

# Knockoff Nets: Stealing Functionality of Black-Box Models (Supplementary material)

Tribhuvanesh Orekondy<sup>1</sup>

Bernt Schiele<sup>1</sup>

Mario Fritz<sup>2</sup>

<sup>1</sup> Max Planck Institute for Informatics    <sup>2</sup> CISPA Helmholtz Center for Information Security  
Saarland Informatics Campus, Germany

## A. Contents

The supplementary material contains:

### A. Contents (this section)

### B. Extended Descriptions

1. Blackbox Models
2. Overlap Between  $P_V$  and  $P_A$
3. Aggregating OpenImages and OpenImages-Faces
4. Additional Implementation Details

### C. Extensions of Existing Results

1. Qualitative Results
2. Sample-efficiency of GT
3. Policies Learnt by adaptive Strategy
4. Reward Ablation

### D. Auxiliary experiments

1. Effect of CNN initialization
2. Seen and Unseen Classes
3. Adaptive Strategy: With/without hierarchy
4. Semi-open World:  $\tau D^2$

## B. Extended Descriptions

In this section, we provide additional detailed descriptions and implementation details.

### B.1. Black-box Models

We supplement Section 5.1 by providing extended descriptions of the blackboxes listed in Table 1 of the main paper. Each blackbox  $F_V$  is trained on one particular image classification dataset.

$P_A$	$P_V$			
	Caltech256 ( $K=256$ )	CUBS200 ( $K=200$ )	Indoor67 ( $K=67$ )	Diabetic5 ( $K=5$ )
ILSVRC ( $Z=1000$ )	108 (42%)	2 (1%)	10 (15%)	0 (0%)
OpenImages ( $Z=601$ )	114 (44%)	1 (0.5%)	4 (6%)	0 (0%)

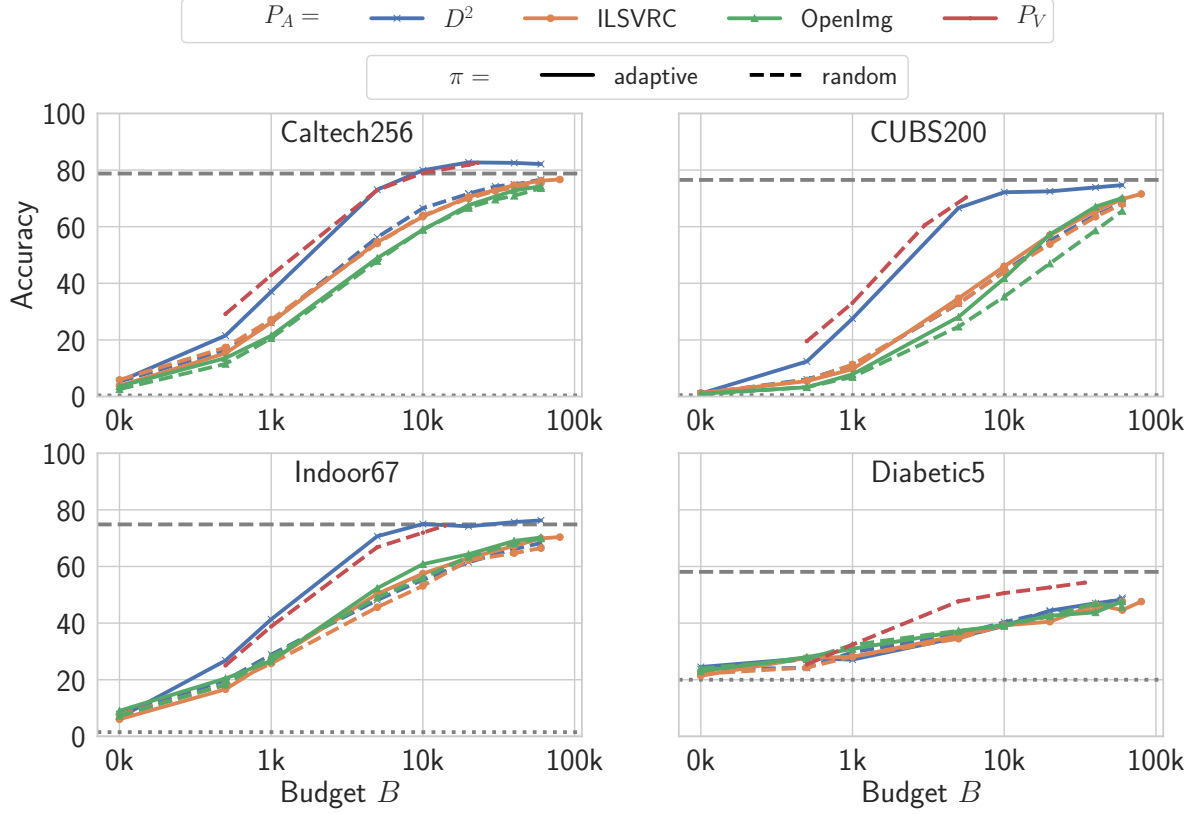
Table S1: Overlap between  $P_A$  and  $P_V$ .

