

# Adaptive Parametric Activation

Konstantinos Panagiotis Alexandridis<sup>1</sup>, Jiankang Deng<sup>1</sup>, Anh Nguyen<sup>2</sup> and Shan Luo<sup>3</sup>

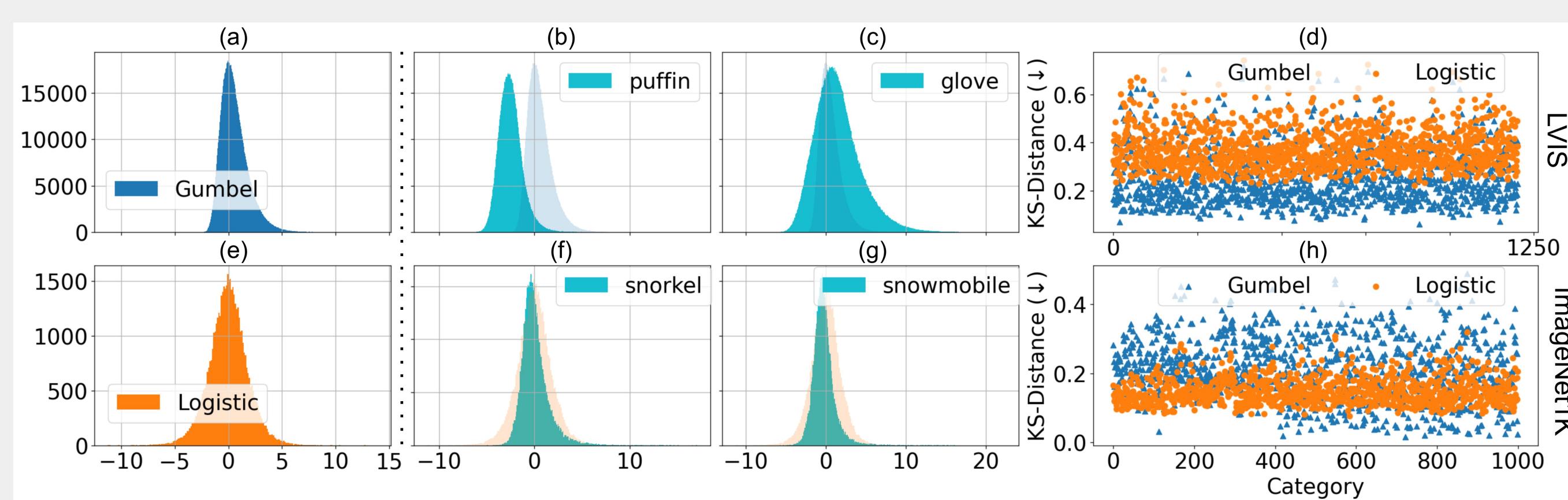
<sup>1</sup>Huawei Noah’s Ark Lab, <sup>2</sup>University of Liverpool, <sup>3</sup>King’s College London



