

## 1 Introduction

- The challenge we aim to solve: **abnormal training examples**. We regard a training example as abnormal unrestrictedly whenever **an observation and its label are semantically unmatched**.



Truck class: The first two images are automobile.



This video is labelled as the person wearing long hair.



This video is labelled as the person wearing black trouser.

Fig 1: Display of abnormal training examples highlighted by red boxes.

The 1st row shows synthetic abnormal examples from corrupted CIFAR-10.

The 2nd and 3rd rows present realistic abnormal examples from video person re-identification benchmark MARS (Zheng et al. 2016).

Three representatives:

- 1) The abnormal images with no person in 3rd row contain no semantic information at all.
- 2) The last abnormal image in 2nd or 3rd row may contain a person that does not belong to any person in the training set.
- 3) We cannot decide the object of interest without any prior when an image contains more than one object, e.g., the 2nd and 3rd last images in 2nd row contain two persons.

- MAE's underfitting problem followed with New Analysis and Interpretation

- **Underfitting observations:** In Table 1, when 40% noise exists, compared with CCE, MAE underfits to clean training data points, thus fitting much fewer abnormal examples.
- **Conclusion:** In Fig 2, MAE emphasises more on uncertain examples, whose probabilities of being classified to its labelled class are 0.5, thus being noise-robust.
- **Interpretation:** According to Fig 2, MAE differentiates samples in a noise-robust way, but its differentiation degree over training examples is too small.

Table 1: Classification accuracy (%) of CCE, MAE, and IMAE on CIFAR-10 (Krizhevsky 2009). 40% of training examples, i.e., the noisy subset, have wrong labels. We test each model's performance on test set, noisy subset and clean subset of training set. The backbone is ResNet56 owning enough capacity (He et al. 2016).

Loss	Test set (Generalisation)	Noisy subset (Noise-ignorance)	Clean subset (Learning ability)
CCE	63.3	75.0	96.2
MAE	66.9	8.1	<b>74.3</b> (worst)
IMAE	<b>81.5</b> (best)	<b>6.5</b> (best)	93.1

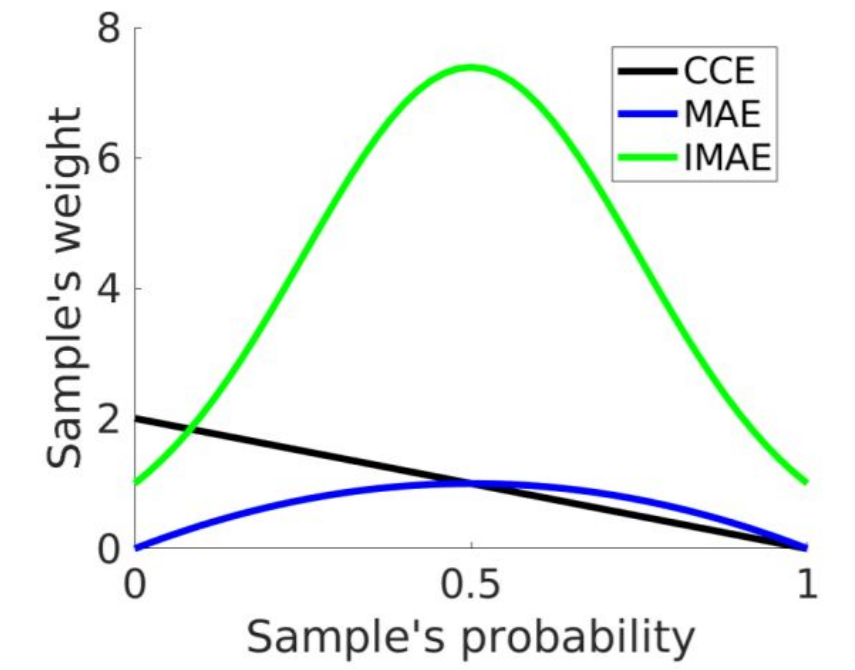


Fig 2: Sample's weight along with sample's probability being classified to its labelled class (input-to-label relevance score).

## 2 Methodology

- New perspectives and new conclusion about MAE.

