

315 **A Proofs**

316 **A.1 Proof of Lemma 2**

317 **Restated Lemma (Lemma 2).** For a balanced target distribution, where $P_t(y) = 1/K$ for all
 318 $y \in [K]$, we have:

$$P_t(y|f_{\text{ft}}(\mathbf{x}), f_{\text{zs}}(\mathbf{x})) = \text{softmax}(f_{\text{ft}}(\mathbf{x}) + f_{\text{zs}}(\mathbf{x}) - \pi_s - \pi_p)(y) \quad (13)$$

319 *Proof.* Denote the output $\mathbf{e} = f_{\text{ft}}(\mathbf{x})$ and $\mathbf{z} = f_{\text{zs}}(\mathbf{x})$. We first use the Bayes Rule to decompose
 320 $P_t(y|\mathbf{e}, \mathbf{z})$ into $P_t(\mathbf{e}, \mathbf{z}|y)$, $P_t(y)$ and $P_t(\mathbf{e}, \mathbf{z})$ in Eq. (14), then rewrite $P_t(\mathbf{e}, \mathbf{z}|y)$ in Eq. (15) accord-
 321 ing to Assumption 1. Focusing on label shift problem [33, 19, 28] where $P(\mathbf{x}|y)$ does not change,
 322 we derive Eq. (16)

$$P_t(y|\mathbf{e}, \mathbf{z}) = \frac{P_t(\mathbf{e}, \mathbf{z}|y)P_t(y)}{P_t(\mathbf{e}, \mathbf{z})} \quad (14)$$

$$= P_t(\mathbf{e}|y)P_t(\mathbf{z}|y) \frac{P_t(y)}{P_t(\mathbf{e}, \mathbf{z})} \quad (15)$$

$$= P_s(\mathbf{e}|y)P_p(\mathbf{z}|y) \frac{P_t(y)}{P_t(\mathbf{e}, \mathbf{z})} \quad (16)$$

$$= \frac{P_s(y|\mathbf{e})P_s(\mathbf{e})}{P_s(y)} \frac{P_p(y|\mathbf{z})P_p(\mathbf{z})}{P_p(y)} \frac{P_t(y)}{P_t(\mathbf{e}, \mathbf{z})} \quad (17)$$

$$= \frac{P_s(y|\mathbf{e})}{P_s(y)} \frac{P_p(y|\mathbf{z})}{P_p(y)} \frac{P_s(\mathbf{e})P_p(\mathbf{z})P_t(y)}{P_t(\mathbf{e}, \mathbf{z})} \quad (18)$$

$$(19)$$

323 Since $P_t(y) = 1/K$ is constant and \mathbf{e}, \mathbf{z} are fixed, we can replace the terms that not rely on y
 324 with a constant C_1 in Eq. (20). Suppose the underlying class-probabilities $P_s(y|\mathbf{e}) \propto \exp(\mathbf{e}_y)$
 325 and $P_p(y|\mathbf{z}) \propto \exp(\mathbf{z}_y)$ for $y \in [K]$. We replace $P_s(y) = \exp(\log P_s(y)) = \exp(\pi_s(y))$ and
 326 $P_p(y) = \exp(\log P_p(y)) = \exp(\pi_p(y))$. Denote the constant C_2 for normalizing P_s and P_p into
 327 probabilities, we get Eq. (21)

$$P_t(y|\mathbf{e}, \mathbf{z}) = \frac{P_s(y|\mathbf{e})}{P_s(y)} \frac{P_p(y|\mathbf{z})}{P_p(y)} C_1 \quad (20)$$

$$= \exp(\mathbf{e} + \mathbf{z} - \pi_s - \pi_p)(y) \frac{C_1}{C_2} \quad (21)$$

328 Because the summation of $P_t(y|\mathbf{e}, \mathbf{z})$ is 1, $\frac{C_1}{C_2} = 1 / \sum_{i \in [K]} \exp(\mathbf{e} + \mathbf{z} - \pi_s - \pi_p)(i)$. Therefore, we
 329 have:

$$P_t(y|f_{\text{ft}}(\mathbf{x}), f_{\text{zs}}(\mathbf{x})) = P_t(y|\mathbf{e}, \mathbf{z}) \quad (22)$$

$$= \frac{\exp(\mathbf{e} + \mathbf{z} - \pi_s - \pi_p)_y}{\sum_{i \in [K]} \exp(\mathbf{e} + \mathbf{z} - \pi_s - \pi_p)_i} \quad (23)$$

$$= \text{softmax}(f_{\text{ft}}(\mathbf{x}) + f_{\text{zs}}(\mathbf{x}) - \pi_s - \pi_p)_y \quad (24)$$

