

# Truncate-Split-Contrast: A Framework for Learning from Mislabeled Videos

Zixiao Wang<sup>1\*</sup>, Junwu Weng<sup>2</sup>, Chun Yuan<sup>1</sup>, Jue Wang<sup>2</sup>

<sup>1</sup> Tsinghua University

<sup>2</sup> Tencent AI Lab

zx-wang19@mails.tsinghua.edu.cn, WE0001WU@e.ntu.edu.sg

yuanc@sz.tsinghua.edu.cn, arphid@gmail.com

## Abstract

Learning with noisy label (LNL) is a classic problem that has been extensively studied for image tasks, but much less for video in the literature. A straightforward migration from images to videos without considering the properties of videos, such as computational cost and redundant information, is not a sound choice. In this paper, we propose two new strategies for video analysis with noisy labels: 1) A lightweight channel selection method dubbed as Channel Truncation for feature-based label noise detection. This method selects the most discriminative channels to split clean and noisy instances in each category; 2) A novel contrastive strategy dubbed as Noise Contrastive Learning, which constructs the relationship between clean and noisy instances to regularize model training. Experiments on three well-known benchmark datasets for video classification show that our proposed truNcatE-split-contrAsT (NEAT) significantly outperforms the existing baselines. By reducing the dimension to 10% of it, our method achieves over 0.4 noise detection F1-score and 5% classification accuracy improvement on Mini-Kinetics dataset under severe noise (symmetric-80%). Thanks to Noise Contrastive Learning, the average classification accuracy improvement on Mini-Kinetics and Sth-Sth-V1 is over 1.6%.

## 1 Introduction

Training deep networks requires large-scale datasets with high-quality human annotations. However, acquiring large-scale clean-annotated data is costly and time-consuming, forcing people to seek low-cost but imprecise alternative labeling. Such labeling inevitably introduces noises: a large number of instances could be annotated with incorrect labels. Recent studies (Zhang et al. 2017; Arpit et al. 2017) have shown that deep neural networks have a high capacity to fit data even under randomly assigned labels, which harms the generalization on unseen data. Therefore, how to train a robust deep learning model in the presence of noisy labels is challenging and is of increasing significance in the industry. To date, the existing LNL approaches mainly focus on image tasks. With the rapidly growing amount of video data on the Internet, designing a noise-robust training strategy for video models becomes imperative. Motivated by the previous success of LNL methods on images, we study the

much less explored problem of applying LNL methods in the video domain.

Depending on whether noisy instances are detected in training, the existing LNL methods can be roughly divided into two types. One is to directly train a noise-robust model in the presence of noisy labels (Patrini et al. 2017; Wang et al. 2019; Ma et al. 2020; Lyu and Tsang 2019; Zhou et al. 2021; Gao, Gouk, and Hospedales 2021). The other one is to explicitly detect the potential noisy instances, and then learn a model by simply excluding them (Huang et al. 2019), or re-using the potential noisy data by estimating the pseudo labels of them (Zhang et al. 2018b; Li, Socher, and Hoi 2019; Li, Xiong, and Hoi 2021; Ortego et al. 2021). This detection strategy is widely adopted in the industry as it not only learns a robust model, but provides a clean dataset as well. Following this strategy, we learn video representations from potentially-mislabeled data with two steps, Noise Detection and Unlabelled Data Utilization.

**Noise Detection.** The *loss-based* method (Han et al. 2018; Zhang et al. 2018b; Huang et al. 2019; Li, Socher, and Hoi 2019) is the common solution for noisy label detection. These methods treat instances with smaller classification losses as clean ones during training, which easily leads to confirmation bias (Li, Socher, and Hoi 2019). Compared with a single loss value, video latent representations naturally contain multi-channel signals, which provide adequate clues for noise detection.

The existing *feature-based* noisy label detection methods (Lee et al. 2018; Han, Luo, and Wang 2019) commonly conduct binary clustering (clean/noisy) on the full-dimensional features before the classification layer. However, we argue that to detect clean/noisy instances, utilizing all channels of a feature learned from classification supervision is not a must. The reason is that the instance feature learned with label supervision may perform well on differentiating categories, but the extra feature dimensions for delicate boundary shaping are not suitable for an unsupervised binary clustering task in label noise detection, especially under the video domain. We find that the extra dimensions are redundant and weaken the performance in noise detection and therefore in final classification. (Sec. 4.4 - Table. 3) Designing a compact network may help reduce the redundancy of feature channels for noise detection, but the limited network capacity will inevitably hurt the classification,

\*This work is done during an internship at Tencent AI Lab.  
Copyright © 2023, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

as these two tasks share the same network and learned features. Channel selection therefore becomes our first choice.

In this paper, we propose a *category-wise* channel selection method, *i.e.* Channel Truncation (CT), for feature-based label noise detection in videos. It evaluates the discriminative ability of each channel of instance representations by simply collecting the temporal statistics across frames of each instance. After sorting all channels based on their noise-discriminative abilities, CT adaptively removes the most confusing ones for each category during training. Afterward, each truncated instance feature is matched with its relevant category clean-prototype to determine whether it is clean or not. This CT for videos can also be simplified to detect noisy label of image data (Appendix).

**Unlabelled Data Utilization.** The previous LNL methods commonly process the detected noise under a semi-supervised learning framework. A pseudo label is assigned to each noisy instance to replace the wrongly annotated one as a supervision signal for model training (Li, Socher, and Hoi 2019; Li, Xiong, and Hoi 2021; Ortego et al. 2021). One drawback of pseudo labeling is that the pseudo label could be ambiguous and unreliable without sophisticated post-processing and data enrichment (Zhang et al. 2018a). This phenomenon is even severe under video domain as the redundant information in videos may mislead the pseudo labeling at the early training stage. In the experiment, we verify that the naive pseudo label has negative impact to the video classification (Sec. 3.2 - Table. 1), while the impact on image is not that severe. Therefore, a new design of unlabelled data utilization in videos is needed. On the other side, little effort has been directed towards enhancing the quality of instance representation when no label is assigned to the detected noise. Compared with guessing the actual label of the noise, the mutual relationship among the instances is relatively easy to estimate after the clean/noise splitting. Inspired by contrastive learning (Oord, Li, and Vinyals 2018; Jing and Tian 2020), we propose a *Noise Contrastive Loss* (NCL) to utilize the unlabelled noise and enlarge the margins among instances from different categories. NCL provides a low-risk contrastive strategy for unlabelled noisy queries, avoiding misguidance from wrong pseudo labels.

