

tailed class distribution







What Makes CLIP More Robust to Long-Tailed Pre-Training Data? **A Controlled Study for Transferable Insights**

Bingchen Zhao² ¹The University of Hong Kong

Yilun Chen³ Jiangmiao Pang³ Xiaojuan Qi¹ ³Shanghai Al Laboratory ²University of Edinburgh

SL under extreme long-tail (or open-world recognition) Sub. Voc is necessary to acquire CLIP knowledge



Figure 8: An extreme case: we train SL models on IN-Caps variants that have tail classes trimmed to only one shot (a & b) or even zero shot (c & d), and evaluate the accuracy on the tail and other classes. • CLIP with a frozen pre-trained text encoder shows superior generalization, which can be acquired by a ***** SL model with ***** fixed class prototypes from CLIP and ***** vocabulary subsampling.

SSL (DINO) on uncurated web data *vs* ImageNet



Figure 10: Transfer learning results of DINO variants pre-trained on LAIONet vs. vanilla DINO trained on ImageNet. Extreme data imbalance makes LAIONet much harder for DINO to learn transferrable representations. The vocabulary subsampling strategy effectively helps DINO alleviate such defects and generally match ImageNet-pretrained performance.









3. Applications to SL and SSL

DINO variants on LAIONet vs. vanilla DINO on ImageNet