330 \square

331 **A.2 Proof of Proposition 2**

332 **Restated Proposition (Proposition 2).** Suppose that the target distribution P_p is class-balanced.
 333 Let $h : \mathbb{R}^K \rightarrow \mathbb{R}^K$ be an arbitrary function that predicts labels using the outputs of the zero-shot
 334 model $f_{\text{zs}}(\mathbf{x})$. Let the derived classifier be denoted as $f_h(\mathbf{x}) = h(f_{\text{zs}}(\mathbf{x}))$. The classifier $f_{\text{zs}} - \pi_p$ is
 335 better than any $f_h(\mathbf{x})$: $\mathcal{R}_t(f_{\text{zs}} - \pi_p) \leq \mathcal{R}_t(f_h(\mathbf{x}))$.

Dataset	Classes	Train size	Test size	Task
ImageNet	1,000	1.28M	50,000	Object-level
CIFAR100	100	50,000	10,000	Object-level
Caltech101	100	4,128	2,465	Object-level
DTD	47	2,820	1,692	Textures
EuroSAT	10	13,500	8,100	Satellite images
FGVCAircraft	100	3,334	3,333	Fine-grained aircraft
Flowers102	102	4,093	2,463	Fine-grained flowers
Food101	101	50,500	30,300	Fine-grained food
OxfordPets	37	2,944	3,669	Fine-grained pets
StanfordCars	196	6,509	8,041	Fine-grained car
SUN397	397	15,880	19,850	Scene-level
UCF101	101	7,639	3,783	Action
ImageNetV2	1,000	-	10,000	Robustness to collocation
ImageNet-Sketch	1000	-	50,889	Robustness to sketch domain
ImageNet-A	200	-	7,500	Robustness to adversarial attack
ImageNet-R	200	-	30,000	Robustness to multi-domains

Table 7: The detailed statistics of datasets for many-shot and few-shot learning.

336 *Proof.* Denote the output $\mathbf{z} = f_{\text{zs}}(\mathbf{x})$. Similar to Eq. (14)-Eq. (24), we have

$$P_t(y|\mathbf{z}) = \frac{P_t(\mathbf{z}|y)P_t(y)}{P_t(\mathbf{z})} \quad (25)$$

$$= \frac{P_p(\mathbf{z}|y)P_t(y)}{P_t(\mathbf{z})} \quad (26)$$

$$= \frac{P_p(y|\mathbf{z})}{P_p(y)} \frac{P_t(y)}{P_t(\mathbf{z})} \quad (27)$$

$$= \exp(\mathbf{z} - \pi_p)(y) / \sum_{i \in [K]} \exp((\mathbf{z} - \pi_p)(i)) \quad (28)$$

$$= \text{softmax}(\mathbf{z} - \pi_p) = \text{softmax}(f_{\text{zs}}(\mathbf{x}) - \pi_p) \quad (29)$$

337 Therefore, we have:

$$\operatorname{argmax}_{y \in \mathcal{Y}} (f_{\text{zs}}(\mathbf{x}) - \pi_p)_y = \operatorname{argmax}_{y \in \mathcal{Y}} \text{softmax}(f_{\text{zs}}(\mathbf{x}) - \pi_p)_y = \operatorname{argmax}_{y \in \mathcal{Y}} P_t(y|f_{\text{zs}}(\mathbf{x})) \quad (30)$$

338 Again, using Lemma 1, any other classifier $f_h(\mathbf{x})$ has higher risk than $f_{\text{zs}}(\mathbf{x}) - \pi_p$, *i.e.*, $\mathcal{R}_t(f_{\text{zs}} - \pi_p) \leq$
339 $\mathcal{R}_t(f_h(\mathbf{x}))$. \square

340 B Experimental Details

341 B.1 Dataset details

342 **Many-shot and few-shot datasets.** For many-shot learning, we use ImageNet, CIFAR100, Stanford-
343 Cars and SUN397 datasets. For few-shot learning, we evaluate models on 15 datasets. The details of
344 each dataset are presented in Table 7.