The framework of *truNcatE-split-contrAsT* (NEAT) is shown in Fig. 1. Each video feature is first *truncated* for clean/noisy instance *splitting*. Then, the detected clean/noisy instances are utilized separately under the supervision of cross entropy and noise *contrastive* loss for model updating. Our main contributions are summarized as follows:

- A lightweight channel selection method for feature-based label noise detection is proposed. It discards the redundant channels to increase the effectiveness and efficiency of noisy/clean instance splitting.
- Noise Contrastive Loss is designed to construct the relationship among instances by referring the estimated clean/noisy splits, and utilizes this relationship to learn visual representations without involving wrong labels.
- To the best of our knowledge, this is the first efficient framework for LNL in video analysis. Extensive experiments show the effectiveness of our method on several video recognition datasets with noisy label settings.

## 2 Related Work

### 2.1 Learning with Noisy Label

LNL on images has been extensively studied in the literature. In this section, we only limit our review to the noise detection methods and the way these methods utilize the detected noisy instances.

**Noise Detection.** There are two main streams of noise detection methods, the loss-based one and the feature-based one. We here mainly discuss feature-based methods. Feature-based methods (Kim et al. 2021) detect noisy samples based on feature similarities between samples. (Lee et al. 2018) takes the cosine similarity between unidentified samples and the prototypes as references to detect noise. (Wu et al. 2020) utilizes  $k$ -nearest neighbor ( $k$ -NN) to build a neighbor graph for each category, treating samples in the dominant sub-component as clean ones. (Ortego et al. 2021) also uses  $k$ -NN and voting to determine whether a sample is clean or not. These works mainly focus on dealing with label noise in image tasks, and so far, no work discusses how things are different in video scenarios. Besides, they refer to the full-dimensional feature for noise detection, which we argue is unnecessary in this task.

**Noise Utilization.** The common strategy for noise utilization is to assign a pseudo label to the detected noisy instance as a supervision label. DivideMix (Li, Socher, and Hoi 2019) estimates the pseudo label from model prediction and applies Mixup (Zhang et al. 2018a) to enhance the reliability of the pseudo label. (Ortego et al. 2021) combines MixUp (Zhang et al. 2018a) and Supervised Contrastive Learning (Khosla et al. 2020) to mitigate the negative impact of noisy labels in representation learning. The label correction in this work is achieved by  $k$ -NN search. Junnan Li *et al.* improve prototypical contrastive learning (PCL) (Li et al. 2020) in (Li, Xiong, and Hoi 2021) by using pseudo-labels to compute class prototypes. This work applies a similar  $k$ -NN-based label correction strategy as (Ortego et al. 2021). The drawback of pseudo labeling is that the corrected label is sometimes not reliable, and label enhancement techniques like sharpening (Li, Socher, and Hoi 2019) are needed. Compared with the pseudo labeling, instead of estimating the unreliable pseudo labels or involving the noisy labels, our proposed NCL utilizes the mutual relationship between the distribution of clean and noisy instances to shape the classification boundary such that the negative impact of estimated labels is mitigated. No complex post-processing and data enrichment is required.

### 2.2 Contrastive Representation Learning

In recent years, the contrastive representation learning methods dominant the existing literature in self-supervised learning (Oord, Li, and Vinyals 2018; Jing and Tian 2020). They optimize similarities of positive (negative) pairs to improve the quality of representations. Techniques like data augmentation (Tian et al. 2020), larger batch size (Chen et al. 2020a; Khosla et al. 2020), network design (Chen et al. 2020b) are always used for better representation learning. Regarding the label noise scenario, the Instance Contrastive Learning (Oord, Li, and Vinyals 2018), the Supervised Contrastive

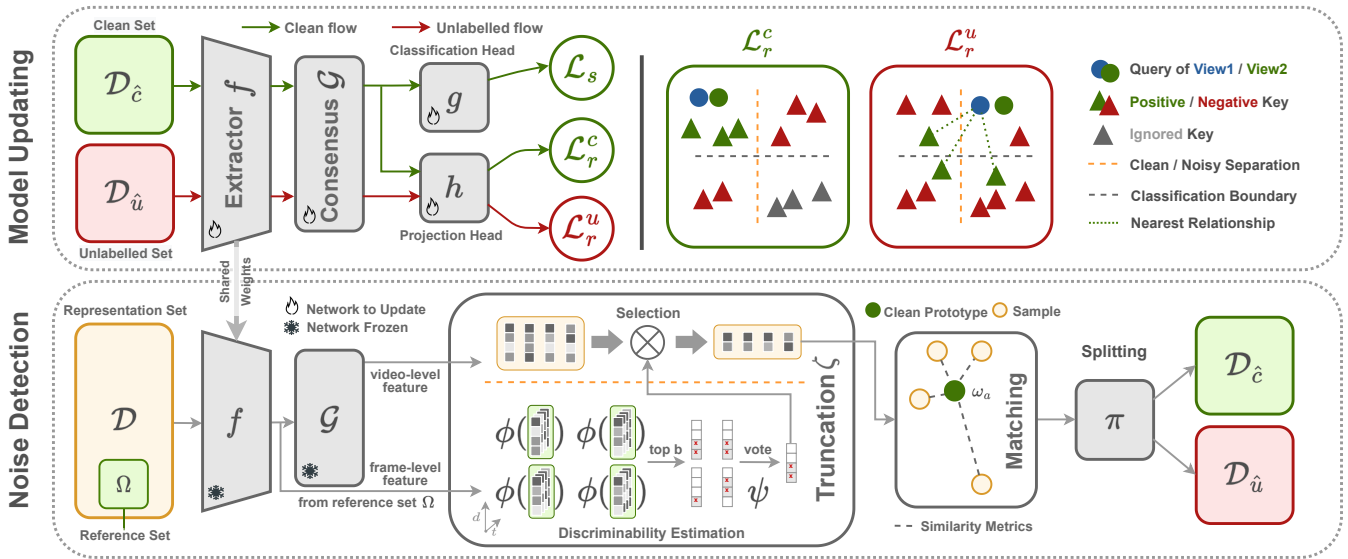


Figure 1: The pipeline of our framework NEAT on the noisy dataset in training. There are two phases in the framework. In the *Noise Detection* phase, the whole training dataset is split-ed into clean/noise clusters by considering the similarities between each pair of dimension-reduced instance representations. In this phase the network is frozen for feature extraction only. During *Model Updating* phase, the detected clean instances are fed into the network for supervised learning, while all instances are utilized in Noise Contrastive Learning for decision boundary shaping. These two phases proceed iteratively. In the first round of model updating, all the instances are utilized for supervised learning.

