



ELSEVIER

Contents lists available at ScienceDirect

## Pattern Recognition

journal homepage: [www.elsevier.com/locate/patcog](http://www.elsevier.com/locate/patcog)

## Towards Non-I.I.D. image classification: A dataset and baselines

Yue He<sup>1</sup>, Zheyang Shen<sup>1</sup>, Peng Cui\*

Lab of Media and Network, Department of Computer Science and Technology, Tsinghua University, Room 9-316, East Main Building, Beijing 100084, PR China

## ARTICLE INFO

## Article history:

Received 14 July 2019

Revised 25 March 2020

Accepted 15 April 2020

Available online xxx

## Keywords:

Non-I.I.D.

Dataset

Context

Bias

ConvNet

Batch balancing

## ABSTRACT

I.I.D.<sup>2</sup> hypothesis between training and testing data is the basis of numerous image classification methods. Such property can hardly be guaranteed in practice where the Non-IIDness is common, causing unstable performances of these models. In literature, however, the Non-I.I.D.<sup>3</sup> image classification problem is largely understudied. A key reason is lacking of a well-designed dataset to support related research. In this paper, we construct and release a Non-I.I.D. image dataset called NICO<sup>4</sup>, which uses contexts to create Non-IIDness consciously. Compared to other datasets, extended analyses prove NICO can support various Non-I.I.D. situations with sufficient flexibility. Meanwhile, we propose a baseline model with ConvNet structure for General Non-I.I.D. image classification, where distribution of testing data is unknown but different from training data. The experimental results demonstrate that NICO can well support the training of ConvNet model from scratch, and a batch balancing module can help ConvNets to perform better in Non-I.I.D. settings.

© 2020 Elsevier Ltd. All rights reserved.

## 1. Introduction

In recent years, machine learning has achieved remarkable progress in a wide range of applications [1–4], mainly owing to the development of deep neural networks [5,6]. One basic hypothesis of machine learning models is that the training and testing data should consist Independent and Identically Distributed (I.I.D.) samples. However, this ideal hypothesis is fragile in real cases where we can hardly impose constraints on the testing data distribution. This implies that the model minimizing empirical error on training data does not necessarily perform well on testing data, leading to the challenge of Non-I.I.D. learning. The problem is more serious when the training samples are not sufficient to approximate the training distribution itself. How to develop Non-I.I.D. learning methods that are robust to distribution shifting is of paramount significance for both academic research and industrial applications.

Benchmark datasets, providing a common ground for competing approaches, are always important to promote the development of a research direction. Take image classification, a prominent learning task, as an example. Its development benefits a lot from the

benchmark datasets, such as PASCAL VOC [7], MSCOCO [8], and ImageNet [9]. In particular, it is the ImageNet, a large-scale and well-structured image dataset, that successfully demonstrates the capability of deep learning and thereafter significantly accelerates the advancement of deep convolutional neural networks. On these datasets, it is easy to establish an I.I.D. image classification setting by random data splitting. But they do not provide an explicit option to simulate a Non-I.I.D. setting. The dataset that can well support the research on Non-I.I.D. image classification is still in vacancy.

In this paper, we construct and release a dataset that is dedicatedly designed for Non-I.I.D. image classification, named NICO (Non-I.I.D. Image dataset with Contexts). The basic idea is to label images with both main concept and contexts. For example, in the category of ‘dog’, images are divided into different contexts such as ‘grass’, ‘car’, ‘beach’, meaning the ‘dog’ is on the grass, in the car, or on the beach respectively. With these contexts, one can easily design an Non-I.I.D. setting by training a model in some contexts and testing it in the other unseen contexts. Meanwhile, the degree of distribution shift can be flexibly controlled by adjusting the proportions of different contexts in training and testing data. Till now, NICO contains 19 classes, 188 contexts and nearly 25,000 images in total. The scale is still increasing, and the current scale has been able to support the training of deep convolution networks from scratch.

The NICO dataset can support, but not limited to, two typical settings of Non-I.I.D. image classification. One is Targeted Non-I.I.D.

\* Corresponding author.

E-mail addresses: [heyue18@mails.tsinghua.edu.cn](mailto:heyue18@mails.tsinghua.edu.cn) (Y. He), [shenzy17@mails.tsinghua.edu.cn](mailto:shenzy17@mails.tsinghua.edu.cn) (Z. Shen), [cui@tsinghua.edu.cn](mailto:cui@tsinghua.edu.cn) (P. Cui).<sup>1</sup> Yue He and Zheyang Shen contributed equally to this work as first authors.<sup>2</sup> I.I.D.: Independent and Identically Distributed<sup>3</sup> Non-I.I.D.: Non-Independent and Identically Distributed<sup>4</sup> NICO: Non-I.I.D. Image dataset with Contexts

image classification, where testing data distribution is known but different from training data distribution. The other is General Non-I.I.D. image classification, where testing data distribution is unknown and different from training data distribution. Apparently, the latter one is much more realistic and challenging. A model learned in one environment could be possibly applied in many other environments. In this case, the robustness of a model in the environments with unknown distribution shift is a highly favorable characteristic. It is especially critical in risk-sensitive applications like medical and security.

Due to the lack of a well-structured and reasonable-scaled dataset, there is still no convolutional neural network model proposed to address the general Non-I.I.D. image classification problem. In this paper, we propose a novel model CNBB<sup>5</sup> (ConvNet with Batch Balancing) as a baseline of exploiting CNN model for general Non-I.I.D. image classification. The experimental results show that the proposed batch balancing mechanism can help a ConvNet model to resist, to some extent, the negative effect brought by Non-IIDness.

In a word, NICO released in this paper is devoted to advance the research about intelligence perception of efficient and robust pattern recognition across diverse environments in visual field. The works that focus on the human cognition, such as causality, always have better interpretability naturally and could design and execute experiments in kinds of Non-I.I.D. settings, compare their performances fairly in NICO. What's more, the ability of adjusting distribution shift controllably can indeed bring more explainability to the models and experiments, especially for deep learning [10,11]. Also the CNBB proposed is a preliminary attempt to introduce causal mechanism into the deep ConvNets.

## 2. Non-I.I.D. image classification

### 2.1. Problem definition

We first give a formal definition of Non-I.I.D. image classification as follow:

**Problem 1. (Non-I.I.D. Image Classification)** Given the training data  $D_{train} = (X_{train}, Y_{train})$ , where  $X_{train} \in \mathbb{R}^{n \times (c \times h \times w)}$  represent the images and  $Y_{train} \in \mathbb{R}^{n \times 1}$  represent the labels. The task is to learn a feature extractor  $g_{\varphi}(\cdot)$  and a classifier  $f_{\theta}(\cdot)$ , so that  $f_{\theta}(g_{\varphi}(\cdot))$  can predict the labels of testing data  $D_{test} = (X_{test}, Y_{test})$  precisely, where  $g_{\varphi}(\cdot) \in \mathbb{R}^{n \times p}$  and  $\psi(D_{train}) \neq \psi(D_{test})$ . Moreover, according to the availability of the prior knowledge on testing data, we further define two different tasks. One is **Targeted Non-I.I.D. Image Classification** where the testing data distribution  $\psi(D_{test})$  is known. The other is **General Non-I.I.D. Image Classification**, which corresponds to a more realistic scenario where the testing data distribution  $\psi(D_{test})$  is unknown.