**Black-box 1: Caltech256 [5].** Caltech-256 is a popular dataset for general object recognition gathered by downloading relevant examples from Google Images and manually screening for quality and errors. The dataset contains 30k images covering 256 common object categories.

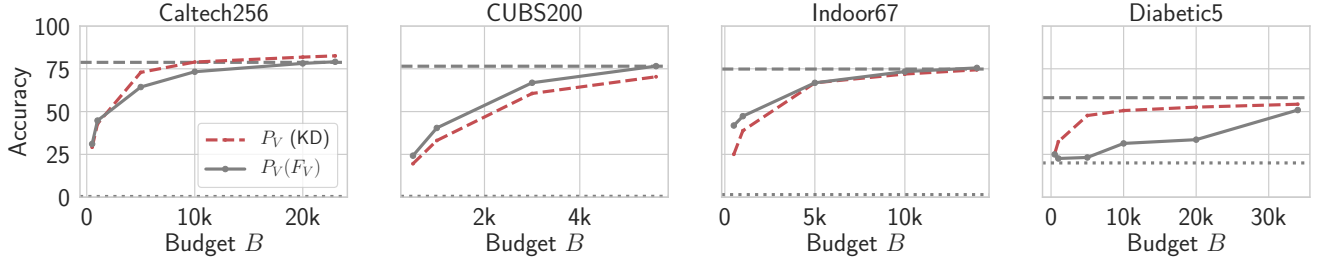
**Black-box 2: CUBS200 [14].** A fine-grained bird-classifier is trained on the CUBS-200-2011 dataset. This dataset contains roughly 30 train and 30 test images for each of 200 species of birds. Due to the low intra-class variance, collecting and annotating images is challenging even for expert bird-watchers.

**Black-box 3: Indoor67 [11].** We introduce another fine-grained task of recognizing 67 types of indoor scenes. This dataset consists of 15.6k images collected from Google Images, Flickr, and LabelMe.

**Black-box 4: Diabetic5 [1].** Diabetic Retinopathy (DR) is a medical eye condition characterized by retinal damage due to diabetes. Cases are typically determined by trained clinicians who look for presence of lesions and vascular abnormalities in digital color photographs of the retina captured using specialized cameras. Recently, a dataset of such 35k retinal image scans was made available as a part of a Kaggle competition [1]. Each image is annotated by a clinician on a scale of 0 (no DR) to 4 (proliferative DR). This highly-specialized biomedical dataset also presents challenges in the form of extreme imbalance (largest class contains  $30\times$  as the smallest one).



**Figure S1: Performance of the knockoff at various budgets.** (Enlarged version of Figure 5) Presented for various choices of adversary’s image distribution ( $P_A$ ) and sampling strategy  $\pi$ . - represents accuracy of blackbox  $F_V$  and .... represents chance-level performance.



**Figure S2: Training on GT vs. KD.** Extension of Figure 5. We compare sample efficiency of first two rows in Table 2: “ $P_V(F_V)$ ” (training with GT data) and “ $P_V(KD)$ ” (training with soft-labels of GT images produced by  $F_V$ )

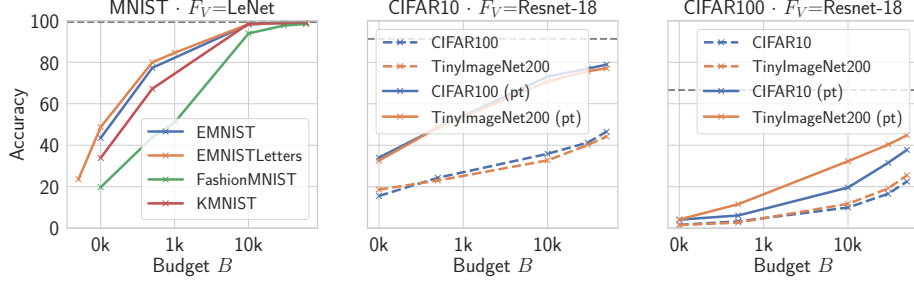
## B.2. Overlap: Open-world

In this section, we supplement Section 5.2.1 in the main paper by providing more details on how overlap was calculated in the open-world scenarios. We manually compute overlap between labels of the blackbox ( $K$ , e.g., 256 Caltech classes) and the adversary’s dataset ( $Z$ , e.g., 1k ILSVRC classes) as:  $100 \times |K \cap Z|/|K|$ . We denote two labels  $k \in K$  and  $z \in Z$  to overlap if: (a) they have the same semantic meaning; or (b)  $z$  is a type of  $k$  e.g.,  $z$  = “maltese dog” and  $k$  = “dog”. The exact numbers are

provided in Table S1. We remark that this is a soft-lower bound. For instance, while ILSVRC contains “Humming-bird” and CUBS-200-2011 contains three distinct species of hummingbirds, this is not counted towards the overlap as the adversary lacks annotated data necessary to discriminate among the three species.