Learning (Khosla et al. 2020) and the Prototypical Contrastive Learning (PCL) (Li et al. 2020) are modified in (Kim et al. 2021), (Ortego et al. 2021) and (Li, Xiong, and Hoi 2021) correspondingly and respectively to fully involve the detected noisy instances.

### 3 Method

Under the video scenario,  $T$  frames are sampled from each video as a clip. The normalized feature of a frame is defined as  $\mathbf{v}$ , which is extracted from a backbone network  $f(\cdot)$ . Afterward, the representation of a clip from a video is defined as  $\mathbf{x} = \mathcal{G}(\mathbf{v}_1, \dots, \mathbf{v}_T)$ , where  $\mathcal{G}$  is a consensus function combining the set  $\{\mathbf{v}_i\}_{i=1}^T$  with normalization. Given a representation dataset  $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^M$  with  $M$  normalized  $\mathbf{x} \in \mathbb{R}^d$ , the goal of a classification task is to find which category  $\mathbf{x}$  belongs to. The annotation of  $\mathbf{x}$  is  $a$ , and it can be also represented as an one-hot vector  $\mathbf{y} \in \{0, 1\}^K$ , in which the  $a$ -th element of  $\mathbf{y}$  is assigned as 1, and the remainings are assigned as 0s. Here  $K$  indicates the number of category. A classification head  $g(\cdot)$  with Softmax operation is defined to predict the probability of  $\mathbf{x}$  belonging to the  $k$ -th category, namely  $p(k|\mathbf{x}) = g(\mathbf{x}, k)$ .

With the noisy label existing, a label  $a$  may be wrongly assigned to a sample not belonging to the  $a$ -th category set. To reduce the negative impact of wrongly annotated samples in model training, our strategy is to first estimate which samples are correctly annotated, *i.e.* clean samples from  $\mathcal{D}_c$ , and which are not, *i.e.* noisy samples from  $\mathcal{D}_u$ , where  $\mathcal{D} = \mathcal{D}_c \cup \mathcal{D}_u$ . The noisy samples are regarded as unlabelled in our framework. This noise detection phase is achieved by a detection function  $\mathcal{D}_\epsilon = \epsilon(\mathcal{D})$ . In the fol-

lowing model updating phase, the labels of the estimated clean samples are directly used in cross entropy loss for label supervision. Meanwhile, the detected noisy samples are involved in model updating as regularization where the impact of wrong labels is ignored. With the notations described above, the model updating process under the noisy label setting can be guided by the loss defined as below,

$$\begin{aligned} \mathcal{L} &= \mathcal{L}_s + \mathcal{L}_r \\ &= - \sum_{\mathbf{x}_i \in \epsilon(\mathcal{D})} \sum_{k=1}^K \mathbf{y}_i(k) \cdot \log g(\mathbf{x}_i, k) + \mathcal{L}_r, \end{aligned} \quad (1)$$

in which the loss  $\mathcal{L}_s$  is defined by cross entropy for label supervision. The detection function  $\epsilon(\cdot)$  filters out the noisy samples in the training set for each category by similarity measure, in which a dimensionality reduction method Channel Truncation, is proposed for the discriminative channel selection. The loss  $\mathcal{L}_r$  is designed as a regularization in model updating involving both estimated clean and noisy samples. We introduce Noise Contrastive Learning in  $\mathcal{L}_r$  to utilize the relationships among clean and noisy instances for decision boundary shaping. The pipeline of our framework is shown in Fig. 1.

#### 3.1 Channel Truncation

Generally, feature-based noise detection methods detect noisy labels by computing the similarity between the query  $\mathbf{x}$  and the clean-prototype  $\mathbf{x}_a$  of category  $a$ . The higher the similarity, the more likely the query  $\mathbf{x}$  to be clean. As full channels are learned to differentiate multiple categories,

they are unnecessary for a much simpler task, *i.e.* differentiating clean/noisy instances in each category. Truncating inessential channels can help detect clean instances better than utilizing all. Therefore, we propose Channel Truncation to truncate the redundant channels and keep the discriminative ones for clean instance detection.

Our method owns a top  $b$  channel selection operation  $\zeta_b(\cdot)$  with a category-level score function  $\psi(\cdot)$ .  $\psi(\cdot)$  evaluates the discriminative ability of each channel in a category, such as  $a$ , by referring to the statistics of its related reference set  $\Omega^a$ . This score function returns a vector with the same dimension as feature  $\mathbf{x}$ , where each element is a score of the corresponding feature channel. Selection operation  $\zeta_b(\cdot)$  picks  $b$  channels with the highest scores of the input representation and returns a dimension-reduced  $\mathbf{w}$  concatenating these top channels. Therefore the truncation function is defined as,

$$\mathbf{w} = \zeta_b(\mathbf{x}, \psi(\Omega^a)). \quad (2)$$

Given truncated-feature  $\mathbf{w}$  and its label  $a$ , the similarity measure between  $\mathbf{w}$  and the class clean-prototype  $\mathbf{w}_a$  is obtained by a similarity function  $\eta_a(\cdot)$ , which is defined as inner product, namely  $\eta_a(\mathbf{w}) = \mathbf{w} \cdot \mathbf{w}_a$ . This similarity indicates how close the instance  $\mathbf{w}$  to the clean cluster of category  $a$ . The class clean-prototype  $\mathbf{w}_a$  is defined as the average  $\mathbf{w}$  of the estimated clean set  $\mathcal{D}_{\hat{c}}$ . We observe that the similarity distribution of the clean and noisy instances gradually becomes a two-peak form during training. Thus to detect the noisy videos, a two-component Gaussian Mixture Model  $\pi(\cdot)$  is utilized to fit the distribution of  $\eta(\mathbf{w})$ . The probability of  $\mathbf{w}$  being noise of the  $a$ -th category is then defined as  $p(\text{noise}|\mathbf{w}; a) = \pi_a(\mathbf{w})$ . Hence, the estimated clean set is obtained by thresholding  $\pi_a(\mathbf{w})$ ,

$$\mathcal{D}_{\hat{c}} = \epsilon(\mathcal{D}) = \{\mathbf{x} \mid \mathbf{x} \in \mathcal{D}, \pi_a(\mathbf{w}) < \xi\}. \quad (3)$$

where  $\xi$  is a threshold and we fix it as 0.5 in the experiments. For simplicity below, we ignore the category index  $a$ . Next, how to design a proper score function  $\psi(\cdot)$ ?

