

Error Slice Discovery via Manifold Compactness

Han Yu¹, Hao Zou¹, Jiashuo Liu¹, Renzhe Xu², Yue He³, Xingxuan Zhang¹, Peng Cui^{1*}

¹Tsinghua University

²Shanghai University of Finance and Economics

³Renmin University of China

yuh21@mails.tsinghua.edu.cn, ahio@163.com, liujiashuo77@gmail.com, xurenzhe@sufe.edu.cn

hy865865@gmail.com, xingxuanzhang@hotmail.com, cuip@tsinghua.edu.cn

Abstract

Despite the great performance of deep learning models in many areas, they still make mistakes and underperform on certain subsets of data, i.e. *error slices*. Given a trained model, it is important to identify its semantically coherent error slices that are easy to interpret, which is referred to as the *error slice discovery* problem. However, there is no proper metric of slice *coherence* without relying on extra information like predefined slice labels. Current evaluation of slice coherence requires access to predefined slices formulated by metadata like attributes or subclasses. Its validity heavily relies on the quality and abundance of metadata, where some possible patterns could be ignored. Besides, current algorithms cannot directly incorporate the constraint of coherence into their optimization objective due to absence of an explicit coherence metric, which could potentially hinder their effectiveness. In this paper, we propose *manifold compactness*, a coherence metric without reliance on extra information by incorporating the data geometry property into its design, and experiments on typical datasets empirically validate the rationality of the metric. Then we develop Manifold Compactness based error Slice Discovery (MCSD), a novel algorithm that directly treats risk and coherence as the optimization objective, and is flexible to be applied to models of various tasks. Extensive experiments on the benchmark and case studies on other typical datasets demonstrate the superiority of MCSD.

1 Introduction

In recent years, with the enhancement of computational power, neural networks have achieved significant progress in numerous tasks (Achiam et al. 2023; Liu et al. 2023a; Tian, Ye, and Doermann 2025). Despite their impressive overall performance, they are far from perfect, and still suffer from performance degradation on some subpopulations (Yang et al. 2023). This substantially hinders their application in risk-sensitive scenarios like medical imaging (Yang et al. 2024), autonomous driving (Chen et al. 2024), etc., where model mistakes may result in catastrophic consequences. Therefore, to avoid the misuse of models, it is a fundamental problem to identify subsets (or slices) where a given model tends to underperform. Moreover, we would like to find coherent interpretable semantic patterns in the underperform-

ing slices. For example, a facial recognition model may underperform in certain demographic groups like elderly females. An autonomous driving system may fail in the face of steep road conditions. Identifying such coherent patterns could help us understand model failures, and we could employ straightforward solutions for improvement like collecting new data (Liu et al. 2023b) or upweighting samples in error slices (Liu et al. 2021a).

Previously, works of *error slice discovery* (d'Eon et al. 2022; Wang et al. 2023) aim for this goal. Despite the emphasis on coherence in error slice discovery, there is no proper metric to assess the coherence of a given slice without additional information like predefined slice labels. On one hand, this impairs the efficacy of the evaluation paradigm of error slice discovery. In the current benchmark (Eyuboglu et al. 2022), with the help of metadata like attributes or subclasses, it predefines slices that are already semantically coherent, and they depict the coherence of a slice discovered by a specific algorithm via the matching degrees between it and the predefined underperforming slices, so as to evaluate the effectiveness of the algorithm. Such practice heavily relies on not only the availability but also the quality of metadata, whose annotations are usually expensive, and may overlook model failure patterns not captured by existing metadata. On the other hand, due to the absence of an explicit coherence metric, current algorithms can only indirectly incorporate the constraint of coherence into their design, e.g. via clustering (Eyuboglu et al. 2022; Wang et al. 2023; Plumb et al. 2023), without treating it as a direct optimization objective. This could potentially impede the development of more effective error slice discovery algorithms.

In this paper, inspired by the data geometry property that high dimensional data tends to lie on a low-dimensional manifold (Belkin and Niyogi 2003; Roweis and Saul 2000; Tenenbaum, Silva, and Langford 2000), we incorporate this property to propose *manifold compactness* as the metric of coherence given a slice, which does not require additional information. We illustrate the validity of the metric by showing that it captures semantic patterns better than depicting coherence via metrics directly calculated in Euclidean space, and is empirically consistent with current evaluation metrics that require predefined slice labels. Then we propose a novel and flexible algorithm named Manifold Compactness based error Slice Discovery (MCSD) that jointly optimizes

*Corresponding author

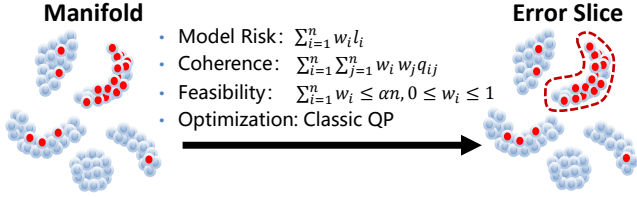


Figure 1: Illustration of MCSD. Blue points are correctly classified by the given trained model, while red ones are wrongly classified. The trained model achieves a good overall accuracy, but exhibit a high error in a certain slice.

the average risk and manifold compactness to identify the error slice. Thus both the risk and coherence, i.e. the desired properties of error slices are explicitly treated as the optimization objective. We illustrate our algorithm in Figure 1. Besides, our algorithm can be directly applied to trained models of different tasks while most error slice discovery methods are restricted to classification only. We provide theoretical analyses of the algorithm. We conduct experiments on dcbench (Eyuboglu et al. 2022) to demonstrate our algorithm’s superiority compared with existing ones. We also provide several case studies on different types of datasets and tasks to showcase the effectiveness and flexibility of our algorithm. Our contributions are summarized below:

- We define manifold compactness as the metric of slice coherence without additional information. We empirically show that it captures semantic patterns well, proving its rationality.
- We propose MCSD, a flexible algorithm that directly incorporates the desired properties of error slices, i.e. risk and coherence, into the optimization objective. It can also be applied to trained models of various tasks.
- We provide theoretical analyses of the algorithm. We conduct experiments on the error slice discovery benchmark to show that our algorithm outperforms existing ones, and we perform diverse case studies to demonstrate the usefulness and flexibility of our algorithm.

2 Problem

Due to space limit, we leave the section of related works in Appendix B. Unless stated otherwise, for random variables, we use uppercase letters, in contrast to a concrete dataset where we use lowercase letters. Consider classic supervised learning. The input variable is denoted as $X \in \mathcal{X}$ and the outcome is denoted as $Y \in \mathcal{Y}$, whose joint distribution is $P(X, Y)$. There exist multiple slices, where j -th slice can be represented as a slice label variable $S^{(j)} \in \{0, 1\}$. For classic supervised learning, the goal is to learn a model $f_\theta : \mathcal{X} \mapsto \mathcal{Y}$ with parameter θ . Denote $\ell : \mathcal{Y} \times \mathcal{Y} \mapsto [0, +\infty]$ as the loss function. Current machine learning algorithms are capable of learning models with a satisfying overall performance, which can be demonstrated via a low risk $\mathbb{E}_P[\ell(f_\theta(X), Y)]$ over the whole population. However, performance degradation could still occur in a certain sub-population or slice. Here we introduce error slice discovery:

Problem 2.1 (Error Slice Discovery). Given a fixed prediction model $f_{\theta_0} : \mathcal{X} \mapsto \mathcal{Y}$ and a validation dataset $\mathcal{D}_{va} = \{(x_i^{va}, y_i^{va})\}_{i=1}^{n_{va}}$, we aim to develop an algorithm \mathcal{A} that takes \mathcal{D}_{va} and f_{θ_0} as input, and learns slicing functions $g_\varphi^{(j)} : \mathcal{X} \times \mathcal{Y} \mapsto \{0, 1\}, 1 \leq j \leq K$. Denote the output of j -th slicing function as \hat{S}_j . We require that the risk in the slice is higher than the population-level risk by a certain threshold: $\mathbb{E}_{X, Y \sim P(X, Y | \hat{S}_j = 1)}[\ell(f_{\theta_0}(X), Y)] > \mathbb{E}_{X, Y \sim P(X, Y)}[\ell(f_{\theta_0}(X), Y)] + \epsilon$, and the discovered slice is as coherent as possible for convenience of interpretation.

The reason why we require an extra validation dataset to implement error slice discovery is that for deep learning models, training data is usually fitted well enough or even nearly perfect. Thus model mistakes on training data carry much less information on models’ generalization ability. This is common practice in previous works (d’Eon et al. 2022; Eyuboglu et al. 2022; Wang et al. 2023). Without ambiguity, we omit the superscript or subscript of “va” for n, x_i, y_i for convenience in the next two sections.

3 Metric

Due to the absence of a proper metric for coherence that is independent of additional information, the current benchmark (Eyuboglu et al. 2022) provides numerous datasets, trained models, and their predefined underperforming slice labels. They employ precision@ k , i.e. the proportion of the top k elements in the discovered slice belonging to the predefined ground-truth error slice as the metric of slice coherence to evaluate error slice discovery algorithms. Although such practice is reasonable to some extent, its effectiveness of evaluation strongly relies on the quality of metadata that composes the underperforming slice labels, which might be not even available under many circumstances.

To eliminate the requirement of predefined slices, we try to propose a new metric of coherence. It is commonly acknowledged that high-dimensional data usually lies on a low-dimensional manifold (Belkin and Niyogi 2003; Roweis and Saul 2000; Tenenbaum, Silva, and Langford 2000). In this case, while direct usage of Euclidean distance cannot properly capture the dissimilarity between data points, the geodesic distance in the metric space of the manifold can. For preliminary justification, here we provide visualization analyses based on different types of dimension-reduction techniques. Among these techniques, PCA mainly preserves pairwise Euclidean distances between data points while t-SNE and UMAP are both manifold learning techniques. In Figure 2, blue dots are correctly classified by the trained model and red dots are wrongly classified. We can see that the visualization of t-SNE and UMAP shows much clearer clustering structures than that of PCA, either having a larger number of clusters or exhibiting larger margins between clusters. This indicates that it could be better to measure coherence in the metric space of a manifold than in the original Euclidean space. Due to space limit, we only present results of the widely adopted facial dataset CelebA (Liu et al. 2015) here, leaving results of other datasets in Appendix A.1, where the same conclusion holds.

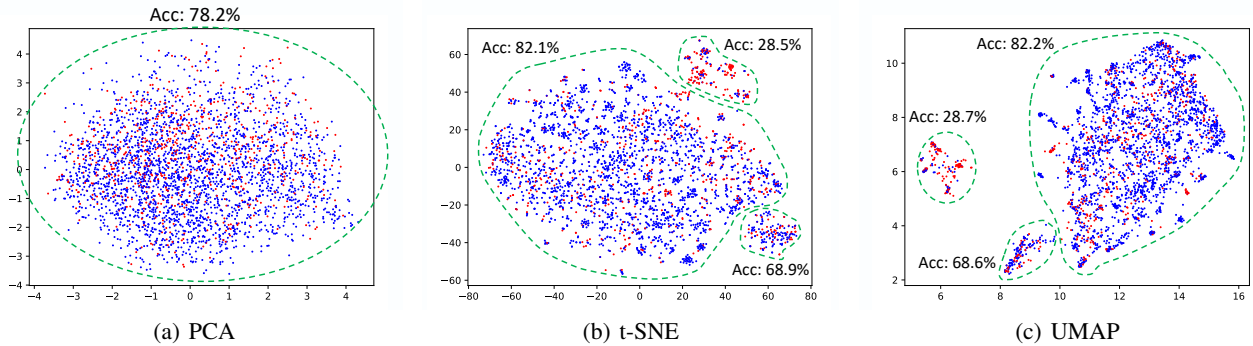


Figure 2: Category “Blond Hair” of CelebA. Visualization of t-SNE and UMAP (manifold-based dimension reduction) shows much clearer clustering structures than that of PCA (mainly preserving Euclidean distances between data points). Thus it could be better to measure coherence in the metric space of a manifold than using metrics directly calculated in Euclidean space.

Therefore, we attempt to define a metric of coherence inside the discovered slice via compactness in the data manifold. In practice, a manifold can be treated as a graph G (Melas-Kyriazi 2020), and we can apply graph learning methods like k-nearest neighbor (kNN) to approximate it (Dann et al. 2022). Given an identified slice $\hat{S} = \{(x_i, y_i) | \hat{s}_i = 1\}$, where \hat{s}_i is the output of the slicing function on i -th sample, we define manifold compactness as:

Definition 3.1 (Manifold Compactness). Consider a given approximation of the data manifold, i.e. a weighted graph $G = (V, E, Q)$. The node set $V = \{v_i\}_{i=1}^n$ corresponds to the dataset $\{(x_i, y_i)\}_{i=1}^n$. The edge set $E = \{e_{ij}\}_{1 \leq i, j \leq n}$, where e_{ij} represents whether node v_i and v_j are connected in the graph G . The weights $Q = \{q_{ij}\}_{1 \leq i, j \leq n}$, where q_{ij} represents the weight of edge e_{ij} . Given a slice \hat{S} , the manifold compactness of it can be defined as:

$$MC(\hat{S}) = \frac{1}{|\hat{S}|} \sum_{(x_i, y_i), (x_j, y_j) \in \hat{S}} q_{ij} \quad (1)$$

This metric is the average weighted degree of nodes of the induced subgraph, whose vertex set corresponds to the slice. The higher it is, the denser or more compact the subgraph is, implying a more coherent slice. Note that when applying this to evaluate multiple slice discovery algorithms, for convenience of comparison, we control the size of \hat{S} for those algorithms to be the same by taking the top αn data points sorted by the slicing function’s prediction probability. Here n is the size of the dataset and $\alpha \in (0, 1]$ is a fixed proportion. The operation of selecting data points with highest prediction probabilities is akin to calculating precision@ k in dcbench (Eyuboglu et al. 2022).

Next, we try to demonstrate the validity and advantages of our proposed coherence metric. A common and representative metric of coherence directly calculated in Euclidean space is variance. Thus we measure variance and manifold compactness respectively on different semantically predefined slices of CelebA (Liu et al. 2015). We use the binary label y to indicate whether the person has blond hair or not,

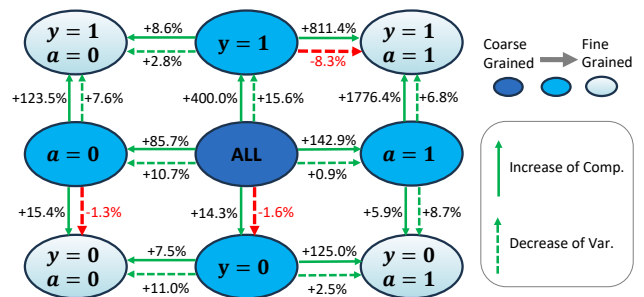


Figure 3: Percentage of increase of manifold compactness (“Comp.”) and decrease of variance (“Var.”) from coarse-grained slices to fine-grained ones in CelebA. For manifold compactness, there is always a positive increase from coarse-grained slices to fine-grained slices. However, in some cases, variance fails to decrease from coarse-grained slices to fine-grained slices as expected, which are marked in red arrows. This could imply that manifold compactness is better at capturing semantic coherence than variance does.

and a to indicate whether the person is male or not. The values of y and a can formulate slices of different granularity. In Figure 3, the most coarse-grained slice is the whole dataset (the darkest circle in the center), the most fine-grained slice is the combination of y and a (the lightest circles in the four corners), and slices of the middle granularity are formulated by either of y and a . Figure 3 shows the percentage of the increase of manifold compactness and the decrease of variance with directed arrows from semantically coarse-grained slices to fine-grained ones. It is intuitive that these digits are supposed to be positive if these two metrics could properly measure semantic coherence. However, for variance, in some cases the value of the more coarse-grained slice is even smaller than the more fine-grained, marked in red arrows. For manifold compactness, there is always a positive increase from semantically coarse-grained slices to fine-grained slices. In this way, we demonstrate that manifold compactness is better at capturing semantic coherence

Algorithm 1: Manifold Compactness based Error Slice Discovery (MCSD)

Input:Validation dataset: $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$.The trained model to be evaluated: $f_{\theta_0} : \mathcal{X} \mapsto \mathcal{Y}$.Size of the slice as a proportion of the dataset: α .Coherence coefficient λ .A pretrained feature extractor: $h_{fe} : \mathcal{X} \mapsto \mathcal{Z}$.**Output:** The identified error slice $\hat{\mathcal{S}}$.**for** $i = 1$ to n **do** Calculate the embedding: $z_i = h_{fe}(x_i)$. Calculate model prediction loss: $l_i = \ell(f_{\theta_0}(x_i), y_i)$.**end for**Establish the kNN graph $G = (V, E, Q)$ based on the embeddings $\{z_i\}_{i=1}^n$.Formulate the quadratic programming problem with variables $\{w_i\}_{i=1}^n$ as Equation (2).