345 **Long-tail datasets.** We use two standard long-tail benchmarks: Places365-LT and ImageNet-LT [29].
346 The skewness of a long-tailed training set is typically represented by imbalanced ratio, which is
347 defined as N_{\max}/N_{\min} . N_{\max} (N_{\min}) denotes the largest (smallest) number of instances per class. A
348 larger imbalanced ratio means a more imbalanced training set. The test sets are divided into three
349 splits: many-shot subset contains classes with > 100 images, medium-shot subset includes classes
350 with ≥ 20 & ≤ 100 images, and few-shot subset covers classes with < 20 images. Details are listed
351 in Table 8.

352 B.2 CLIP zero-shot

353 We use prompt ensembling of 80 prompts provided by CLIP [45] for ImageNet, CIFAR100, and Cal-
354 tech101 to improve performance, *i.e.*, averaging the text embedding of many captions, *e.g.*, “a photo

Dataset	Size of all classes	Size of many classes	Size of medium classes	Size of few classes	Size of training samples	Imbalanced ratio
Places365-LT	365	131	163	71	62.5K	996
ImageNet-LT	1000	385	479	136	186K	256

Table 8: Details of long-tailed datasets.

of a $\{c_k\}$.” and “an image of a $\{c_k\}$.”. For OxfordPets, StanfordCars, Flowers102, Food101, FGV-CAircraft, EuroSAT, UCF101, DTD and SUN397, we use the pre-defined prompt from CoOp [52].

B.3 Fine-tuned models

End-to-end and linear probe fine-tuning. We follow WiSE-FT [45] to implement fine-tuning. We initialize the classifier with the zero-shot classifier and the output of the image encoder Φ_v is normalized during fine-tuning. We fine-tune for a total of 10 epochs using AdamW [30] optimizer with default hyper-parameters $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$ and weight decay 0.1. We choose a batch size of 512. We use the same data augmentation and cosine-annealing learning rate schedule as [45].

B.4 Prompt tuning.

Prompt tuning like CoOp [52] automates prompt engineering by learning the prompt given few samples from downstream tasks. CoOp provides two options of prompt design: unified prompt that is shared among all classes and class-specific prompt that is different for each class. In this paper, we adopt the class-specific prompt design as the fine-tuned model to implement GLA. In specific, given the word embedding \mathbf{t}_k^0 initialized by zero-shot prompts, we aim to learn a collection of class-specific word embedding $\mathbf{R} = \{\mathbf{r}_k\}_{k=1}^K$, such that the text input $\mathbf{t}_k = \mathbf{t}_k^0 + \mathbf{r}_k$ minimizes the empirical risk: $\mathbf{R}^* = \operatorname{argmin}_{\mathbf{R}} \mathbb{E}_{\mathbf{x}, y} [y \neq \operatorname{argmax}_i f(x; \mathbf{R})_i]$.

We adhere CoOp to use CLIP ResNet-50 as image encoder for few-shot classification. The word embedding \mathbf{R} is initialized from zeros. For the m few-shot classification setting (where $m \in \{1, 2, 4, 8, 16\}$), we randomly sample m training and m validation points from the respective full datasets. For all few-shot datasets except ImageNet, the training epoch is set to 200 for 16/8 shots, 100 for 4/2 shots, and 50 for 1 shot. For ImageNet, the epoch is set to 50 for all shots. We fine-tune the prompt with SGD optimizer decayed by the cosine annealing rule. The base initial learning rate and batch size are set to 10^{-4} and 32. When given an m -shot sample setting, we increase the learning rate and batch size by m times simultaneously to accelerate the training speed.

B.5 Estimation of the class prior

To estimate the log-probability of the pre-training distribution $\hat{\pi}_s = \log \mathbf{q}$, we utilize the optimization toolkit Cooper [13] from <https://github.com/cooper-org/cooper>. \mathbf{q} is initialized as a uniformed distribution, $\mathbf{q}(y) = \frac{1}{K}$ for all $y \in [K]$. We use the standard SGD as the primal and dual optimizers for 2000 steps.