In order to intuitively quantify the degree of distribution shift between  $\psi(D_{train})$  and  $\psi(D_{test})$ , we define the Non-I.I.D. Index as follow:

**Definition 1. Non-I.I.D. Index (NI)** Given a feature extractor  $g_{\varphi}(\cdot)$  and a class  $C$ , the degree of distribution shift between training data  $D_{train}^C$  and testing data  $D_{test}^C$  is defined as:

$$NI(C) = \left\| \frac{\overline{g_{\varphi}(X_{train}^C)} - \overline{g_{\varphi}(X_{test}^C)}}{\sigma(g_{\varphi}(X^C))} \right\|_2,$$

where  $X^C = X_{train}^C \cup X_{test}^C$ ,  $\overline{(\cdot)}$  represents the first order moment,  $\sigma(\cdot)$  is the std used to normalize the scale of features and  $\|\cdot\|_2$  represents the 2-norm.

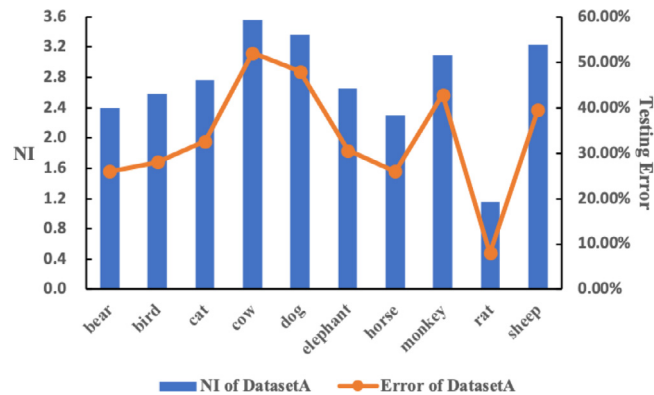


Fig. 1. NI (represented by the bar-type) and testing error (represented by the curve-type) of each class in Dataset A.

### 2.2. Existence of Non-IIDness

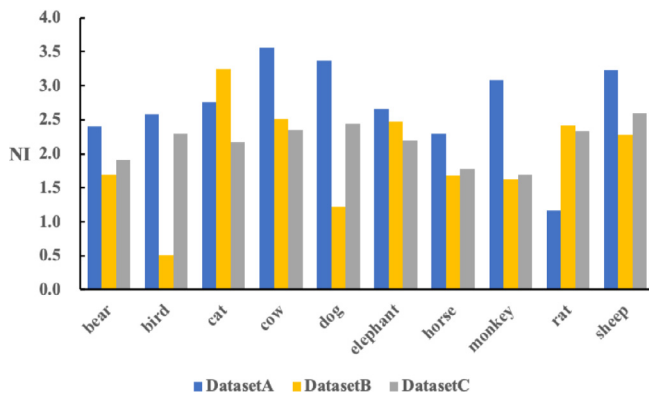
In real cases, the I.I.D. hypothesis can never be strictly satisfied, meaning that Non-IIDness ubiquitously exists in previous datasets [12]. Here we take ImageNet as an example. ImageNet is in a hierarchical structure, where each class (e.g. dog) contains multiple subclasses (e.g. different kinds of dogs). For each subclass, it provides training and testing (validation) subsets of images. To verify the Non-IIDness in ImageNet, we select 10 common animal classes (e.g. dog, cat) and construct a new dataset using 10 instantiated subclasses (e.g. Labrador, Persian), each randomly drawn from those classes. Using the training and testing subsets, we train and evaluate a ConvNet on image classification task. The structure of the ConvNet used in this paper is similar to AlexNet (details seen in **Appendix**), and we take the last FC layer of the ConvNet as the feature extractor  $g_{\varphi}$ . Note that model structure is used in all subsequent analysis (including on NICO) for fair comparison, and thus selected by trading-off performance and required training data scale. But as a base model with sufficient learning capacity, the specific model structure does not affect the conclusions. We repeat this collection procedure for 3 times, obtain 3 new datasets (*Dataset A*, *Dataset B* and *Dataset C*) and calculate the NI and testing error for each class respectively. As an example, we plot the results of *Dataset A* in Fig. 1. We can find that:

- NI is above zero for all classes, which implies the Non-IIDness between training and testing data is ubiquitous even in large-scale datasets like ImageNet.
- Different classes have different NI values and higher NI value corresponds to higher testing error.

The strong correlation between NI and testing error can be further proved by their high pearson correlation coefficients ( $r = 0.95$ ) and small  $p\_value$  ( $2e-15$ ).

One may argue that the numerical value of NI is conditioned on feature extractor and could not be compared cross different datasets due to the supervised learning. In fact, we only use it to analyse the trend of distribution bias by some intervention between training and testing subsets from the one data source. To learn the feature extractor and use it to compute NI in the same dataset could guarantee that the change of NI is only caused by the predefined specific intervention, which increases more controllability and explainability to the corresponding experiments. Otherwise the unknown external disturb would be drawn in. In later paragraph, we use NI to make an empirical analysis on the new dataset we construct to prove that NICO can support various Non-I.I.D. situations flexibly and consciously.

<sup>5</sup> CNBB: ConvNet with Batch Balancing.



**Fig. 2.** NI of each class in 3 different datasets constructed from ImageNet. Different datasets instantiate the same classes with different subclasses.

### 2.3. Limitations of existing datasets

Throughout the development of computer vision research, benchmark datasets have always played a critical role on both providing a common ground for algorithm evaluation and driving new directions. Specifically, for image classification task, we can enumerate several milestone datasets such as PASCAL VOC, MSCOCO and ImageNet. However, existing benchmark datasets cannot well support the Non-I.I.D. image classification. First of all, despite the manifested Non-I.I.D. in ImageNet and other datasets, as shown in Fig. 1, the overall degree of distribution shift between training and testing data for each class is relatively small, making these datasets less challenging from the angle of Non-I.I.D. image classification. More importantly, there is no explicit way to control the degree of distribution shift between training and testing data in the existing datasets. As illustrated in Fig. 2, if we instantiate the same class with different subclasses in ImageNet and obtain 3 datasets with identical structure, the NI of a given class is fairly unstable across different datasets. Without a controllable way to simulate different levels of Non-I.I.D.ness, competing approaches cannot be evaluated fairly and systematically on those datasets. Those said, a dataset that is dedicatedly designed for Non-I.I.D. image classification beyond the above limitations is demanded.

## 3. The NICO dataset

In this section, we introduce the properties and collection process of the dataset, followed by preliminary empirical results in different Non-I.I.D. settings supported by this dataset.

### 3.1. Context for Non-I.I.D. images

The essential idea of generating Non-I.I.D. images is to enrich the labels of an image with both conceptual and contextual labels. Different from previous datasets that only label an image with the major concept (e.g. dog), we also label the concrete context (e.g. on grass) that the concept appears in. Then it is easy to simulate a Non-I.I.D. setting by training and testing the model of a concept with different contexts. A good model for Non-I.I.D. image classification is expected to perform well in both training contexts and testing contexts.