## B.3. Dataset Aggregation

All datasets used in the paper (except OpenImages) have been used in the form made publicly available by the authors. We use a subset of OpenImages due to storage con-



**Figure S3: Training with non-ImageNet initializations of knockoff models.** Shown for various choices of blackboxes  $F_V$  (subplots) and adversary’s image distribution  $P_A$  (lines). All victim blackbox models are trained from scratch; test accuracy indicated by ---. All knockoff models are either trained from scratch, or pretrained on the corresponding  $P_A$  task (suffixed with ‘(pt)’).

straints imposed by its massive size (9M images). The description to obtain these subsets are provided below.

**OpenImages.** We retrieve 2k images for each of the 600 OpenImages [8] “boxable” categories, resulting in 554k unique images.  $\sim 19$ k images are removed for either being corrupt or representing Flickr’s placeholder for unavailable images. This results in a total of 535k unique images.

**OpenImages-Faces.** We download all images (422k) from OpenImages [8] with label “/m/0dzct: Human face” using the OID tool [13]. The bounding box annotations are used to crop faces (plus a margin of 25%) containing at least  $180 \times 180$  pixels. We restrict to at most 5 faces per image to maintain diversity between train/test splits. This results in a total of 98k faces images.

#### B.4. Additional Implementation Details

In this section, we provide implementation details to supplement discussions in the main paper.

**Input Transformations.** While training the blackbox models  $F_V$  we augment training data by applying input transformations: random  $224 \times 224$  crops and horizontal flips. This is followed by performing normalizing the image using standard Imagenet mean and standard deviation values. While training the knockoff model  $F_A$  and for evaluation, we resize the image to  $256 \times 256$ , obtain a  $224 \times 224$  center crop and normalize as before.

**Training  $F_V = \text{Diabetic5}$ .** We train this model using a learning rate of 0.01 (while this is 0.1 for the other models) and a weighted loss. Due to the extreme imbalance between classes of the dataset, we weigh each class as follows. Let  $n_k$  denote the number of images belonging to class  $k$  and let  $n_{\min} = \min_k n_k$ . We weigh the loss for each class  $k$  as  $n_{\min}/n_k$ . From our experiments with weighted loss, we found approximately 8% absolute improvement in overall accuracy on the test set. However, the training of knockoffs of all blackboxes are identical in all aspects, including a non-weighted loss irrespective of the victim blackbox targeted.

**Creating ILSVRC Hierarchy.** We represent the 1k labels of ILSVRC as a hierarchy (Figure 4b) in the form: root node “entity”  $\rightarrow N$  coarse nodes  $\rightarrow$  1k leaf nodes. We obtain  $N$  (30 in our case) coarse labels as follows: (i) a 2048-d mean feature vector representation per 1k labels is obtained using an Imagenet-pretrained ResNet; (ii) we cluster the 1k features into  $N$  clusters using scikit-learn’s [10] implementation of agglomerative clustering; (iii) we obtain semantic labels per cluster (i.e., coarse node) by finding the common parent in the Imagenet semantic hierarchy.

**Adaptive Strategy.** Recall from Section 6, we train the knockoff in two phases: (a) *Online*: during transfer set construction; followed by (b) *Offline*: the model is retrained using transfer set obtained thus far. In phase (a), we train  $F_A$  with SGD (with 0.5 momentum) with a learning rate of 0.0005 and batch size of 4 (i.e., 4 images sampled at each  $t$ ). In phase (b), we train the knockoff  $F_A$  from scratch on the transfer set using SGD (with 0.5 momentum) for 100 epochs with learning rate of 0.01 decayed by a factor of 0.1 every 60 epochs. We used  $\Delta=25$ .

### C. Extensions of Existing Results









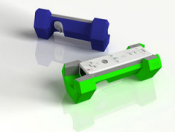









In this section, we present extensions of existing results discussed in the main paper.

#### C.1. Qualitative Results
















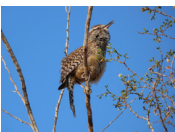


Qualitative results to supplement Figure 6 are provided in Figures S4-S7. Each row in the figures correspond to an output class of the blackbox whose images the knockoff has never encountered before. Images in the “transfer set” column were randomly sampled from ILSVRC [4, 12]. In contrast, images in the “test set” belong to the victim’s test set (Caltech256, CUBS-200-2011, etc.).

#### C.2. Sample Efficiency: Training Knockoffs on GT

We extend Figure 5 in the main paper to include training on the same ground-truth data used to train the blackboxes.








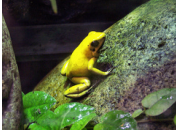










Transfer Set			Test Set		
					
Buddha: 0.65 Boxing glove: 0.07 Jesus Christ: 0.034	Buddha: 0.49 Cake: 0.06 Minotaur: 0.05	Buddha: 0.28 Minotaur: 0.09 Dog: 0.09	<u>Buddha: 0.68</u> Tombstone: 0.07 Elephant: 0.04	<u>Buddha: 0.5</u> Minaret: 0.08 Tower Pisa: 0.05	People: 0.26 <u>Buddha: 0.25</u> T-shirt: 0.14
					