B.6 Long-tail learning baselines and training details

We compared with 5 long-tailed classification methods:

1. Standard ERM: We learn the model by standard empirical risk minimization on the long-tailed data.
2. Learnable Weight Scaling (LWS) [22]: We first learn the model by standard ERM, then fix the model and learn to re-scale the magnitude of the classifier using class-balanced sampling.
3. Logit Adjustment (LA) [33]: We first learn the model by standard ERM, then compensates the long-tailed distribution by subtracting a class-dependent offset to the model outputs.
4. Balanced Softmax (BS) [40] modifies the Softmax cross-entropy loss which explicitly accommodate the label distribution shift during optimization.

Model	Source	Target	
	ViT-B/32	ViT-B/16	ViT-L/14
Original zero-shot model $f_{zs}(\mathbf{x})$	63.4	68.8	75.6
Debiased zero-shot model $f_{zs}(\mathbf{x}) - \hat{\pi}_p$	65.4	69.3	76.3

Table 9: Estimated π_p is transferable across different backbones. $\hat{\pi}_p$ is estimated using CLIP ViT-B/32.

395 5. BALLAD [31] first fine-tunes the vision-language models via contrastive loss on long-tailed
396 data, then freezes the backbone and finally employs an adapter to enhance the representations
397 of tail classes with re-sampling strategies.

398 For all combinations of the fine-tuning baselines and long-tailed learning methods, visual backbones
399 are initialized from CLIP-ResNet-50 and all classifiers are initialized by feeding prompt with class
400 names to the text encoder. We use SGD for all experiments with a momentum of 0.9 for 50 epochs
401 with batch size of 512. The initial learning rate is set to 1.6×10^{-3} which is decayed by the cosine
402 annealing rule. To mitigate explosive gradients, we use the warmup learning rate equals to 10^{-5}
403 during the first epoch. For the sake of fairness in comparison, all hyper-parameters of baselines are
404 carefully searched using grid search on the validation set.

405 C Additional Experiments

406 C.1 Estimated label distribution is transferable

407 The estimated $\hat{\pi}_p$ should be transferable across different zero-shot models if they are trained on the
408 same pre-training dataset. To confirm this, we estimate π_p using CLIP ViT-B/32 based zero-shot
409 model, and use it to debias zero-shot models based on CLIP ViT-B/16 and ViT-L/14. Results are
410 shown in Table 9, where our debiased zero-shot models based on CLIP ViT-B/16 and ViT-L/14 using
411 $\hat{\pi}_p$ estimated from ViT-B/32 show clear performance gains over original zero-shot models.

412 C.2 Few-shot learning accuracy

413 We provide mean and standard deviation in Table 10 in for {1, 2, 4, 8, 16} shots on all 11 few-shot
414 learning datasets.

Dataset	1 shot	2 shots	4 shots	8 shots	16 shots
ImageNet	61.65 ± 0.15	62.64 ± 0.01	63.32 ± 0.07	64.51 ± 0.09	65.61 ± 0.03
Caltech101	89.08 ± 0.09	90.25 ± 0.25	90.98 ± 0.43	91.90 ± 0.21	92.58 ± 0.42
OxfordPets	87.79 ± 0.15	87.86 ± 0.21	88.22 ± 0.21	88.09 ± 0.27	89.53 ± 0.16
StanfordCars	60.00 ± 0.14	63.10 ± 0.42	66.25 ± 0.19	69.87 ± 0.09	73.95 ± 0.11
Flowers102	73.45 ± 0.60	81.00 ± 0.46	88.31 ± 0.65	92.89 ± 0.46	95.41 ± 0.32
Food101	78.41 ± 0.07	78.62 ± 0.07	78.68 ± 0.06	78.85 ± 0.19	79.54 ± 0.47
FGVCAircraft	20.22 ± 0.59	22.09 ± 0.37	24.65 ± 0.85	28.23 ± 0.44	31.99 ± 0.50
SUN397	64.29 ± 0.19	66.32 ± 0.16	68.01 ± 0.08	69.99 ± 0.18	71.64 ± 0.21
DTD	47.38 ± 1.23	50.75 ± 1.46	56.90 ± 0.20	62.73 ± 0.80	65.78 ± 0.49
EuroSAT	56.50 ± 1.34	67.26 ± 3.58	72.40 ± 2.43	77.59 ± 1.84	84.93 ± 1.89
UCF101	65.32 ± 0.17	68.42 ± 0.81	70.88 ± 0.50	74.23 ± 0.24	76.07 ± 0.03

Table 10: GLA Accuracy (%) with standard deviation of few-shot learning on 11 datasets.