Employ Gurobi to solve the problem in Equation (2).

for $i = 1$ to n **do** $\hat{s}_i = 1$ **if** $w_i > \alpha$ -Quantile of $\{w_i\}_{i=1}^n$ **else** 0.**end for****return:** $\hat{\mathcal{S}} = \{(x_i, y_i) | \hat{s}_i = 1\}$

than variance does. Still due to the space limit, we only provide results of CelebA here, and leave detailed values and results of other datasets, along with comparisons against other metrics directly calculated in Euclidean space like Median Absolute Deviation (MedianAD) in Appendix A.1, where we reach the same conclusion. Besides, in Table 1 of Section 5, we have also empirically shown that the rank order of the four methods according to precision metrics is generally the same as that of manifold compactness. Since the precision metrics are based on predefined slice labels with semantic meanings, it implies that our proposed coherence metric could capture semantic patterns well and is appropriate for evaluation of slice discovery algorithms even when predefined slice labels are absent.

4 Algorithm

We introduce Manifold Compactness based error Slice Discovery (MCSD), a novel error slice discovery algorithm that incorporates the data geometry property by taking manifold compactness into account. In this way, the metrics of both risk and coherence can be treated as the explicit objective of optimization, thus better enabling the identified error slice to exhibit consistent and easy-understanding semantic meanings. The detailed algorithm is described in Algorithm 1. It is worth noting that although we mainly focus on the identified worst-performing slice for convenience of analyses and comparison, our algorithm could discover more error slices by removing the first discovered slice from the validation dataset and applying our algorithm repeatedly to the rest of the dataset for more error slices. Related experiments and analyses are included in Appendix A.2.

First, we approximate the data manifold via a graph. To facilitate the graph learning approach, we obtain the embeddings of the dataset via a pretrained feature extractor (Rad-

ford et al. 2021), i.e. $z_i = h_{fe}(x_i)$, which follows previous works of error slice discovery (Eyuboglu et al. 2022; Wang et al. 2023). Then we construct a kNN graph $G = (V, E, Q)$ based on the embeddings $\{z_i\}_{i=1}^n$, which is a widely adopted manifold learning approach (Zemel and Carreira-Perpiñán 2004; Pedronette, Gonçalves, and Guilherme 2018; Dann et al. 2022). In the graph G , the edge weight $q_{ij} = 1$ if z_j is among the k nearest neighbors of z_i , or else $q_{ij} = 0$.

For the convenience of optimization, instead of direct hard selection, we assign a sample weight w_i for each (x_i, y_i) , which is the variable to be optimized and is restricted in the range $[0, 1]$. We theoretically prove the equivalence between hard selection and sample weight optimization in Appendix A.3. Considering the model risk, we employ the weighted average mean of loss $\sum_{i=1}^n w_i l_i$ as our optimization objective, where $l_i = \ell(f_{\theta_0}(x_i), y_i)$ is the model prediction loss of i th sample given f_{θ_0} . Considering coherence, we adopt manifold compactness in Definition 3.1 as the optimization objective, i.e. $\sum_{i=1}^n \sum_{j=1}^n w_i w_j q_{ij}$. We add these two objectives with a hyperparameter λ . We also restrict the size of the identified slice to be no more than a proportion α of the dataset. Thus we formulate the optimization problem as a quadratic programming (QP) problem:

$$\begin{aligned} \max_{\{w_i\}_{i=1}^n} & \sum_{i=1}^n w_i l_i + \lambda \sum_{i=1}^n \sum_{j=1}^n w_i w_j q_{ij} \\ \text{s.t.} & \sum_{i=1}^n w_i \leq \alpha n \\ & 0 \leq w_i \leq 1, \quad \forall 1 \leq i \leq n \end{aligned} \quad (2)$$

The above QP problem can be solved by classic optimization algorithms or powerful mathematical optimization solvers like Gurobi (Gurobi Optimization 2021). After solving for the proper sample weights $\{w_i\}_{i=1}^n$, we select the top αn samples sorted by the weights as the error slice $\hat{\mathcal{S}}$. Note that in most previous algorithms' workflow, they require prediction probabilities as input (Eyuboglu et al. 2022; Plumb et al. 2023; Wang et al. 2023), thus only applicable to classification, while our algorithm takes prediction loss as input, naturally more flexible and applicable to various tasks.

5 Experiments

In this section, we conduct extensive experiments to demonstrate the validity of our proposed metric and the advantages of our algorithm MCSD compared with previous methods. For quantitative results, we conduct experiments on the error slice discovery benchmark *dcbench* (Eyuboglu et al. 2022). Besides, we conduct experiments on other types of datasets like classification for medical images (Irvin et al. 2019), object detection for driving (Yu et al. 2020), and detection of toxic comments (Borkan et al. 2019), which showcase the great potential of our algorithm to be applied to various tasks. The baselines we compare with are Spotlight (d'Eon et al. 2022), Domino (Eyuboglu et al. 2022), and PlaneSpot (Plumb et al. 2023). More experimental details are included in Appendix A.4.

For evaluation, we compute manifold compactness as the main coherence metric along with the average performance

Table 1: Results of dcbench. We mark the best method in bold type and underline the second-best in terms of each metric. ‘‘Comp.’’ means ‘‘Manifold Compactness’’. ‘‘Corr.’’ means ‘‘Correlation’’. ‘‘ \uparrow ’’ indicates that higher is better. ‘‘%’’ indicates the digits are percentage values.

Metric	Precision@10 (%) \uparrow			Precision@25 (%) \uparrow			Average Precision (%) \uparrow			Manifold Comp. \uparrow		
Method	Corr.	Rare	Noisy	Corr.	Rare	Noisy	Corr.	Rare	Noisy	Corr.	Rare	Noisy
Spotlight	32.3	28.7	43.2	32.2	26.4	40.9	28.9	16.4	22.7	<u>4.78</u>	2.67	4.20
Domino	<u>36.2</u>	<u>52.5</u>	<u>51.7</u>	<u>33.8</u>	<u>52.3</u>	<u>50.0</u>	<u>29.9</u>	<u>37.7</u>	<u>31.3</u>	4.14	<u>4.06</u>	<u>5.53</u>
PlaneSpot	26.1	18.1	29.4	22.3	18.1	27.8	21.8	14.3	18.8	2.93	1.59	3.30
MCSD	47.4	61.1	60.6	45.6	59.8	57.4	40.3	52.4	38.4	6.22	7.81	8.71

of the given model f_{θ_0} on the identified slice \hat{S} . For classification, the performance metric is accuracy. For object detection, it is Average Precision (AP). Note that there are two aspects of evaluation simultaneously. In this case, we put more emphasis on coherence than performance, since we only require the performance of the identified slice to be low to a certain degree but expect it to be as coherent as possible for the benefits of interpretation. This is similar to dcbench (Eyuboglu et al. 2022) where coherence outweighs performance and is chosen as the main evaluation metric.

For running time comparison and related analyses of our method and the baselines, we leave results in Appendix A.5. For the choice and analyses of hyperparameters, we leave them in Appendix A.6. For the improvement of the original models utilizing the discovered error slices, we leave results in Appendix A.7. Due to space limit, we put more examples including those of Spotlight and PlaneSpot in Appendix A.11, about 20 images for each identified slice.

5.1 Benchmark results: Dcbench

Dcbench (Eyuboglu et al. 2022) offers 886 publicly available settings for error slice discovery. Each setting consists of a trained ResNet-18 (He et al. 2016), a validation dataset and a test dataset, both with predefined underperforming slice labels. The validation dataset and its error slice labels are taken as input of slice discovery methods, while the test dataset and its error slice labels are used for evaluation. There are three types of slices in dcbench: correlation slices, rare slices, and noisy label slices. The correlation slices are generated from CelebA (Liu et al. 2015), while the other two types of slices are generated from ImageNet (Deng et al. 2009). More details are included in Appendix A.8. In terms of evaluation metrics, we employ precision@ k and average precision following dcbench’s practice, where precision@ k is the proportion of samples with top k highest probabilities output by the learned slicing function that belongs to the predefined underperforming slice, and average precision is calculated based on precision@ k with different values of k . We also calculate manifold compactness as Definition 3.1. For all these metrics, a higher value indicates higher coherence of the identified slice, thus implying a more effective algorithm capable of error slice discovery.

Effectiveness of our method Table 1 shows that MCSD outperforms other methods across all three types of error slices in precision@10, precision@25, average precision,

and manifold compactness. This greatly exhibits the strengths of our method compared with existing ones in error slice discovery. Among the baselines, Domino consistently ranks 2nd, also showing a fair performance.

Validity of our metric It is also worth noting that the proposed metric manifold compactness shows a strong consistency with other metrics. Table 1 shows that the rank order of the four methods based on precision metrics is usually MCSD, Domino, Spotlight, PlaneSpot, the same as the rank order based on manifold compactness, except for the correlation slice where the rank order of Domino and Spotlight switches. While other metrics require access to predefined labels of underperforming slices, our metric does not. This demonstrates the validity and advantages of our proposed manifold compactness when measuring coherence and evaluating error slice discovery algorithms.

5.2 Case Study: CelebA

CelebA (Liu et al. 2015) is a large facial dataset of 202,599 images, each with annotations of 40 binary attributes. In the setting of subpopulation shift, it is the most widely adopted dataset since it is easy to generate spurious correlations between two specific attributes by downsampling the dataset (Yang et al. 2023; Sagawa et al. 2019; Liu et al. 2021a). Different from settings in dcbench, in this case study we follow Sagawa et al. (2019) to treat the binary label of blond hair as the target of prediction and directly use the whole dataset of CelebA (Liu et al. 2015) without downsampling, thus closer to the real scenario. In terms of implementation details, we employ the default data split provided by CelebA and follow the training process of ERM in (Sagawa et al. 2019) to train a ResNet-50. We apply error slice discovery algorithms on both categories respectively, thus taking advantage of outcome labels that are known during slice discovery. We also illustrate results of directly selecting top αn_{te} samples sorted by prediction losses.

From Table 2, we can see that for both categories of CelebA, our algorithm identifies the most coherent underperforming slice in terms of manifold compactness, where higher is better. Although it ranks 2nd for the category of blond hair in terms of accuracy, where lower is better, for the task of error slice discovery, we put more emphasis on coherence since we want the identified slices to be interpretable, and we only require the performance of the slice to be lower than a threshold compared with the overall perfor-

Table 2: Results on CelebA and the overall accuracy of the trained model. “Acc.” means “Accuracy”. “Comp.” means “Manifold Compactness”. “↑” indicates higher is better, while “↓” indicates lower is better. We mark the best method in bold type and underline the second-best. “%” indicates the digits are percentage values.

Blond Hair?	Yes		No	
Method	Acc. (%) ↓	Comp. ↑	Acc. (%) ↓	Comp. ↑
Spotlight	26.3	5.71	65.9	3.35
Domino	34.6	6.07	82.1	3.58
PlaneSpot	68.4	2.92	93.6	1.13
MCS D	33.8	8.09	75.7	5.54
Overall	76.4	-	98.2	-

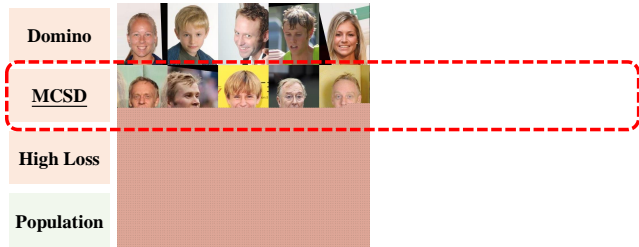


Figure 4: Images randomly sampled from slices of CelebA. Left five columns are results of the category “Blond Hair”. Right five columns are results of the category “Not Blond Hair”. We can see that MCS D is capable of finding error slices that are more coherent than others.

mance, as stated in Section 3. In Figure 4, left five columns and right five columns are from two categories separately. Four rows correspond to randomly sampled images from different sources: the error slice that Domino identifies, the error slice that MCS D identifies, top αn_{te} samples sorted by the loss, and all samples of the corresponding category. We can see that images from the error slice identified by MCS D obviously exhibit more coherent characteristics than others.

For the category of blond hair, images in the row of MCS D are all faces of males, conforming to the intuition that models may learn the spurious correlation between blond hair and female, and could be inclined to make mistakes in subgroups like males with blond hair in the row of MCS D. Although more than half of the images for Domino in the blond hair category are also males, its coherence is much smaller than that of MCS D, making it hard for humans to interpret the failure pattern when compared with images of the whole population. Besides, in the third row, when simply taking account of the prediction loss to select risky samples, it is also difficult to extract the common pattern. For the category of not blond hair, although both Domino and sorting-by-loss can extract the pattern of faces being female with brown hair or blond hair (label noise), MCS D identifies more detailed common characteristics that faces in the images are not only female, but bear vintage styles like in the 20th century, which also constitute a riskier slice than Domino in terms of accuracy. It is also worth noting that MCS D achieves a higher manifold compactness than

Domino in Table 2, consistent with that the identified slice of MCS D exhibits more coherent semantics in Figure 4, further confirming the rationality of our coherence metric.

Table 3: Results on CheXpert and overall accuracy of the trained model. “Acc.” means “Accuracy”. “Comp.” means “Manifold Compactness”. “↑” indicates higher is better, while “↓” indicates lower is better. We mark the best method in bold type and underline the second-best. “%” indicates the digits are percentage values.

Ill?	Yes		No	
Method	Acc. (%) ↓	Comp. ↑	Acc. (%) ↓	Comp. ↑
Spotlight	19.5	2.10	64.9	4.70
Domino	31.5	1.53	88.4	2.82
PlaneSpot	42.8	3.66	69.5	3.17
MCS D	31.5	4.70	63.3	4.87
Overall	45.5	-	91.0	-

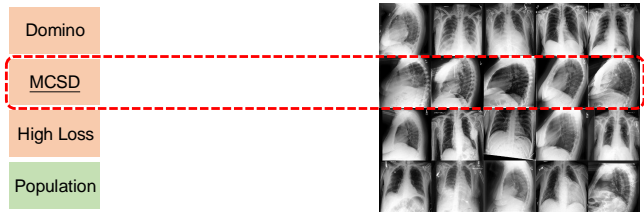


Figure 5: Images randomly sampled from slices of CheXpert. Left five columns are results of the category “Ill”. Right five columns are results of the category “Healthy”. We can see that MCS D is capable of finding error slices that are more coherent than others.

5.3 Case Study: CheXpert

To demonstrate the effectiveness of our algorithm on other types of data, we conduct experiments on a medical imaging dataset, i.e. CheXpert (Irvin et al. 2019), where the task is to predict whether patients are ill or not based on their chest X-ray images. It contains 224,316 images from 65,240 patients. We follow the data split and training process of Yang et al. (2023) to train a ResNet-50. Still, we apply algorithms to images of ill and healthy patients respectively.

In Table 3, we can see that MCS D still achieves highest manifold compactness and relatively low slice accuracy in terms of the discovered error slice for both ill and healthy patients. In Figure 5, for ill patients, images sampled from the error slice discovered by MCS D are all taken from the frontal view, while there are different views for images sampled from other sources. For healthy patients, images corresponding to MCS D are all taken from the left lateral view, while other rows constitute images from different views, making it difficult to extract the common risky pattern. These results showcase MCS D’s usefulness in medical imaging, which is a highly risk-sensitive task and deserves more attention for error slice discovery and failure pattern interpretation. Besides, the consistency of the order of coherence for MCS D and Domino in Table 3 and Figure 5 also confirms the rationality of our proposed coherence metric.

5.4 Case Study: BDD100K

Table 4: Results of algorithms on BDD100K for two categories, along with the overall AP of the trained model. “Comp.” means “Manifold Compactness”. “↑” indicates that higher is better, while “↓” indicates that lower is better. We mark the best method in bold type. “%” indicates that the digits are percentage values.

Category	Pedestrian		Traffic Light	
Method	AP (%) ↓	Comp. ↑	AP (%) ↓	Comp. ↑
Spotlight	57.3	2.05	46.3	2.61
MCSD	53.8	6.60	57.3	4.78
Overall	71.4	-	69.2	-

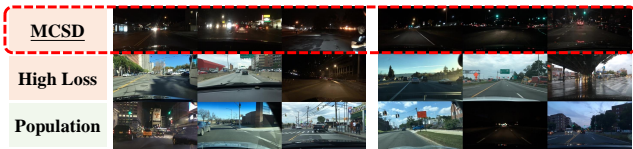


Figure 6: Images randomly sampled from slices of BDD100K. Left three columns are results of the category “Pedestrian”. Right three columns are results of the category “Traffic Light”. It shows that MCSD can find error slices that are more coherent than others.