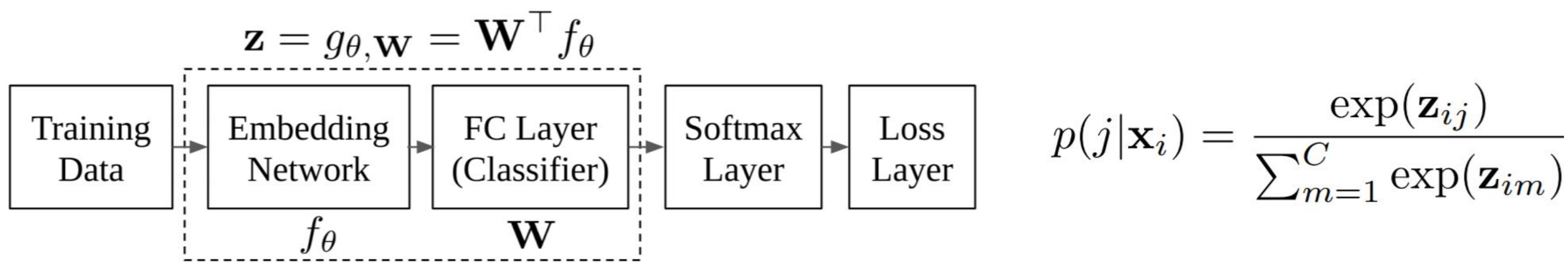


Fig 3: Pipeline of a deep feature embedding network. The output of softmax layer is interpreted as classification probabilities.

$$L_{CCE}(\mathbf{X}; f_{\theta}, \mathbf{W}) = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C q(j|\mathbf{x}_i) \log p(j|\mathbf{x}_i) \quad L_{MAE}(\mathbf{X}; f_{\theta}, \mathbf{W}) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C |p(j|\mathbf{x}_i) - q(j|\mathbf{x}_i)|$$

$$= -\frac{1}{N} \sum_{i=1}^N \log p(y_i|\mathbf{x}_i). \quad = \frac{2}{N} \sum_{i=1}^N (1 - p(y_i|\mathbf{x}_i)),$$

$$\frac{\partial L_{CCE}(\mathbf{x}_i)}{\partial p(j|\mathbf{x}_i)} = \begin{cases} -p(y_i|\mathbf{x}_i)^{-1}, & j = y_i \\ 0, & j \neq y_i \end{cases}$$

$$\frac{\partial L_{MAE}(\mathbf{x}_i)}{\partial p(j|\mathbf{x}_i)} = \begin{cases} -2, & j = y_i \\ 0, & j \neq y_i \end{cases}$$

*Gradient magnitude w.r.t. probabilities*--Prior conclusion in Zhang and Sabuncu 2018: CCE is sensitive to abnormal data points while MAE treat all examples equally, thus being robust.

$$\frac{\partial L_{CCE}(\mathbf{x}_i)}{\partial \mathbf{z}_{ij}} = \begin{cases} p(y_i|\mathbf{x}_i) - 1, & j = y_i \\ p(j|\mathbf{x}_i), & j \neq y_i \end{cases} \quad \frac{\partial L_{MAE}(\mathbf{x}_i)}{\partial \mathbf{z}_{ij}} = \begin{cases} 2p(y_i|\mathbf{x}_i)(p(y_i|\mathbf{x}_i) - 1), & j = y_i \\ 2p(y_i|\mathbf{x}_i)p(j|\mathbf{x}_i), & j \neq y_i \end{cases}$$

$$w_{CCE}(\mathbf{x}_i) = \left\| \frac{\partial L_{CCE}(\mathbf{x}_i)}{\partial \mathbf{z}_i} \right\|_1 = 2(1 - p(y_i|\mathbf{x}_i))$$

$$w_{MAE}(\mathbf{x}_i) = \left\| \frac{\partial L_{MAE}(\mathbf{x}_i)}{\partial \mathbf{z}_i} \right\|_1 = 4p(y_i|\mathbf{x}_i)(1 - p(y_i|\mathbf{x}_i))$$

*Gradient magnitude w.r.t. logits*--Our conclusion is illustrated in Fig 2.

- Interpretation of MAE's underfitting?