**Oracle Selection** Ideally, when the true splits of correctly and wrongly annotated instances are known beforehand, we take all the clean training data from the category  $a$  as reference set  $\Omega$ , and the channel discriminative ability can be measured by the within-/ between-cluster variance. Following the Fisher Discriminant Analysis, for a specific category  $a$ , the oracle score function  $\psi_o(\cdot)$  is then defined as:

$$\psi_o = \frac{(\boldsymbol{\mu}_c - \boldsymbol{\mu}_u)^2}{\boldsymbol{\sigma}_c^2 + \boldsymbol{\sigma}_u^2} \quad (4)$$

$$\boldsymbol{\mu}_c = \frac{1}{|\Omega|} \sum_{\mathbf{x}_i \in \Omega} \mathbf{x}_i, \quad \boldsymbol{\sigma}_c^2 = \frac{1}{|\Omega|} \sum_{\mathbf{x}_i \in \Omega} (\mathbf{x}_i - \boldsymbol{\mu}_c)^2,$$

where  $\boldsymbol{\mu}_c \in \mathbb{R}^d$  and  $\boldsymbol{\sigma}_c \in \mathbb{R}^d$  are the average values and standard deviation of all  $\mathbf{x}_i$  in set  $\Omega$ , respectively. All the operations defined here in the equation are element-wise. Similarly,  $\boldsymbol{\mu}_u$  and  $\boldsymbol{\sigma}_u$  are the statistics of the unlabelled (mislabelled) set. The higher the score, the more discriminative the corresponding channels. As the true splits of correctly and wrongly annotated instances are unknown in training, the oracle score in Eq.4 *cannot* be used as the selection criterion

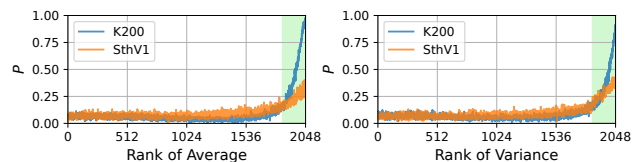


Figure 2: The statistical relation between oracle score and amplitudes/variance on K200 and SthV1 under symmetric-40% noise setting at fifth epoch. Each recorded point bears the coordinates  $(r, p)$ .  $r$  is the ranking of the corresponding statistics. The higher the  $r$ , the larger the amplitudes/variance.  $p$  is the probability of the relevant channels picked by the top  $b$  oracle selection. The top  $b$  amplitudes/variance area is filled with green. ( $b = 200$ )

of channels. However, this measurement can be treated as an evaluation metric of practicable  $\psi(\cdot)$ .

**Proposed Selection** A video, in essence, consists of both the scene and motion semantics (Wang and Hoai 2018; Choi et al. 2019; Weng et al. 2020; Wang et al. 2021). The information of global scene appearance generally remains almost unchanged in a video, while the motion-relevant information varies on the temporal domain. To well differentiate videos, it is critical to distill both signals. We therefore utilize the temporal *average* and *variance* to roughly search the semantics-intense channels. As expected, we experimentally discover that channels with larger temporal *average* and *variance* of amplitudes tend to achieve higher oracle scores, as shown in Fig. 2. Thus, we propose to estimate the channel discriminative ability by introducing an instance-level score function  $\phi(\cdot)$  and defining category-level score function  $\psi(\cdot)$  as a histogram counting the top  $b$  instance-level score events of each channel in a certain category. The function  $\phi(\cdot)$  is designed as temporal statistics of  $\{\mathbf{v}_t\}_{t=1}^T$  from  $G$ , which is lightweight and suitable for the video scenario. Here we have two versions of  $\phi(\cdot)$ , the average pooling  $\phi_{ave}$ , and the temporal variation  $\phi_{var}$  of the representation channel. They are defined as,

$$\phi_{ave}(\mathbf{v}_1, \dots, \mathbf{v}_T) = \sum_t \mathbf{v}_t / T, \quad (5)$$

$$\phi_{var}(\mathbf{v}_1, \dots, \mathbf{v}_T) = \sum_t (\mathbf{v}_t - \phi_{ave})^2 / T.$$

We set  $\Omega$  as the clean set  $\mathcal{D}_{\hat{c}}$  from the last epoch and initialize  $\Omega$  as the dataset  $\mathcal{D}$  in the first training epoch. The *average* responses extract the similar signal among  $T$  frames. The *variance* of the amplitude distills the dissimilar signal across the frames. Each operation here is element-wise. It is observed that the selected channels are closely related to the semantics, *e.g.* scene and motion, of a video clip in the specific category, about which we will show more visualizations in the Appendix. Moreover, a further detailed analysis of different temporal statistics will be discussed in Sec.4.

### 3.2 Noise Contrastive Learning

With the estimated clean and unlabelled splits, namely  $\mathcal{D}_{\hat{c}}$  and  $\mathcal{D}_{\hat{u}}$ , the motivation of Noise Contrastive Learning (NCL) is to fully utilize the estimated noisy samples

in model updating, and further enlarge the margins among samples from different categories.

To reach this goal, contrastive learning (Oord, Li, and Vinyals 2018) is involved in NCL to enforce the consistency within each clean cluster and enlarge the dissimilarity between clean and noisy clusters. We first randomly sample two clips from a video to represent two different views of this video. Both view clips consist of  $T$  frames following a certain sampling strategy. In NCL, we take one clip representation of a video as a query and set the remaining clips from the same and other videos as keys. The selection of positive and negative keys varies from clean query to noisy query. For a query representation  $\mathbf{x}_i \in \mathcal{D}_{\hat{c}}$  from *clean* cluster of category  $a_i$ , its sets of positive and negative keys  $\mathcal{P}_{\hat{c}}^i$ ,  $\mathcal{N}_{\hat{c}}^i$  can be defined as,

$$\begin{aligned} \mathcal{P}_{\hat{c}}^i &= \{\mathbf{x}_j | \mathbf{x}_j \in \mathcal{D}_{\hat{c}}, i \neq j, a_i = a_j\}, \\ \mathcal{N}_{\hat{c}}^i &= \{\mathbf{x}_j | \mathbf{x}_j \in \mathcal{D}_{\hat{u}}, a_i = a_j\} \cup \{\mathbf{x}_l | \mathbf{x}_l \in \mathcal{D}_{\hat{c}}, a_i \neq a_l\}, \end{aligned} \quad (6)$$

