
Calibrated Ensembles: A Simple Way to Mitigate ID-OOD Accuracy Tradeoffs

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We often see undesirable tradeoffs in robust machine learning where out-of-
2 distribution (OOD) accuracy is at odds with in-distribution (ID) accuracy. A “robust”
3 classifier obtained via specialized techniques like removing spurious features has
4 better OOD but worse ID accuracy compared to a “standard” classifier trained via
5 vanilla ERM. On six distribution shift datasets, we find that simply ensembling a
6 standard and a robust model is a strong baseline—we match the ID accuracy of a
7 standard model with only a small drop in OOD accuracy compared to the robust
8 model. However, calibrating these models in-distribution surprisingly improves
9 the OOD accuracy of the ensemble and eliminates the tradeoff and we achieve the
10 best of both ID and OOD accuracy over the original models.

11 1 Introduction

12 Machine learning models typically suffer large drops in accuracy in the presence of distribution
13 shift where the test distribution is different from the training distribution. As ML systems are widely
14 deployed, it is important for models to have good “out-of-distribution” (OOD) accuracy. There has
15 been a lot of research interest in tackling this robustness problem under various settings such as
16 robustness to spurious correlations (1; 2; 3), domain generalization (4; 5), robustness to demographic
17 shifts (6; 7) among others. Almost universally across these different settings, an unfortunate tradeoff
18 arises. Robustness interventions typically improve the OOD accuracy but simultaneously cause a
19 drop in the “in-distribution” (ID) accuracy on new test points from the original distribution.

20 This tradeoff is a major hurdle in using the multitude of proposed robustness interventions. In practice,
21 most inputs are likely to be ID, so it is unsatisfactory to use a “robust” model that has high OOD
22 performance but performs less accurately on these majority ID points. On the other hand, “standard
23 models” (trained without robustness interventions) fail catastrophically in the presence of even small
24 shifts, and it can be highly dangerous to use a standard model even if OOD points are rare. In this
25 work, we ask *is there a general strategy by which we can achieve high accuracy both in-distribution
26 and out-of-distribution and mitigate tradeoffs arising in robustness?*

27 We consider four benchmark datasets (DomainNet, CIFAR → STL, ImageNet → ImageNet-R, and
28 BREEDS-Entity-30) and two real world satellite remote sensing datasets (Landcover and Cropland),
29 that have been used in prior work on robustness. Our work spans different types of robustness
30 interventions (projecting out spurious correlations, zero-shot language prompting, freezing pretrained
31 features), data modalities (image and time series data), and model architectures (vision transformers,
32 deep convolutional networks, time series convolution). Averaged across these datasets, robustness
33 interventions increase OOD accuracy from 63% to 74%, but decrease ID accuracy from 88% to 85%.

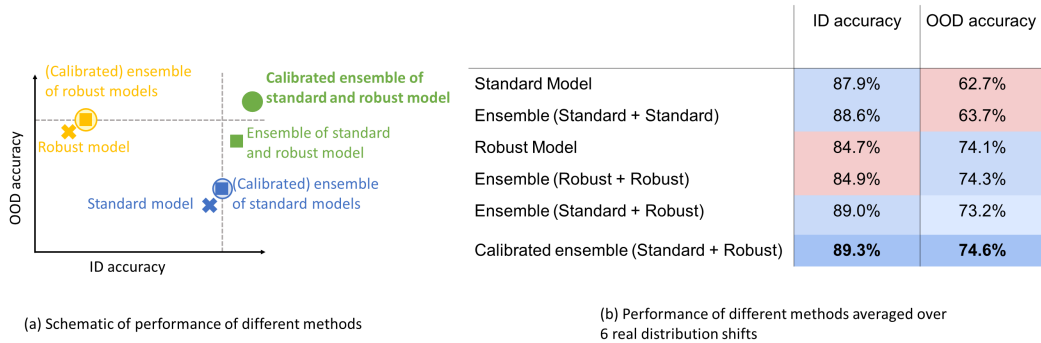


Figure 1: In many settings, we have a ‘standard’ model that performs better in-distribution, and a ‘robust’ model that performs better out-of-distribution. Simply ensembling these two models (e.g., by adding their probabilities), gets better ID accuracy than the standard and robust models, and closes most of the OOD gap. Calibrating the models in-distribution (no access to OOD data) before ensembling them leads to further improvements. Note that ensembling two standard or two robust models does not close the gap and only leads to small improvements.