MAE's differentiation degree over data points is too small, which means the relative contribution of one example versus another is not recognised well and the majority contribute almost equally. Therefore, MAE generally underfits training data especially when noise rate is high, as justified in Table 1.

- How to solve it?

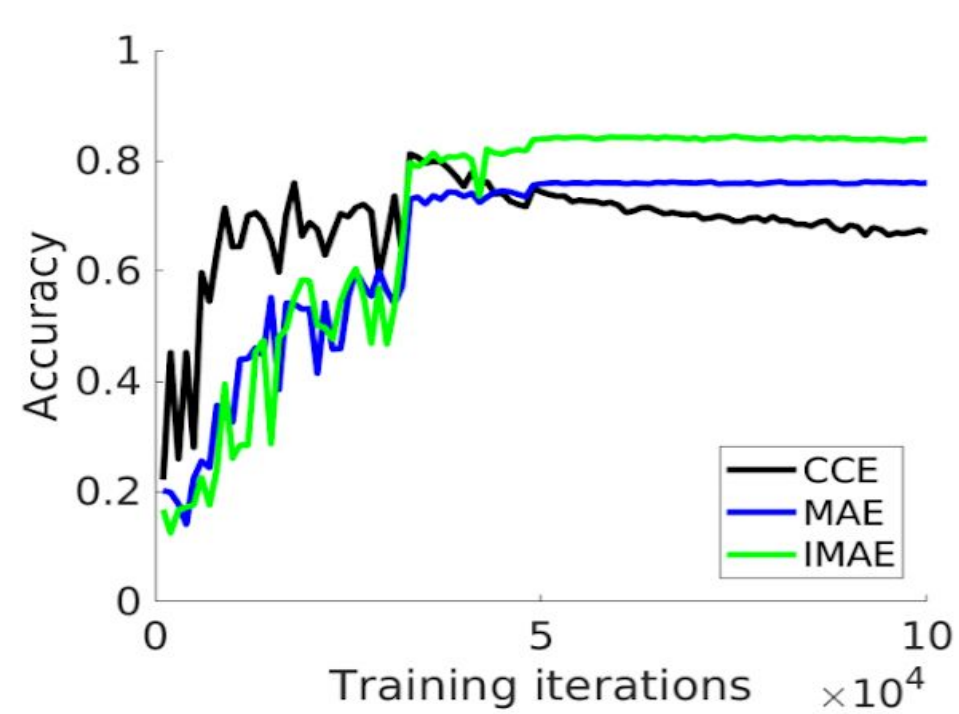
$$w_{IMAE}(\mathbf{x}_i) = \exp(T \cdot p(y_i|\mathbf{x}_i)(1 - p(y_i|\mathbf{x}_i)))$$

$$\frac{\partial L_{IMAE}(\mathbf{x}_i)}{\partial \mathbf{z}_i} = \frac{\partial L_{MAE}(\mathbf{x}_i)}{\partial \mathbf{z}_i} \cdot \frac{w_{IMAE}(\mathbf{x}_i)}{w_{MAE}(\mathbf{x}_i)}$$