In NICO, we mainly incorporate two kinds of contexts. One is the attributes of a concept (or object), such as color, action, and shape. Some examples of 'context + concept' pairs include *white bear*, *climbing monkey* and *double decker* etc. The other kind of con-

texts is the background or scene of a concept. The examples of 'context + concept' pairs include *cat on snow*, *horse aside people* and *airplane in sunrise* etc. Samples of different contexts in the NICO dataset are shown in Fig. 3.

In order to provide more flexible Non-I.I.D. settings, we tend to select the contexts that occur in multiple concepts. Then for a given concept, a context may occur in both positive samples and negative samples (that are sampled from other concepts). This provides another flexibility to let a context included in training positive samples appear or do not appear in training negative samples, which will yield different Non-I.I.D. settings.

There are some related datasets that also supply contextual information in addition to major concepts, such as NUSWIDE [13] and MSCOCO [8]. NUSWIDE dataset and its extended version [14] focus on social image understanding including the tasks of tag completion, image retrieval and so on, launching the deep models like DCB [15], WDMF [16] and WDM [17]. MSCOCO promotes the researches of various detection and segmentation a lot. However, none of these datasets are towards Non-I.I.D. image classification specifically. That is to say, one cannot build and adjust the shift of distribution to meeting various Non-I.I.D. settings well and conveniently. For example an image always has multiple and overlapping tags for one category in NUSWIDE. So it's hard to divide different data distributions controllably, especially for the compositional bias setting below. And the sample size of each context (tag) per class is quite imbalanced in other datasets. That says only NICO can well support kinds of Non-I.I.D. researches about robust and explainable machine learning [18–21].

### 3.2. Data collection and statistics

Referring to ImageNet, MSCOCO and other classical datasets [22,23], we first confirm two superclasses: *Animal* and *Vehicle*. For each superclass, we select classes from the 272 candidates in MSCOCO, with the criterion that the selected classes in a superclass should have large inter-class differences. For context selection, we exploit YFCC100m [24] browser<sup>6</sup> and first derive the frequently co-occurred tag list for a given concept (i.e. class label). We then filter out the tags that occur in only a few concepts. Finally, we manually screen all tags and select the ones that are consistent with our definition of contexts (i.e. object attributes or backgrounds and scenes).

After obtaining the conceptual and contextual tags, we concatenate a given conceptual tag and each of its contextual tags to form a query, input the query into the API of Google and Bing image search, and collect the top-ranked images as candidates. Finally, in the phase of screening, we select images into the final dataset according to the following criteria:

- The content of an image should correctly reflect its concept and context.
- Given a class, the number of images in each context should be adequate and as balance as possible across contexts.

Note that we do not conduct image registration or filtering by object centralization, so that the selected images are more realistic and in wild than those in ImageNet.

The NICO dataset will be continuously updated and expanded. Till now, there are two superclasses: *Animal* and *Vehicle*, with 10 classes for *Animal* and 9 classes for *vehicle*. Each class has 9 or 10 contexts. The average size of contexts per class ranges from 83 to 215, and the average size of classes is about 1300 images, which is similar to ImageNet. In total, there are 25,000 images in the NICO dataset. As NICO is in a hierarchical structure, it is easy to be ex-



**Fig. 3.** Samples with contexts in NICO. Images in the first row are dogs of *Animal*, assigned to different contexts below it. The second and third row correspond to horse of *Animal* and boat of *Vehicle* respectively.

**Table 1**  
Data size of each class in NICO.

<i>Animal</i>	Data size	<i>Vehicle</i>	Data size
Bear	1609	Airplane	930
Bird	1590	Bicycle	1639
Cat	1479	Boat	2156
Cow	1192	Bus	1009
Dog	1624	Car	1026
Elephant	1178	Helicopter	1351
Horse	1258	Motorcycle	1542
Monkey	1117	Train	750
Rat	846	Truck	1000
Sheep	918		

panded. More statistics on NICO is reported in Table 1. The dataset can be downloaded through the link.<sup>7</sup>

### 3.3. Supported Non-I.I.D. settings

By dividing a class into different contexts, NICO provides the flexibility of simulating Non-I.I.D. settings in different levels. To name a few, here we list 4 typical settings.

**Setting 1. Minimum bias.** Given a class, we can ignore the contexts, and randomly split all images of the class into training and testing subsets as positive samples. Then we can randomly sample images belonging to other classes into training and testing subsets as negative samples. In this setting, the way of random sampling leads to minimum distribution shift between training and testing distributions in the dataset, which simulates a nearly i.i.d. scenario.

**Setting 2. Proportional bias.** Given a class, when sampling positive samples, we use all contexts for both training and testing, but the percentage of each context is different in training and testing subsets. For example, we can let one context take the majority in training data while taking minority in testing, which is consistent with the natural phenomena that visual concepts follow a power law distribution[25]. The negative sampling process is the same as Setting 1. In this setting, the level of distribution shift can be tuned by adjusting the proportion difference between training and testing subsets for each context.

**Setting 3. Compositional bias.** Given a class, not every testing context that the positive samples belong to appears in training subset simultaneously. Such a setting is quite common in real scene, because available datasets could not contain all the potential contexts in nature due to the limitations of sampling time and space. Intuitively, the distribution shift from observed contexts to unseen contexts is usually large. The less number of testing contexts observed in training generally leads to the higher distribution shift. A more radical distribution shift can be further achieved by combining compositional bias and proportional bias.

**Setting 4. Adversarial bias.** Given a class, the positive sampling process is the same as Setting 3. For negative sampling, we tend to select the negative samples from the contexts that have not been (or have been) included in positive training samples to form the negative training (or testing) subset. In this way, the distribution shifting is even higher than Setting 3, and the existing classification model developed under i.i.d. assumption is more prone to be confused.

The above 4 settings are designed to generate Non-I.I.D. training and testing subsets. Under each setting, we can conduct either Targeted or General Non-I.I.D. image classification by assuming the distribution of testing subset is known or unknown.

### 3.4. Empirical analysis

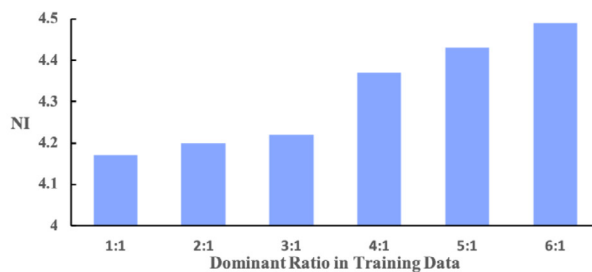
To verify the effectiveness of NICO in supporting Non-I.I.D. image classification, we conduct a series of empirical analysis. It is worth noting that, in each setting, only the distribution of training or testing data changes, while the structure of ConvNet and the size of training data keep the same.