Compared with most previous algorithms (Eyuboglu et al. 2022; Wang et al. 2023; Plumb et al. 2023) requiring prediction probabilities as a part of input and are only designed for classification tasks, our algorithm is flexible to be employed in various tasks since it takes prediction losses as input. To illustrate its benefits of extending to other tasks, we conduct a case study on BDD100K (Yu et al. 2020), a large-scale dataset composed of driving scenes with abundant annotations. It includes ten tasks, of which we investigate object detection in our paper. The number of images in its object detection task is 79,863, which we split into train, validation, and test datasets with the ratio 2:1:1. We train a YOLOv7 (Wang, Bochkovskiy, and Liao 2023) and try to identify coherent error slices for it. We employ Average Precision (AP) as the performance metric that is widely adopted in detection tasks. Of the 13 categories in the task, we select 2 categories with a relatively high overall performance and a large sample size, i.e. pedestrian and traffic light. We apply our algorithm MCSD for each of them respectively. Note that we do not compare with Domino or PlaneSpot since neither of them is applicable to tasks other than classification.

Table 4 shows that MCSD successfully identifies error slices whose AP are lower than those of overall for both categories, and whose coherence is higher than that of Spotlight in terms of manifold compactness. In Figure 6, each row corresponds to five images randomly sampled from a given source. Left three columns correspond to the category of pedestrians, while right three columns correspond to the category of traffic lights. For both pedestrians and traffic lights, samples from the source of MCSD are coherent in

Table 5: Results on CivilComments, along with overall accuracy of the trained model. “Comp.” means “Manifold Compactness”. “↑” indicates that higher is better, while “↓” indicates that lower is better. We mark the best method in bold type. “%” indicates that the digits are percentage values.

Toxic?	Yes		No	
Method	Acc. (%) ↓	Comp. ↑	Acc. (%) ↓	Comp. ↑
Spotlight	48.6	5.10	91.0	7.33
Domino	56.1	<u>5.98</u>	<u>87.9</u>	6.55
PlaneSpot	46.3	1.65	96.5	2.99
MCSD	25.2	8.56	60.8	7.67
Overall	61.2	-	90.9	-

that they are all taken at night. This conforms to the intuition that it is more difficult to recognize and locate objects when the light is poor. However, directly sampling from the high-loss images can hardly exhibit any common patterns. This reveals the potential of MCSD to extend to other tasks.

5.5 Case Study: CivilComments

In addition to experiments on visual tasks, to demonstrate the applicability of our method to other types of data, we conduct experiments on CivilComments (Borkan et al. 2019), a text dataset of 244,436 comments included in popular distribution shift benchmarks (Yang et al. 2023; Koh et al. 2021). Its task is to predict whether a given comment is toxic or not. We follow the data split and training process of Yang et al. (2023) to train a BERT_{base}. We apply algorithms to toxic and non-toxic comments respectively. In Table 5, we can see that MCSD identifies slices of the lowest accuracy and highest manifold compactness in both categories. We also list two parts of comments that are respectively sampled from the slice identified by applying MCSD to the “toxic” category and from all comments of “toxic” category in Appendix A.9 (**Warning**: many of these comments are severely offensive or sensitive), where each part contains 10 comments. We employ ChatGPT to tell the main difference between the two parts of comments and the reply is “Part 1 is characterized by detailed, historical, and ethical discussions with a critical stance on conservatism and a defense of marginalized groups”. We further check and confirm that comments in part 1, i.e. the slice identified by our method, mostly present a positive attitude towards minority groups in terms of gender, race, or religion. This implies that the model tends to treat comments with excessively positive attitudes towards minority groups as non-toxic, some of which are actually toxic. These results demonstrate our method’s usefulness in text data.

6 Conclusion

In this paper, inspired by the data geometry property, to better evaluate error slice discovery methods, we propose manifold compactness as a metric of slice coherence, which does not rely on predefined underperforming slice labels. Furthermore, with the help of explicit metrics for risk and coherence, we develop an algorithm that directly incorporates risk

and coherence into the optimization objective, which is also flexible to be applied to different scenarios. We conduct experiments on a benchmark and perform multiple case studies to demonstrate both the validity of our proposed metric and the superiority of our algorithm, along with the potential to be flexibly extended to various types of tasks.

Acknowledgements

This work was supported by Tsinghua-Toyota Joint Research Fund, NSFC (No. 62425206, 62141607), and Beijing Municipal Science and Technology Project (No. Z241100004224009). Peng Cui is the corresponding author. All opinions in this paper are those of the authors and do not necessarily reflect the views of the funding agencies.

A Appendix

A.1 More results of preliminary studies of manifold compactness

In this part, we provide more experimental results that demonstrate the validity and advantages of our proposed coherence metric, i.e. manifold compactness. In Section 3 we only present results of CelebA, while here we also present results on other datasets like CheXpert and BDD100K.

Visualization Analyses We provide visualization results of different dimension-reduction methods: PCA, t-SNE, and UMAP, where PCA mainly preserves pairwise Euclidean distances between data points while t-SNE and UMAP are both manifold learning techniques. We employ features extracted by the image encoder of CLIP-ViT-B/32 as input of the dimension-reduction methods. Thus the original dimension (dimension of features extracted by the image encoder of CLIP-ViT-B/32) is 512 and the reduced dimension is 2 for convenience of visualization. In Figure 7 and 8, blue dots are correctly classified by the trained model and red dots are wrongly classified. In Figure 9, the color is brighter when the loss is higher. All three visualizations illustrate that t-SNE and UMAP show much clearer clustering structures than PCA, either showing a larger number of clusters or exhibiting larger margins between clusters. Such results indicate that it is better to measure coherence in the metric space of a manifold instead of using metrics directly calculated in Euclidean space.

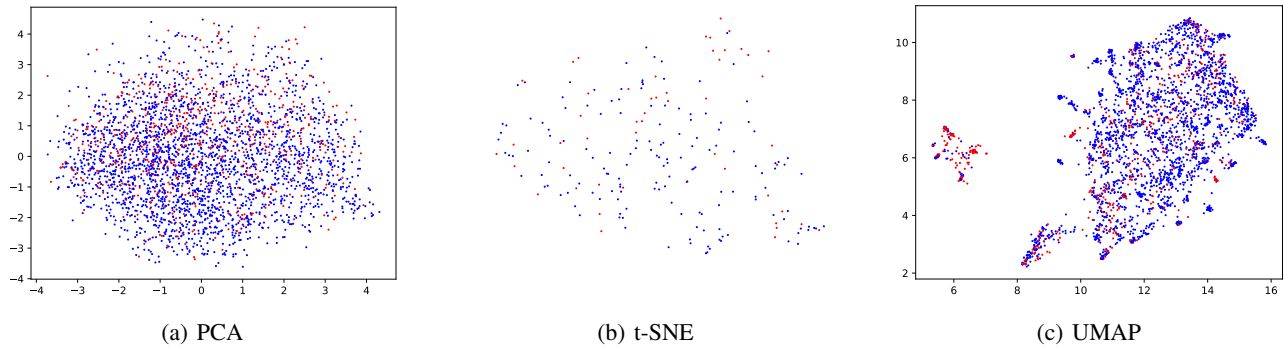


Figure 7: Visualization: Category “blond hair” of CelebA.

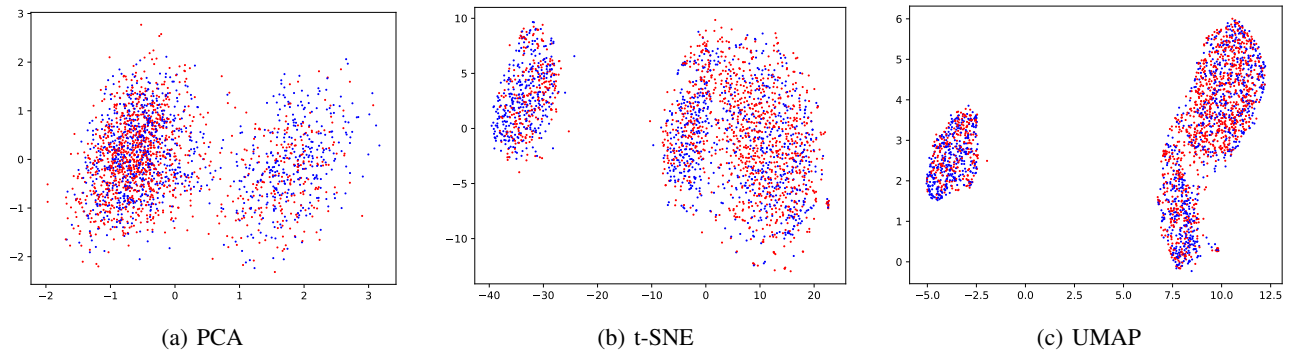


Figure 8: Visualization: Category “ill” of CheXpert.

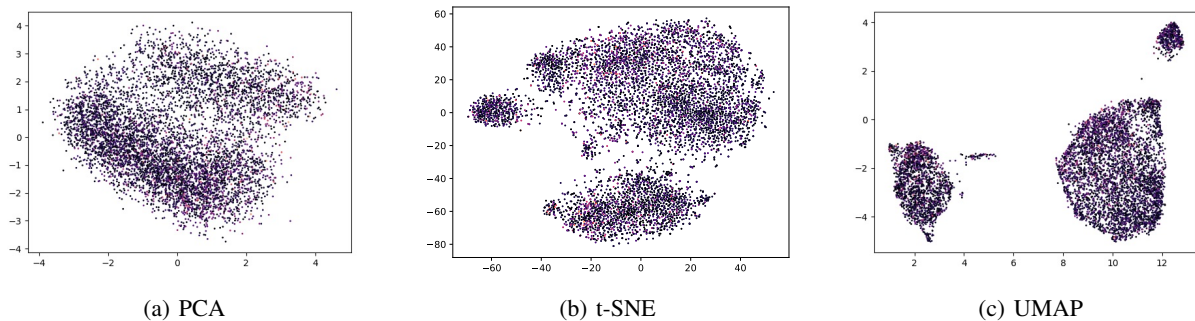


Figure 9: Visualization: Category “pedestrian” of BDD100K.

Comparison with Variance We compare manifold compactness with variance, a common and representative metric of coherence directly calculated in the Euclidean space, on different semantically predefined slices. Note that since the slices are not of the same size, to compare manifold compactness of different slices properly, for each given slice we randomly sample a subset of size 150 with 20 times, and average the manifold compactness of 20 subsets as the manifold compactness of the given slice. For CelebA, we use the binary label y to indicate whether the person has blond hair or not, and a to indicate whether the person is male or not. From Table 6, we can see that for manifold compactness, its value of the more fine-grained slice, i.e. the more coherent slice, is larger than the more coarse-grained slice. For example, the manifold compactness of $y = 1 \& a = 0$ is 0.38, larger than that of $y = 1$ (the value is 0.35) or $a = 0$ (the value is 0.13). Such a relationship holds for every pair of slices. However, in terms of variance, for example, variance of $y = 1 \& a = 1$ is 39.6, larger than that of $y = 1$ whose value is 36.2, which is contrary to our expectation that variance of the more fine-grained slice is smaller than that of the more coarse-grained slice. For CheXpert, we use the binary label y to indicate whether the person is ill or not, and a to indicate whether the person is male or not. We also find that for manifold compactness, its value of the more fine-grained slice, i.e. the more coherent slice, is larger than the more coarse-grained slice, while the value of variance is not consistent with the granularity of the slice. We also additionally compare with other metrics that are directly calculated in Euclidean distance, including Mean Absolute Deviation (MeanAD), Median Absolute Deviation (MedianAD), and Interquartile Range (IQR). We find that they exhibit similar phenomena to variance, i.e. the metric value of the more coarse-grained slice is sometime even smaller than that of the more fine-grained slice, which contradicts our expectation. Thus we demonstrate that manifold compactness is better at capturing semantic coherence than variance does.

Table 6: Comparing manifold compactness with metrics directly calculated in Euclidean space.

Dataset	CelebA					CheXpert				
Slice	Comp.	Var.	MeanAD	MedianAD	IQR	Comp.	Var.	MeanAD	MedianAD	IQR
All	0.07	42.9	113.9	96.2	192.9	0.07	9.4	42.3	34.7	69.2
$y = 1$	0.35	36.2	102.7	86.5	173.1	0.12	10.1	42.1	34.6	69.3
$y = 0$	0.08	43.6	114.6	97.0	194.3	0.07	9.4	42.3	34.6	69.1
$a = 1$	0.17	42.5	114.2	96.4	193.1	0.08	9.8	41.3	33.8	67.5
$a = 0$	0.13	38.3	106.9	90.1	180.4	0.08	9.0	43.2	35.4	70.7
$y = 1, a = 1$	3.19	39.6	107.3	90.4	181.6	0.15	9.5	40.8	33.6	67.4
$y = 1, a = 0$	0.38	35.2	101.1	85.4	171.0	0.17	10.1	43.1	35.4	71.2
$y = 0, a = 1$	0.18	42.5	114.1	96.3	193.0	0.09	9.9	41.3	33.8	67.4
$y = 0, a = 0$	0.15	38.8	107.1	90.2	180.7	0.09	8.6	43.2	35.4	70.6

A.2 Showcase for multiple error slices

In this part, we compare both the worst slice and the second worst slice discovered by our algorithm MCS D and the previous SOTA algorithm Domino. For MCS D, we remove the first error slice from the validation dataset and apply our algorithm again to the rest of the validation dataset to acquire the second error slice. For Domino, we select the slice with the highest and second highest prediction error in the validation dataset as the worst and second worst slice. Results of the blond hair category of CelebA are shown in Figure 10. We find that MCS D is also capable of identifying a coherent slice where faces are female with vintage styles, similar to the error slice also identified by MCS D in Figure 15, while only the pattern of female can be captured in the second worst slice identified by Domino.

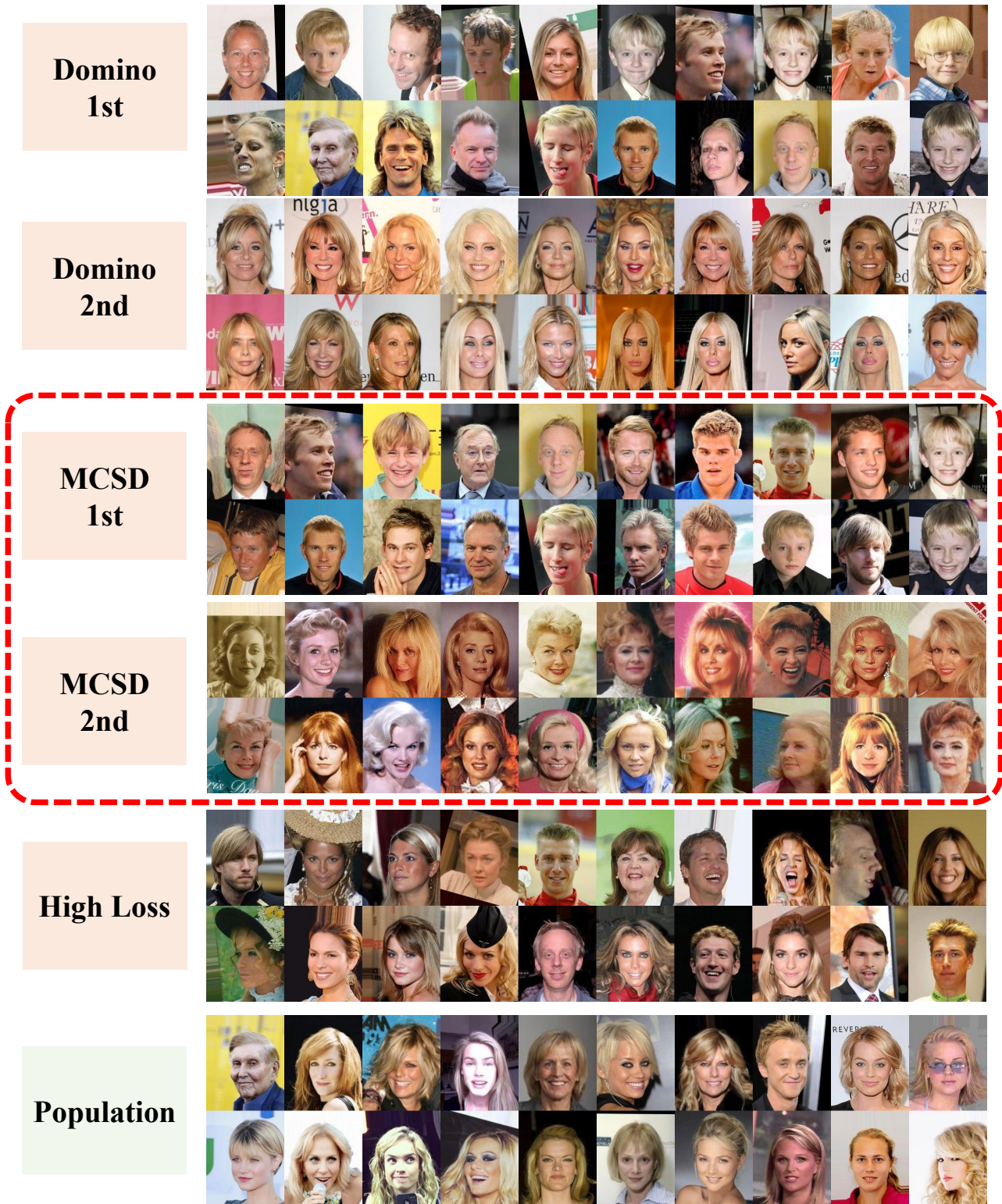


Figure 10: Showcase of multiple error slices for each algorithm on the category “Blond Hair” of CelebA.