34 We first explore the natural strategy of ensembling the standard and robust models—concretely, we
 35 add the probabilities of each model to obtain a prediction with the hope that when the two models
 36 conflict, the more confident model (with larger probability) dictates the final prediction. We find that
 37 this surprisingly simple baseline already performs quite well—on average across all our datasets, this
 38 closes 80% of the gap between the OOD of standard models, while outperforming both models ID.
 39 However, vanilla ensembling still leaves a gap as it underperforms the robust model OOD.

40 We find that simply calibrating both models ID (adjusting their predicted confidence to match their
 41 accuracy, on *in-distribution* data) before ensembling them closes this gap. *Calibrated ensembles get*
 42 *an average accuracy of 89.3% ID and 74.6% OOD, and outperform both the standard and robust*
 43 *model, ID and OOD.* The other method in the literature to alleviate robustness induced tradeoffs is
 44 self-training that uses large amount of unlabeled data (8; 9; 10). On the two remote-sensing datasets
 45 with additional unlabeled data, we find that calibrated ensembles match self-training on these datasets
 46 despite its simplicity and without requiring any unlabeled data.

47 While our method is intuitive, it is intriguing that it works so well because ensembling seems to rely
 48 on good uncertainty estimates while it is common wisdom that uncertainty estimates of deep networks
 49 are unreliable OOD even after calibrating in-distribution (11). Indeed, on the six datasets we test on,
 50 the models fare poorly on standard uncertainty metrics OOD, even after calibration. The expected
 51 calibration error of the standard model across all datasets is 12%. Even the relative confidences of
 52 the models can be incorrect—on the remote sensing dataset (Landcover), the standard model is on
 53 average 6% more confident in its OOD predictions than the robust model, even though the standard
 54 model is less accurate OOD. But at the granularity of individual points, calibrated ensembles are able
 55 to combine predictions effectively and achieve high ID and OOD accuracy.

56 2 Setup

57 Consider a K -class classification task, where the goal is to predict targets $y \in [K]$ from inputs $x \in \mathbb{R}^d$.

58 **Models:** A model $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$ takes an input $x \in \mathbb{R}^d$ and outputs $f(x) \in \mathbb{R}^k$ where $f(x)_i$ denotes the
 59 model’s confidence that the output is $y = i$. The model predicts the label $\hat{y} = \operatorname{argmax}_k f(x)_k$. The
 60 confidences can be converted into probabilities using a softmax function:

61 **Data:** Let P_{id} and P_{ood} denote the underlying distribution of (x, y) pairs in-distribution and
 62 out-of-distribution, respectively. We have a validation set $\{(x_i^{\text{val}}, y_i^{\text{val}})\}_{i=1}^{n_{\text{val}}} \sim P_{\text{id}}$ used for early
 63 stopping and calibration, a held-out in-distribution test set $\{(x_i^{\text{test}}, y_i^{\text{test}})\}_{i=1}^{n_{\text{test}}} \sim P_{\text{id}}$, and a held-out
 64 out-of-distribution test set $\{(x_i^{\text{ood}}, y_i^{\text{ood}})\}_{i=1}^{n_{\text{ood}}} \sim P_{\text{ood}}$. All methods can use the ID validation set for

Table 1: *In-distribution (ID)* accuracies for the standard model, robust model, and calibrated ensembling, across six datasets. Calibrated ensembling matches or outperforms the better model in all cases, and on average outperforms both the standard and robust models.

	Ent30	DomNet	CIFAR10	Land	Crop	ImNet
Standard	93.6 (0.2)	83.9 (1.0)	97.4 (0.1)	76.9 (0.3)	95.3 (0.0)	80.5 (-)
Robust	90.7 (0.2)	89.2 (0.1)	92.0 (0.0)	72.7 (0.2)	95.1 (0.1)	68.4 (-)
Cal ensemble	93.7 (0.1)	91.2 (0.6)	97.2 (0.1)	77.1 (0.2)	95.6 (0.1)	81.1 (-)

65 tuning hyperparameters, early stopping, and calibration. The ID and OOD test set are only used for
 66 evaluation. We evaluate a model f on the average accuracy on the ID and OOD test sets.