respectively. By this way, the estimated clean instances in the same category are forced to be close to each other, and the ones from different categories are pushed away from one another. When the query representation  $\mathbf{x}_i \in \mathcal{D}_{\hat{u}}$  is from *unlabelled* cluster, its sets of keys  $\mathcal{P}_{\hat{u}}^i$  and  $\mathcal{N}_{\hat{u}}^i$  are defined as,

$$\begin{aligned} \mathcal{P}_{\hat{u}}^i &= \{\mathbf{x}_j | \mathbf{x}_j \in \text{NN}(\mathbf{x}_i)\} \cup \{\tilde{\mathbf{x}}_i\}, \\ \mathcal{N}_{\hat{u}}^i &= \mathcal{D} - \mathcal{P}_{\hat{u}}^i \cup \{\mathbf{x}_i\}, \end{aligned} \quad (7)$$

in which  $\mathbf{x}_i$  and  $\tilde{\mathbf{x}}_i$  are the two views of a query video. The positive keys of the noisy query also come from the nearest neighbours. Specifically, we retrieve the top similar keys of a query  $\mathbf{x}_i$  by  $k$ -NN function  $\text{NN}(\mathbf{x})$  which returns the  $B$  nearest neighbors of input instance in the set  $\mathcal{D} - \{\mathbf{x}\}$ . With the proposed positive and negative sets  $\mathcal{P}_{\hat{c}}, \mathcal{P}_{\hat{u}}, \mathcal{N}_{\hat{c}}$  and  $\mathcal{N}_{\hat{u}}$ , we extend the InfoNCE (Oord, Li, and Vinyals 2018) loss to two noisy contrastive losses  $\mathcal{L}_r^c$  and  $\mathcal{L}_r^u$  to force the clean and noisy clusters to be apart from each other:

$$\mathcal{L}_r = \frac{1}{|\mathcal{D}_{\hat{c}}|} \sum_{\mathbf{x}_q \in \mathcal{D}_{\hat{c}}} \mathcal{L}_r^c(\mathbf{x}_q) + \frac{1}{|\mathcal{D}_{\hat{u}}|} \sum_{\mathbf{x}_q \in \mathcal{D}_{\hat{u}}} \mathcal{L}_r^u(\mathbf{x}_q) \quad (8)$$

where  $\mathcal{L}_r^c = \gamma(\mathcal{P}_{\hat{c}}, \mathcal{N}_{\hat{c}})$  and  $\mathcal{L}_r^u = \gamma(\mathcal{P}_{\hat{u}}, \mathcal{N}_{\hat{u}})$ .  $\gamma(\cdot)$  is defined as:

$$\gamma(\mathcal{P}, \mathcal{N}) = -\frac{1}{|\mathcal{P}|} \sum_{\mathbf{x}_+ \in \mathcal{P}} \log \frac{\exp(\mathbf{z}_q \cdot \mathbf{z}_+ / \tau)}{\sum_{\mathbf{x}_j \in \mathcal{P} \cup \mathcal{N}} \exp(\mathbf{z}_q \cdot \mathbf{z}_j / \tau)}, \quad (9)$$

in which  $\mathbf{z} \in \mathbb{R}^{\hat{d}}$  is the normalized representation mapped by projection head  $h(\cdot)$ , *i.e.*  $\mathbf{z} = h(\mathbf{x})$ , and  $\hat{d} < d$ . The symbol “ $\cdot$ ” denotes the inner (dot) product between two vectors.  $\tau$  is the temperature scaling.

Compared with assigning pseudo labels to the estimated unlabelled instances (Li, Socher, and Hoi 2019), the proposed NCL is low risk as it does not involve the label of the estimated noisy samples, which may severely mislead the model updating. Even though the true labels of the noisy samples are ignored in label supervision of the model updating, the *mutual exclusion* between the clean and noisy instances in each category can still well shape the decision

Noise Type Noise Ratio	Symmetric				Asymmetric		
	20%	40%	60%	80%	10%	20%	40%
GCE(2018)	53.1	49.6	42.1	23.4	54.0	52.0	41.4
SCE(2019)	64.5	57.8	48.1	27.9	67.2	62.0	46.9
TopoFilter(2020)	61.4	55.5	37.7	14.9	65.7	63.6	55.5
Co-teaching(2018)	61.0	60.9	56.9	32.5	60.6	60.2	46.9
M-correction(2019)	66.7	62.3	54.8	40.1	65.5	62.1	52.9
CT-all	68.4	64.2	58.3	43.3	69.1	67.4	51.9
CT-var	69.2	<b>67.1</b>	<b>61.1</b>	<b>48.4</b>	<b>70.0</b>	68.0	55.9
CT-ave	<b>69.4</b>	66.9	61.0	48.1	69.8	<b>68.6</b>	<b>58.4</b>
CT-oracle	70.4	67.7	61.5	49.6	70.6	70.0	58.9
Clean Only	70.5	68.3	64.9	58.8	70.9	70.3	68.6
DivideMix(2019)	69.4	65.9	60.7	46.5	68.2	67.5	53.1
CT-var + PL-H	68.1	65.6	58.3	37.9	66.5	65.8	54.6
CT-var + PL-S	68.2	66.0	60.8	45.4	67.1	66.8	55.8
CT-var + PL-K	67.7	65.0	60.1	47.4	66.4	65.7	55.0
CT-var + CL	69.4	66.6	60.4	12.2	68.3	67.6	54.7
CT-var + SCL	70.3	67.6	61.4	46.9	70.2	68.1	57.0
CT-var + NCL	<b>70.9</b>	<b>68.6</b>	<b>63.4</b>	<b>49.9</b>	<b>70.5</b>	<b>69.5</b>	<b>59.2</b>

Table 1: Testing Accuracies (%) on Mini-Kinetics Dataset.

boundary indirectly. In each category, most of the detected noisy instances do not share the same label as the clean ones. Therefore when expanding the gap between clean and noisy clusters, the decision margin between each category pair is simultaneously enlarged. The NCL and CT nourish each other. The model regularized by NCL provides a good representation for CT to differentiate clean and noisy instances better. Meanwhile, a well splitting of clean / unlabelled instances by CT supplies reasonable separation of positive and negative key sets for Noise Contrastive Learning.