#### 3.4.1. Minimum bias setting

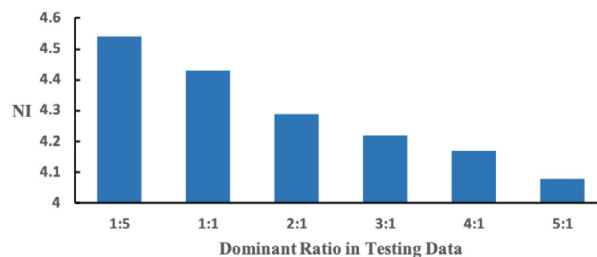
In this setting, we randomly sample 8000 images for training and 2000 images for testing from *Animal* and *Vehicle* superclasses respectively. The average testing accuracy and *NI* over all the classes are 49.6%, 3.85 for *Animal* superclass and 63.0%, 3.20 for *Vehicle* superclass. We can find that *NI* in NICO is much higher than *NI* in ImageNet even if there is no explicit bias (due to random sampling) when we construct the training and testing subsets. This is because the images in NICO are typically non-iconic images with rich contextual information and non-canonical viewpoints, which is more challenging from the perspective of image classification.

<sup>6</sup> <http://www.yfcc100m.org/>

<sup>7</sup> <https://www.dropbox.com/sh/8mouawi5guaupyb/AAD4fdySrA6fn3PgSmhKwFgva?dl=0>



(a) Average  $NI$  over all classes in *Animal* superclass with respect to various dominant ratio of training data, while the dominant ratio of testing data is fixed to 1:1 (uniform sampling).



(b) Average  $NI$  over all classes in *Animal* superclass with respect to various dominant ratio of testing data, while the dominant ratio of training data is fixed to 5:1.

**Fig. 4.**  $NI$  in proportional bias setting.

### 3.4.2. Proportional bias setting

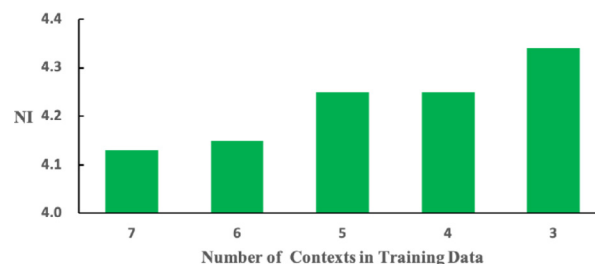
In this setting, we let all the contexts appear in both training and testing data, and randomly select one dominant context in training data (or testing data) for each class in *Animal* superclass. Such experimental settings comply with the natural phenomena that a majority of visual contexts are rare except a few common ones [25]. Specifically, we define the dominant ratio as follow:

$$\text{Dominant Ratio} = \frac{N_{\text{dominant}}}{N_{\text{minor}}},$$

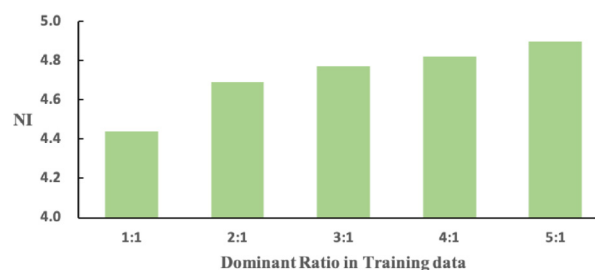
where  $N_{\text{dominant}}$  refers to the sample size of the dominant context and  $N_{\text{minor}}$  refers to the average size of other contexts where we uniformly sample other contexts. We conduct two experiments where either dominant ratio of training data or testing data is fixed, and vary the other one. We plot the results in Fig. 4(a) and (b). From the figures, we can clearly find a consistent pattern that the  $NI$  becomes higher as the discrepancy between dominant ratio of training data and testing data becomes larger. As a result, by tuning the dominant ratio of training data (or testing data), we can easily simulate different extents of distribution shift as we want.

### 3.4.3. Compositional bias setting

Compared to proportional bias setting, compositional bias setting simulates a condition where the knowledge obtained from training data is insufficient to characterize the whole distribution. To do so, we choose a subset of contexts for a given class when constructing the training data and testing the model with all the contexts. By varying the number of contexts observed in training data, we can simulate different extents of information loss and distribution shift. From Fig. 5, we can find that the  $NI$  consistently decreases when we could observe more contexts in training data. A more radical distribution shift can be achieved by combining the notion of proportional bias and compositional bias. Given a particular class in *Vehicle* superclass, we choose 7 contexts for training and the other 3 contexts for testing, and further let one context dominate the training data. By doing so, we can obtain a more severe Non-I.I.D. condition between training and testing data than previous two settings, as illustrated by the results from Fig. 6.



**Fig. 5.**  $NI$  in compositional bias setting: average  $NI$  over all classes in *Vehicle* superclass with respect to the number of contexts used in training data.



**Fig. 6.**  $NI$  in the combined setting of compositional bias and proportional bias: average  $NI$  over all classes in *Vehicle* superclass with respect to various dominant ratio of training data, where contexts in testing data is totally unseen in training.

### 3.4.4. Adversarial bias setting

Given a target class, we define a context as confounding context if it only appears in the negative samples of training data and positive samples of testing data. In this experiment, we choose four classes in *Animal* superclass as target classes and report the  $NI$  w.r.t various number of confounding contexts in Fig. 7. The experimental results indicate that the number of confounding contexts has consistent influence on the  $NI$  of different classes. Given any target class, we can simulate a more harsh distribution shift and further confuse the ConvNet by adding more confounding contexts.

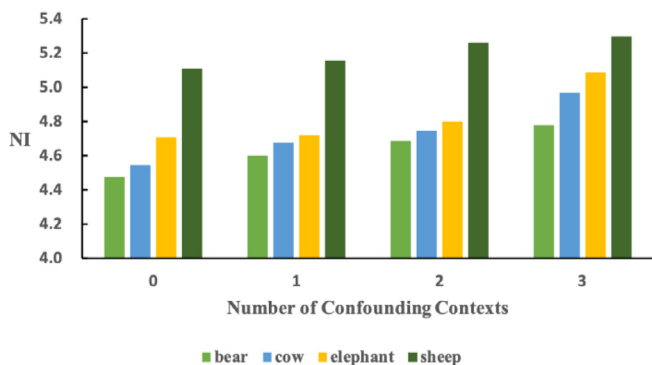


Fig. 7. NI in the adversarial bias setting: NI of target class with respect to the number of confounding contexts.

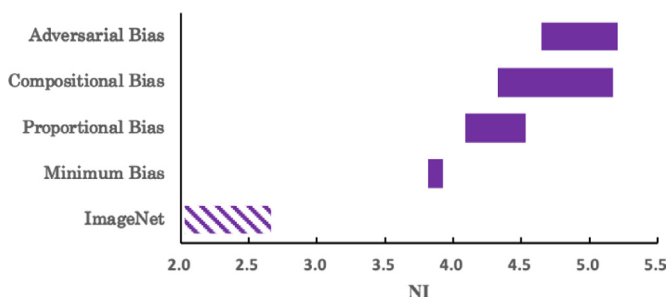


Fig. 8. Range of average NI over Animal superclass for different settings supported in NICO.

Finally, we show the range of NI in different Non-I.I.D. settings in Fig. 8. We can see the level of NI in NICO is significantly higher than ImageNet, and there is an obvious ascending trend from Minimum Bias to Adversarial Bias settings.

The intention of releasing NICO is to appeal more attention and promote the research about the intrinsic mechanism of stable and robust learning cross various environments. Intuitively, such mechanisms should be consistent with human cognitive habits, such as causality, possibly lighting up the way to the “strong AI” [26]. However some skills that could still improve the performance to some extent do not reveal the essential law of intelligence. For example, a human would not find out the bounding box firstly before object classification and such localization methods still fail to recognise when the attributes of object, like the color of a bear in NICO, change in testing environments. Our position is to forbid these methods to be applied in NICO for pursuing high scores only and not approve the corresponding results.

