

## 456 A Description of Baselines

- 457 • Empirical risk minimization (ERM) [49] minimizes the sum of errors across domains and  
458 examples.
- 459 • Invariant risk minimization (IRM) [1] learns a feature representation such that the optimal  
460 linear classifier on top of that representation matches across domains. For WILDS datasets,  
461 we pull baseline performance from [51]. For ISIC, we use the implementation from [16].
- 462 • Deep CORAL [45] penalizes differences in the means and covariances of the feature  
463 distributions (i.e., the distribution of last layer activations in a neural network) for each  
464 domain. For WILDS datasets, we pull baseline performance from [51]. For ISIC, we use  
465 the implementation from [16].
- 466 • Fish targets domain generalization by maximizing the inner product between gradients from  
467 different domains. For WILDS datasets, we pull baseline performance from the original  
468 paper. For ISIC, we use the implementation from [16].
- 469 • LISA augments the set of training data by randomly performing two types of mixup-  
470 style [64] interpolations: intra-label (same label, different domain) and inter-label (same  
471 domain, different label). For WILDS datasets, we pull baseline performance from the  
472 original paper, and train their implementation on the ISIC dataset.
- 473 • CLOvE [51] finds an invariant classifier by enforcing the classifier to be calibrated across  
474 all training domains. While the original paper proposes several model variants leveraging  
475 this idea, we report their best-performing variant, which starts with a trained CORAL  
476 model and finetunes the weights using a regularized cross-entropy loss. The regularizer  
477 aggregates Maximum Mean Calibration Error (MMCE) [27] over all training domains. For  
478 WILDS datasets, we pull baseline performance from [51]. As their implementation is not  
479 publicly-available, we implement it for ISIC.

## 480 B Description of Datasets

481 Representative examples of the 3 datasets are shown in Fig. 4.

- 482 • **Camelyon-17.** We use Camelyon-17 from the WILDS benchmark [4, 23], which provides  
483 450,000 lymph-node scans sampled from 5 hospitals. Camelyon-17 is a medical image  
484 classification task where the input  $x$  is a  $96 \times 96$  image and the label  $y$  is whether there  
485 exists tumor tissue in the image. The environment denotes the hospital that the patch was  
486 taken from. The training dataset is drawn from the first 3 hospitals, while out-of-distribution  
487 validation and out-of-distribution test datasets are sampled from the 4th hospital and 5th  
488 hospital, respectively.
- 489 • **ISIC.** The melanoma dataset is from the International Skin Imaging Collaboration (ISIC)  
490 archive<sup>7</sup>. Data from the archive are collected by different organizations at different points in  
491 time [7, 8, 9, 17, 41, 43, 47]. There are about 70k data samples in total. In particular, the  
492 resized input image  $x$  is a  $224 \times 224$  image and a binary target label  $y$  denotes whether the  
493 image exhibit is melanoma or not. The environment is the hospital from which the image  
494 was collected<sup>8</sup>. We follow a similar setup to Camelyon-17. The training dataset is drawn  
495 from the first 3 hospitals, while out-of-distribution validation and out-of-distribution test  
496 datasets are sampled from the 4th hospital and 5th hospital, respectively. For preprocessing,  
497 we filter out datapoints that are not specifically categorized as "benign" or "malignant" (e.g.  
498 "indeterminate"). The OOD validation dataset is from the "Barcelona1" site indicator and  
499 the OOD test dataset is the "Vienna1" site indicator.
- 500 • **FMoW.** The FMoW dataset is from the WILDS benchmark [6, 23], a satellite image classi-  
501 fication task which includes 62 classes and 80 domains (16 years x 5 regions). Concretely,  
502 the input  $x$  is a  $224 \times 224$  RGB satellite image, the label  $y$  is one of the 62 building or land  
503 use categories, and the environment represents the year that the image was taken as well

<sup>7</sup><https://www.isic-archive.com>

<sup>8</sup>Hospitals are Hospital Clinic of Barcelona, Medical University of Vienna, University of Queensland Diamantina Institute, Memorial Sloan Kettering Cancer Center, University of Sydney Melanoma Diagnostic Centre, and University of Pittsburgh Medical Center.

504  
505  
506  
507

as its corresponding geographical region – Africa, the Americas, Oceania, Asia, or Europe. The train/test/validation splits are based on the time when the images are taken. Specifically, images taken before 2013 are used as the training set. Images taken between 2013 and 2015 are used as the validation set. Images taken after 2015 are used for testing.