A.3 Theoretical analyses

In this subsection, we conduct theoretical analyses on our optimization objective, i.e. Equation (2). First we theoretically prove that the objective not only explicitly considers the manifold compactness inside the identified slice, but also implicitly considers the separability between samples in and out of the identified slice. Then we prove that the optimization objective, where optimized variables are continuous, is equivalent to the discrete version of sample selection with an appropriate assumption. This confirms the validity of our transformation of the problem from the discrete version into the continuous version for the convenience of optimization.

First, we prove a lemma for convenience of later theoretical analyses.

Lemma A.1. *For the inequality constraint $\sum_{i=1}^n w_i \leq \alpha n$ in Equation (2), the equality can be achieved for the solution of Equation (2).*

Proof. Denote the solution of Equation (2) as $w_1^*, w_2^*, \dots, w_n^*$. Assume $\sum_{i=1}^n w_i^* < \alpha n$. Since $\sum_{i=1}^n w_i^* < \alpha n \leq n$, there exists at least one sample weight satisfying $w_k^* < 1$. Let $w'_k = \min\{w_k^* + \alpha n - \sum_{i=1}^n w_i^*, 1\}$. We can see that all constraints in Equation (2) are still be satisfied. However, since $w'_k > w_k^*$, the objective has become larger than before. Thus $w_1^*, w_2^*, \dots, w_n^*$ are not a solution for Equation (2). Since the initial assumption leads to a contradiction, we prove that for the solution of Equation (2), we have $\sum_{i=1}^n w_i^* = \alpha n$.

Next, we prove that Equation (2) also implicitly takes the separability between samples in and out of the identified slice by proving that it is equivalent to an objective that explicitly takes the separability into account.

Proposition A.2. *Maximizing $\sum_{i=1}^n w_i l_i + \lambda \sum_{i=1}^n \sum_{j=1}^n w_i w_j q_{ij}$ under the constraints in Equation (2) is equivalent to maximizing $\sum_{i=1}^n w_i l_i + \lambda_1 \sum_{i=1}^n \sum_{j=1}^n w_i w_j q_{ij} - \lambda_2 \sum_{i=1}^n \sum_{j=1}^n w_i (1 - w_j) q_{ij}$ under the same constraints, where $\lambda = \lambda_1 + \lambda_2$*

Proof. Note that since $\{q_{ij}\}_{1 \leq i, j \leq n}$ corresponds to a kNN graph, we have:

$$\sum_{j=1}^n q_{ij} = k, \forall 1 \leq i \leq n \quad (3)$$

Combined with Lemma A.1, we have:

$$\begin{aligned} & \sum_{i=1}^n w_i l_i + \lambda_1 \sum_{i=1}^n \sum_{j=1}^n w_i w_j q_{ij} - \lambda_2 \sum_{i=1}^n \sum_{j=1}^n w_i (1 - w_j) q_{ij} \\ &= \sum_{i=1}^n w_i l_i + (\lambda_1 + \lambda_2) \sum_{i=1}^n \sum_{j=1}^n w_i w_j q_{ij} - \lambda_2 \sum_{i=1}^n \sum_{j=1}^n w_i q_{ij} \\ &= \sum_{i=1}^n w_i l_i + (\lambda_1 + \lambda_2) \sum_{i=1}^n \sum_{j=1}^n w_i w_j q_{ij} - \lambda_2 \sum_{i=1}^n w_i k \\ &= \sum_{i=1}^n w_i l_i + \lambda \sum_{i=1}^n \sum_{j=1}^n w_i w_j q_{ij} - \lambda_2 \alpha n k \end{aligned} \quad (4)$$

Since $\lambda_2 \alpha n k$ is constant, we have proved the equivalence. Note that the other objective has an extra term of $\sum_{i=1}^n \sum_{j=1}^n w_i (1 - w_j) q_{ij}$, which exactly represents the separability between samples in and out of the identified slice.

Finally, we prove that Equation (2) is equivalent to the original discrete version of sample selection with proper assumptions.

Proposition A.3. *Assume there exists an ordering of sample index $\{r_i\}_{1 \leq i \leq n}$ satisfying that $q_{r_{i-1}, r_i} = q_{r_i, r_{i-1}} = 0$ and $\alpha \cdot n \in \mathbb{N}^+$. Then there exists a solution of Equation (2) $\mathbf{w}^* = \{w_1^*, w_2^*, \dots, w_n^*\}$ such that $w_i^* \in \{0, 1\}, \forall 1 \leq i \leq n$.*

Proof. Define $J(\mathbf{w}) = \sum_{i=1}^n w_i l_i + \lambda \sum_{i=1}^n \sum_{j=1}^n w_i w_j q_{ij}$. We denotes an optimal solution \mathbf{w}' achieving the optimal value of $J(\mathbf{w})$. Then we can find an optimal solution \mathbf{w}^* satisfying $w_i^* \in \{0, 1\}, \forall 1 \leq i \leq n$ and $J(\mathbf{w}^*) = J(\mathbf{w}')$.

We initialize $\mathbf{w}^{(0)} = \mathbf{w}'$. Then we sweep a variable i from 1 to $n - 1$. For each iteration with $1 \leq j \leq n - 1$, we generate a new weight vector $\mathbf{w}^{(j)}$ by the following process. We can prove that $\mathbf{w}^{(j)}$ is one solution of Equation (2) and $w_{r_i}^{(j)} \in \{0, 1\}, \forall 1 \leq i \leq j$ by mathematical induction, which is already satisfied for $j = 0$.

Firstly, we assign $w_{r_i}^{(j)} = w_{r_i}^{(j-1)}$ for $1 \leq i \leq j - 1$ and $j + 2 \leq i \leq n$ and denote $C = w_{r_j}^{(j-1)} + w_{r_{j+1}}^{(j-1)} \in [0, 2]$.

If $w_{r_j}^{(j-1)} \in \{0, 1\}$, we assign $w_{r_j}^{(j)} = w_{r_j}^{(j-1)}$ and $w_{r_{j+1}}^{(j)} = w_{r_{j+1}}^{(j-1)}$.

Otherwise, we reformulate the function $J(\mathbf{w}^{(j)})$ as following (for the sake of brevity, we omit the superscript of (j)):

$$\begin{aligned}
J(\mathbf{w}^{(j)}) &= w_{r_j} l_{r_j} + w_{r_{j+1}} l_{r_{j+1}} + \sum_{i \notin \{r_j, r_{j+1}\}} w_i l_i + \lambda \sum_{i \notin \{r_j, r_{j+1}\}} \sum_{s \notin \{r_j, r_{j+1}\}} w_i w_s q_{is} \\
&+ \lambda \left(w_{r_j} \sum_{i \neq r_j} (q_{i, r_j} + q_{r_j, i}) + w_{r_{j+1}} \sum_{i \neq r_{j+1}} (q_{i, r_{j+1}} + q_{r_{j+1}, i}) + w_{r_j} w_{r_{j+1}} (q_{r_j, r_{j+1}} + q_{r_{j+1}, r_j}) \right) \\
&= w_{r_j} l_{r_j} + (C - w_{r_j}) l_{r_{j+1}} + \sum_{i \notin \{r_j, r_{j+1}\}} w_i l_i + \lambda \sum_{i \notin \{r_j, r_{j+1}\}} \sum_{s \notin \{r_j, r_{j+1}\}} w_i w_s q_{is} \\
&+ \lambda \left(w_{r_j} \sum_{i \neq r_j} (q_{i, r_j} + q_{r_j, i}) + (C - w_{r_j}) \sum_{i \neq r_{j+1}} (q_{i, r_{j+1}} + q_{r_{j+1}, i}) + w_{r_j} (C - w_{r_j}) (q_{r_j, r_{j+1}} + q_{r_{j+1}, r_j}) \right)
\end{aligned} \tag{5}$$

We can see that $J(\mathbf{w}^{(j)})$ is a quadratic or linear function with respect to $w_{r_j}^{(j)}$. Since $q_{r_j, r_{j+1}} = q_{r_{j+1}, r_j} = 0$, $J(\mathbf{w})$ becomes a linear function of $w_{r_j}^{(j)}$.

Because setting $w_{r_j}^{(j)} = w_{r_j}^{(j-1)} \in (0, 1)$ is a global minimum, the coefficient of $J(\mathbf{w})$ with respect to $w_{r_j}^{(j)}$ equals zero. Thus $J(\mathbf{w})$ is constant with respect to $w_{r_j}^{(j)}$. Therefore, we set the value of $w_{r_j}^{(j)}$ and $w_{r_{j+1}}^{(j)}$ as the following two rules:

- If $1 \leq C < 2$, we assign $w_{r_j}^{(j)} = 1$ and $w_{r_{j+1}}^{(j)} = C - 1$
- If $0 < C < 1$, we assign $w_{r_j}^{(j)} = 0$ and $w_{r_{j+1}}^{(j)} = C$

It is obvious that $J(\mathbf{w}^{(j)}) = J(\mathbf{w}^{(j-1)})$. Therefore, $\mathbf{w}^{(j)}$ also achieves the optimal value for Equation (2). Since $w_{r_i}^{(j)} = w_{r_i}^{(j-1)} \in \{0, 1\}, \forall 1 \leq i < j$ and $w_{r_j}^{(j)} \in \{0, 1\}$, we conclude that $w_{r_i}^{(j)} = w_{r_i}^{(j-1)} \in \{0, 1\}, \forall 1 \leq i \leq j$.

Finally, we obtain $\mathbf{w}^{(n-1)}$ where we have $w_{r_i}^{(n-1)} \in \{0, 1\}, \forall 1 \leq i \leq n - 1$. Since the sum of $\mathbf{w}^{(n-1)}$ is an integer αn , $w_{r_n}^{(n-1)}$ is also an integer. According to the construction of $w_{r_n}^{(n-1)}$ in the above two rules, it can be found that $0 \leq w_{r_n}^{(n-1)} \leq 1$. Therefore, we have $w_i^{(n-1)} \in \{0, 1\}, \forall 1 \leq i \leq n$.

The pursued solution of \mathbf{w}^* can be obtained by setting $\mathbf{w}^* = \mathbf{w}^{(n-1)}$.

Remark Since $n \gg k$, it is likely that the constructed graph is extremely sparse. Therefore, it is easy to find a sample ordering that the contiguous samples are not connected, which means that our assumption is satisfied. This proposition proves the equivalence between our continuous optimization formulation to the discrete sample selection.

A.4 Experimental details and related analyses

Details of baselines are listed below:

- Spotlight (d'Eon et al. 2022): It learns a point in the embedding space as the risky centroid, and chooses the closest points to the centroid as the error slice.
- Domino (Eyuboglu et al. 2022): It develops an error-aware Gaussian mixture model (GMM) by incorporating predictions into the modeling process of GMM.
- PlaneSpot (Plumb et al. 2023): It combines the prediction confidence and the reduced two-dimensional representation together as the input of a GMM.

Note that in all our experiments, we apply algorithms to the validation dataset $\{x_i^{\text{va}}, y_i^{\text{va}}\}_{i=1}^{n_{\text{va}}}$ to obtain the slicing function g_φ , and then employ g_φ on the test dataset $\{x_i^{\text{te}}, y_i^{\text{te}}\}_{i=1}^{n_{\text{te}}}$ to acquire the prediction probability of each test sample belonging to the error slice. We choose the top αn_{te} samples from the test dataset sorted by the prediction probabilities as the error slice \hat{S} , and calculate evaluation metrics based on it. As for our method that outputs the error slice of the validation dataset instead of a slicing function, to compare with other methods, we additionally train a binary MLP classifier $g_\varphi : \mathcal{Z} \mapsto [0, 1]$ on top of embeddings, by treating samples in the error slice as positive examples, and treat the rest as negative ones. However, it is worth noting that our method is effective at error slice discovery without this additional slicing function, which is illustrated in Appendix A.4. Meanwhile, we explain the choice of the feature extractor h_{f_e} and conduct analyses in Appendix A.4.

Ablation of the slicing function In Algorithm 1, after acquiring the desired samples, we additionally train an MLP as the slicing function. This is because the standard evaluation process of error slice discovery, an error slice discovery method is applied to validation data to obtain the slicing function, and then the slicing function is applied to test data to calculate evaluation metrics and conduct case analyses. Such practice is also adopted by dcbench (Eyuboglu et al. 2022), the benchmark of error slice discovery. Thus we follow this practice by additionally train a slicing function after we obtain the desired samples, so that we can compare fairly with previous methods. However, even without training this slicing function, our method can still

produce meaningful results of case studies. We show examples randomly sampled from optimized results of Equation (2) in Figure 11 and 12. We can see that there is still high coherence in the identified slices.



Figure 11: Examples of the category “blond hair” of CelebA directly sampling from optimized results of Equation (2).



Figure 12: Examples of the category “not blond hair” of CelebA directly sampling from optimized results of Equation (2).

Choice of feature extractors Following previous works of error slice discovery (Eyuboglu et al. 2022; Wang et al. 2023), for image data, we employ the image encoder of CLIP with a backbone of ViT-B/32 to extract embeddings of images for error slice discovery algorithms in our main experiments. For text data, we employ pretrained BERT_{base} to extract embeddings. To show that our method is flexible in the choice of the feature extractor h_{fe} , we conduct additional experiments on CelebA by changing CLIP-ViT-B/32 to CLIP-ResNet50 and ImageNet-supervised-pretrained ResNet50. Table 7 shows that whatever the pretrained feature extractor is, MCSD consistently identifies slices of low accuracy and outperforms other methods in terms of manifold compactness. It is worth noting that this conclusion is valid even for MCSD with ResNet50, which is generally considered as a weaker pretrained feature extractor than ViT-B/32 employed by baselines. In Figure 13, we can see that MCSD with different pretrained feature extractors truly identifies coherent error slices for the blond hair category of CelebA. As for the practice of using pretrained feature extractors, it is acceptable and generally adopted in previous works of error slice discovery (Eyuboglu et al. 2022; Jain et al. 2023; Plumb et al. 2023).

Table 7: Experiments using different pretrained feature extractors.

Blond Hair?	Yes		No	
	Acc. (%) ↓	Comp. ↑	Acc. (%) ↓	Comp. ↑
Spotlight	26.3	5.71	65.9	3.35
Domino	34.6	6.07	82.1	3.58
PlaneSpot	68.4	2.92	93.6	1.13
MCSD(CLIP-ViT-B/32)	33.8	8.09	75.7	5.54
MCSD(CLIP-ResNet50)	29.3	8.77	71.7	5.38
MCSD(Supervised-ResNet50)	32.8	7.22	<u>67.0</u>	4.75
Overall	76.4	-	98.2	-



Figure 13: Left five images are sampled from the slice identified by MCS D (CLIP-ResNet50). Right five images are sampled from the slice identified by MCS D (Supervised-ResNet50).

A.5 Time comparison

Since the optimization process of our method formulates a non-convex quadratic programming problem and we employ Gurobi optimizer to solve it, it is hard to analyze the time complexity. However, we directly provide a running time comparison. Here we report the running time of different methods on CelebA. The two rows of results correspond to the two categories of CelebA, where the time is measured in seconds. The time consumption of solely constructing the kNN graph is also listed in the last column. We can see that although our method MCS D requires longer running time, the time cost is still generally low and acceptable. Furthermore, in terms of scalability to very large datasets, it is worth noting that our method only requires a validation dataset to work. The validation dataset is essentially a subset sampled from the whole dataset, whose size is much smaller than that of the whole dataset. For example, in CelebA, the validation data size is only 19,867, about 1/10 of the whole dataset size of 202,599. This indicates that for a very large dataset, we could sample a small and appropriate proportion of the whole dataset, and it would be possible for our method to be still effective when being applied to the subset. Besides, the construction of the kNN graph is fast and only takes up a small proportion of running time, which is not the bottleneck of time consumption.