#### 4. General Non-I.I.D. image classification

In this section, we propose a novel model for General Non-I.I.D. image classification.

In the literature of Non-I.I.D. image classification, most previous methods are proposed for Targeted Non-I.I.D. image classification. Domain adaptation and covariate shift methods [27–29] are proposed to match distributions, transform feature space or learn invariant features between training data and testing data. These methods can achieve good performances but are less feasible in practice due to the fact that they need prior knowledge on testing data distribution. On the other hand, several methods are proposed to liberalize the need of testing data information in Targeted Non-I.I.D. image classification. For example, domain generalization methods [30,31] only use training data to learn a domain-agnostic model or invariant representations. However, these methods about transfer learning [32,33] require the training data has multiple do-

main and we know which domain each sample belongs to. Moreover, the performance of these methods is highly dependent on the diversity of training data.

Recently, growing attention has been paid on General Non-I.I.D. learning. In the literature of causality [34], an ideal model to resolve selection bias is to make policy based on causal variables, which keep stable across different domains [35]. Popular methods based on observational data to estimate the causal effect of a treatment on the outcome include propensity score matching [36,37], markov blankets [38] and confounder balancing [39,40] and etc. [41]. Lately [42] leverages causality for predictive modeling. By performing global confounder balancing, one can accurately identify the stable features that are insensitive to unknown distribution shift for prediction. Shen et al. [43] proposes a causally regularized logistic regression called CRLR<sup>8</sup> for General Non-I.I.D. image classification and achieves good performance in a relatively small dataset. Other literatures, such as RSNMF [44], expect to learn robust image representations with the help of sparse coding technology. However, due to the lack of well-structured and reasonable-scaled dataset, these methods cannot leverage the powerful deep representation learning techniques (e.g. ConvNets) and therefore are not favourable for large-scale image classification tasks.

In this work, with the help of NICO, we extend the notion of global confounder balancing into ConvNet, and propose a novel model called CNBB, ConvNet with Batch Balancing.

##### 4.1. ConvNet with batch balancing

The key idea in CRLR is global confounder balancing, which successively sets each feature as treatment variable, and learns an optimal set of sample weights that can balance the distribution of treated and control groups for any treatment variable. Thereafter, the correlations among features will be disentangled and their true effects on class label can be more accurately estimated.

To introduce the notion of global confounder balancing into deep learning, we mainly face two challenges:

- Confounder balancing methods assume features to be in binary form, while we generally have continuous features in ConvNet.
- For global confounder balancing, we need to learn a new set of sample weights for all the training samples in one iteration.

This is not feasible for ConvNet where we cannot feed all the training data into the model at once.

To overcome these challenges, we introduce a quantization loss for feature binarization and propose a batch confounder balancing method. Specifically, given a batch of training images, we define the quantization loss as follows:

$$Loss_q = - \sum_{i=1}^n \|g_\varphi(x_i)\|_2^2, \quad (1)$$

where  $n$  refers to the batch size,  $x_i$  refers to the  $i$ th sample in a batch and  $g_\varphi$  refers to the feature extractor (here we use the last FC layer in ConvNet as  $g_\varphi$ ). By minimizing  $Loss_q$ , we can amplify the feature activated by tanh function from  $(-1, 1)$  to approach to  $\{-1, 1\}$ .

Following the CRLR, we successively regard each feature as treatment, calculate the balancing loss of confounders and sum it over all the features globally. Formally, we solve the batch

<sup>8</sup> CRLR: Causally Regularize Logistic Regression.

confounder balancing problem as follows:

$$\min_W Lossb = \sum_{j=1}^p \left\| \frac{g_\varphi(X)_{-j}^T \cdot (W \odot I_j)}{W^T \cdot I_j} - \frac{g_\varphi(X)_{-j}^T \cdot (W \odot (1 - I_j))}{W^T \cdot (1 - I_j)} \right\|_2^2 + \alpha \|W\|_2^2 \quad \text{s.t.} \quad \sum_{i=1}^n W_i = 1, W \geq 0, \quad (2)$$

where  $W$  represents sample weights,  $I_j$  means the  $j$ th column of  $I$ , and  $I_{ij}$  refers to the treatment status of sample  $i$  when setting feature  $j$  as treatment variable, and  $\|W\|_2^2$  can reduce the variance of weights to prevent the weights from overfitting outlier samples. Different from CRLR, we define the confounder balancing loss w.r.t. a batch of training samples instead of the whole training samples. Moreover, the sample weights and model parameters are jointly optimized through a supervised way in CRLR, while in CNBB we first fix the model parameters (a.k.a. representation) and learn the sample weights  $W$  through an unsupervised way.

As far as we have learnt an optimal set of sample weights for a batch which can balance the confounder distribution, then we combine the weighted softmax loss and quantization loss and propose our CNBB model:

$$\min_{\theta, \varphi} Lossp = - \sum_{i=1}^n w_i \ln(f_\theta(g_\varphi(x_i)) \cdot y_i) + \lambda Lossq, \quad (3)$$

where  $f_\theta$  refers to softmax layer and  $\lambda$  is a trade-off parameter between classification and quantization.

Algorithm 1 gives the complete steps of the batch balancing method and Fig. 9 illustrates it intuitively.

#### 4.2. Experiments on NICO

In this section, we evaluate the proposed ConvNet with batch balancing (CNBB) in the task of General Non-I.I.D. image classification based on NICO.

##### 4.2.1. Experimental settings

One should note that only category labels are available in our environment which is the most common in real and the most fundamental image classification task. Although some methods, like DCE [15], attempt to learn refined semantics with rich weakly-supervised information, contexts of class are actually hidden variables here. For fair comparison, we choose a typical structure of CNN and CNN with batch normalization [45] (CNN+BN) as baselines. The latter is a popular method in deep learning to improve

#### Algorithm 1 ConvNets with batch balancing (CNBB).

**Input:** Train dataset  $D_{train} = \{(x_i, y_i) | i = 1, \dots, n\}$   
**Output:** Non-linear parameters  $\theta$  and  $\varphi$   
 Initialize  $\theta^{(0)}, \varphi^{(0)}$  and  $t_1 \leftarrow 0$   
**repeat**  
   Sample batch of images  $\{(x_1, y_1), \dots, (x_m, y_m)\}$   
   Extract image features  $\{g_{\varphi^{(t_1)}}(x_1), \dots, g_{\varphi^{(t_1)}}(x_m)\}$   
   Calculate indicator matrix  $I$  of features  
   Initialize sample weights  $W^{(0)}$  and  $t_2 \leftarrow 0$   
   **repeat**  
     Optimize  $W^{(t_2+1)}$  to minimize  $Lossb$  in Eq. 2  
      $t_2 \leftarrow t_2 + 1$   
   **until**  $Lossb$  converges or  $t_2$  reaches maximum  
   Predict  $\{f_{\theta^{(t_1)}}(g_{\varphi^{(t_1)}}(x_1)), \dots, f_{\theta^{(t_1)}}(g_{\varphi^{(t_1)}}(x_m))\}$   
   Optimize  $\theta^{(t_1+1)}$  and  $\varphi^{(t_1+1)}$  to minimize  $Lossp$  in Eq. 3  
    $t_1 \leftarrow t_1 + 1$   
**until**  $Lossp$  converges or  $t_1$  reaches maximum  
**return:**  $\theta$  and  $\varphi$