415 References

- 416 [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson,
417 Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual
418 language model for few-shot learning. *NeurIPS*, 2022.
- 419 [2] Martin Arjovsky. *Out of distribution generalization in machine learning*. PhD thesis, New York
420 University, 2020.
- 421 [3] Eric Bauer and Ron Kohavi. An empirical comparison of voting classification algorithms:
422 Bagging, boosting, and variants. *Machine learning*, 1999.
- 423 [4] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von
424 Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the
425 opportunities and risks of foundation models. *Technical Report*, 2021.
- 426 [5] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative
427 components with random forests. In *ECCV*, 2014.
- 428 [6] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
429 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
430 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.
431 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz
432 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec
433 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In
434 *NeurIPS*, 2020.
- 435 [7] Guangyi Chen, Weiran Yao, Xiangchen Song, Xinyue Li, Yongming Rao, and Kun Zhang.
436 Prompt learning with optimal transport for vision-language models. *ICLR*, 2023.
- 437 [8] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi.
438 Describing textures in the wild. In *CVPR*, 2014.
- 439 [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
440 hierarchical image database. In *CVPR*, 2009.
- 441 [10] Thomas G Dietterich. Ensemble methods in machine learning. In *Multiple Classifier Systems:
442 First International Workshop*, 2000.
- 443 [11] Li Fei-Fei, Rob Fergus, and Pietro Perona. Learning generative visual models from few training
444 examples: An incremental bayesian approach tested on 101 object categories. In *CVPRW*, 2004.
- 445 [12] Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning
446 and an application to boosting. *Journal of computer and system sciences*, 1997.
- 447 [13] Jose Gallego-Posada and Juan Ramirez. Cooper: a toolkit for Lagrangian-based constrained
448 optimization. <https://github.com/cooper-org/cooper>, 2022.
- 449 [14] Lars Kai Hansen and Peter Salamon. Neural network ensembles. *IEEE transactions on pattern
450 analysis and machine intelligence*, 12(10):993–1001, 1990.
- 451 [15] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel
452 dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal
453 of Selected Topics in Applied Earth Observations and Remote Sensing*, 2019.
- 454 [16] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo,
455 Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A
456 critical analysis of out-of-distribution generalization. In *CVPR*, 2021.
- 457 [17] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural
458 adversarial examples. In *CVPR*, 2021.
- 459 [18] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network.
460 *NeurIPS Workshop*, 2014.

- 461 [19] Youngkyu Hong, Seungju Han, Kwanghee Choi, Seokjun Seo, Beomsu Kim, and Buru Chang.
462 Disentangling label distribution for long-tailed visual recognition. In *CVPR*, 2021.
- 463 [20] Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon
464 Wilson. Averaging weights leads to wider optima and better generalization. *UAI*, 2018.
- 465 [21] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan
466 Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning
467 with noisy text supervision. In *ICML*, 2021.
- 468 [22] Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and
469 Yannis Kalantidis. Decoupling representation and classifier for long-tailed recognition. In *ICLR*,
470 2020.
- 471 [23] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for
472 fine-grained categorization. In *CVPRW*, 2013.
- 473 [24] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- 474 [25] Ananya Kumar, Tengyu Ma, Percy Liang, and Aditi Raghunathan. Calibrated ensembles can
475 mitigate accuracy tradeoffs under distribution shift. In *UAI*. PMLR, 2022.
- 476 [26] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable
477 predictive uncertainty estimation using deep ensembles. *NeurIPS*, 2017.
- 478 [27] Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. Ensemble distillation for robust
479 model fusion in federated learning. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and
480 H. Lin, editors, *NeurIPS*, 2020.
- 481 [28] Zachary Lipton, Yu-Xiang Wang, and Alexander Smola. Detecting and correcting for label shift
482 with black box predictors. In *ICML*, 2018.
- 483 [29] Ziwei Liu, Zhongqi Miao, Xiaohang Zhan, Jiayun Wang, Boqing Gong, and Stella X Yu.
484 Large-scale long-tailed recognition in an open world. In *CVPR*, 2019.
- 485 [30] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019.
- 486 [31] Teli Ma, Shijie Geng, Mengmeng Wang, Jing Shao, Jiasen Lu, Hongsheng Li, Peng Gao, and
487 Yu Qiao. A simple long-tailed recognition baseline via vision-language model. *arXiv preprint*
488 *arXiv:2111.14745*, 2021.
- 489 [32] Subhansu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-
490 grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- 491 [33] Aditya Krishna Menon, Sadeep Jayasumana, Ankit Singh Rawat, Himanshu Jain, Andreas Veit,
492 and Sanjiv Kumar. Long-tail learning via logit adjustment. In *ICLR*, 2021.
- 493 [34] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large
494 number of classes. In *Indian Conference on Computer Vision, Graphics & Image Processing*,
495 2008.
- 496 [35] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In
497 *CVPR*, 2012.
- 498 [36] Hieu Pham, Zihang Dai, Golnaz Ghiasi, Kenji Kawaguchi, Hanxiao Liu, Adams Wei Yu, Jiahui
499 Yu, Yi-Ting Chen, Minh-Thang Luong, Yonghui Wu, et al. Combined scaling for zero-shot
500 transfer learning. *arXiv preprint arXiv:2111.10050*, 2021.
- 501 [37] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
502 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
503 models from natural language supervision. In *ICML*, 2021.
- 504 [38] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet
505 classifiers generalize to imagenet? In *ICML*, 2019.