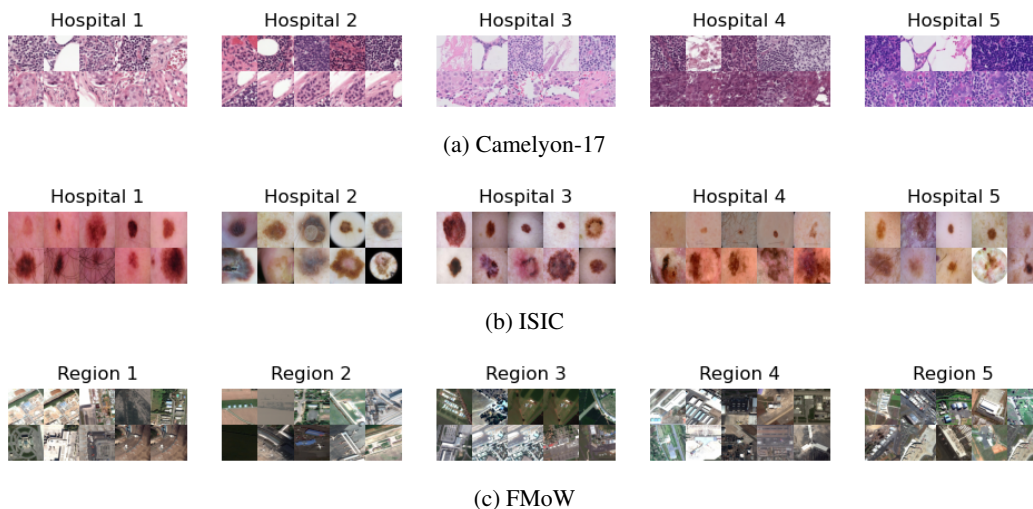


Figure 4: Representative images for datasets, separated by domain. Each row depicts a separate class. For FMoW, for simplicity, we show 2 classes out of 62 and only images before 2013.

508 **C Hyperparameter Details**

509 Table 3 shows hyperparameter settings for all datasets, where NW-specific hyperparameters are below  
 510 the midline. For all models, we use pretrained ImageNet weights. For  $\lambda$ , we perform a grid-search  
 511 over the values  $\{0.01, 0.1, 1\}$ . Fig. 5 depicts  $NW_e^B$  performance vs  $N_c$  for Camelyon-17 and ISIC  
 512 datasets. We find that performance is relatively insensitive to  $N_c$  above  $\sim 5$  examples per class.

Table 3: Hyperparameter settings for various datasets.

<i>Hyperparameter</i>	<i>Camelyon-17</i>	<i>ISIC</i>	<i>FMoW</i>
Learning rate	1e-4	5e-5	1e-4
Weight decay	1e-4	0	1e-2
Scheduler	None	None	StepLR
Batch size	32	8	8
Architecture	DenseNet-121	ResNet-50	DenseNet-121
Optimizer	SGD	Adam	Adam
Maximum Epoch	10	5	60
$N_q$	8	8	8
$N_c$	8	8	1
$N_s$	$N_c \times 2 = 16$	$N_c \times 2 = 16$	$N_c \times 62 = 62$
$\lambda$	0.01	0.01	0.1
$k$	3	3	3

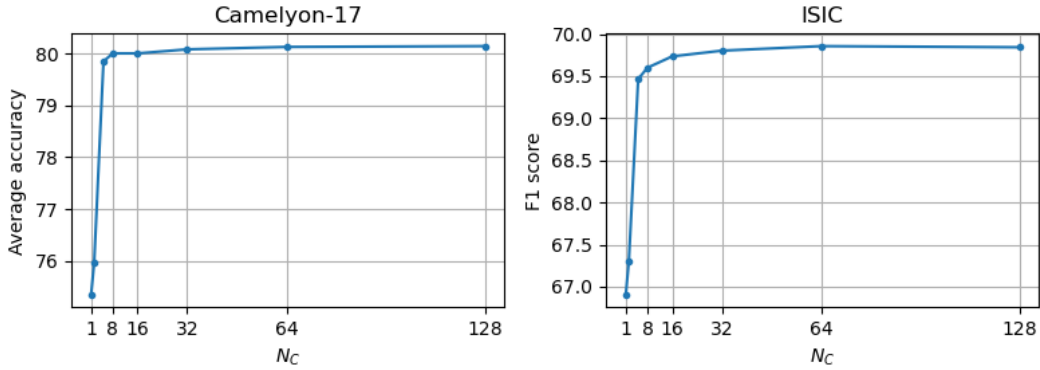


Figure 5:  $NW_e^B$  performance vs  $N_c$  for Camelyon-17 and ISIC datasets. Full mode. Performance is relatively insensitive to  $N_c$  above  $\sim 5$  examples per class.

513 **D Table of Runtimes**

514 Table 4 shows approximate runtimes for various datasets during training and inference. All experiments are performed on a GPU.

Table 4: Approximate runtimes for various algorithms. Training time is time to complete maximum epochs as specified in Table 3, and does not include validation. Inference time is time to evaluate the entire test set. Averaged over all training runs.