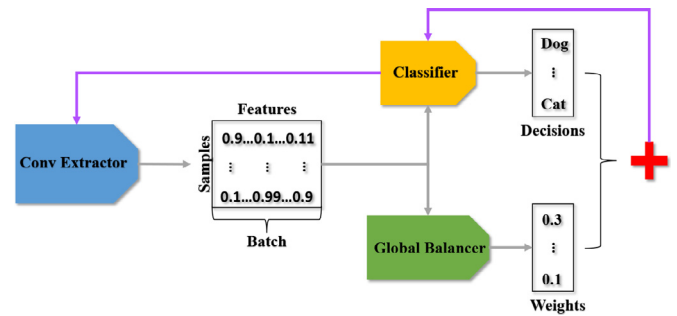


Fig. 9. Info flow in CNBB. The gray and purple lines refer to the forward and backward processes respectively. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

Table 2

Performances of different methods on test accuracy (%) for proportional bias in *Animal* superclass.

Exp2	1 : 5	1 : 1	2 : 1	3 : 1	4 : 1
CNN	37.17	37.80	41.46	42.50	43.23
CNN+BN	38.70	<b>39.60</b>	41.64	42.00	43.85
CNBB	<b>39.06</b>	<b>39.60</b>	<b>42.12</b>	<b>43.33</b>	<b>44.15</b>

Table 3

Performances of different methods on test accuracy (%) for compositional bias in *Vehicle* superclass.

Exp3	3	4	5	6	7
CNN	40.61	42.32	43.34	44.03	44.03
CNN+BN	<b>41.98</b>	38.85	43.12	44.71	44.31
CNBB	41.41	<b>43.34</b>	<b>44.54</b>	<b>45.96</b>	<b>45.16</b>

the generalization ability of CNN by normalizing the scale of activations. All the methods are implemented using PyTorch and optimized by stochastic gradient descent.

We design four experiments according to the supported Non-I.I.D. settings of NICO in Section 3.3

- Minimum bias (Exp 1): In this experiment, we randomly sample 8000 images for training and 2000 images for testing.
- Proportional bias (Exp 2): In this experiment, we fix the dominant ratio of training data to 5:1, and vary the dominant ratio of testing data from 1:5 to 4:1.
- Compositional bias (Exp 3): In this experiment, we vary the number of contexts observed in training data from 3 to 7 while let all the contexts appear in testing data.
- Combined Proportional & Compositional bias (Exp 4): To simulate a more harsh condition, for each class, we randomly select 7 contexts for training and the other 3 contexts for testing. Furthermore, we vary the dominant ratio of training data from 1:1 to 5:1 while fix the dominant ratio of testing data to 1:1.

##### 4.2.2. Experimental results

We calculate the average testing accuracy of all the methods for each experiment. First of all, CNBB is comparable with CNN in the minimum bias setting, with a slightly higher accuracy (49.94% v.s. 49.60%), and CNN+BN performs worst (46.48%). For the other three experiments with explicit distribution shift between training data and testing data, CNBB outperforms the other baselines at almost every setting, as shown in Tables 2–4, indicating its effectiveness in Non-I.I.D. image classification. Note that the performance of CNN with batch normalization is relatively unstable compared to original CNN across different experiments. It is mainly because, in the General Non-I.I.D. setting, the agnostic distribution shift between training and testing data cannot be effectively normalized only based on the training data. Comparatively, the batch balanc-

**Table 4**

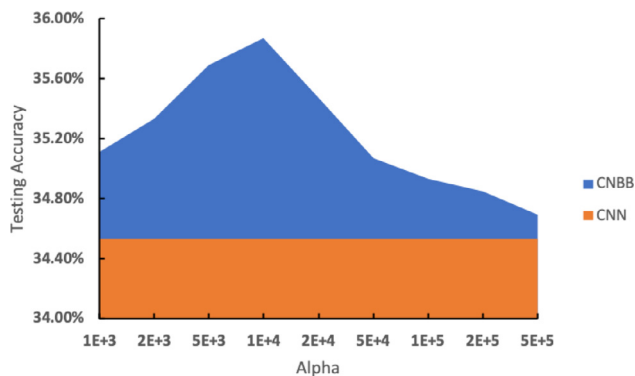
Performances of different methods of test accuracy (%) for combined proportional & compositional bias in *Vehicle* super-class.

Exp4	1 : 1	2 : 1	3 : 1	4 : 1	5 : 1
CNN	37.07	35.20	34.53	34.13	33.73
CNN + BN	33.87	32.93	31.20	30.93	30.67
CNBB	<b>38.98</b>	<b>36.89</b>	<b>35.87</b>	<b>35.33</b>	<b>35.02</b>

**Table 5**

The range of  $NI$  with respect to the average improvement of performance to CNN.

Experiment	Improvement	$NI$
Exp1	0.33%	3.81–3.93
Exp2	1.22%	4.17–4.53
Exp3	1.22%	4.13–4.34
Exp4	1.49%	4.44–4.90



**Fig. 10.** Parameter sensitivity analysis of Exp4. Testing accuracy with respect to the trade-off parameter  $\lambda$  in Eq. (2) while we set dominant ratio of training data to 3:1. The blue area represents the improvement of CNBB against CNN. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

ing module can enable CNBB to identify more stable features and therefore resist the negative effect brought by distribution shift to some extent.

We further summarize the improvement of CNBB over the best baseline in different experiments. From Table 5, we can clearly find that with the discrepancy between the training and testing data getting larger (indicated by higher  $NI$ ), CNBB gains larger improvement over baselines, which demonstrates the advantage of our method in more challenging Non-I.I.D. settings.

Finally, we analyze the hyperparameter  $\alpha$ .  $\alpha$  eventually plays the role of trading-off the valid sample size and degree of batch balancing. In theory, when  $\alpha$  is extremely large, the weights of samples tend to be uniform, resulting in a largest valid sample size. When  $\alpha$  is zero, the algorithm tends to converge to a situation where sample weights concentrate on only a few images, although leading to an optimal batch balancing. Both of valid sample size and degree of batch balancing are critical for the performances of Non-I.I.D. image classification. As in Eq. (2), we tune the hyperparameter  $\alpha$  with 9 values (1e3 to 5e5) in all the experiments. Taking the case where training dominant ratio is 3:1 in Table 4 as an example, a convex hull is clear in Fig. 10. Along with the increasing  $\alpha$ , the gain of CNBB will tend to vanish eventually. The results fully demonstrate the effectiveness of batch balancing module.

## 5. Conclusion and future works

In this paper, we introduce a new dataset NICO for promoting the research on Non-I.I.D. image classification. To the best of

our knowledge, NICO is the first well-structured Non-I.I.D. image dataset with reasonable scale to support the training of ConvNets. By incorporating the idea of context, NICO can provide various Non-I.I.D. settings and create different levels of Non-IIDness consciously. We also propose a simple baseline model with ConvNet structure for General Non-I.I.D. image classification problem, where testing data bear agnostic distribution shift from training data. Empirical results clearly demonstrate the capability of NICO on training the ConvNets and the superiority of the proposed model in various Non-I.I.D. settings.