Floppy-disk: 0.77 Socks: 0.04 Mattress: 0.03	Floppy-disk: 0.46 iPod: 0.18 CD: 0.06	Floppy-disk: 0.18 Flashlight: 0.08 Pez-dispenser: 0.07	<u>Floppy-disk: 0.74</u> Necktie: 0.04 Video Projec.: 0.02	<u>Floppy-disk: 0.61</u> Socks: 0.02 Teddy bear: 0.02	iPod: 0.49 <u>Floppy-disk: 0.08</u> CD: 0.04
					
Tomato: 0.96 Grapes: 0.01 Boxing glove: 0.00	Tomato: 0.43 Welding mask: 0.12 Bowling ball: 0.01	Tomato: 0.27 Dice: 0.18 Stained glass: 0.01	<u>Tomato: 0.94</u> Boxing glove: 0.04 Bowling ball: 0.00	<u>Tomato: 0.41</u> Cactus: 0.07 Iguana: 0.6	Mushroom: 0.17 <u>Tomato: 0.10</u> Birdbath: 0.07

**Figure S4: Qualitative results: Caltech256.** Extends Figure 6 in the main paper. GT labels are underlined, **correct** knockoff top-1 predictions in green and **incorrect** in red.




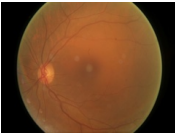
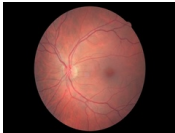






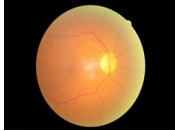



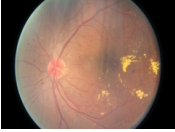


Transfer Set			Test Set		
					
Gadwall: 0.65 Nighthawk: 0.06 Horned Lark: 0.05	Gadwall: 0.31 B. Swallow: 0.14 N. Flicker: 0.12	Gadwall: 0.14 Chuck. Widow: 0.13 Swain. Warbler: 0.1	<u>Gadwall: 0.95</u> Mallard: 0.01 Rb. Merganser: 0.00	<u>Gadwall: 0.44</u> Mallard: 0.15 Rb. Merganser: 0.11	Pom. Jaeger: 0.16 Black Tern: 0.11 Herring Gull: 0.07
					
Lin. Sparrow: 0.82 Ovenbird: 0.07 House Sparrow: 0.03	Lin. Sparrow: 0.49 Mockingbird: 0.18 N. Waterthrush: 0.07	Lin. Sparrow: 0.32 Song Sparrow: 0.06 Tree Sparrow: 0.05	<u>Lin. Sparrow: 0.64</u> Hen. Sparrow: 0.05 Clay c. Sparrow: 0.03	<u>Lin. Sparrow: 0.50</u> Song Sparrow: 0.24 Hen. Sparrow: 0.04	Hen. Sparrow: 0.28 Ovenbird: 0.11 <u>Lin. Sparrow: 0.10</u>
					
Cactus Wren: 0.95 W. Meadowlark: 0.02 Lin. Sparrow: 0.00	Cactus Wren: 0.88 Geococcyx: 0.04 W. Meadowlark: 0.03	Cactus Wren: 0.33 N. Flicker: 0.28 Lin. Sparrow: 0.07	<u>Cactus Wren: 0.86</u> Rock Wren: 0.02 R. Blackbird: 0.01	<u>Cactus Wren: 0.82</u> N. Flicker: 0.02 Geococcyx: 0.02	Geococcyx: 0.25 <u>Cactus Wren: 0.20</u> Nighthawk: 0.08

**Figure S5: Qualitative results: CUBS200.** Extends Figure 6 in the main paper.

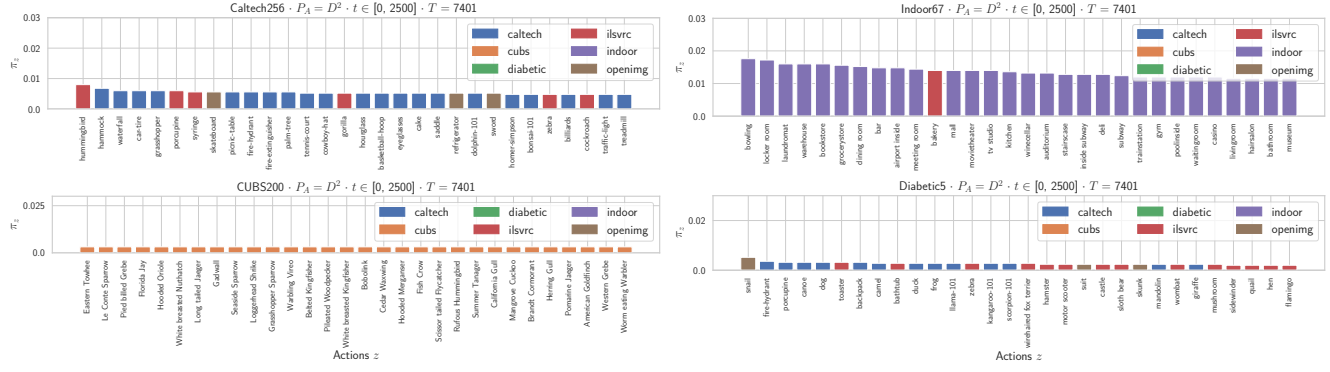


Transfer Set			Test Set		
					
Casino: 0.79 Deli: 0.03 Jewel. Shop: 0.03	Casino: 0.38 Airport ins.: 0.35 Bowling: 0.05	Casino: 0.18 Bar: 0.18 Movie theater: 0.13	<u>Casino: 0.99</u> Deli: 0.00 Toystore: 0.00	<u>Casino: 0.88</u> Toystore: 0.08 Bar: 0.01	Restaurant: <u>0.46</u> Bar: 0.24 Airport ins.: 0.07
					