Table 8: Time comparison measured in seconds.

Blond Hair?	Data Size	Spotlight	Domino	PlaneSpot	MCS D	kNN graph
Yes	3,056	16.8	1.3	7.1	39.1	2.0
No	16,811	93.0	26.3	45.3	171.0	13.5

A.6 Hyperparameter selection and analyses

Hyperparameter selection For the hyperparameter of coherence coefficient λ , we fix it as 1 for experiments on dcbench. For the case studies, we set the search space of λ as $\{0.5, 0.8, 1.0, 1.5, 2.0, 2.5, 3.0\}$. We split the validation dataset into two halves, apply our algorithms on one half to obtain a slicing function, apply the slicing function on the other half, and calculate the average performance and manifold compactness of the discovered slice. We choose λ that maximizes manifold compactness under the condition that the slice performance is significantly lower than the overall performance, where the threshold can be customized for different tasks. In our experiments, we set it as 15 percent point for accuracy in terms of image classification, and 10 percent point for average precision (AP) in terms of object detection. For the hyperparameter of slice size α , in our experiments we set it as 0.05 when the size of the validation dataset is smaller than 5,000, and set it as 0.01 otherwise. For building kNN graphs, we fix $k = 10$ in our experiments.

Hyperparameter analyses In this part, we conduct hyperparameter analyses on the category “Blond Hair” of CelebA for the coherence coefficient λ , the size α , and the number of neighbors k when building the kNN graph. From Table 9, we find that both the accuracy and manifold compactness are best when λ and α are in a moderate range, neither too large nor too small. This implies the importance of the balance between pursuing high error and high coherence, which could be achieved by the tuning strategy mentioned in Appendix A.6. This also implies the importance of appropriately controlling α , i.e. the size of the slice, which is set according to experience in our implementation. Its selection is left for future work.

For k , we initially find that $k = 10$ works well and thus fix it. In Table 9 where the manifold compactness of other values has been rescaled to the case of $k = 10$, we can see that the accuracy of the identified slice is generally low compared with the overall accuracy of the blond hair category (76.4%), and the compactness is high when $10 \leq k \leq 30$. Although $k = 15$ is slightly better than $k = 10$ in terms of compactness, it is still appropriate to select $k = 10$ since it is computationally more efficient.

A.7 Performance improvement via utilization of the discovered error slices

We conduct experiments to show performance improvement that the identified error slices could bring via data collection guided by the interpretable characteristics of identified error slices, following the practice of non-algorithmic interventions of Liu et al. (2023b). For example, for a given trained image classification model on CelebA, the identified error slice exhibits characteristics of blond hair male, then we could be guided to collect specific data of the targeted characteristics of blond hair male and add

Table 9: Hyperparameter analyses on the category ‘‘Blond Hair’’ of CelebA for the coherence coefficient λ , the size α , and the number of neighbors k . ‘‘ \uparrow ’’ indicates that higher is better, while ‘‘ \downarrow ’’ indicates that lower is better. We mark the best method in bold type and underline the second-best. ‘‘%’’ indicates that the digits are percentage values.

λ	Acc. (%) \downarrow	Comp. \uparrow	α	Acc. (%) \downarrow	Comp. \uparrow	k	Acc. (%) \downarrow	Comp. \uparrow
0	27.8	2.94	0.005	46.2	1.00	3	21.1	4.16
0.5	19.6	3.36	0.01	19.2	3.31	5	20.3	5.01
0.8	18.1	3.71	0.03	22.8	5.84	10	33.8	8.09
1.0	<u>19.6</u>	4.85	0.05	<u>33.8</u>	8.09	15	42.1	8.13
1.5	<u>30.1</u>	7.14	0.1	48.9	<u>8.07</u>	20	41.4	7.60
2.0	33.8	8.09	0.15	49.4	7.60	30	36.8	7.98
2.5	39.9	7.93	0.2	56.6	7.40	50	36.8	6.73
3.0	45.1	<u>7.99</u>	0.3	67.7	7.54	100	34.2	5.66

to the training data, which is a non-algorithmic intervention and a straightforward and practical way of improving performance of the original model after interpreting characteristics of the identified error slice. Here we compare the results of guided data collection and random data collection. To simulate the guided data collection process, for CelebA, since the identified error slice for the category of blond hair is male and there are extra annotations of sex, we add the images annotated as blond hair male in validation data to training data. Since the identified error slice for not blond hair category is female bearing vintage styles, and there are no related attribute annotations, we directly add the images of the identified slice to the training data. For CheXpert, the identified error slices are from the frontal view for ill patients and from the left lateral view for healthy patients, and CheXpert has annotations of views, so we add the corresponding images in validation data to training data. To simulate the random data collection process, we randomly sample the same number of images from validation data and add to training data for each dataset. Then we retrain the model three times with varying random seeds.

Table 10: Performance of different data collection strategies. ‘‘ \uparrow ’’ indicates that higher is better. We mark the best strategy in bold type. ‘‘%’’ indicates that the digits are percentage values.

CelebA	Average Acc. (%) \uparrow	Worst Group Acc. (%) \uparrow	CheXpert	Average Acc. (%) \uparrow	Worst Group Acc. (%) \uparrow
Original	95.3	37.8	Original	86.7	40.3
Random	95.3 \pm 0.4	42.4 \pm 2.2	Random	88.1 \pm 0.2	50.0 \pm 1.1
Guided	95.4\pm0.5	59.8\pm3.0	Guided	88.7\pm0.8	70.1\pm1.7

Here worst group accuracy is defined following a distribution shift benchmark (Yang et al. 2023), where CelebA and CheXpert are divided into groups according to annotated attributes, and worst group accuracy is an important metric. From Table 10, we can see that guided data collection outperforms the original model and random data collection in both metrics, especially in worst group accuracy. This illustrates that our method is beneficial to performance improvement in practical applications.

A.8 Benchmark details

Dcbench (Eyuboglu et al. 2022) offers a large number of settings for the task of error slice discovery. Each setting consists of a trained ResNet-18 (He et al. 2016), a validation dataset and a test dataset, both with labels of predefined underperforming slices. The validation dataset and its error slice labels are taken as the input of slice discovery methods, while the test dataset and its error slice labels are used for evaluation.

There are 886 settings publicly available in the official repository of dcbench¹, comprising three types of slices: correlation slices, rare slices, and noisy label slices. The correlation slices are generated from CelebA (Liu et al. 2015), a facial dataset with abundant binary facial attributes like whether the person wears lipstick. Correlation slices include 520 settings. They bear resemblance to subpopulation shift (Yang et al. 2023), where a subgroup is predefined as the minor group by generating spurious correlations between two attributes when sampling training data. That subgroup also tends to be the underperforming group after training. The other two types of slices are generated from ImageNet (Deng et al. 2009), which has a hierarchical class structure. Rare slices include 118 settings constructed by controlling the proportion of a predefined subclass to be small. Noisy label slices include 248 settings formulated by adding label noise to a predefined subclass. Although many settings comprise more than one predefined error slice, we check and find that the given model actually achieves even better performance on a number of slices than the corresponding overall performance. For accurate and convenient evaluation, we select the worst-performing slice of the given model for each setting.

¹<https://github.com/data-centric-ai/dcbench>

A.9 Examples from CivilComments

Warning: Many of these comments are severely offensive or sensitive

Here in Table 11 we list two parts of comments that are respectively sampled from the slice identified by applying MCSD to the “toxic” category and from all comments of “toxic” category. Since some comments are too long, we do not list the complete comments but additionally list the id of these comments in the dataset for convenience of checking. We check the complete comments and confirm that comments belonging to the error slice identified by MCSD mostly exhibit a positive attitude towards minority groups in terms of gender, race, or religion. This implies that the model tends to treat comments with positive attitudes towards minority groups as non-toxic, while some of these comments are also offensive and toxic.

Table 11: Comments and their id that are respectively sampled from the slice identified by applying MCSD to the “toxic” category and from all comments of “toxic” category. (**Warning: Many of these comments are severely offensive or sensitive**)

Slice	Content	Id
MCSD	The Kingdom of Hawai’i has a long, proud history of being the most diverse and inclusive nation...	5054686
	You have it a bit twisted, TomZ. You say “ ...there is evidence that Judge supported same-sex acts ...	5977193
	scuppers: Go outside. Seriously. You are so invested in this narrative that you’re completely losing...	5865832
	It’s not racist to shun people who believe apostates and blasphemers against Islam should be killed...	6155392
	So, libs have as a leader a person with , IQ less than 70, who can not make a full statement without...	5316528
	Wow! The US Catholic Church is learning only this year - in 2017 - that racism is rife in the country...	6320767
	Who is asking for special accommodations here? Transgender people who just want to exist and live...	5665161
	Poor analogy. Both the KKK and the Blacks are Christians. So cross burning is racial not religious....	348794
	Brian Griffin quoted Mencken: “The common man’s a fool”. Peter Griffin is proof.	691891
Wow.... Trump isn’t a very deep thinker, and neither is anyone who supports this stupid rule. First...	5438617	
Population	Asian countries for Asians. Black countries for Blacks. but White countries for everybody? That’s genocide.	6299730
	The rapist was a Stanford student; his victim was not. The judge was a Stanford alumnus. We’re looking...	343947
	I wonder if Trump would be in favour of hot black women kneeling?	6020870
	Plato condemned homosexual relationships as contrary to nature. What are you smoking and where can ...	5614294
	Black Pride = being black and proud Gay Pride = being gay and proud White Pride = NAZI!	5815448
	That was Brennan, under Obama, not our good American, Christian president. You really should not make...	6250038
	So when it’s a pretty white woman murdered by her boyfriend, the ANCWL pickets outside the courtroom...	5763897
	pnw mike, you are right! hillary is a liar. One metric comes from independent fact-checking website...	520428
Reminds me of an old Don Rickles joke. “Why do jewish men die before their wives?.....Because they ...	5215731	
When do see the piece on the worlds most annoying Catholics? Buddhist? Muslims?	375375	

A.10 Code

Code is available at <https://github.com/h-yu16/MCSD>.

A.11 More examples for case studies

In this part, we provide more examples for the case studies of visual tasks in our main paper. For CelebA (Figure 14 and 15) and CheXpert (Figure 16 and 17), we randomly sample 20 images from each slice and put 10 images in a row. For BDD100K (from Figure 18 to 25), we randomly sample 18 images for each slice and put 3 images in a row for clearer presentation. We also draw the predicted bounding box with red color and the ground truth bounding box with yellow color. Experimental findings are basically the same as those in our main paper. MCSD still consistently identifies coherent slices in these three cases. Note that in CheXpert, previous algorithms like Spotlight and PlaneSpot are also able to identify coherent slices, illustrating a certain degree of their effectiveness in error slice discovery.

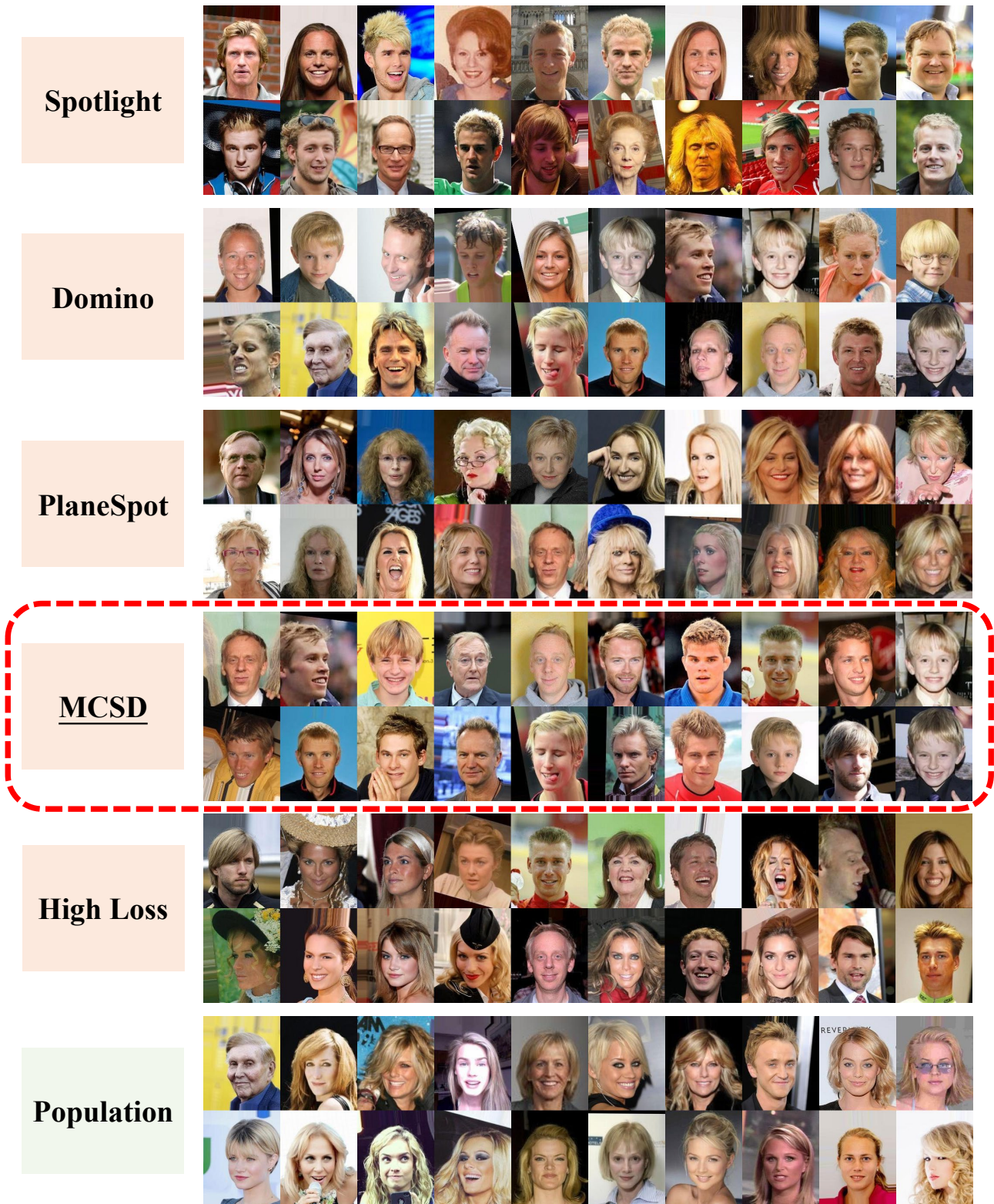


Figure 14: More examples of the category “Blond Hair” of CelebA.