Our future works will focus on the followings. Firstly, both quality and quantity of NICO continue to be improved. Orthogonal contexts, denoised images and proper area ratio of objects will be explored to make NICO more controllable to tune bias and response to the Non-I.I.D. uniquely. And we will expand the scale of dataset from all the levels for adequate demands. Secondly, more settings about different forms of Non-I.I.D. are expected to be exploited. So other visual concepts may be added to NICO if needed and the ways of using NICO to meet new settings will be given in detail. Thirdly, more effective models will be designed for addressing problems in different settings of Non-I.I.D. image classification.

## Declaration of Competing Interest

None.

## Acknowledgments

This work was supported in part by National Key R&D Program of China (No. 2018AAA0102004), National Natural Science Foundation of China (Nos. U1936219, 61772304, 61531006, U1611461), Beijing Academy of Artificial Intelligence (BAAI), and a grant from the Institute for Guo Qiang, Tsinghua University.

## Appendix A

**Table A.1**

Basic structure of ConvNet used in this paper.

Structure of ConvNet		
Layer	Filter	height & width
input	3	(64 * 64)
conv	64	(64 * 64)
relu		
maxpool	64	(32 * 32)
conv	128	(32 * 32)
relu		
maxpool	128	(16 * 16)
conv	256	(16 * 16)
relu		
maxpool	256	(8 * 8)
conv	512	(8 * 8)
relu		
maxpool	512	(4 * 4)
conv	1024	(4 * 4)
relu		
maxpool	1024	(2 * 2)
fc	512	1
relu		
fc	50	1
tanh		
fc	10/9	1
softmax		



**Table A.2**Data size of each context for every class in *Animal* superclass.

<i>Animal</i>										
Bear	black 245	brown 220	eating grass 133	in forest 243	in water 169	lying 217	on ground 97	on snow 111	on tree 70	white 104
Bird	eating 187	flying 203	in cage 90	in hand 94	in water 81	on branch 239	on grass 242	on ground 276	on shoulder 77	standing 101
Cat	at home 274	eating 270	in cage 109	in river 141	in street 177	in water 50	on grass 140	on snow 137	on tree 50	walking 131
Cow	aside people 56	at home 77	eating 147	in forest 131	in river 139	lying 162	on grass 147	on snow 135	spotter 75	standing 123
Dog	at home 92	eating 264	in cage 122	in street 87	in water 139	lying 143	on beach 280	on grass 158	on snow 238	running 101
Elephant	eating 122	in circus 114	in forest 160	in river 178	in street 90	in zoo 162	lying 69	on grass 103	on snow 69	standing 111
Horse	aside people 124	at home 86	in forest 146	in river 73	in street 77	lying 141	on beach 165	on grass 165	on snow 138	running 143
Monkey	climbing 88	eating 168	in cage 77	in forest 140	in water 118	on beach 50	on grass 106	on snow 102	sitting 168	walking 100
Rat	at home 126	eating 169	in cage 57	in forest 85	in hole 50	in water 85	lying 50	on grass 124	on snow 50	running 50
Sheep	aside people 50	at sunset 66	eating 116	in forest 95	in water 71	lying 109	on grass 132	on road 111	on snow 87	walking 81

**Table A.3**Data size of each context for every class in *Vehicle* superclass.

<i>Vehicle</i>										
Airplane	around cloud 87	aside mountain 76	at airport 153	at night 76	in city 55	in sunrise 70	on beach 104	on grass 53	taking off 128	with pilot 128
Bicycle	in garage 143	in street 113	in sunset 134	on beach 131	on grass 219	on road 125	on snow 163	shared 225	velodrome 220	with people 166
Boat	at wharf 219	cross bridge 190	in city 194	in river 265	in sunset 196	on beach 168	sailboat 252	with people 143	wooden 248	yacht 281
Bus	aside traffic light 35	aside tree 165	at station 95	at yard 74	double decker 221	in city 199	on bridge 45	on snow 124	with people 51	
Car	at park 80	in city 149	in sunset 89	on beach 102	on booth 112	on bridge 36	on road 146	on snow 184	on track 89	with people 39
Helicopter	aside mountain 165	at heliport 185	in city 69	in forest 124	in sunset 160	on beach 107	on grass 147	on sea 156	on snow 180	with people 58
Motorcycle	in city 194	in garage 148	in street 173	in sunset 157	on beach 122	on grass 99	on road 162	on snow 134	on track 185	with people 168
Train	aside mountain 63	at station 158	cross tunnel 36	in forest 100	in sunset 94	on beach 46	on bridge 54	on snow 129	subway 70	
Truck	aside mountain 62	in city 77	in forest 91	in race 134	in sunset 155	on beach 97	on bridge 44	on grass 78	on road 145	on snow 117

## References

- [1] S. Ren, K. He, R. Girshick, J. Sun, (2015). Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems* (pp. 91–99).
- [2] Y. Ma, Y. He, F. Ding, S. Hu, J. Li, X. Liu, *Progressive generative hashing for image retrieval*, In *IJCAI*, 2018, pp. 871–877.
- [3] J. Li, X. Liu, M. Zhang, D. Wang, Spatio-temporal deformable 3D convnets with attention for action recognition, *Pattern Recognit.* (2019), doi:10.1016/j.patcog.2019.107037.
- [4] A. Liu, X. Liu, J. Fan, Y. Ma, A. Zhang, H. Xie, D. Tao, Perceptual-sensitive GAN for generating adversarial patches, *AAAI*, 2019.
- [5] A. Krizhevsky, I. Sutskever, G.E. Hinton, Imagenet classification with deep convolutional neural networks, in: F. Pereira, C.J.C. Burges, L. Bottou, K.Q. Weinberger (Eds.), *Advances in Neural Information Processing Systems*, 25, 2012, pp. 1097–1105.
- [6] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [7] M. Everingham, S.M.A. Eslami, L.V. Gool, C.K.I. Williams, J. Winn, A. Zisserman, The pascal visual object classes challenge: aretrospective, *Int. J. Comput. Vis.* 111 (1) (2015) 98–136.
- [8] T.Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollr, C.L. Zitnick, (2014, September). Microsoft coco: Common objects in context. In *European conference on computer vision* (pp. 740–755). Springer, Cham.
- [9] J. Deng, W. Dong, R. Socher, L.J. Li, K. Li, F.F. Li, Imagenet: a large-scale hierarchical image database, in: *IEEE Conference on Computer Vision & Pattern Recognition*, 2009.
- [10] J. Kauffmann, K.-R. Müller, G. Montavon, Towards explaining anomalies: a deep taylor decomposition of one-class models, *Pattern Recognit.* 101 (2020) 107198.
- [11] B. Zhou, D. Bau, A. Oliva, A. Torralba, Interpreting deep visual representations via network dissection, *IEEE Trans. Pattern Anal. Mach. Intell.* 41 (9) (2018) 2131–2145.
- [12] A. Torralba, A.A. Efros, (2011, June). Unbiased look at dataset bias. In *CVPR 2011* (pp. 1521–1528). IEEE.
- [13] T.S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, Y. Zheng, Nus-wide: a real-world web image database from national university of singapore, in: *Acm International Conference on Image & Video Retrieval*, 2009.
- [14] J. Tang, X. Shu, G.J. Qi, Z. Li, M. Wang, S. Yan, R. Jain, Tri-clustered tensor completion for social-aware image tag refinement, *IEEE transactions on pattern analysis and machine intelligence* 39 (8) (2016) 1662–1674.
- [15] Z. Li, J. Tang, T. Mei, Deep collaborative embedding for social image understanding, *IEEE transactions on pattern analysis and machine intelligence* 41 (9) (2018) 2070–2083.
- [16] Z. Li, J. Tang, Weakly supervised deep matrix factorization for social image understanding, *IEEE Trans. Image Process.* 26 (1) (2016) 276–288.
- [17] Z. Li, J. Tang, Weakly supervised deep metric learning for community-contributed image retrieval, *IEEE Transactions on Multimedia* 17 (11) (2015) 1989–1999.
- [18] Y. Shi, Y. Han, Q. Zhang, X. Kuang, Adaptive iterative attack towards explainable adversarial robustness, *Pattern Recognit.* 105 (2020) 107309.
- [19] D. Liu, L. Zhang, T. Luo, L. Tao, Y. Wu, Towards interpretable and robust hand detection via pixel-wise prediction, *arXiv preprint: 2001.04163* (2020).
- [20] M. Wu, M.C. Hughes, S. Parbhoo, M. Zazzi, V. Roth, F. Doshi-Velez, Beyond sparsity: tree regularization of deep models for interpretability, in: *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [21] S. Muñoz Romero, A. Gorostiaga, C. Soguero-Ruiz, I. Mora-Jiménez, J.L. Rojo-Álvarez, Informative variable identifier: expanding interpretability in feature selection, *Pattern Recognit.* 98 (2020) 107077.
- [22] A. Krizhevsky, G. Hinton, et al., Learning multiple layers of features from tiny images, Technical Report, Citeseer, 2009.