Ins. Subway: 0.87 Video store: 0.11 Dental office: 0.01	Ins. Subway: 0.59 Florist: 0.12 Greenhouse: 0.06	Ins. Subway: 0.37 Airport ins.: 0.12 Train station: 0.08	<u>Ins. Subway: 0.96</u> Casino: 0.02 Museum: 0.01	Ins. Subway: <u>0.86</u> Train station: 0.07 Subway: 0.05	Corridor: <u>0.45</u> Ins. Subway: 0.11 Bar: 0.09
					
Prison cell: 0.52 Elevator: 0.20 Airport ins.: 0.11	Prison cell: 0.23 Museum: 0.19 Nursery: 0.17	Prison cell: 0.21 Museum: 0.12 Airport ins.: 0.11	<u>Prison cell: 0.83</u> Kitchen: 0.03 Locker room: 0.03	<u>Prison cell: 0.52</u> Subway: 0.08 Nursery: 0.07	Wine cellar: <u>0.31</u> <u>Prison cell: 0.17</u> Staircase: 0.09

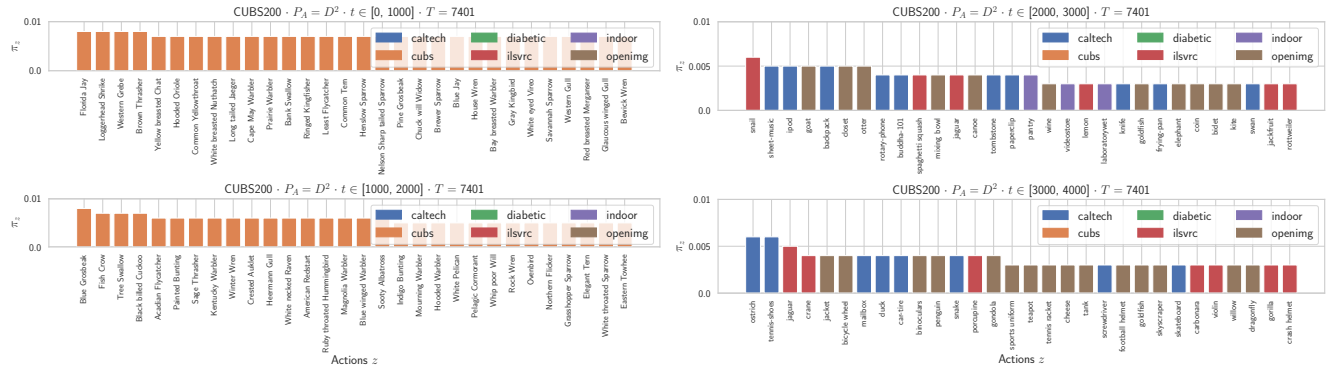
**Figure S6: Qualitative results: Indoor67.** Extends Figure 6 in the main paper. GT labels are underlined, correct top-1 knockoff predictions in green and incorrect in red.

Transfer Set			Test Set		
					
No DR: 0.73 Proliferative: 0.12 Moderate: 0.08	No DR: 0.48 Mild: 0.36 Moderate: 0.16	No DR: 0.30 Moderate: 0.29 Proliferative: 0.28	<u>No DR: 0.50</u> Mild: 0.33 Moderate: 0.13	<u>No DR: 0.36</u> Mild: 0.33 Moderate: 0.28	Mild: <u>0.53</u> <u>No DR: 0.43</u> Moderate: 0.03
					
Moderate: 0.69 No DR: 0.31 Mild: 0.01	Moderate: 0.63 No DR: 0.15 Severe: 0.13	Moderate: 0.35 Mild: 0.32 No DR: 0.22	<u>Moderate: 0.48</u> Mild: 0.316 No DR: 0.21	<u>Moderate: 0.35</u> Mild: 0.31 No DR: 0.23	No DR: <u>0.36</u> Mild: 0.33 <u>Moderate: 0.26</u>
					
Severe: 0.73 Proliferative: 0.23 Moderate: 0.04	Severe: 0.70 Proliferative: 0.30 Moderate: 0.00	Severe: 0.53 Mild: 0.16 Moderate: 0.15	<u>Severe: 0.57</u> Moderate: 0.23 Proliferative: 0.19	<u>Severe: 0.41</u> Proliferative: 0.29 Moderate: 0.24	Moderate: <u>0.62</u> <u>Severe: 0.16</u> Mild: 0.13