67 3 Methods and Datasets

68 **Calibrated ensembles.** Given two models f_1 and f_2 , calibrated ensembles first calibrate each model
 69 using temperature scaling (12) with the cross-entropy loss l on the *in-distribution* validation data:

$$T_j = \operatorname{argmin}_T \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} l\left(\frac{f_j(x_i^{\text{val}})}{T}, y_i^{\text{val}}\right) \quad \text{for } j \in \{1, 2\} \quad (3.1)$$

70 We then ensemble the two models by adding up the probabilities that they predict (13):

$$\hat{p} = \frac{1}{2} \left(\operatorname{softmax}\left(\frac{f_1(x)}{T_1}\right) + \operatorname{softmax}\left(\frac{f_2(x)}{T_2}\right) \right) \quad (3.2)$$

71 **Other ensembles.** As a baseline, we consider a *tuned ensemble*: outputting a weighted average of
 72 the standard and robust model’s probabilities, where the weight $\alpha \in [0, 1]$ is tuned to maximize accuracy
 73 on the in-distribution set.

$$\hat{p} = \alpha \operatorname{softmax}(f_{\text{std}}(x)) + (1 - \alpha) \operatorname{softmax}(f_{\text{rob}}(x)) \quad (3.3)$$

74 In *vanilla ensembling* the weight α is set to 1/2. We also considered other ways of combing the model
 75 outputs (adding logits vs probabilities) and found that the results were similar.

76 **Datasets** We run experiments spanning three different types of robustness interventions: projecting
 77 out spurious metadata, zero-shot language prompting in CLIP, and freezing pretrained features. These
 78 experiments span multiple model architectures (vision transformers, deep convolutional networks, time
 79 series convolution) and data modalities (image and time series data), and include two real world remote
 80 sensing datasets used in prior work studying ID-OOD tradeoffs (9). See Appendix A for more details.

81 4 Experiments

82 **Strong ID and OOD accuracy:** Calibrating and then ensembling a standard and a robust model, gets
 83 the best of both worlds, typically outperforming the standard and robust model both ID (Table 1) and
 84 OOD (Table 2). Averaged across the datasets, calibrated ensembles get 89.3% ID (vs 87.9% for the
 85 standard model and 84.7% for the robust model) and 74.6% OOD (vs 74.1% for the robust model
 86 and 62.7% for the standard model).

87 **Competitive with self-training:** The remote sensing datasets have lots of unlabeled data so prior
 88 work (9) uses self-training on these datasets to mitigate the ID-OOD accuracy tradeoff. Table 3 shows
 89 that calibrated-ensembles match or outperform self-training on both datasets, both ID and OOD. We
 90 believe this is an interesting result because calibrated ensembling is a simple method and does not
 91 need additional unlabeled data.

92 **Calibration is important:** We find that a strong baseline of tuning the ensemble weights on ID data has
 93 lower accuracy than calibrated ensembles OOD (Table 4; calibrated ensembles average: 74.6%, tuned

Table 2: *Out-of-distribution (OOD)* accuracies across six datasets. Calibrated ensembling matches or outperforms the better model in 4/6 cases, and on average outperforms both the standard and robust models. For the remaining two datasets, DomainNet and ImageNet-R, calibrated ensembles close 96% and 93% of the gap between the standard and robust model.

	Ent30	DomNet	STL	Land	Crop	ImNet-R
Standard	60.7 (0.1)	55.3 (0.4)	82.4 (0.3)	55.7 (1.1)	85.6 (5.8)	36.2 (-)
Robust	63.2 (1.1)	87.2 (0.1)	85.1 (0.2)	60.4 (1.1)	89.8 (0.4)	59.1 (-)
Cal ensemble	64.8 (0.5)	85.9 (0.2)	87.3 (0.2)	60.8 (0.8)	91.3 (0.8)	57.4 (-)

Table 3: Calibrated ensembles are competitive with self-training (9), which requires unlabeled data.

	Cropland		Landcover	
	ID Acc	OOD Acc	ID Acc	OOD Acc
Standard model	95.3 (0.0)	85.6	76.9	55.7
Robust model	95.1 (0.1)	89.8	72.7	60.4
Self-training	95.3 (0.2)	90.6 (0.6)	77.0 (0.4)	61.0 (0.7)
Cal ensembling	95.6 (0.1)	91.0 (0.8)	77.1 (0.2)	60.8 (0.8)

Table 4: OOD accuracies: calibrated ensembles outperform vanilla ensembles and even tuned ensembles where the combination weights are tuned to maximize in-distribution accuracy. Averaged across the datasets, calibrated ensembles get an OOD accuracy of 74.6%, while tuned ensembles get an accuracy of 71.3%. The in-distribution accuracies of the methods are very close (within 0.2% of each other).

