



Task	Source val	DEV	Entropy	BNM	EnsV	CODA
DomainNet26	real → sketch	0.5	4.6	3.1	<b>0.4</b>	1.0
	real → clip	4.5	49.3	0.3	6.5	2.8
	real → paint	1.3	34.7	2.3	2.0	<b>0.2</b>
	sketch → real	4.7	6.3	9.1	8.5	<b>0.1</b>
	sketch → clip	2.4	13.4	5.9	6.4	<b>3.0</b>
	sketch → paint	3.8	3.5	3.2	3.2	<b>0.1</b>
	clip → real	1.5	6.3	5.6	5.6	<b>0.3</b>
	clip → sketch	2.8	4.9	5.1	4.9	<b>2.2</b>
	clip → paint	3.1	7.8	4.2	4.2	<b>1.2</b>
	paint → real	0.0	36.6	3.3	3.3	<b>0.1</b>
WILDS	paint → sketch	0.7	26.1	4.3	4.3	<b>0.9</b>
	paint → clip	2.9	16.2	6.3	6.3	<b>0.4</b>
MSV	iwildcam	9.8	-	-	-	<b>0.9</b>
	camelyon	<b>4.3</b>	-	-	-	13.1
	fmow	0.8	-	-	-	0.9
	civilcomments	-	-	-	-	<b>3.8</b>
GLUE	cifar10-low	-	-	-	-	<b>2.3</b>
	cifar10-high	-	-	-	-	<b>4.9</b>
	pacs	-	-	-	-	<b>0.4</b>
	cola	-	-	-	-	5.3
	mnli	-	-	-	-	<b>3.0</b>
	qnli	-	-	-	-	<b>3.3</b>
	qqp	-	-	-	-	<b>1.1</b>
RTE	rte	-	-	-	-	<b>14.8</b>
	sst2	-	-	-	-	<b>3.6</b>
	mrpc	-	-	-	-	<b>1.0</b>
	mrcp	-	-	-	-	<b>1.0</b>

**Table 5. Unsupervised model selection results.** We report regret at step 0 (lower is better) for all methods on all tasks. Best method for each task is in bold. CODA matches or exceeds state-of-the-art performance on 20 out of 26 tasks. Note that because models/predictions for MSV and GLUE are black-box/one-hot (as in ModelSelector [45]), the only comparison we are able to make is to EnsV.

idation accuracy on “source” data, *i.e.* using a validation set from the same distribution as the training set.

**Target entropy** [43] Shannon entropy is computed on prediction scores on the test set. Model selection is performed by selecting the model with the lowest entropy, *i.e.* the model that is most confident in its test predictions.

**Deep embedded validation** [70] Computes a classification loss for each source validation sample, and weights each loss based on a computed probability that the sample “belongs to” the test domain. This probability comes from a separate domain classifier trained on source and test data.

**Batch nuclear norm** [11] Performs singular value decomposition on the prediction matrix. Model selection is performed by selecting the model with the minimum nuclear norm (*i.e.* minimum sum of singular values).

**EnsV** [20] All models in the hypothesis set are ensembled. To perform model selection, accuracy is estimated for each model with respect to the ensemble’s predictions.

### 9.3. Unaggregated results on all tasks

In Fig. 8 we visualize regret and cumulative regret at every step for all baselines on all tasks.