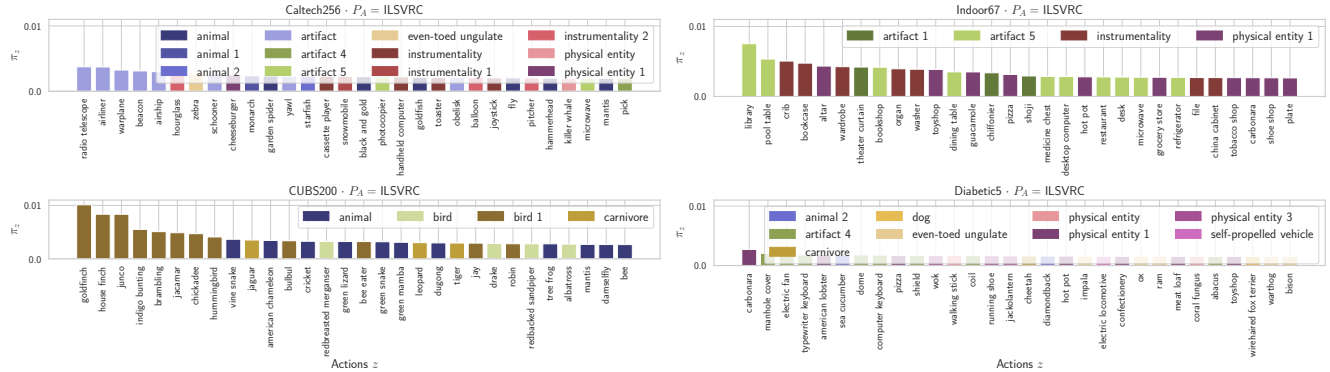
**Figure S7: Qualitative results: Diabetic5.** Extends Figure 6 in the main paper.



(a) Closed world.



(b) Closed world. Analyzing policy over time  $t$  for CUBS200.

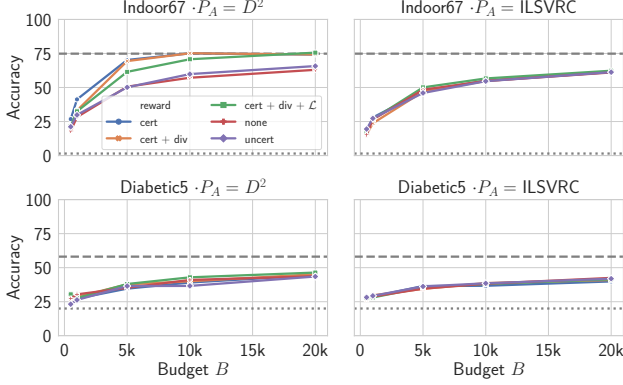


(c) Open world.

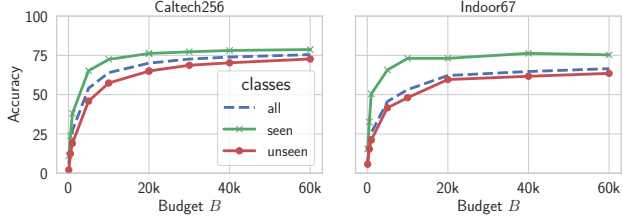
**Figure S8: Policies learnt by adaptive strategy. Supplements Figure 7 in the main paper.**

This extension “ $P_V(F_V)$ ” is illustrated in Figure S2, displayed alongside KD approach. The figure represents the sample-efficiency of the first two rows of Table 2. Here we observe: (i) comparable performance in all but one case (Diabetic5, discussed shortly) indicating KD is an effective approach to train knockoffs; (ii) we find KD achieve

better performance in Caltech256 and Diabetic5 due to regularizing effect of training on soft-labels [6] on an imbalanced dataset.



**Figure S9: Reward Ablation.** Supplements Figure 8 in the main paper.



**Figure S10: Per class evaluation.** Per-class evaluation split into seen and unseen classes.

### C.3. Policies learnt by Adaptive

We inspected the policy  $\pi$  learnt by the adaptive strategy in Section 6.1. In this section, we provide policies over all blackboxes in the closed- and open-world setting. Figures S8a and S8c display probabilities of each action  $z \in Z$  at  $t = 2500$ .

Since the distribution of rewards is non-stationary, we visualize the policy over time in Figure S8b for CUBS200 in a closed-world setup. From this figure, we observe an evolution where: (i) at early stages ( $t \in [0, 2000]$ ), the approach samples (without replacement) images that overlaps with the victim’s train data; and (ii) at later stages ( $t \in [2000, 4000]$ ), since the overlapping images have been exhausted, the approach explores related images from other datasets e.g., “ostrich”, “jaguar”.

### C.4. Reward Ablation

The reward ablation experiment (Figure 8 in the main paper) for the remaining datasets are provided in Figure S9. We make similar observations as before for Indoor67. However, since  $F_V = \text{Diabetic5}$  demonstrates confident predictions in all images, we find little-to-no improvement for knockoffs of this victim model.

## D. Auxiliary Experiments

In this section, we present experiments to supplement existing results in the main paper.

### D.1. Effect of CNN Initialization

In our experiments (Section 6), the victim and knockoff models are initialized with ImageNet pretrained weights<sup>1</sup>, a de facto when training CNNs with a limited amount of data. In this section, we study influence of different initializations of the victim and adversary models.