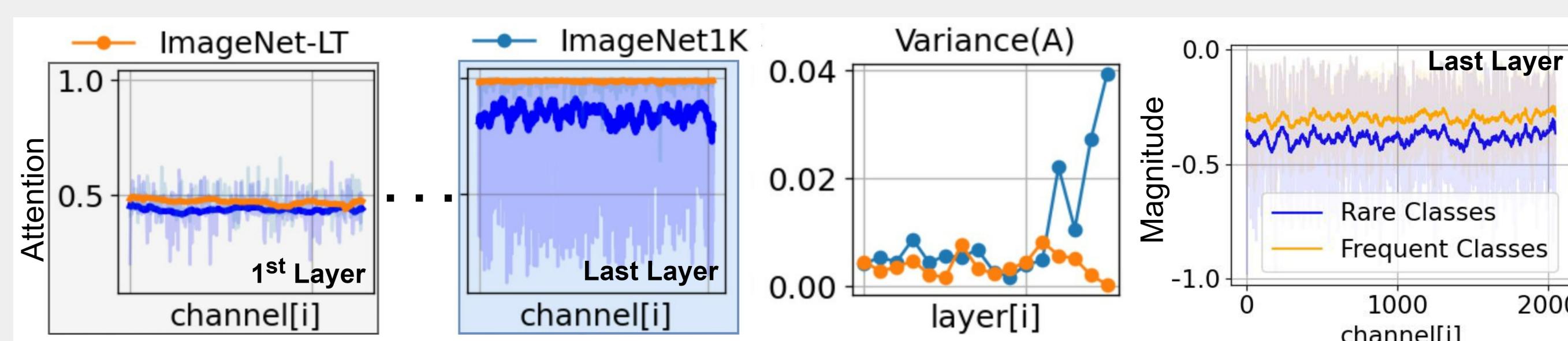
## SUMMARY

1. We study the importance of the activation function in balanced and imbalanced classification problems;
2. We propose the novel Adaptive Parametric Activation (APA) function that unifies most common activation functions under a single formula;
3. We have validated the efficacy of APA in balanced and imbalanced benchmarks surpassing the SOTA;

## MOTIVATION



The degree of class imbalance affects the shape of the learned logit distributions. When training with the imbalanced LVIS dataset (top), the logit distribution of the classes are closer to the Gumbel distribution, when measured with the Kolmogorov-Smirnov (KS) distance. In contrast, when training with the ImageNet1K (bottom), the logit distributions are closer to the Logistic distribution.



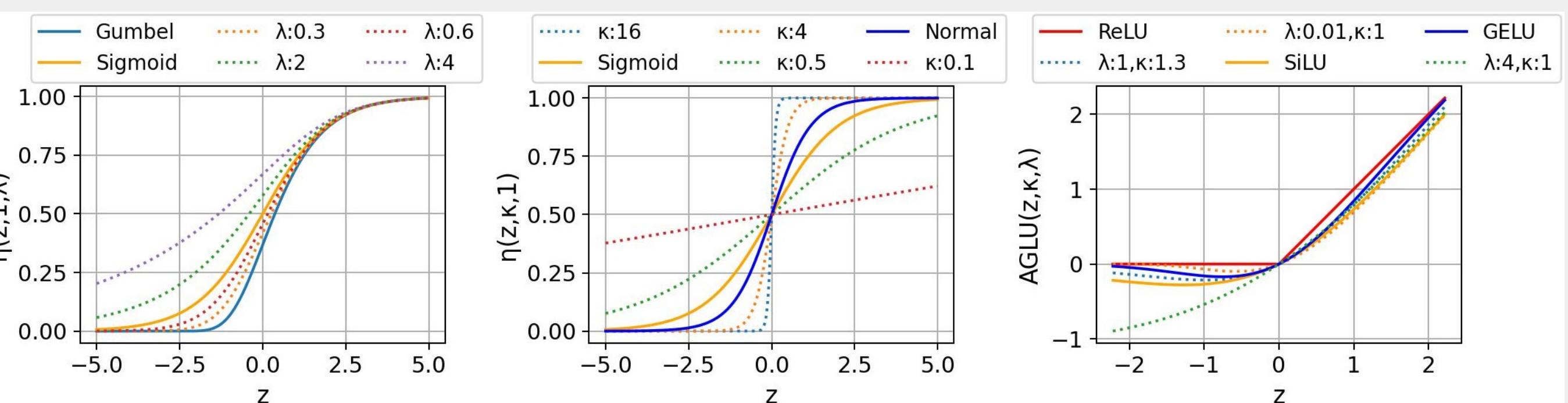
Class imbalance affects the channel attention signals of the model’s most semantic layers. When training with ImageNet-LT, the attention signal becomes all-ones and the variance goes to zero. This shows that the semantic attention is biased to the frequent classes, because the activation magnitude is smaller for the rare classes than the frequent classes.