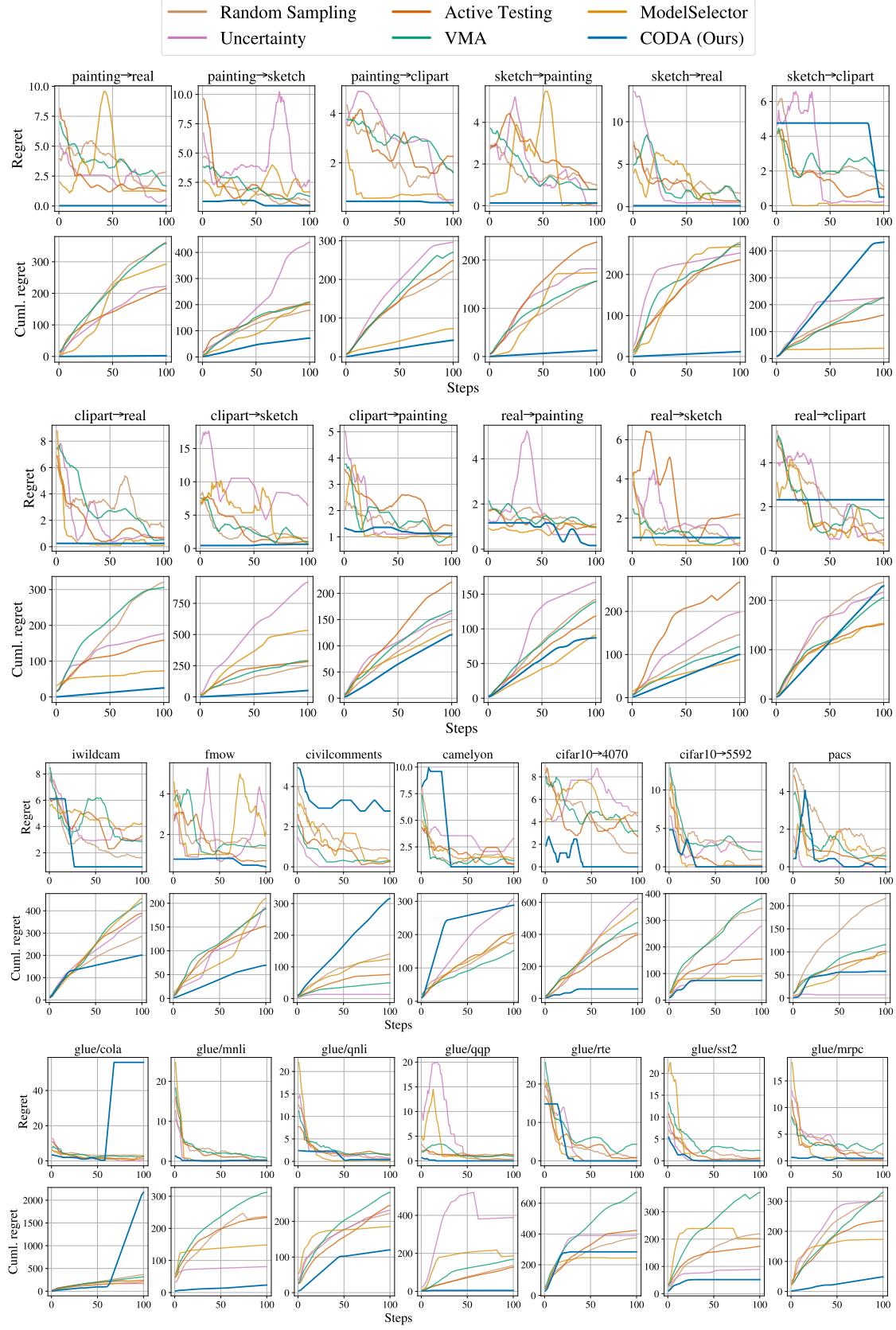


Figure 8. Results on all benchmarks.

## 10. Implementation details

### 10.1. Dawid-Skene data generating process

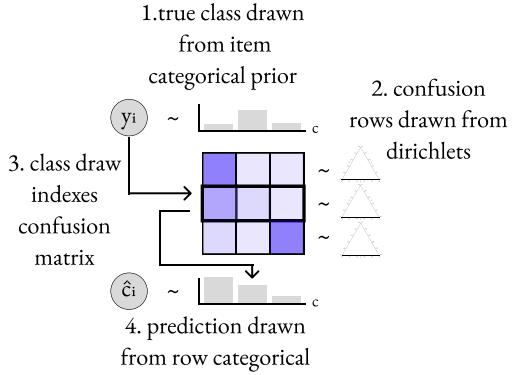


Figure 9. A visual depiction of the Dawid-Skene [14] data generating process that we adapt to active model selection. See Sec. 4.1 of the main paper for more details.

We visualize the data generating process of the Bayesian implementation of the Dawid-Skene model in Fig. 9, as described in Sec. 4.1. We repeat the text here for easy reference.

The data generating process proceeds as follows:

1. Each data point’s true class label  $y_i$  is drawn randomly from per-data-point prior distributions over which class that data point could be,  $y_i \sim \text{Cat}(\pi(x_i))$ .
2. Each row of the classifier’s confusion matrix is drawn randomly from per-row distributions,  $M_{k,c} \sim \theta_{k,c}$ , where  $\theta_{k,c}$  is the prior distribution over what the row of the confusion matrix could be. To accommodate Bayesian updates, we initialize each  $\theta_{k,c}$  to be a Dirichlet prior.
3. The sampled true class indexes into the corresponding row of the classifier’s confusion matrix,  $M_{k,y_i}$ .
4. The classifier’s prediction for that data point is sampled from the distribution over that row’s cells,  $\hat{c}_{k,i} \sim \text{Cat}(M_{k,y_i})$ .