	<i>Algorithm</i>	<i>Camelyon-17</i>	<i>ISIC</i>	<i>FMoW</i>
Training	ERM	7 hr	1 hr	22 hr
	NW	14 hr	2 hr	40 hr
Inference	ERM	10 min	2 min	10 min
	NW, Random	15 min	3 min	20 min
	NW, Full	2 hr	15 min	1 hr
	NW, Ensemble	2 hr	15 min	1 hr
	NW, Cluster	2.2 hr	17 min	1.1 hr
	NW, Probe	10 min	2 min	10 min

515

516 **E Imbalanced ISIC Experiments**

517 As ISIC exhibits significant label imbalance where the positive class is much less represented than  
 518 the negative class (see Fig. 6), we experiment with an NW variant without support-set label balancing  
 519 as an ablation. For both variants, we set  $N_c = 8$ . To train the imbalanced variant, we sample a  
 520 mini-batch support set by first sampling one image per class (to guarantee both classes are represented  
 521 in the support at least once), and then sampling the rest of the images randomly from the dataset.

522 To characterize the performance of both variants across class imbalances, we change the prevalence  
 523 of  $y = 0$  in the test set by removing negative class images until the desired prevalence is achieved.  
 524 Note the default prevalence is  $\sim 0.85$ . Then, we compute the accuracy over this manipulated test set  
 525 (note, prior work has shown that F1-score is not a good metric for comparing classifiers with different  
 526 label imbalances [3]).

527 Fig 7 shows results. We observe that at high prevalence where the proportion of negative classes  
 528 matches in training and test, the imbalanced variant outperforms the balanced variant, whereas the  
 529 opposite is true for low prevalence. This makes sense because at high prevalence, the class imbalance  
 530 is similar for the training and test domains; thus, a model which overpredicts the negative class  
 531 is usually right. On the other hand, this fails in test sets where the prevalence is flipped (i.e. low  
 532 prevalence). These results suggest that NW<sup>B</sup> is a more robust classifier in the presence of label shift.

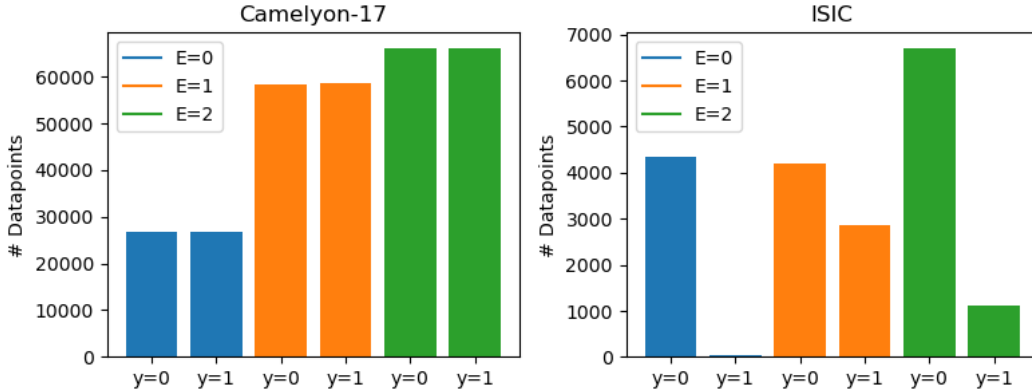


Figure 6: Number of datapoints separated by class for Camelyon-17 and ISIC datasets. There is significant label imbalance for the ISIC dataset.

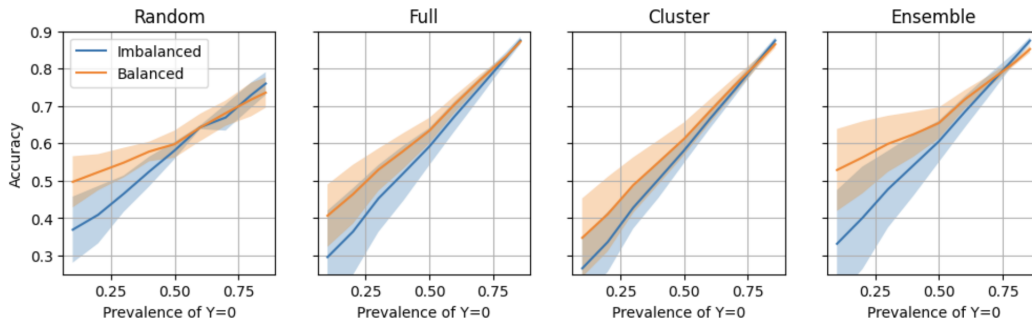


Figure 7: Accuracy of NW (imbalanced) and NW<sup>B</sup> (balanced) models over varying prevalence of  $y = 0$  for ISIC dataset. At low prevalence where the prevalence differs the most from training domains, we observe that model performance is higher for NW<sup>B</sup>. The default prevalence is  $6705/7818 = 0.8576$ , which is the right-most value in the x-axis.