Full Track dataset improves mAP by 4.2



# Baseline: Cascade R-CNN with ResNext101 64x4d pretrained on COCO

mAP on validation set of Tiny Track

pretrained on Full Track

## Tiny track experiments

Other tricks improve mAP by 5.3 



#### hul ByteDance

## Our final results

Validation set (Full tra

Test set (Full track

Validation set (Tiny tr

Test set (Tiny track

	mAP	
ack)	34.5	
<b>k)</b>	31.1	
rack)	34.8	
k)	27.4	



## Experiment details

#### **Basic setting**

- Cascade R-CNN with 3 stages  $\bigcirc$
- FPN  $\bigcirc$
- Deformable convolution  $\bigcirc$

#### Backbones

- ResNeXt 101 64x4d / 32x8d  $\bigcirc$
- SENet154  $\bigcirc$
- Resnet152  $\bigcirc$

#### **Training Pipeline and settings**

- ImageNet pre-train  $\rightarrow$  COCO pre-train for 12 epochs  $\bigcirc$
- $\bigcirc$
- Tiny Track: fine-tuning for 10 epochs (Ir 0.1 for 4 epochs, 0.01 for 6 epochs)  $\bigcirc$
- Batch size: 80 (2 imgs/GPU \* 40 GPUs)  $\bigcirc$

Full Track: training for 20 epochs (Ir 0.1 for 6 epochs, 0.01 for 10 epochs, 0.001 for 4 epochs)

## Conclusion

- **Data distribution matters** 
  - Long tail distribution greatly degrades the overall performance  $\bigcirc$
- Expert helps general model
  - Expert model can improve APs for long tail classes  $\bigcirc$
- General model also helps expert
  - Large data pre-training helps the learning of long tail classes  $\bigcirc$
- Long tail problem for object detection has not been solved

### **by ByteDance**

## We are hiring!

areas (@Beijing, Shanghai, Shenzhen):

Machine learning, natural language processing, computer vision, speech recognition and synthesis, and distributed systems.

Email:lab-hr@bytedance.com



# THANKS.

# **ByteDance**