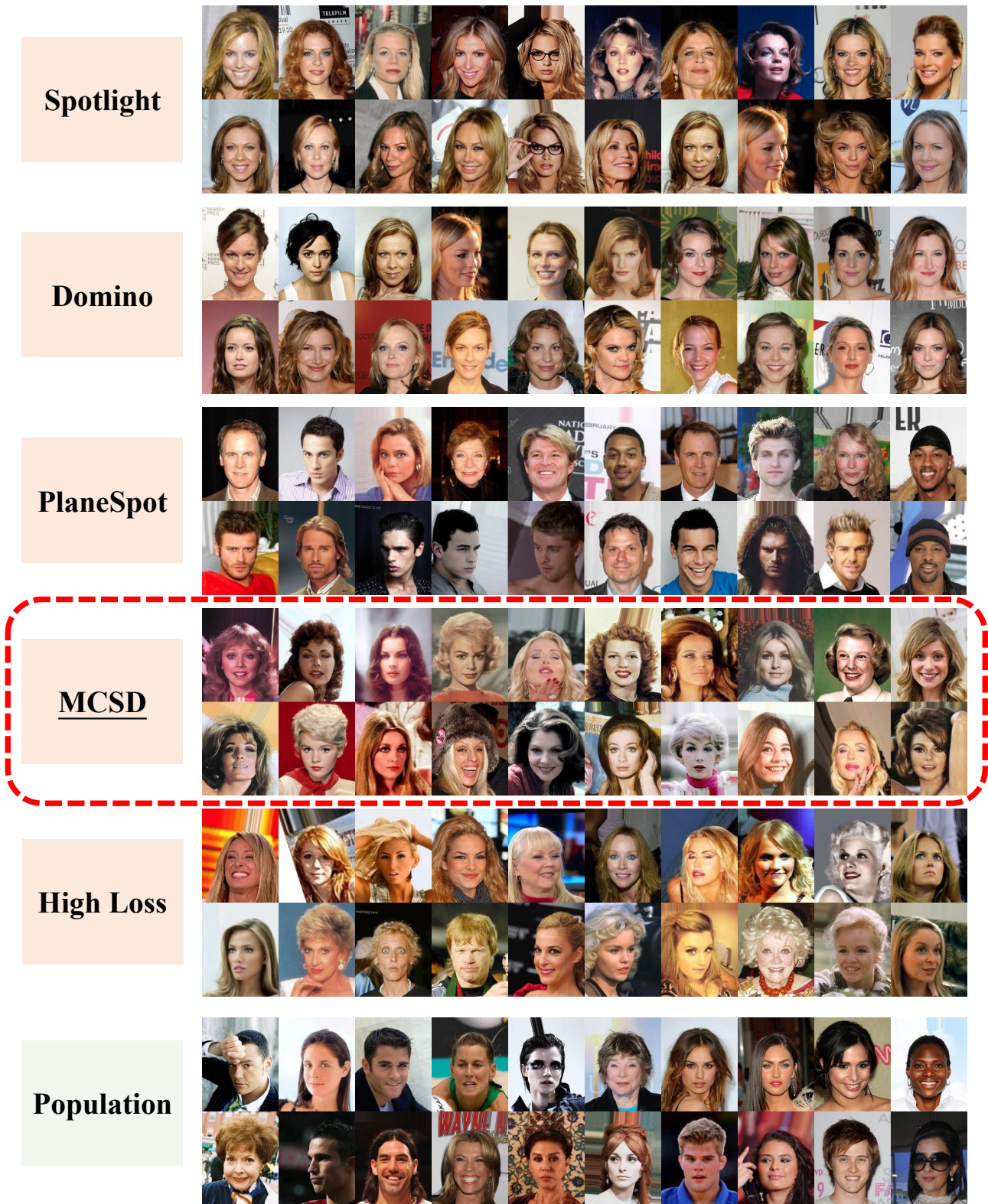
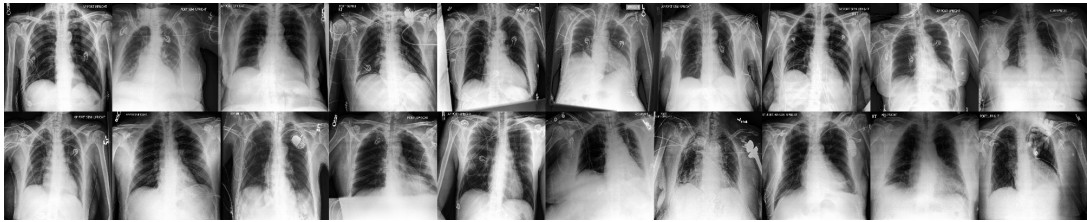
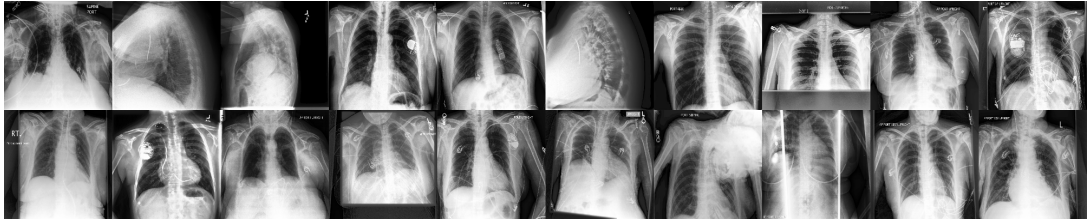


Figure 15: More examples of the category “Not Blond Hair” of CelebA.

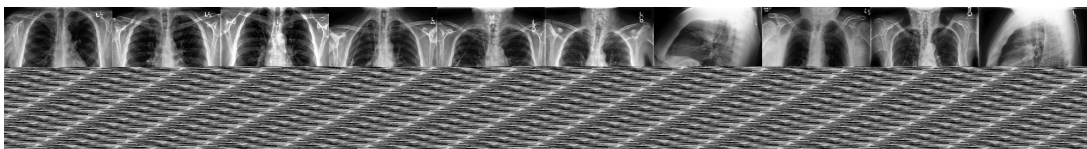
Spotlight



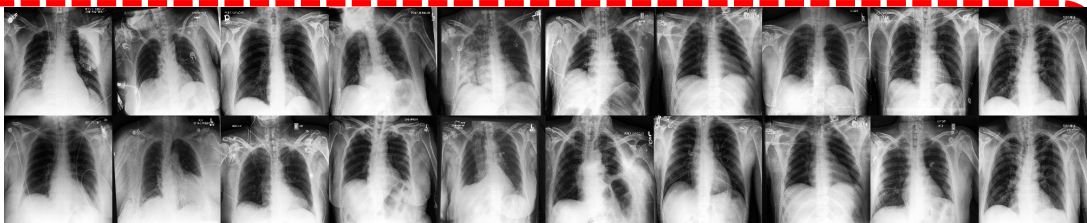
Domino



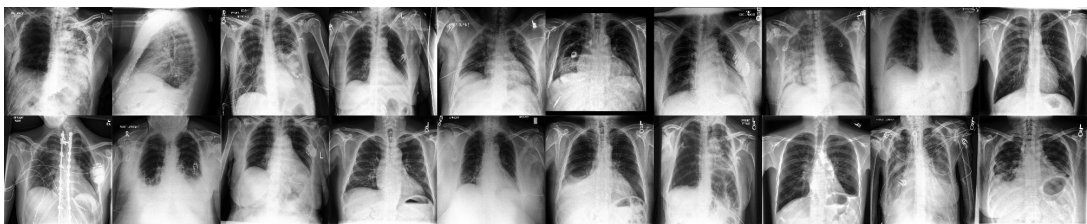
PlaneSpot



MCSD



High Loss



Population

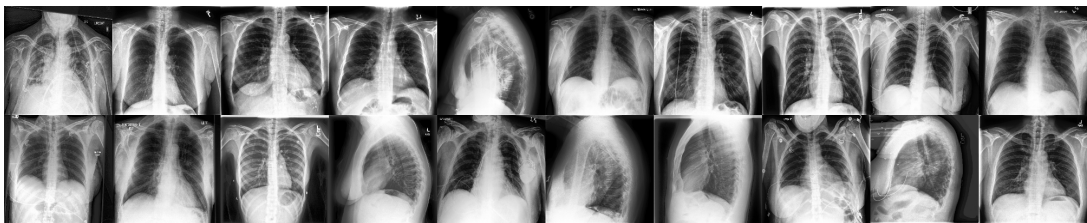


Figure 16: More examples of the category “ill” of CheXpert.

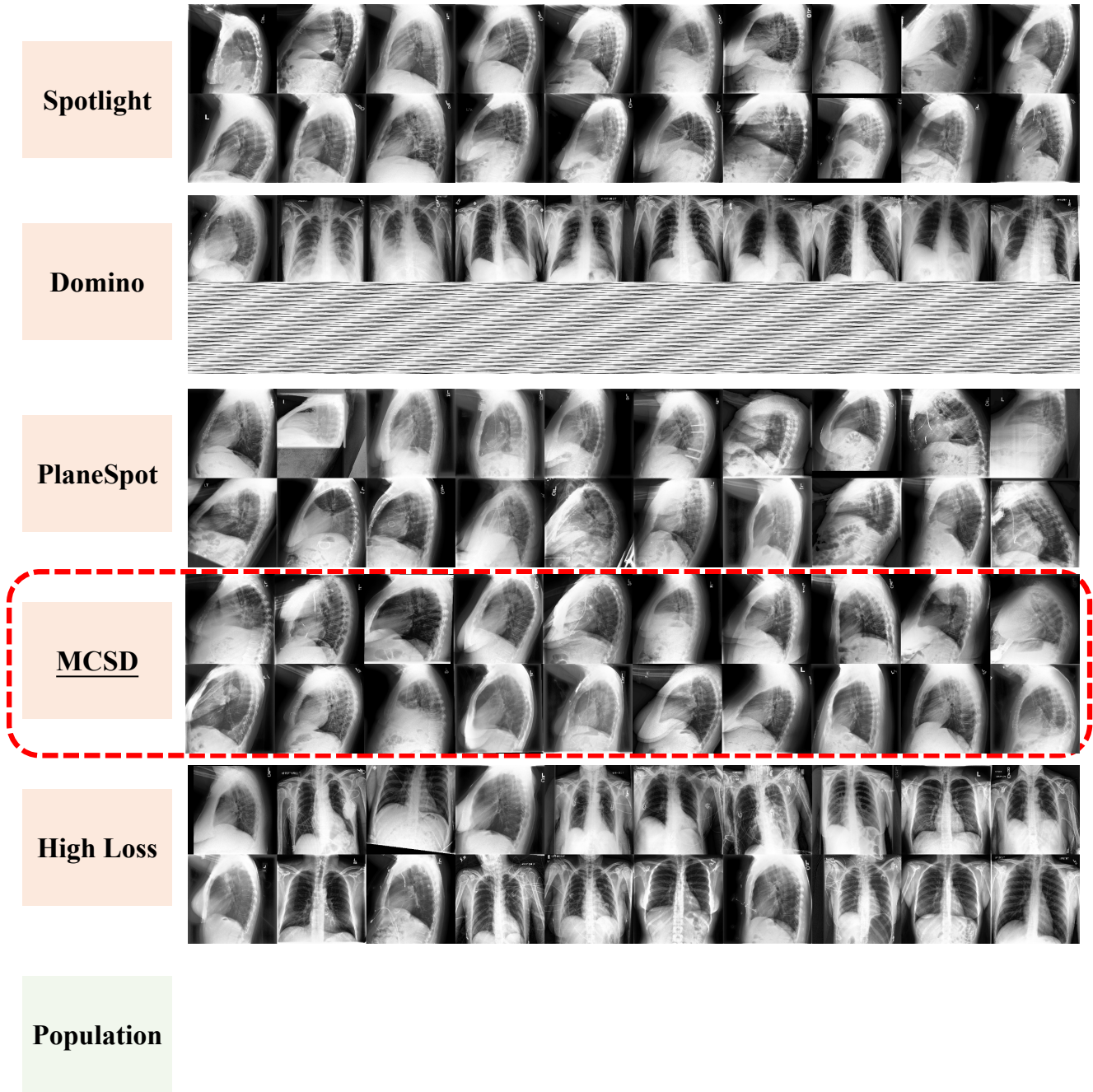


Figure 17: More examples of the category “healthy” of CheXpert.



Figure 18: More examples of the category “Pedestrian” of BDD100K via Spotlight.

Figure 19: More examples of the category “Pedestrian” of BDD100K via MCSD.

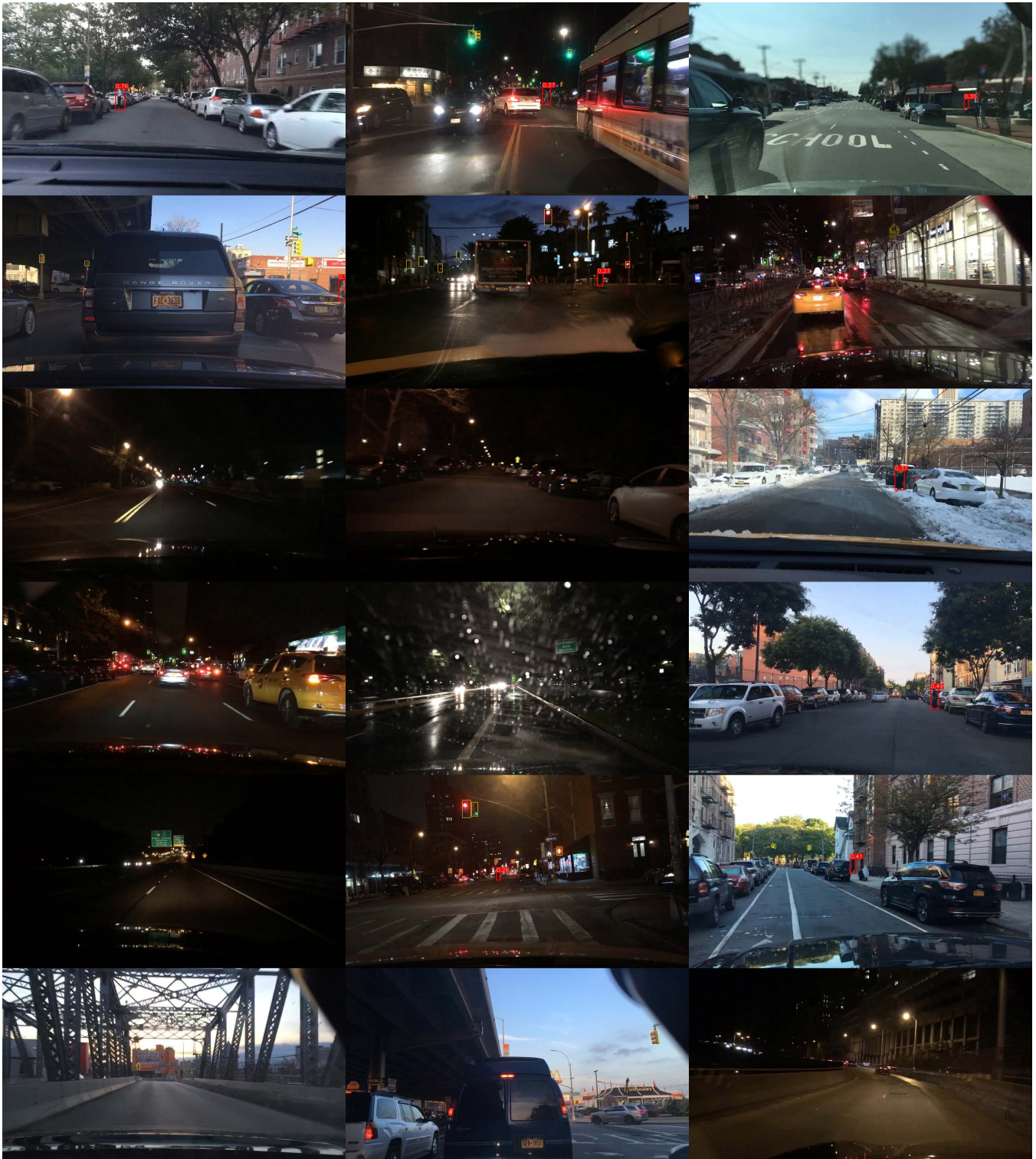


Figure 20: More examples of the category “Pedestrian” of BDD100K sampling from high loss images.

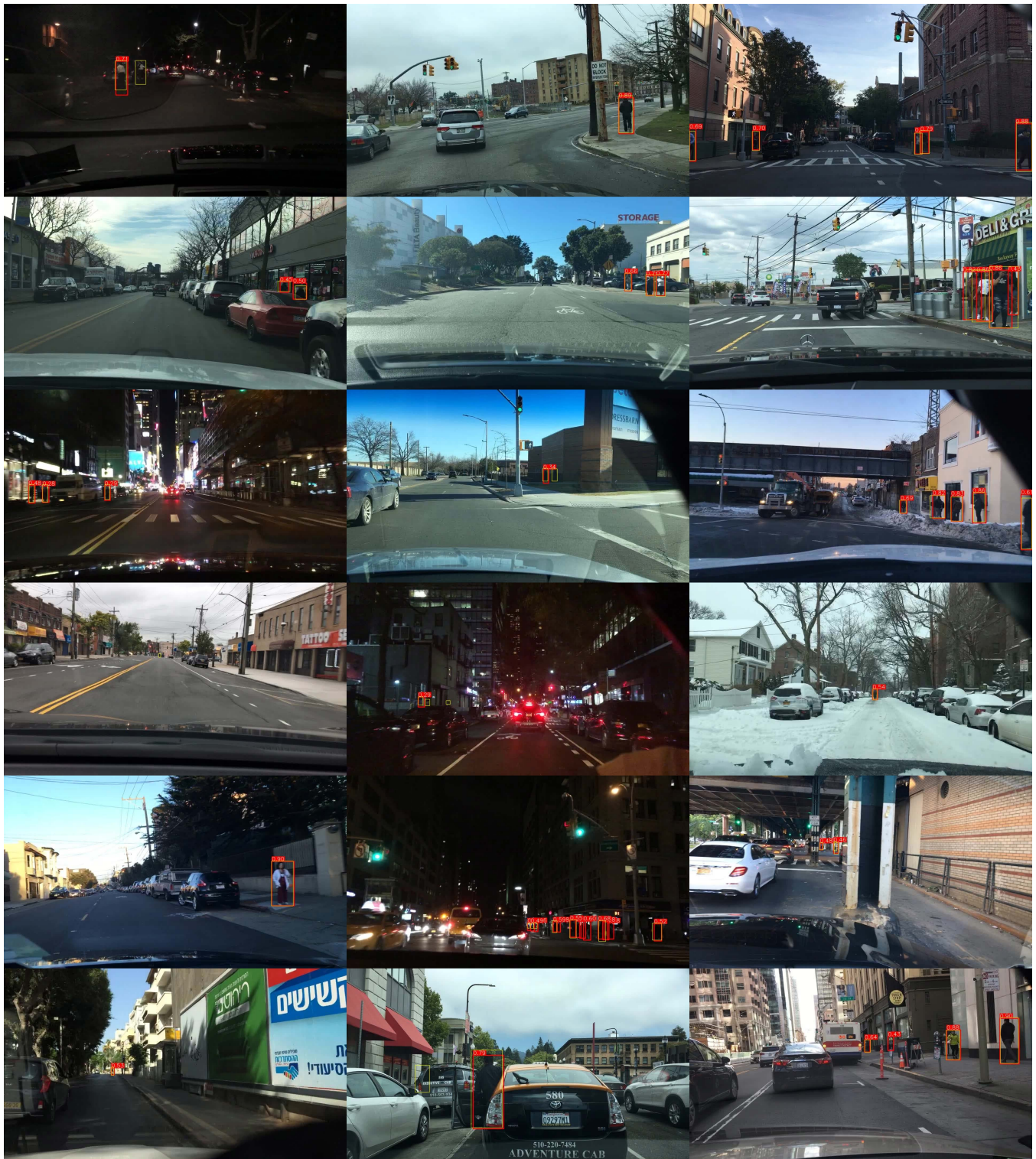


Figure 21: More examples of the category “Pedestrian” of BDD100K sampling from the whole population.