	Ent30	DomNet	STL	Land	Crop	ImNet-R
Vanilla	64.6 (0.4)	78.7 (1.3)	87.2 (0.2)	59.5 (1.0)	90.9 (0.2)	58.0 (-)
Tuned ID	64.6 (0.6)	86.3 (0.6)	85.7 (0.9)	58.7 (1.2)	87.3 (5.7)	45.4 (-)
Calibrated ID	64.8 (0.5)	85.9 (0.2)	87.3 (0.2)	60.8 (0.8)	91.3 (0.8)	57.4 (-)

94 ensembles: 71.3%) and calibrated ensembles only have a 0.2% drop in ID accuracy relative to tuned
 95 ensembles. Naturally, we expect the tuned ensemble to do the best ID since its weights are tailored
 96 for ID—what is surprising is that the calibrated ensembles do so much better OOD without using any
 97 OOD data either. Calibrated ensembles outperform vanilla ensembles both ID and OOD as well.

98 **Standard ensembles do not mitigate tradeoffs:** As we show in Figure 1, simply ensembling two
 99 standard models or two robust models (even with calibration) does not mitigate ID-accuracy tradeoffs.

100 5 Related Works and Discussion

101 **Calibration.** Calibration (in-distribution) has been widely studied (14; 12; 15; 16; 17; 18). Ovadia
 102 et al. (11) show that if we calibrate a model ID, it still has poor uncertainties OOD.

103 **Ensembling.** Typically ensemble members are trained on the same data with a different random
 104 seed (13) or augmentation (19)—in these settings prior work has shown that calibration does not
 105 help (20; 11). Indeed, calibration has minimal effect when we ensemble two standard, or two robust mod-
 106 els, but leads to clear improvements when we combine two very different models (standard and robust).

107 **Mitigating ID-OOD tradeoffs.** Prior work self-trains on large amounts of unlabeled data to mitigate
 108 ID-OOD tradeoffs (8; 9; 10). In concurrent work, Wortsman et al. (21) show on ImageNet and variants
 109 (e.g. ImageNet-R) that there *exists* a way to ensemble a CLIP zero-shot and fine-tuned model to get
 110 good ID and OOD accuracy—but this might require OOD data which is not available. In fact, we show
 111 that the natural way to *learn* how to weight ensemble members—selecting the weights to optimize
 112 in-distribution accuracy—does not mitigate the ID-OOD gap, but calibrated ensembles do.

References

- 113
- 114 [1] Christina Heinze-Deml and Nicolai Meinshausen. Conditional variance penalties and domain
115 shift robustness. *arXiv preprint arXiv:1710.11469*, 2017.
- 116 [2] Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally
117 robust neural networks for group shifts: On the importance of regularization for worst-case
118 generalization. In *International Conference on Learning Representations (ICLR)*, 2020.
- 119 [3] Yining Chen, Colin Wei, Ananya Kumar, and Tengyu Ma. Self-training avoids using spurious fea-
120 tures under domain shift. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- 121 [4] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk
122 minimization. *arXiv preprint arXiv:1907.02893*, 2019.
- 123 [5] Baochen Sun and Kate Saenko. Deep coral: Correlation alignment for deep domain adaptation.
124 In *European Conference on Computer Vision (ECCV)*, 2016.
- 125 [6] Tatsunori B. Hashimoto, Megha Srivastava, Hongseok Namkoong, and Percy Liang. Fairness
126 without demographics in repeated loss minimization. In *International Conference on Machine
127 Learning (ICML)*, 2018.
- 128 [7] John Duchi, Tatsunori Hashimoto, and Hongseok Namkoong. Distributionally robust losses
129 against mixture covariate shifts. [https://cs.stanford.edu/~thashim/assets/
130 publications/condrisk.pdf](https://cs.stanford.edu/~thashim/assets/publications/condrisk.pdf), 2019.
- 131 [8] Aditi Raghunathan, Sang Michael Xie, Fanny Yang, John C. Duchi, and Percy Liang.
132 Understanding and mitigating the tradeoff between robustness and accuracy. In *International
133 Conference on Machine Learning (ICML)*, 2020.
- 134 [9] Sang Michael Xie, Ananya Kumar, Robert Jones, Fereshte Khani, Tengyu Ma, and Percy Liang.
135 In-N-out: Pre-training and self-training using auxiliary information for out-of-distribution
136 robustness. In *International Conference on Learning Representations (ICLR)*, 2021.
- 137 [10] Fereshte Khani and Percy Liang. Removing spurious features can hurt accuracy and affect groups
138 disproportionately. In *ACM Conference on Fairness, Accountability, and Transparency (FAccT)*,
139 2021.
- 140 [11] Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, D Sculley, Sebastian Nowozin, Joshua V.
141 Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty?
142 evaluating predictive uncertainty under dataset shift. In *Advances in Neural Information
143 Processing Systems (NeurIPS)*, 2019.
- 144 [12] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural
145 networks. In *International Conference on Machine Learning (ICML)*, pages 1321–1330, 2017.
- 146 [13] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable
147 predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information
148 Processing Systems (NeurIPS)*, 2017.
- 149 [14] Mahdi Pakdaman Naeini, Gregory F. Cooper, and Milos Hauskrecht. Binary classifier calibration:
150 Non-parametric approach. *arXiv*, 2014.
- 151 [15] Ananya Kumar, Percy Liang, and Tengyu Ma. Verified uncertainty calibration. In *Advances
152 in Neural Information Processing Systems (NeurIPS)*, 2019.
- 153 [16] Allan H Murphy. A new vector partition of the probability score. *Journal of Applied Meteorology*,
154 12(4):595–600, 1973.
- 155 [17] Morris H. DeGroot and Stephen E. Fienberg. The comparison and evaluation of forecasters.
156 *Journal of the Royal Statistical Society. Series D (The Statistician)*, 32:12–22, 1983.