### 10.2. Computing $P_{\text{Best}}$

We illustrate visually the computation of  $P_{\text{Best}}$  from Sec. 4.3 in Fig. 10 in the simplified case of two models. To compute the probability that Model 1 is best, we integrate over all possible accuracy values that Model 1 *could* have. For every possible accuracy, we compute the probability that Model 1 has that accuracy (defined by its PDF  $f$ ), multiplied that the probability Model 2 has accuracy *less* than that value (defined by its CDF  $F$ ).

## 11. Data and model details

We provide more details about the datasets and models in our benchmarking suite in Tab. 6, Tab. 7, and Tab. 8.

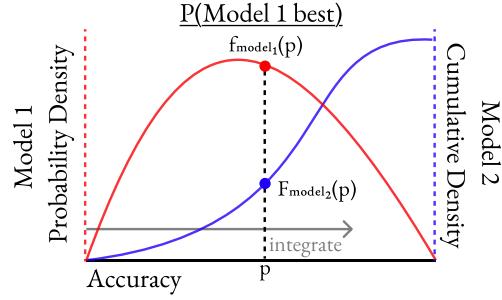


Figure 10. A visual depiction of the integration technique used to construct our distribution over which model is best,  $P_{\text{Best}}$ . See Sec. 10.2 of the supplemental and Sec. 4.3 of the main paper for more details.

DomainNet126				
Task	Num. Classes	Num. Checkpoints	Test set size	
sketch_painting	126	10 algs 20 epochs	= 200	3086
sketch_real	126	10 algs 20 epochs	= 200	20939
sketch_clipart	126	10 algs 20 epochs	= 200	5611
clipart_painting	126	10 algs 20 epochs	= 200	3086
clipart_sketch	126	10 algs 20 epochs	= 200	7313
clipart_real	126	10 algs 20 epochs	= 200	20939
painting_sketch	126	10 algs 20 epochs	= 200	7313
painting_real	126	10 algs 20 epochs	= 200	20939
painting_clipart	126	10 algs 20 epochs	= 200	5611
real_painting	126	10 algs 20 epochs	= 200	3086
real_sketch	126	10 algs 20 epochs	= 200	7313
real_clipart	126	10 algs 20 epochs	= 200	5611

Table 6. **Dataset details: DomainNet126.** For each transfer task, we first train a “source-only” model on the source domain. We then train 10 UDA models for each transfer task (source domain → target domain) using the Powerful Benchmark codebase [43].

<b>WILDS</b>					
Task	Num. Classes	Num. Checkpoints	Test set size	Regret w/ val.	
RxRx1	1139	4 algs 18 epochs = 72	34,432	0.1	
Amazon	5	4 algs 3 epochs = 12	100,050	0.1	
CivilComments	2	4 algs 5 epochs = 20	133,782	N/A	
fMoW	62	4 algs 60 epochs = 240	22,108	0.8	
iWildCam	182	4 algs 12 epochs = 48	42,791	9.8	
Camelyon17	2	4 algs 10 epochs = 40	85,054	4.3	

Table 7. **Dataset details: WILDS.** We show metrics for all classification datasets in WILDS where we could perform model selection. In our main experiments, we only use the benchmarks where near-perfect model selection can be performed trivially using the default in-distribution validation set. We train all models ourselves using the public code from Koh et al. [30].

<b>MSV</b>				
Dataset	Num. Classes	Num. Checkpoints	Test set size	
CIFAR10-High	10	80	10000	
CIFAR10-Low	10	80	10000	
PACS	7	30	9991	

<b>GLUE</b>				
Dataset	Num. Classes	Num. Checkpoints	Test set size	
CoLA	2	109	1043	
MNLI	3	82	9815	
QNLI	2	90	5463	
QQP	2	101	40430	
RTE	2	87	277	
SST2	2	97	872	
MRPC	2	95	408	

Table 8. **Dataset details: MSV and GLUE.** All model checkpoints sourced directly from Okanovic et al. [45].