- 506 [39] William J Reed. The pareto, zipf and other power laws. *Economics letters*, 2001.
- 507 [40] Jiawei Ren, Cunjun Yu, Xiao Ma, Haiyu Zhao, Shuai Yi, et al. Balanced meta-softmax for
508 long-tailed visual recognition. *NeurIPS*, 2020.
- 509 [41] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human
510 actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
- 511 [42] Kaihua Tang, Jianqiang Huang, and Hanwang Zhang. Long-tailed classification by keeping the
512 good and removing the bad momentum causal effect. *NeurIPS*, 2020.
- 513 [43] Rohan Taori, Achal Dave, Vaishaal Shankar, Nicholas Carlini, Benjamin Recht, and Ludwig
514 Schmidt. Measuring robustness to natural distribution shifts in image classification. *NeurIPS*,
515 2020.
- 516 [44] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global represen-
517 tations by penalizing local predictive power. *NeurIPS*, 2019.
- 518 [45] Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca
519 Roelofs, Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, et al.
520 Robust fine-tuning of zero-shot models. In *CVPR*, 2022. [https://arxiv.org/abs/2109-](https://arxiv.org/abs/2109.01903)
521 [01903](https://arxiv.org/abs/2109.01903).
- 522 [46] Tong Wu, Ziwei Liu, Qingqiu Huang, Yu Wang, and Dahua Lin. Adversarial robustness under
523 long-tailed distribution. *CVPR*, 2021.
- 524 [47] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database:
525 Large-scale scene recognition from abbey to zoo. In *CVPR*, 2010.
- 526 [48] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui
527 Wu. Coca: Contrastive captioners are image-text foundation models. *TMLR*, 2022.
- 528 [49] Renrui Zhang, Rongyao Fang, Peng Gao, Wei Zhang, Kunchang Li, Jifeng Dai, Yu Qiao, and
529 Hongsheng Li. Tip-adapter: Training-free clip-adapter for better vision-language modeling.
530 *ECCV*, 2022.
- 531 [50] Songyang Zhang, Zeming Li, Shipeng Yan, Xuming He, and Jian Sun. Distribution alignment:
532 A unified framework for long-tail visual recognition. In *CVPR*, 2021.
- 533 [51] Tong Zhang. Statistical behavior and consistency of classification methods based on convex
534 risk minimization. *The Annals of Statistics*, 2004.
- 535 [52] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for
536 vision-language models. *IJCV*, 2022.
- 537 [53] Beier Zhu, Yulei Niu, Yucheng Han, Yue Wu, and Hanwang Zhang. Prompt-aligned gradient
538 for prompt tuning. *arXiv preprint arXiv:2205.14865*, 2022.
- 539 [54] Beier Zhu, Yulei Niu, Saeil Lee, Minhoe Hur, and Hanwang Zhang. Debaised fine-tuning for
540 vision-language models by prompt regularization. *AAAI*, 2023.