## ADAPTIVE PARAMETRIC ACTIVATION

To boost the activation signals under long-tail learning, we propose the Adaptive Parametric Activation (APA):

$$\eta_{ad}(z, \kappa, \lambda) = (\lambda \exp(-\kappa z) + 1)^{\frac{1}{-\lambda}}. \quad (1)$$

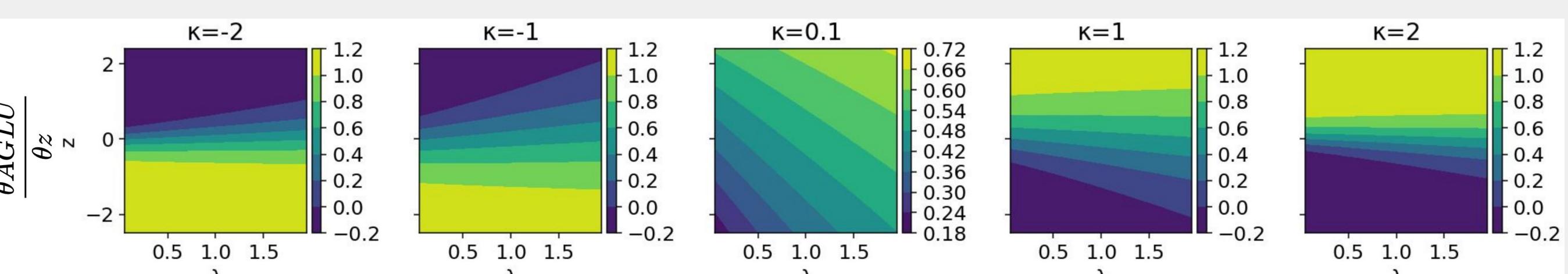
APA adapts its activation rate, according to the input’s distribution using two learnable parameters  $\kappa$  and  $\lambda$ .



The Adaptive Generalised Linear Unit (AGLU) is:

$$AGLU(z, \kappa, \lambda) = z \cdot \eta_{ad}(z, \kappa, \lambda) \quad (2)$$

The  $\kappa$  parameter controls the RELU-ness and the  $\lambda$  parameter controls the leakage, making AGLU flexible and expressive, as shown by it’s derivative with respect to  $z$ .



AGLU unifies many common activation functions. For example, if  $\kappa = \lambda = 1$  then AGLU becomes SiLU, if  $\kappa = 1.702$  and  $\lambda = 1$ , it approximates GELU, if  $\kappa \rightarrow \infty$ , then it approximates ReLU and if  $\lambda \rightarrow \infty$ , then AGLU becomes the identity function.

Name	Formula	Range
RELU	$\eta(z) = \max(0, z)$	$(0, \infty)$
GELU	$\eta(z) = z \sigma(1.702z)$	$(-0.17, \infty)$
SiLU	$\eta(z) = z \sigma(z)$	$(-0.28, \infty)$
Mish	$\eta(z) = z \tanh(\ln(1 + \exp(z)))$	$(-0.31, \infty)$
PRELU	$\eta(z, \kappa) = \max(0, z) + \kappa \min(0, z)$	$(-\infty, \infty)$
ELU	$\eta(z, \kappa) = \max(0, z) + \kappa(\exp(\min(0, z)) - 1)$	$(-\kappa, \infty)$
AGLU (ours)	$\eta(z, \kappa, \lambda) = z \cdot (\lambda \exp(-\kappa z) + 1)^{\frac{1}{-\lambda}}$	$(-\infty, \infty)$

## RESULTS

Results on long-tailed image classification, using Squeeze and Excite (SE) ResNets (R).  $APA^*$  shows that both APA and AGLU are used inside the SE model.

Method	Acc	Method	Acc	Activations	Acc
SE-X50	56.7	SE-R50	71.3	ReLU	57.4
RIDE(4E)	56.8	BCL	71.8	GELU	57.5
ResLT	56.1	ResLT	70.5	SiLU	57.1
$APA^*$ (ours)	<b>59.1</b>	$APA^*$ (ours)	<b>74.8</b>	AGLU	<b>57.9</b>
(a)ImageNet-LT		(b)iNaturalist18		(c)ImageNet-LT	