## 4 Experiments

### 4.1 Settings and Datasets

**Model** In this section, We choose TSM-ResNet50 (Lin, Gan, and Han 2019) with ImageNet-pretraining as our backbone for all experiment settings. All training hyper-parameters are adapted from the original work (Lin, Gan, and Han 2019) and detailed in Appendix. Following previous work, we choose  $\tau = 0.1$  for all experiments in this paper. A memory bank with the size of  $16 \times K$  is added for all contrastive learning loss to provide ample training positive-negative pairs. All baselines in this paper are re-implemented by open-source codes.

**Datasets** We conduct experiments on three large-scale video classification datasets, Kinetics-400 (K400), Mini-kinetics (K200), and Something-Something-V1 (SthV1).

- **K400** This is a large-scale action recognition dataset with 400 categories, and it contains  $\sim 240k$  videos for training, and  $\sim 20k$  (50 per category) for validation. Note that K400 is not a balanced dataset. The amount of training samples in each category are in the range of [250, 1000].

- **K200** This is a balanced subset of K400 with 200 categories. In this dataset, each category contains 400 videos for training and 25 videos for validation.

- **SthV1** Something-Something is a challenging dataset focusing on temporal reasoning. In this dataset, the scenes and

Dataset Noise Type Noise Ratio	Kinetics						Something V1					
	Symmetric			Asymmetric			Symmetric			Asymmetric		
	40%	60%	80%	10%	20%	40%	40%	60%	80%	10%	20%	40%
TopoFilter(2020)	49.1	40.4	21.5	56.1	55.1	47.8	21.4	4.3	1.2	35.8	34.5	26.9
Co-teaching(2018)	53.9	51.4	31.3	51.4	51.4	41.7	23.6	14.5	4.5	24.0	24.6	18.5
M-correction(2019)	57.3	51.2	40.3	54.8	53.9	47.1	24.1	15.1	4.5	36.6	35.2	22.0
CT-ave	<b>60.3</b>	<b>56.6</b>	<b>46.3</b>	<b>62.9</b>	<b>61.3</b>	48.0	<b>36.3</b>	26.2	4.6	41.5	37.8	28.3
CT-var	60.1	55.8	45.7	62.8	61.0	<b>48.4</b>	36.2	<b>26.7</b>	<b>4.8</b>	<b>41.6</b>	<b>38.7</b>	<b>28.7</b>
Clean Only	61.6	58.8	54.3	70.3	68.0	61.7	40.7	36.0	24.8	44.0	43.5	40.6
CT*-var	61.1	56.7	48.6	63.6	62.5	57.8	36.8	27.0	5.8	41.7	40.2	32.8
CT-var+CL	60.0	55.3	7.6	61.7	59.5	45.2	31.8	1.6	1.1	37.8	36.2	27.1
CT-var+SCL	60.5	56.5	46.5	63.0	60.9	48.8	37.6	28.4	2.9	41.5	40.2	29.4
CT-var+NCL	<b>61.2</b>	<b>57.2</b>	<b>46.9</b>	<b>63.3</b>	<b>61.5</b>	<b>49.1</b>	<b>38.3</b>	<b>30.1</b>	<b>5.1</b>	<b>41.8</b>	<b>40.8</b>	<b>30.0</b>

Table 2: Testing Accuracies (%) on Kinetics and Something V1 Dataset.

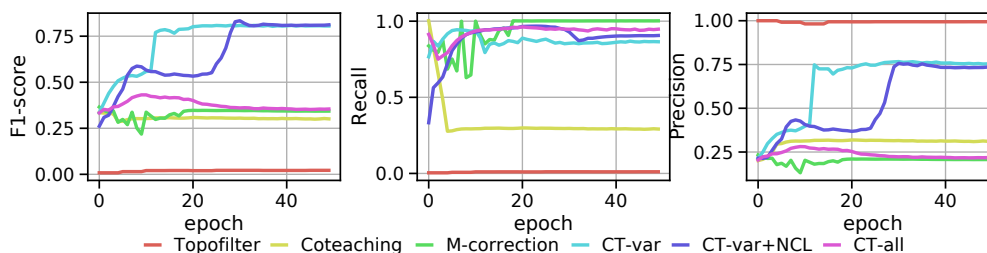


Figure 3: F1-Score/Precision/Recall in K200 dataset with 80% symmetric label noise.

objects in each category are various, which strongly requires the classification model to focus on the *temporal dynamics* among video frames. SthV1 is an unbalanced dataset with 174 categories. It contains  $\sim 86k$  videos for training and  $\sim 12k$  for validation.

**Noise Types** There are two kinds of label noise in the experiment: symmetric and asymmetric noise. Hyper-parameter  $\rho$  controls the noise ratio.

- **Symmetric Noise:** Each noisy sample in the training set is independently and uniformly assigned to a random label other than its true label.

- **Asymmetric Noise:** The samples in one category can only be assigned to a specific category other than the true label.

## 4.2 Label Noise Detection

We compare the Channel Truncation (CT) with multiple well-known methods for LNL. GCE (Zhang and Sabuncu 2018) and SCE (Wang et al. 2019) are robust loss functions. We also involve Clean Only baseline in the experiment as reference, which is trained with the correctly annotated samples. TopoFilter (Wu et al. 2020) is a feature-based method that detects noise by an L2-norm-based clustering. Co-teaching (Han et al. 2018) is a loss-based method with a dual-network design. The noise ratio  $\rho$  is assumed to be already known in Co-teaching. M-correction (Arazo et al. 2019) is also a loss-based method with beta mixture model for clean/noisy instance separation. DivideMix (Li, Socher, and Hoi 2019) applies Mix-up (Zhang et al. 2018a) on both labeled and unlabeled samples to enrich the dataset. As they

are designed for images, we here re-implement video versions of them for comparison.

All methods share the same network architecture (TSM-ResNet50), training hyper-parameters, noise settings, and train/test splits of datasets for a fair comparison. We set  $b$  to 200 in all settings. Since each query is expected to have 16 true positive keys in the memory bank with size  $16 \times K$ , the number of nearest keys  $B$  is set as 16. The reference set  $\Omega$  is initialized with the whole dataset  $\mathcal{D}$  except for Asymmetric-40%. In the setting of Asymmetric-40%, we use  $k$ -means for binary clustering and regard the cluster with the dominant instances as the initialization of  $\Omega$ . Note that we choose  $\phi_{var}$  as the score function for all experiment and denote our method with  $\phi_{var}$  as CT-var. As an additional comparison, CT-all takes all the channels for noise detection, *i.e.*  $b = 2048$ . As a tight upper bound of CT, CT-oracle is implemented by utilizing oracle selection for channel selection.