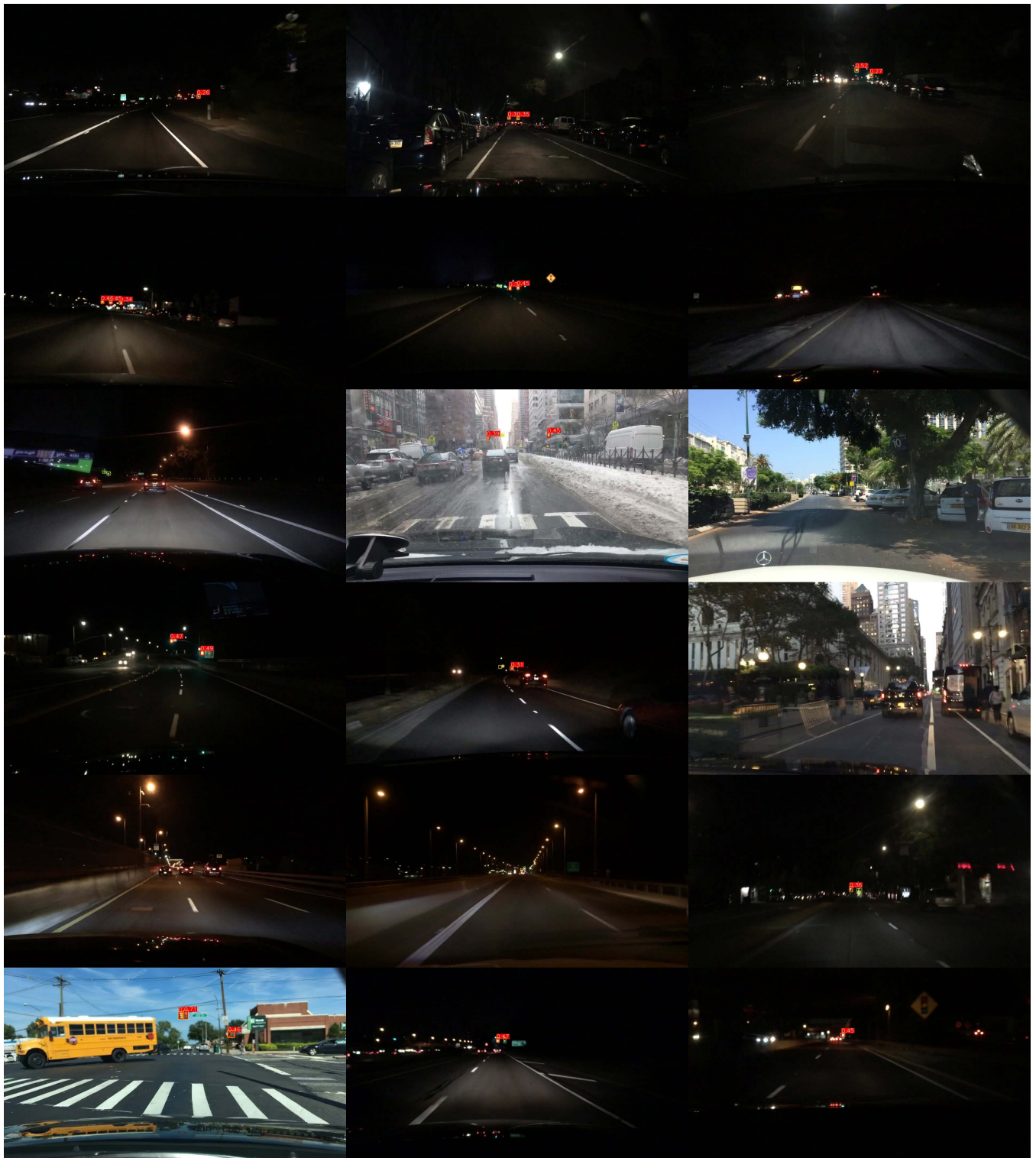


Figure 22: More examples of the category “Traffic Light” of BDD100K via Spotlight.

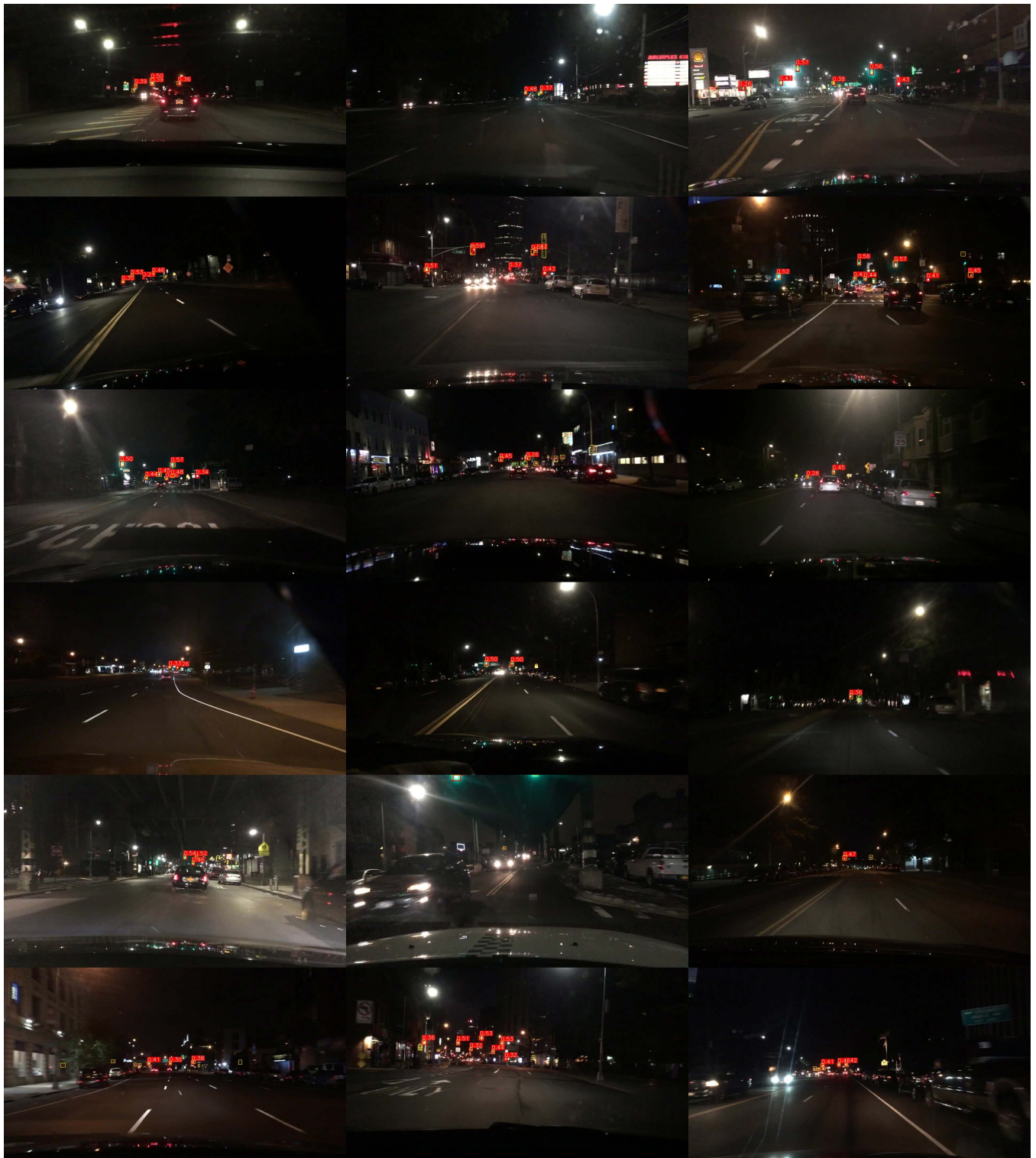


Figure 23: More examples of the category “Traffic Light” of BDD100K via MCSD.

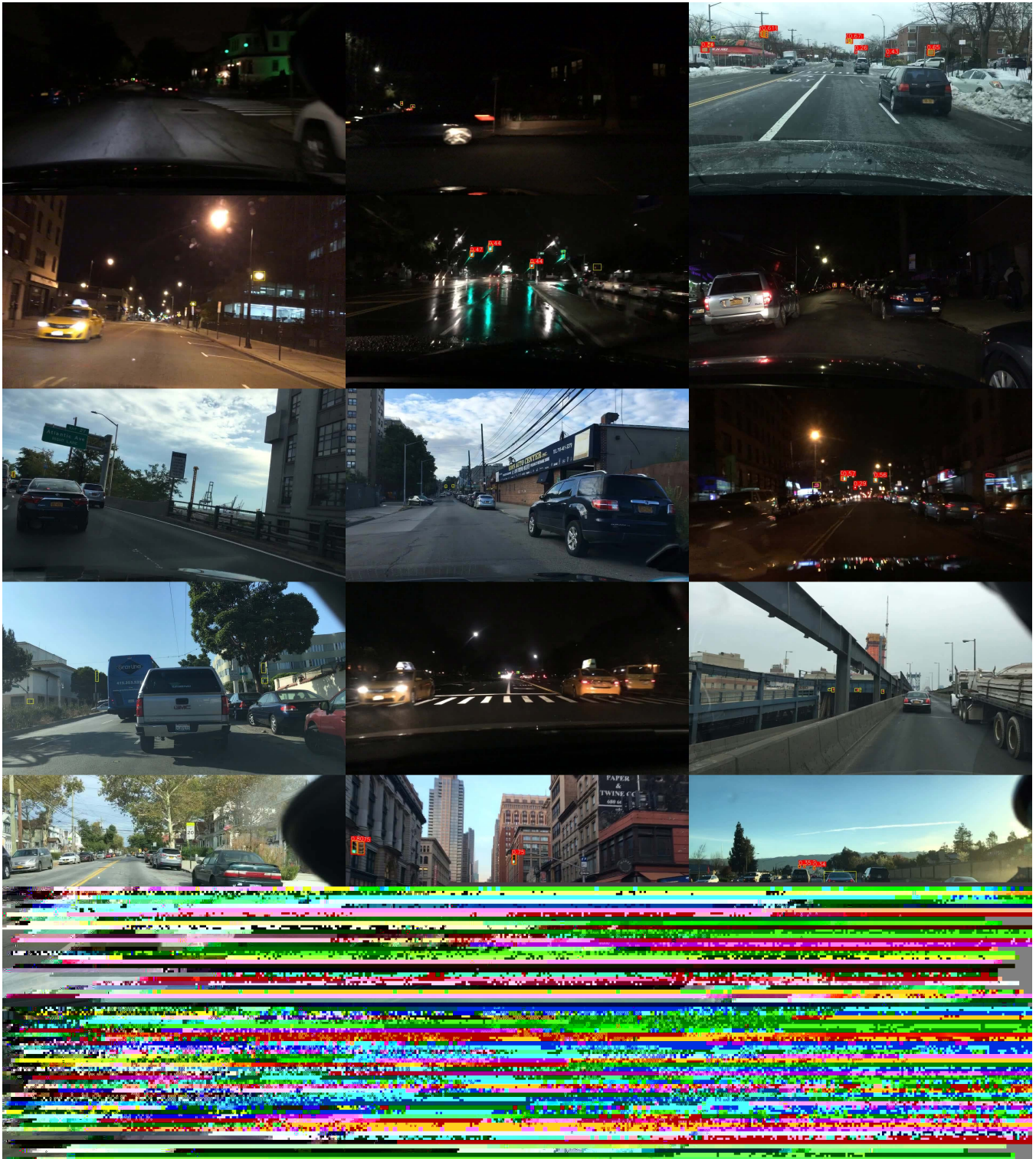


Figure 24: More examples of the category “Traffic Light” of BDD100K sampling from high loss images.

Figure 25: More examples of the category “Traffic Light” of BDD100K sampling from the whole population.

B Related Work

Subpopulation Shift It is widely acknowledged that models tend to make systematic mistakes on some subpopulations, which leads to the problem of subpopulation shift (Yang et al. 2023). To guarantee the worst subpopulation performance, some generate pseudo environment labels (Creager, Jacobsen, and Zemel 2021; Nam et al. 2021) and then apply existing invariant learning methods (Arjovsky et al. 2019; Krueger et al. 2021). Others take advantage of importance weighting to upweight the minority group or worst group like GroupDRO (Sagawa et al. 2019) and JTT (Liu et al. 2021a), or to combine with Mixup (Zhang et al. 2018) for more benefits (Han et al. 2022). A more recent work points out that current methods for subpopulation shifts heavily rely on the availability of group labels during model selection, and even simple data balancing techniques can achieve competitive performance (Idrissi et al. 2022). For better comparisons between algorithms and promotion of future algorithm development, Yang et al. (Yang et al. 2023) establish a comprehensive benchmark across various types of datasets.

OOD Generalization As a superset of subpopulation shift, Out-of-Distribution (OOD) generalization aims to address distribution shift more generally. There exists multiple branches of works. Domain generalization (Huang et al. 2020; Zhou et al. 2021; Zhang et al. 2023b,a) and invariant learning (Liu et al. 2021b, 2025, 2024; Creager, Jacobsen, and Zemel 2021) try to learn relationships that are generalizable to unseen distributions via utilizing multiple training environments. Distributionally robust optimization (DRO) (Duchi and Namkoong 2021; Volpi et al. 2018; Sinha, Namkoong, and Duchi 2018) aims to optimize the worst distribution centered around the training distribution. Stable learning (Kuang et al. 2020; Shen et al. 2020; Yu et al. 2023, 2025a) decorrelates covariates to mitigate the usage of spurious correlations.

Error Slice Discovery Instead of designing algorithms to improve generalization, error slice discovery has also attracted much attention recently. It is more flexible in that it can be followed by either non-algorithmic interventions like collecting more data for error slices, or algorithmic interventions like upweighting data belonging to error slices. There are mainly two paradigms for the process of error slice discovery. The first paradigm, also the more traditional practice, separates error slice discovery and later interpretation via case analyses or with the help of multi-modal models. Spotlight (d’Eon et al. 2022) attempts to learn a centroid in the representation space and employ the distance to this centroid as the error degree. InfEmbed (Wang et al. 2023) employs the influence of training samples on each test sample as embeddings used for clustering. PlaneSpot (Plumb et al. 2023) concatenates model prediction probability with dimension-reduced representation for clustering. All of them interpret the identified slices via case analyses directly. Meanwhile, the error-aware Gaussian mixture algorithm Domino (Eyuboglu et al. 2022) is followed by finding the best match between candidate text descriptions and the discovered slice in the representation space of multi-modal models like CLIP (Radford et al. 2021). This paradigm has a relatively high requirement for the coherence of identified error slices so that they can be interpreted. The second paradigm incorporates the discovery and interpretation of error slices together. HiBug (Chen, Li, and Xu 2023) and PRIME (Rezaei et al. 2024) divide the whole population of data into subgroups through proposing appropriate attributes and conducting zero-shot classification for these attributes using pretrained multi-modal models, and then directly calculate average performance for subgroups to identify the risky ones. The obtained subgroups are naturally interpretable via the combination of the attribute pseudo labels. Two recent works (Wiles, Albuquerque, and Goyal 2022; Gao et al. 2023) generate data from diffusion models (Rombach et al. 2022) before identifying error slices, avoiding the need of an extra validation dataset for slice discovery. It is obvious that the second paradigm heavily relies on the quality of proposed attributes and the capability of pretrained multi-modal models. Besides, it is worth noting that error slice discovery is closely related to the evaluation of OOD generalization (Wu et al. 2024; Yu et al. 2024a; Blanchet et al. 2024; Yu et al. 2024b, 2025b) since it is equivalent to evaluating models under subpopulation shift.