- 157 [18] Tilmann Gneiting and Adrian E. Raftery. Weather forecasting with ensemble methods. *Science*,
158 310, 2005.
- 159 [19] Asa Cooper Stickland and Iain Murray. Diverse ensembles improve calibration. *arXiv*, 2020.
- 160 [20] Xixin Wu and M. Gales. Should ensemble members be calibrated? *arXiv*, 2021.
- 161 [21] Mitchell Wortsman, Gabriel Ilharco, Mike Li, Jong Wook Kim, Hannaneh Hajishirzi, Ali Farhadi,
162 Hongseok Namkoong, and Ludwig Schmidt. Robust fine-tuning of zero-shot models. *arXiv*
163 *preprint arXiv:2109.01903*, 2021.
- 164 [22] Sherrie Wang, William Chen, Sang Michael Xie, George Azzari, and David B. Lobell. Weakly
165 supervised deep learning for segmentation of remote sensing imagery. *Remote Sensing*, 12, 2020.
- 166 [23] E. Vermote. MOD09A1 MODIS/terra surface reflectance 8-day L3 global 500m SIN grid V006.
167 <https://doi.org/10.5067/MODIS/MOD09A1.006>, 2015.
- 168 [24] Pall Oskar Gislason, Jon Atli Benediktsson, and Johannes R. Sveinsson. Random forests for
169 land cover classification. *Pattern Recognition Letters*, 27(4):294–300, 2006.
- 170 [25] Marc Rußwurm, Sherrie Wang, Marco Korner, and David Lobell. Meta-learning for few-shot
171 land cover classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
172 *Pattern Recognition Workshops*, pages 200–201, 2020.
- 173 [26] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
174 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.
175 Learning transferable visual models from natural language supervision. In *International*
176 *Conference on Machine Learning (ICML)*, volume 139, pages 8748–8763, 2021.
- 177 [27] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo,
178 Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin
179 Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization.
180 *arXiv preprint arXiv:2006.16241*, 2020.
- 181 [28] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
182 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. ImageNet large scale visual
183 recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.
- 184 [29] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation.
185 In *Association for Computational Linguistics (ACL)*, 2021.
- 186 [30] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe,
187 Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning
188 for NLP. *arXiv*, 2019.
- 189 [31] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment
190 matching for multi-source domain adaptation. In *International Conference on Computer Vision*
191 *(ICCV)*, 2019.
- 192 [32] Shuhan Tan, Xingchao Peng, and Kate Saenko. Class-imbalanced domain adaptation: An
193 empirical odyssey. *arXiv preprint arXiv:1910.10320*, 2020.
- 194 [33] Geoff French, Michal Mackiewicz, and Mark Fisher. Self-ensembling for visual domain
195 adaptation. In *International Conference on Learning Representations*, 2018.
- 196 [34] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report,
197 University of Toronto, 2009.
- 198 [35] Adam Coates, Andrew Ng, and Honlak Lee. An analysis of single-layer networks in unsupervised
199 feature learning. In *Proceedings of the Fourteenth International Conference on Artificial*
200 *Intelligence and Statistics*, volume 15, pages 215–223, 2011.