To achieve reasonable performance in our limited data setting, we perform the following experiments on comparatively smaller models and datasets. We choose three victim blackboxes (all trained after random initialization) using the following datasets: MNIST [9], CIFAR10 [7], and CIFAR100 [7]. We train a LeNet-like model<sup>2</sup> for MNIST, and Resnet-18 models for CIFAR-10 and CIFAR-100.

While we use the same blackbox model architecture for the knockoff, we either randomly initialize them or pre-train them on a different task. Consequently, in the following experiments, both the victim and knockoff have different initializations. We repeat our experiment using random policy (Section 4.1.1) and using as the query set  $P_A$ : (a) when  $P_V = \text{MNIST}$ : EMNIST [3] (superset of MNIST containing alpha numeric characters [A-Z, a-z, 0-9]), EMNISTLetters ([A-Z, a-z]), FashionMNIST [15] (fashion items spanning 10 classes e.g., trouser, coat) and KMNIST [2] (Japanese Hiragana characters spanning 10 classes); (b) when  $P_V = \text{CIFAR10}$ : CIFAR100 [7] and TinyImageNet200<sup>3</sup> (subset of ImageNet with 500 images per each of 200 classes); and (c) when  $P_V = \text{CIFAR100}$ : CIFAR10 and TinyImageNet200. Note that the output classes between CIFAR10 and CIFAR100 are disjoint.

From Figure S3, we observe: (i) model stealing is possible even when the knockoffs are randomly initialized. For instance, when stealing MNIST, we recover  $0.98 \times$  victim accuracy across all choices of  $P_A$ ; (ii) pretraining the knockoff model – even on a different task – improves sample efficiency of model stealing attacks e.g., when  $F_V = \text{CIFAR10-resnet18}$ , querying images from  $P_V$  improves the knockoff accuracy after 50k queries from 46.5% to 78.9%.

### D.2. Seen and Unseen classes

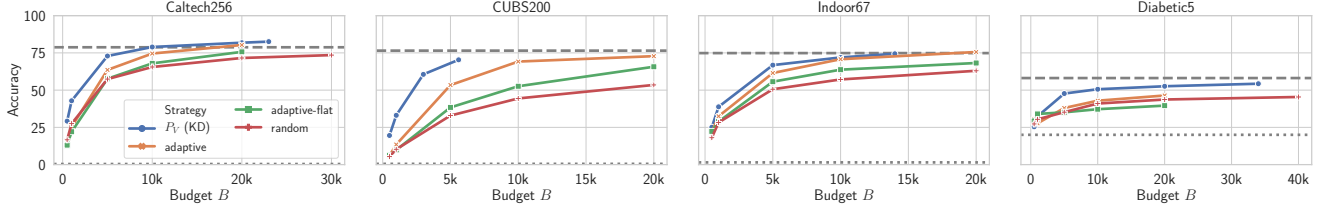
We now discuss evaluation to supplement Section 5.2.1 and Section 6.1.

In Section 6.1, we highlighted strong performance of the knockoff even among classes that were never encountered (see Table S1 for exact numbers) during training. To elaborate, we split the blackbox output classes into “seen” and “unseen” categories and present mean per-class accuracies in Figure S10. Although we find better performance on

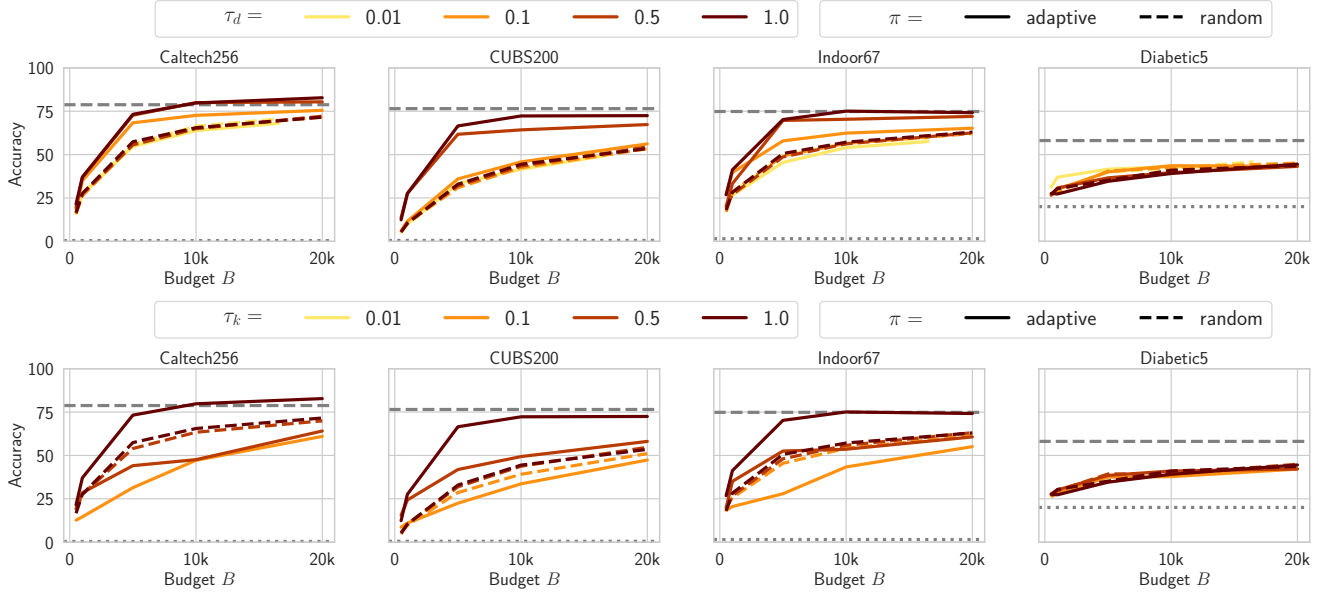
<sup>1</sup>Alternatives for ImageNet pretrained models across a wide range of architectures were not available at the time of writing