Error Prediction Another branch of works sharing a similar goal with error slice discovery is error prediction (or performance prediction). Although they are also able to find slices with high error, they focus on predicting the overall error rate given an unlabeled test dataset, and measure the effectiveness of error prediction methods via the gap between the predicted performance and the ground-truth one. Moreover, they do not emphasize the coherence and interpretability of error slices. Currently, there are several ways for error prediction. Some employ model output properties on the given test data like model confidence (Garg et al. 2021; Guillory et al. 2021), neighborhood smoothness (Ng et al. 2022), prediction dispersity (Deng et al. 2023), invariance under transformations (Deng, Gould, and Zheng 2021), etc. Inspired by domain adaptation (Long et al. 2015; Ben-David et al. 2010), some make use of distribution discrepancy between training data and unlabeled test data (Deng and Zheng 2021; Yu et al. 2022; Lu et al. 2023). Others utilize model disagreement between two models identically trained except random initialization and batch order during training (Jiang et al. 2021; Baek et al. 2022; Chen et al. 2021; Kirsch and Gal 2022), which exhibits SOTA performance in the error prediction task (Trivedi, Koutra, and Thiagarajan 2023).

References

- Achiam, J.; Adler, S.; Agarwal, S.; Ahmad, L.; Akkaya, I.; Aleman, F. L.; Almeida, D.; Altenschmidt, J.; Altman, S.; Anadkat, S.; et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Arjovsky, M.; Bottou, L.; Gulrajani, I.; and Lopez-Paz, D. 2019. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*.
- Baek, C.; Jiang, Y.; Raghunathan, A.; and Kolter, J. Z. 2022. Agreement-on-the-line: Predicting the performance of neural networks under distribution shift. *Advances in Neural Information Processing Systems*, 35: 19274–19289.

Belkin, M.; and Niyogi, P. 2003. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural computation*, 15(6): 1373–1396.

Ben-David, S.; Blitzer, J.; Crammer, K.; Kulesza, A.; Pereira, F.; and Vaughan, J. W. 2010. A theory of learning from different domains. *Machine learning*, 79(1-2): 151–175.

Blanchet, J.; Cui, P.; Li, J.; and Liu, J. 2024. Stability evaluation through distributional perturbation analysis. In *Forty-first International Conference on Machine Learning*.

Borkan, D.; Dixon, L.; Sorensen, J.; Thain, N.; and Vasserman, L. 2019. Nuanced metrics for measuring unintended bias with real data for text classification. In *Companion proceedings of the 2019 world wide web conference*, 491–500.

Chen, J.; Liu, F.; Avci, B.; Wu, X.; Liang, Y.; and Jha, S. 2021. Detecting errors and estimating accuracy on unlabeled data with self-training ensembles. *Advances in Neural Information Processing Systems*, 34: 14980–14992.

Chen, L.; Wu, P.; Chitta, K.; Jaeger, B.; Geiger, A.; and Li, H. 2024. End-to-end autonomous driving: Challenges and frontiers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.

Chen, M.; Li, Y.; and Xu, Q. 2023. HiBug: On Human-Interpretable Model Debug. *Advances in Neural Information Processing Systems*, 36.

Creager, E.; Jacobsen, J.-H.; and Zemel, R. 2021. Environment inference for invariant learning. In *International Conference on Machine Learning*, 2189–2200. PMLR.

Dann, E.; Henderson, N. C.; Teichmann, S. A.; Morgan, M. D.; and Marioni, J. C. 2022. Differential abundance testing on single-cell data using k-nearest neighbor graphs. *Nature Biotechnology*, 40(2): 245–253.

Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, 248–255. Ieee.

Deng, W.; Gould, S.; and Zheng, L. 2021. What does rotation prediction tell us about classifier accuracy under varying testing environments? In *International Conference on Machine Learning*, 2579–2589. PMLR.

Deng, W.; Suh, Y.; Gould, S.; and Zheng, L. 2023. Confidence and Dispersity Speak: Characterising Prediction Matrix for Unsupervised Accuracy Estimation. *arXiv preprint arXiv:2302.01094*.

Deng, W.; and Zheng, L. 2021. Are labels always necessary for classifier accuracy evaluation? In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 15069–15078.

d’Eon, G.; d’Eon, J.; Wright, J. R.; and Leyton-Brown, K. 2022. The spotlight: A general method for discovering systematic errors in deep learning models. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, 1962–1981.

Duchi, J. C.; and Namkoong, H. 2021. Learning models with uniform performance via distributionally robust optimization. *The Annals of Statistics*, 49(3): 1378–1406.

Eyuboglu, S.; Varma, M.; Saab, K. K.; Delbrouck, J.-B.; Lee-Messer, C.; Dunnmon, J.; Zou, J.; and Re, C. 2022. Domino: Discovering Systematic Errors with Cross-Modal Embeddings. In *International Conference on Learning Representations*.

Gao, I.; Ilharco, G.; Lundberg, S.; and Ribeiro, M. T. 2023. Adaptive testing of computer vision models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 4003–4014.

Garg, S.; Balakrishnan, S.; Lipton, Z. C.; Neyshabur, B.; and Sedghi, H. 2021. Leveraging unlabeled data to predict out-of-distribution performance. In *International Conference on Learning Representations*.

Guillory, D.; Shankar, V.; Ebrahimi, S.; Darrell, T.; and Schmidt, L. 2021. Predicting with confidence on unseen distributions. In *Proceedings of the IEEE/CVF international conference on computer vision*, 1134–1144.

Gurobi Optimization, L. 2021. Gurobi optimizer reference manual.

Han, Z.; Liang, Z.; Yang, F.; Liu, L.; Li, L.; Bian, Y.; Zhao, P.; Wu, B.; Zhang, C.; and Yao, J. 2022. Umix: Improving importance weighting for subpopulation shift via uncertainty-aware mixup. *Advances in Neural Information Processing Systems*, 35: 37704–37718.

He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.

Huang, Z.; Wang, H.; Xing, E. P.; and Huang, D. 2020. Self-challenging improves cross-domain generalization. In *European conference on computer vision*, 124–140. Springer.

Idrissi, B. Y.; Arjovsky, M.; Pezeshki, M.; and Lopez-Paz, D. 2022. Simple data balancing achieves competitive worst-group-accuracy. In *Conference on Causal Learning and Reasoning*, 336–351. PMLR.

Irvin, J.; Rajpurkar, P.; Ko, M.; Yu, Y.; Ciurea-Ilcus, S.; Chute, C.; Marklund, H.; Haghighi, B.; Ball, R.; Shpanskaya, K.; et al. 2019. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, 590–597.

Jain, S.; Lawrence, H.; Moitra, A.; and Madry, A. 2023. Distilling Model Failures as Directions in Latent Space. In *The Eleventh International Conference on Learning Representations*.

Jiang, Y.; Nagarajan, V.; Baek, C.; and Kolter, J. Z. 2021. Assessing Generalization of SGD via Disagreement. In *International Conference on Learning Representations*.

Kirsch, A.; and Gal, Y. 2022. A Note on "Assessing Generalization of SGD via Disagreement". *Transactions on Machine Learning Research*.

Koh, P. W.; Sagawa, S.; Marklund, H.; Xie, S. M.; Zhang, M.; Balsubramani, A.; Hu, W.; Yasunaga, M.; Phillips, R. L.; Gao, I.; et al. 2021. Wilds: A benchmark of in-the-wild distribution shifts. In *International conference on machine learning*, 5637–5664. PMLR.

Krueger, D.; Caballero, E.; Jacobsen, J.-H.; Zhang, A.; Binas, J.; Zhang, D.; Le Priol, R.; and Courville, A. 2021. Out-of-distribution generalization via risk extrapolation (rex). In *International Conference on Machine Learning*, 5815–5826. PMLR.

Kuang, K.; Xiong, R.; Cui, P.; Athey, S.; and Li, B. 2020. Stable prediction with model misspecification and agnostic distribution shift. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, 4485–4492.

Liu, E. Z.; Haghgoo, B.; Chen, A. S.; Raghunathan, A.; Koh, P. W.; Sagawa, S.; Liang, P.; and Finn, C. 2021a. Just train twice: Improving group robustness without training group information. In *International Conference on Machine Learning*, 6781–6792. PMLR.

Liu, H.; Li, C.; Wu, Q.; and Lee, Y. J. 2023a. Visual instruction tuning. *Advances in neural information processing systems*, 36: 34892–34916.

Liu, J.; Hu, Z.; Cui, P.; Li, B.; and Shen, Z. 2021b. Heterogeneous risk minimization. In *International Conference on Machine Learning*, 6804–6814. PMLR.

Liu, J.; Wang, T.; Cui, P.; and Namkoong, H. 2023b. On the need for a language describing distribution shifts: Illustrations on tabular datasets. *Advances in Neural Information Processing Systems*, 36.

Liu, S.; Fan, C.; Cheng, K.; Wang, Y.; Cui, P.; Sun, Y.; and Liu, Z. 2024. Inductive meta-path learning for schema-complex heterogeneous information networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.

Liu, S.; He, Y.; Wang, H.; Yang, W.; Wang, Y.; Cui, P.; and Liu, Z. 2025. Environment Inference for Learning Generalizable Dynamical System. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.

Liu, Z.; Luo, P.; Wang, X.; and Tang, X. 2015. Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, 3730–3738.

Long, M.; Cao, Y.; Wang, J.; and Jordan, M. 2015. Learning transferable features with deep adaptation networks. In *International conference on machine learning*, 97–105. PMLR.

Lu, Y.; Qin, Y.; Zhai, R.; Shen, A.; Chen, K.; Wang, Z.; Kolouri, S.; Stepputtis, S.; Campbell, J.; and Sycara, K. 2023. Characterizing Out-of-Distribution Error via Optimal Transport. *arXiv preprint arXiv:2305.15640*.

Melas-Kyriazi, L. 2020. The mathematical foundations of manifold learning. *arXiv preprint arXiv:2011.01307*.

Nam, J.; Kim, J.; Lee, J.; and Shin, J. 2021. Spread Spurious Attribute: Improving Worst-group Accuracy with Spurious Attribute Estimation. In *International Conference on Learning Representations*.

Ng, N.; Cho, K.; Hulkund, N.; and Ghassemi, M. 2022. Predicting out-of-domain generalization with local manifold smoothness. *arXiv preprint arXiv:2207.02093*.

Pedronette, D. C. G.; Gonçalves, F. M. F.; and Guilherme, I. R. 2018. Unsupervised manifold learning through reciprocal kNN graph and Connected Components for image retrieval tasks. *Pattern Recognition*, 75: 161–174.

Plumb, G.; Johnson, N.; Cabrera, A.; and Talwalkar, A. 2023. Towards a More Rigorous Science of Blindspot Discovery in Image Classification Models. *Transactions on Machine Learning Research*.

Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763. PMLR.

Rezaei, K.; Saberi, M.; Moayeri, M.; and Feizi, S. 2024. PRIME: Prioritizing Interpretability in Failure Mode Extraction. In *The Twelfth International Conference on Learning Representations*.

Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10684–10695.

Roweis, S. T.; and Saul, L. K. 2000. Nonlinear dimensionality reduction by locally linear embedding. *science*, 290(5500): 2323–2326.

Sagawa, S.; Koh, P. W.; Hashimoto, T. B.; and Liang, P. 2019. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv preprint arXiv:1911.08731*.

Shen, Z.; Cui, P.; Zhang, T.; and Kunag, K. 2020. Stable learning via sample reweighting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, 5692–5699.

Sinha, A.; Namkoong, H.; and Duchi, J. 2018. Certifying Some Distributional Robustness with Principled Adversarial Training. In *International Conference on Learning Representations*.

Tenenbaum, J. B.; Silva, V. d.; and Langford, J. C. 2000. A global geometric framework for nonlinear dimensionality reduction. *science*, 290(5500): 2319–2323.

Tian, Y.; Ye, Q.; and Doermann, D. 2025. Yolov12: Attention-centric real-time object detectors. *arXiv preprint arXiv:2502.12524*.

Trivedi, P.; Koutra, D.; and Thiagarajan, J. J. 2023. A Closer Look At Scoring Functions And Generalization Prediction. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 1–5. IEEE.

Volpi, R.; Namkoong, H.; Sener, O.; Duchi, J. C.; Murino, V.; and Savarese, S. 2018. Generalizing to unseen domains via adversarial data augmentation. *Advances in neural information processing systems*, 31.

Wang, C.-Y.; Bochkovskiy, A.; and Liao, H.-Y. M. 2023. YOLOv7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 7464–7475.

Wang, F.; Adebayo, J.; Tan, S.; Garcia-Olano, D.; and Kokhlikyan, N. 2023. Error Discovery by Clustering Influence Embeddings. *Advances in Neural Information Processing Systems*, 36.

Wiles, O.; Albuquerque, I.; and Goyal, S. 2022. Discovering bugs in vision models using off-the-shelf image generation and captioning. *arXiv preprint arXiv:2208.08831*.

Wu, J.; Liu, J.; Cui, P.; and Wu, S. Z. 2024. Bridging multicalibration and out-of-distribution generalization beyond covariate shift. *Advances in Neural Information Processing Systems*, 37: 73036–73078.

Yang, Y.; Zhang, H.; Gichoya, J. W.; Katabi, D.; and Ghassemi, M. 2024. The limits of fair medical imaging AI in real-world generalization. *Nature Medicine*, 30(10): 2838–2848.

Yang, Y.; Zhang, H.; Katabi, D.; and Ghassemi, M. 2023. Change is Hard: A Closer Look at Subpopulation Shift. In *International Conference on Machine Learning*, 39584–39622. PMLR.

Yu, F.; Chen, H.; Wang, X.; Xian, W.; Chen, Y.; Liu, F.; Madhavan, V.; and Darrell, T. 2020. Bdd100k: A diverse driving dataset for heterogeneous multitask learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2636–2645.

Yu, H.; Cui, P.; He, Y.; Shen, Z.; Lin, Y.; Xu, R.; and Zhang, X. 2023. Stable Learning via Sparse Variable Independence. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, 10998–11006.

Yu, H.; He, Y.; Xu, R.; Li, D.; Zhang, J.; Zou, W.; and Cui, P. 2025a. Sample Weight Averaging for Stable Prediction. *arXiv preprint arXiv:2502.07414*.

Yu, H.; Li, K.; Li, D.; He, Y.; Zhang, X.; and Cui, P. 2025b. ODP-Bench: Benchmarking Out-of-Distribution Performance Prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 1846–1858.

Yu, H.; Liu, J.; Zhang, X.; Wu, J.; and Cui, P. 2024a. A survey on evaluation of out-of-distribution generalization. *arXiv preprint arXiv:2403.01874*.

Yu, H.; Zhang, X.; Xu, R.; Liu, J.; He, Y.; and Cui, P. 2024b. Rethinking the evaluation protocol of domain generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 21897–21908.

Yu, Y.; Yang, Z.; Wei, A.; Ma, Y.; and Steinhardt, J. 2022. Predicting out-of-distribution error with the projection norm. In *International Conference on Machine Learning*, 25721–25746. PMLR.

Zemel, R.; and Carreira-Perpiñán, M. 2004. Proximity graphs for clustering and manifold learning. *Advances in neural information processing systems*, 17.

Zhang, H.; Cisse, M.; Dauphin, Y. N.; and Lopez-Paz, D. 2018. mixup: Beyond Empirical Risk Minimization. In *International Conference on Learning Representations*.

Zhang, X.; Xu, R.; Yu, H.; Dong, Y.; Tian, P.; and Cui, P. 2023a. Flatness-aware minimization for domain generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 5189–5202.

Zhang, X.; Xu, R.; Yu, H.; Zou, H.; and Cui, P. 2023b. Gradient norm aware minimization seeks first-order flatness and improves generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 20247–20257.

Zhou, K.; Yang, Y.; Qiao, Y.; and Xiang, T. 2021. Domain generalization with mixstyle. *arXiv preprint arXiv:2104.02008*.