201 [36] Xinlei Chen, Haoqi Fan, Ross B. Girshick, and Kaiming He. Improved baselines with momentum
202 contrastive learning. *arXiv*, 2020.

203 [37] Shibani Santurkar, Dimitris Tsipras, and Aleksander Madry. Breeds: Benchmarks for
204 subpopulation shift. *arXiv*, 2020.

205 A Details on datasets

206 A.1 Spurious metadata

207 We run experiments on two remote sensing datasets used in prior work studying ID-OOD tradeoffs (9).
208 These datasets consist of a core input x (image data or time series data) and metadata z (e.g., location,
209 meteorological climate data). The metadata is spuriously correlated with the target—using the
210 metadata to predict labels improves accuracy in-distribution (ID), but hurts accuracy out-of-distribution.
211 Xie et al. (9) consider a standard model that takes in both the core inputs and metadata to predict the
212 target, and a robust model that only takes in the core inputs and does some additional pretraining. They
213 call these the ‘aux-in’ and ‘aux-out’ models respectively.

214 **Cropland.** The goal is to predict whether a satellite image is of a cropland or not. The core input
215 x is an RGB satellite image, and the metadata z consists of location coordinates and vegetation bands.
216 The original dataset is from Wang et al. (22), and we use U-net model checkpoints from Xie et al. (9).

217 **Landcover.** The goal is to predict the land type from satellite data at a given location. Here, the
218 core input x is a time series measured by NASA’s MODIS satellite (23), and z is climate data (e.g.,
219 temperature) at that location. The dataset is from Gislason et al. (24); Rußwurm et al. (25). We use
220 model checkpoints from Xie et al. (9) where they use 1D convolutions for time series data.

221 A.2 Zero-shot language prompting

222 Radford et al. (26) (CLIP) pretrain a model on a large multi-modal language and vision dataset.
223 The model can then predict the label of an image by comparing the image embedding, with the
224 language embedding for prompts such as ‘photo of an apple’ or ‘photo of a banana’. They show that
225 this zero-shot language prompting approach can be much more accurate out-of-distribution than the
226 traditional method of fine-tuning the entire model.

227 **ImageNet → ImageNet-R.** We use a CLIP vision transformer, specifically a ViT-B/16, which is the
228 best publicly available model. The robust model uses language prompts to make zero-shot predictions
229 on ImageNet-Renditions (27), a dataset containing cartoon, graffiti, video game, etc, renditions of
230 ImageNet classes. The standard model initializes with weights from the CLIP model, and fine-tunes
231 on ImageNet (28) training data for 10 epochs, before making predictions on ImageNet-R. We note
232 that the robust model gets 10% lower accuracy ID (on ImageNet validation examples), but gets 20%
233 higher accuracy OOD (on ImageNet-R test examples)

234 A.3 Freezing pretrained features

235 When adapting a pretrained model to an ID dataset, typically all the model parameters are fine-tuned.
236 Recent work looks at ‘lightweight’ fine-tuning, where only parts of the model are adapted—this can
237 often do better OOD even though the ID performance is worse (29; 30). We consider three distribution
238 shift datasets where the standard model starts from a pretrained initialization and fine-tunes all
239 parameters on an ID dataset, and the robust model only learns the top linear ‘head’ layer.

240 **DomainNet.** A standard domain adaptation dataset (31). Here, our ID dataset contains ‘sketch’
241 images (e.g., drawings of apples, elephants, etc), and the OOD dataset contains ‘real’ photos of the
242 same categories. We use the version of the dataset from Tan et al. (32). We start from a CLIP pretrained
243 ResNet50 and either fine-tune (to get a standard model) or train the head layer (to get a robust model).

244 **CIFAR-10** → **STL**. Another standard domain adaptation dataset (33), where the ID is CIFAR-
245 10 (34), and the OOD is STL (35). We start from a ResNet50 pretrained on unlabeled ImageNet
246 examples using MoCo-v2 (36) and either fine-tune (to get a standard model) or train the head layer
247 (to get a robust model).

248 **Living-17**. Part of the BREEDS benchmark (37), here the goal is to classify an image as one of 17
249 animal categories such as ‘bear’—the ID dataset contains images of black bears and sloth bears and
250 the OOD dataset has images of brown bears and polar bears. We start from a ResNet50 pretrained
251 on unlabeled ImageNet examples using MoCo-v2 (36) and either fine-tune (to get a standard model)
252 or train the head layer (to get a robust model).