<sup>2</sup><https://github.com/pytorch/examples/blob/master/mnist/main.py>

<sup>3</sup><https://tiny-imagenet.herokuapp.com/>



**Figure S11: Hierarchy.** Evaluating adaptive with and without hierarchy using  $P_A = D^2$ . -- represents accuracy of blackbox  $F_V$  and . . . represents chance-level performance.



**Figure S12: Semi-open world:  $\tau_d$  and  $\tau_k$ .**

classes seen while training the knockoff, performance of unseen classes is remarkably high, with the knockoff achieving  $>70\%$  performance in both cases.

### D.3. Adaptive: With and without hierarchy

The adaptive strategy presented in Section 4.1.2 uses a hierarchy discussed in Section 5.2.2. As a result, we approached this as a hierarchical multi-armed bandit problem. Now, we present an alternate approach adaptive-flat, without the hierarchy. This is simply a multi-armed bandit problem with  $|Z|$  arms (actions).

Figure S11 illustrates the performance of these approaches using  $P_A = D^2$  ( $|Z| = 2129$ ) and rewards  $\{\text{certainty, diversity, loss}\}$ . We observe adaptive consistently outperforms adaptive-flat. For instance, in CUBS200, adaptive is  $2\times$  more sample-efficient to reach accuracy of 50%. We find the hierarchy helps the adversary (agent) better navigate the large action space.

### D.4. Semi-open World

The closed-world experiments ( $P_A = D^2$ ) presented in Section 6.1 and discussed in Section 5.2.1 assumed access to the image universe. Thereby, the overlap between  $P_A$  and  $P_V$  was 100%. Now, we present an intermediate overlap scenario **semi-open world** by parameterizing the overlap as: (i)  $\tau_d$ : The overlap between *images*  $P_A$  and  $P_V$  is  $100 \times \tau_d$ ; and (ii)  $\tau_k$ : The overlap between *labels*  $K$  and  $Z$  is  $100 \times \tau_k$ . In both these cases  $\tau_d, \tau_k \in (0, 1]$  represents the fraction of  $P_A$  used.  $\tau_d = \tau_k = 1$  depicts the closed-world scenario discussed in Section 6.1.

From Figure S12, we observe: (i) the random strategy is unaffected in the semi-open world scenario, displaying comparable performance for all values of  $\tau_d$  and  $\tau_k$ ; (ii)  $\tau_d$ : knockoff obtained using adaptive obtains strong performance even with low overlap e.g., a difference of at most 3% performance in Caltech256 even at  $\tau_d = 0.1$ ; (iii)  $\tau_k$ : although the adaptive strategy is minimally affected in few cases (e.g., CUBS200), we find the performance drop



due to a pure exploitation (certainty) that is used. We observed recovery in performance by using all rewards indicating exploration goals (diversity, loss) are necessary when transitioning to an open-world scenario.

## References

- [1] Eyepacs. <https://www.kaggle.com/c/diabetic-retinopathy-detection>. Accessed: 2018-11-08. 1
- [2] Tarin Clanuwat, Mikel Bober-Irizar, Asanobu Kitamoto, Alex Lamb, Kazuaki Yamamoto, and David Ha. Deep learning for classical japanese literature, 2018. 7
- [3] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and André van Schaik. Emnist: an extension of mnist to handwritten letters. *arXiv preprint arXiv:1702.05373*, 2017. 7
- [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 3
- [5] Gregory Griffin, Alex Holub, and Pietro Perona. Caltech-256 object category dataset. 2007. 1
- [6] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv:1503.02531*, 2015. 6
- [7] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, Cite-seer, 2009. 7
- [8] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Tom Duerig, et al. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *arXiv:1811.00982*, 2018. 3
- [9] Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. *AT&T Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2:18, 2010. 7
- [10] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011. 3
- [11] Ariadna Quattoni and Antonio Torralba. Recognizing indoor scenes. In *CVPR*, 2009. 1
- [12] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 2015. 3
- [13] Angelo Vittorio. Toolkit to download and visualize single or multiple classes from the huge open images v4 dataset. [https://github.com/EsVM/0IDv4\\_ToolKit](https://github.com/EsVM/0IDv4_ToolKit), 2018. 3
- [14] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011. 1
- [15] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017. 7