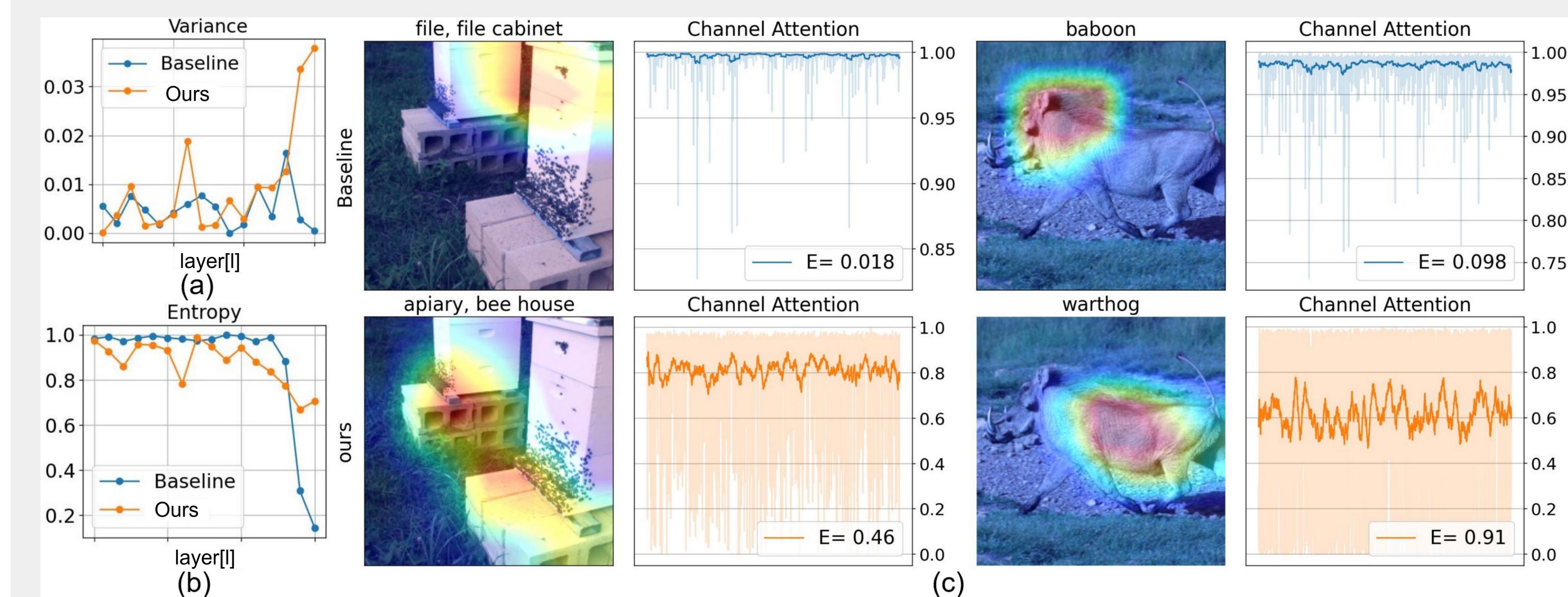
Results on LVISv1 using Mask-RCNN-SE-ResNet50, our  $APA^*$  increases the rare class performance  $AP^r$ .

Method	$AP$	$AP^r$	$AP^c$	$AP^f$	$AP^b$
GOL	28.2	20.6	28.9	30.8	28.1
GOL w/ $APA^*$ (ours)	<b>29.1</b>	<b>21.6</b>	<b>29.6</b>	<b>31.7</b>	<b>29.0</b>

$APA^*$  generalises well to balanced benchmarks like COCO, V3DET and ImageNet1K, using ResNet50.

Method	$AP^b$	$AP^m$	Method	$AP^b$	Method	top-1	
MRCNN	39.2	35.4	CRCNN	31.6	w/ SE	77.5	
w/ SE	40.5	36.9	w/ SE	33.3	w/ $APA^*$	<b>78.7</b>	
w/ $APA^*$	<b>41.2</b>	<b>37.6</b>	w/ $APA^*$	<b>35.4</b>	(a)COCO	(b)V3Det	(c)ImageNet1k

$APA^*$  increases the variance of the channel attention signals, thus it debiases the model. Also, it increases the channel attention’s entropy in the semantic layers of the model, allowing the model to better recognise the rare classes like *apiary* and *warthog*.



Contact: alex.kostas@gmail.com