Table. 1 and 2 shows the best testing accuracy across all epochs. The best results are marked as **bold**. On the K200, K400, and SthV1 datasets, our method greatly outperforms other LNL baselines under different levels of symmetric/asymmetric label noise.

Due to the unbalanced number of samples in different categories, under the asymmetric noise with large  $\rho$ , some categories' number of noisy instances may be larger than that of clean ones. The clean/noisy instances in these unbalanced datasets are theoretically indistinguishable if extra information is not provided. To tackle this issue, we keep a balanced instance set, where each category contains *ten* samples with

$b$	100	200	400	1600	2048 (all)
CT- <i>var</i>	60.6	<b>61.1</b>	60.8	59.0	58.3
CT- <i>var</i> +NCL	62.9	<b>63.4</b>	63.3	62.0	61.6

Table 3: Testing Accuracies (%) on K200 with 60% Symmetric Noise and Different Number of Kept Channels  $b$ .

Score	Top	Middle	Bottom
CT- <i>var</i>	<b>48.4</b>	34.6	25.5
CT- <i>ave</i>	<b>48.1</b>	35.4	20.9

Table 4: Test accuracy (%) on K200 with 80% Symmetric Noise. Channels with higher score benefit more to model performance.

correct labels. By taking this small clean sets as  $\Omega$ , CT\* achieves a substantial improvement in these unbalanced situations.

Fig. 3 shows the F1-score, Precision, and Recall of noise detection during training on K200 with symmetric-80% noise. Our method beats the baselines in almost every training epoch in F1-Score. We notice that most methods have a similar F1-score at the very beginning, and no one has an obvious advantage over others when the model is weak. It takes about ten epochs for CT-*var* to warm up before noise detection ability leaps. The warm-up duration of CT-*var*+NCL is much larger than that of CT-*var*, which may be caused by the randomly-initialized projection head  $h(\cdot)$ . Thanks to NCL, our method tends to discover more clean instances leading to higher recall but a little bit lower precision than CT-*var*. The improved classification accuracy shows that this Recall-Precision trade-off is beneficial.

### 4.3 Regularization on Unlabelled Noisy Sample

To fully utilize the unlabelled noisy set  $\mathcal{D}_{\hat{u}}$ , our Noise Contrastive Loss (NCL) provides a low-risk regularization in video classification scenario with noisy labels. We compare NCL with two famous contrastive losses: standard Contrastive Learning (CL) loss and Supervised Contrastive Learning (SCL) (Khosla et al. 2020) loss. CL loss involves all instances in training. As SCL requires the label supervisions, instances from  $\mathcal{D}_{\hat{u}}$  are ignored in our SCL implementation. Training settings are shared for all methods, and the only difference is the contrastive strategy, which rules how positive/negative keys are selected. The results shown in Table 1 and Table 2 indicate that NCL can stably improve model performance. We also implement three commonly used pseudo-label-based (PL) strategies, in which a pseudo label is assigned to each unlabelled noisy instance. PL-H (Zhang et al. 2021)/PL-S (Li, Socher, and Hoi 2019) generates a hard/soft label for unlabelled instances according to model prediction. PL-K (Ortego et al. 2021) utilizes  $k$ -NN on the reference set  $\Omega$  to predict pseudo label for instances in  $\mathcal{D}_{\hat{u}}$ . Here there are 64 nearest neighbors considered. We compare NCL with PL on the K200 dataset as shown in Table 1. As we mentioned in Sec.3.2, a pseudo label without

Score Function	$\phi_{ave}$	$\phi_{var}$
Mini-Kinetics	<b>94.9</b>	77.9
Something V1	45.5	<b>50.3</b>

Table 5: Comparison of Score Function on Different Datasets with 40% Symmetric Noise (larger is better).

robust post-processing is less reliable when the generalization of a model is weak. Tremendous wrong pseudo labels will dramatically degrade the supervision quality and thus hurt model prediction.

### 4.4 Ablation Studies

**Parameter Sensitivity** The hyper-parameter  $b$  controls how many dimensions are kept after channel truncation. As shown in Table 3, CT performs better when  $b$  is close to 200.  $b = 2048$  means no channels is truncated.

**Channel Selection** We select the top, middle and bottom ranked  $b$  channels, respectively, for noise detection and evaluate their impacts on final classification. As shown in Table 4, model performance dramatically drops if we only keep  $b$  channels with low scores, which proves the significance of channel selection.

**Score Function** We expect the proposed selections can keep as many the most discriminative channels for each category as the oracle selection does. To verify the ability of our proposed selections, we quantify the similarity of channel selection between the proposed selections and the oracle selection by the intersection size of selected channels by these two criteria at fifth epoch in model updating. The larger the size of the intersection, and the better the score function of the proposed selection. We report the evaluation in Table 5. On the scene-based dataset K200,  $\psi_{ave}$  shares the most selected channels with the Oracle Selection  $\psi_o$ . While on the temporal reasoning dataset SthV1, channels selected by the  $\psi_{var}$  is more discriminative. This observation keeps consistent with their performance in testing accuracy.

**Extension to Images** We also extend  $\phi(\cdot)$  to an image version. Details can be found in Appendix.

## 5 Discussion and Conclusion

In this work, we propose an efficient framework for learning with noisy labels in video classification. This framework includes a feature-based noise detection method, the Channel Truncation, and a regularization approach, the Noise Contrastive Loss. The proposed CT adaptively selects discriminative channels for clean instances filtering in each model updating round, achieving good noise detection performance even with severe interference of noisy labels. The novel NCL learning strategy fully utilizes the detected noisy instances to assist model updating in reshaping the classification boundaries among categories. Extensive experiments validate the effectiveness of our framework NEAT in both noise detection and classification under various settings.

## Acknowledgements

This work was partially supported by the National Key R&D Program of China (2022YFB4701400/4701402), SZSTC Grant (JCYJ20190809172201639, WDZC2020082020065-5001), and Shenzhen Key Laboratory (ZDSYS2021062309-2001004).