- [23] A. Kuznetsova, H. Rom, N. Alldrin, J. Uijlings, I. Krasin, J. Pont-Tuset, S. Kamali, S. Popov, M. Mallocci, T. Duerig, (2020). The open images dataset v4. *International Journal of Computer Vision*, 1–26.
- [24] B. Thomee, D.A. Shamma, G. Friedland, B. Elizalde, K. Ni, D. Poland, D. Borth, L.-J. Li, Yfcc100m: The new data in multimedia research, arXiv preprint: arXiv:1503.01817 (2015).
- [25] A. Clauset, C.R. Shalizi, M.E.J. Newman, Power-law distributions in empirical data, *SIAM Rev.* 51 (4) (2009) 661–703.
- [26] J. Pearl, D. Mackenzie, *The Book of Why: The New Science of Cause and Effect*, Basic Books, 2018.
- [27] M. Long, H. Zhu, J. Wang, M.I. Jordan, Deep transfer learning with joint adaptation networks, in: *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, JMLR. org, 2017, pp. 2208–2217.
- [28] E. Sangineto, G. Zen, E. Ricci, N. Sebe, We are not all equal: personalizing models for facial expression analysis with transductive parameter transfer, in: *Proceedings of the 22nd ACM International Conference on Multimedia*, ACM, 2014, pp. 357–366.
- [29] C. Deng, X. Liu, C. Li, D. Tao, Active multi-kernel domain adaptation for hyperspectral image classification, *Pattern Recognit.* 77 (2018) 306–315.
- [30] M. Ghifary, W. Bastiaan Kleijn, M. Zhang, D. Balduzzi, Domain generalization for object recognition with multi-task autoencoders, in: *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 2551–2559.
- [31] K. Muandet, D. Balduzzi, B. Schölkopf, Domain generalization via invariant feature representation, in: *International Conference on Machine Learning*, 2013, pp. 10–18.
- [32] S.J. Pan, Q. Yang, A survey on transfer learning, *IEEE Trans. Knowl. Data Eng.* 22 (10) (2010) 1345–1359.
- [33] C. Deng, Y. Xue, X. Liu, C. Li, D. Tao, Active transfer learning network: a unified deep joint spectralspatial feature learning model for hyperspectral image classification, *IEEE Trans. Geosci. Remote Sens.* 57 (3) (2019) 1741–1754.
- [34] J. Pearl, *Causality: Models, Reasoning and Inference*, 29, Springer.
- [35] P.R. Rosenbaum, D.B. Rubin, The central role of the propensity score in observational studies for causal effects, *Biometrika* 70 (1) (1983) 41–55.
- [36] H. Bang, J.M. Robins, Doubly robust estimation in missing data and causal inference models, *Biometrics* 61 (4) (2005) 962–973.
- [37] P.C. Austin, An introduction to propensity score methods for reducing the effects of confounding in observational studies, *Multivar. Behav. Res.* 46 (3) (2011) 399–424.
- [38] J.-P. Pellet, A. Elisseeff, Using Markov blankets for causal structure learning, *J. Mach. Learn. Res.* 9 (Jul) (2008) 1295–1342.
- [39] J. Hainmueller, Entropy balancing for causal effects: a multivariate reweighting method to produce balanced samples in observational studies, *Political Anal.* 20 (1) (2012) 25–46.
- [40] K. Kuang, P. Cui, B. Li, M. Jiang, S. Yang, Estimating treatment effect in the wild via differentiated confounder balancing, in: *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ACM, 2017, pp. 265–274.
- [41] F. Li, K.L. Morgan, A.M. Zaslavsky, Balancing covariates via propensity score weighting, *J. Am. Stat. Assoc.* 113 (521) (2018) 390–400.
- [42] K. Kuang, P. Cui, S. Athey, R. Xiong, B. Li, Stable prediction across unknown environments, in: *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, ACM, 2018, pp. 1617–1626.
- [43] Z. Shen, P. Cui, K. Kuang, B. Li, P. Chen, Causally regularized learning with agnostic data selection bias, in: *2018 ACM Multimedia Conference on Multimedia Conference*, ACM, 2018, pp. 411–419.
- [44] Z. Li, J. Tang, X. He, Robust structured nonnegative matrix factorization for image representation, *IEEE Trans. Neural Netw. Learn. Syst.* PP (99) (2017) 1–14.
- [45] S. Ioffe, C. Szegedy, Batch normalization: accelerating deep network training by reducing internal covariate shift, arXiv preprint: arXiv:1502.03167 (2015).

**Yue He** is a Ph.D. candidate from Lab of Media and Network, Department of Computer Science and Technology, Tsinghua University. His research interests include causal inference and deep learning.

**Zheyang Shen** is a Ph.D. candidate from Lab of Media and Network, Department of Computer Science and Technology, Tsinghua University. His research interests include causal inference and stable learning.

**Peng Cui** is an associate professor with tenure in Department of Computer Science and Technology, Tsinghua University. His research interests include network representation learning, human behavioral modeling, and social-sensed multimedia computing. His contact address is Room 9-316, East Main Building, Tsinghua University, Beijing 100084, P.R.China.