## References

- Arazo, E.; Ortego, D.; Albert, P.; O'Connor, N.; and McGuinness, K. 2019. Unsupervised Label Noise Modeling and Loss Correction. In *ICML*.
- Arpit, D.; Jastrzebski, S.; Ballas, N.; Krueger, D.; Bengio, E.; Kanwal, M. S.; Maharaj, T.; Fischer, A.; Courville, A.; Bengio, Y.; et al. 2017. A closer look at memorization in deep networks. In *ICML*.
- Chen, T.; Kornblith, S.; Norouzi, M.; and Hinton, G. 2020a. A simple framework for contrastive learning of visual representations. In *ICML*. PMLR.
- Chen, X.; Fan, H.; Girshick, R.; and He, K. 2020b. Improved baselines with momentum contrastive learning. arXiv:2003.04297.
- Choi, J.; Gao, C.; Messou, J. C.; and Huang, J.-B. 2019. Why Can't I Dance in the Mall? Learning to Mitigate Scene Bias in Action Recognition. *NeurIPS*.
- Gao, B.; Gouk, H.; and Hospedales, T. M. 2021. Searching for Robustness: Loss Learning for Noisy Classification Tasks. arXiv:2103.00243.
- Han, B.; Yao, Q.; Yu, X.; Niu, G.; Xu, M.; Hu, W.; Tsang, I.; and Sugiyama, M. 2018. Co-teaching: Robust training of deep neural networks with extremely noisy labels. In *NeurIPS*.
- Han, J.; Luo, P.; and Wang, X. 2019. Deep self-learning from noisy labels. In *ICCV*.
- Huang, J.; Qu, L.; Jia, R.; and Zhao, B. 2019. O2u-net: A simple noisy label detection approach for deep neural networks. In *ICCV*.
- Jing, L.; and Tian, Y. 2020. Self-supervised visual feature learning with deep neural networks: A survey. *T-PAMI*.
- Khosla, P.; Teterwak, P.; Wang, C.; Sarna, A.; Tian, Y.; Isola, P.; Maschinot, A.; Liu, C.; and Krishnan, D. 2020. Supervised Contrastive Learning. *NeurIPS*, 33.
- Kim, C. D.; Jeong, J.; Moon, S.; and Kim, G. 2021. Continual Learning on Noisy Data Streams via Self-Purified Replay. In *ICCV*.
- Lee, K.-H.; He, X.; Zhang, L.; and Yang, L. 2018. Cleannet: Transfer learning for scalable image classifier training with label noise. In *CVPR*.
- Li, J.; Socher, R.; and Hoi, S. C. 2019. DivideMix: Learning with Noisy Labels as Semi-supervised Learning. In *ICLR*.
- Li, J.; Xiong, C.; and Hoi, S. C. 2021. Learning From Noisy Data With Robust Representation Learning. In *ICCV*.
- Li, J.; Zhou, P.; Xiong, C.; and Hoi, S. 2020. Prototypical Contrastive Learning of Unsupervised Representations. In *ICLR*.
- Lin, J.; Gan, C.; and Han, S. 2019. Tsm: Temporal shift module for efficient video understanding. In *ICCV*.
- Lyu, Y.; and Tsang, I. W. 2019. Curriculum loss: Robust learning and generalization against label corruption. In *ICML*.
- Ma, X.; Huang, H.; Wang, Y.; Romano, S.; Erfani, S.; and Bailey, J. 2020. Normalized Loss Functions for Deep Learning with Noisy Labels. arXiv:2006.13554.
- Oord, A. v. d.; Li, Y.; and Vinyals, O. 2018. Representation learning with contrastive predictive coding. arXiv:1807.03748.
- Ortego, D.; Arazo, E.; Albert, P.; O'Connor, N. E.; and McGuinness, K. 2021. Multi-Objective Interpolation Training for Robustness to Label Noise. In *CVPR*.
- Patrini, G.; Rozza, A.; Krishna Menon, A.; Nock, R.; and Qu, L. 2017. Making deep neural networks robust to label noise: A loss correction approach. In *CVPR*.
- Tian, Y.; Sun, C.; Poole, B.; Krishnan, D.; Schmid, C.; and Isola, P. 2020. What makes for good views for contrastive learning? *NeurIPS*.
- Wang, J.; Gao, Y.; Li, K.; Lin, Y.; Ma, A. J.; Cheng, H.; Peng, P.; Huang, F.; Ji, R.; and Sun, X. 2021. Removing the Background by Adding the Background: Towards Background Robust Self-supervised Video Representation Learning. In *CVPR*.
- Wang, Y.; and Hoai, M. 2018. Pulling actions out of context: Explicit separation for effective combination. In *CVPR*.
- Wang, Y.; Ma, X.; Chen, Z.; Luo, Y.; Yi, J.; and Bailey, J. 2019. Symmetric cross entropy for robust learning with noisy labels. In *ICCV*.
- Weng, J.; Luo, D.; Wang, Y.; Tai, Y.; Wang, C.; Li, J.; Huang, F.; Jiang, X.; and Yuan, J. 2020. Temporal distinct representation learning for action recognition. In *ECCV*. Springer.
- Wu, P.; Zheng, S.; Goswami, M.; Metaxas, D.; and Chen, C. 2020. A Topological Filter for Learning with Label Noise. *NeurIPS*, 33.
- Zhang, B.; Wang, Y.; Hou, W.; Wu, H.; Wang, J.; Okumura, M.; and Shinozaki, T. 2021. Flexmatch: Boosting semi-supervised learning with curriculum pseudo labeling. arXiv:2110.08263.
- Zhang, C.; Bengio, S.; Hardt, M.; Recht, B.; and Vinyals, O. 2017. Understanding deep learning requires rethinking generalization. *ICLR*.
- Zhang, H.; Cisse, M.; Dauphin, Y. N.; and Lopez-Paz, D. 2018a. mixup: Beyond Empirical Risk Minimization. In *ICLR*.
- Zhang, J.; Wang, Z.; Meng, J.; Tan, Y.-P.; and Yuan, J. 2018b. Boosting positive and unlabeled learning for anomaly detection with multi-features. *T-MM*.
- Zhang, Z.; and Sabuncu, M. 2018. Generalized cross entropy loss for training deep neural networks with noisy labels. In *NeurIPS*.
- Zhou, X.; Liu, X.; Wang, C.; Zhai, D.; Jiang, J.; and Ji, X. 2021. Learning with Noisy Labels via Sparse Regularization. In *ICCV*.