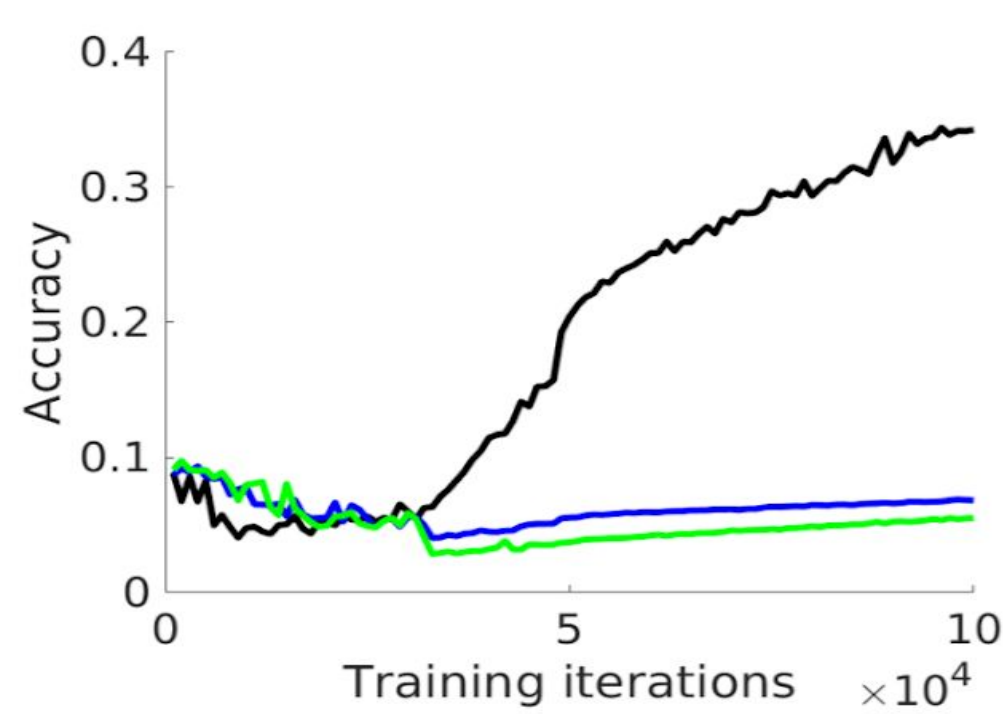
$$\Rightarrow \left\| \frac{\partial L_{IMAE}(\mathbf{x}_i)}{\partial \mathbf{z}_i} \right\|_1 = w_{IMAE}(\mathbf{x}_i).$$

- Being the same as MAE, IMAE is a symmetric loss since it does not revise the loss computation of MAE. That is to say, IMAE's loss values are bounded.
- Being more flexible and practical, IMAE has a hyper-parameter T to control the weighting variance over training samples.
- The theoretical analysis on MAE's noise-robustness based on loss values (Ghosh, Kumar, and Sastry 2017) can be directly applied to IMAE.

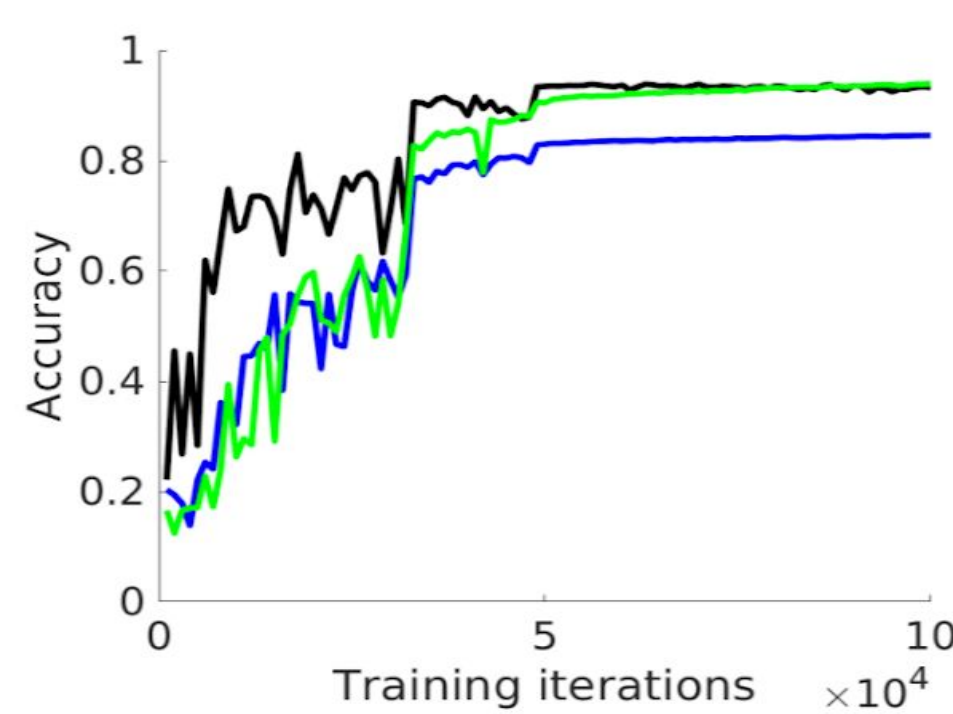
## 3 Experiments on CIFAR-10 with symmetric label noise (For more experiments, please see <https://arxiv.org/pdf/1903.12141.pdf>)



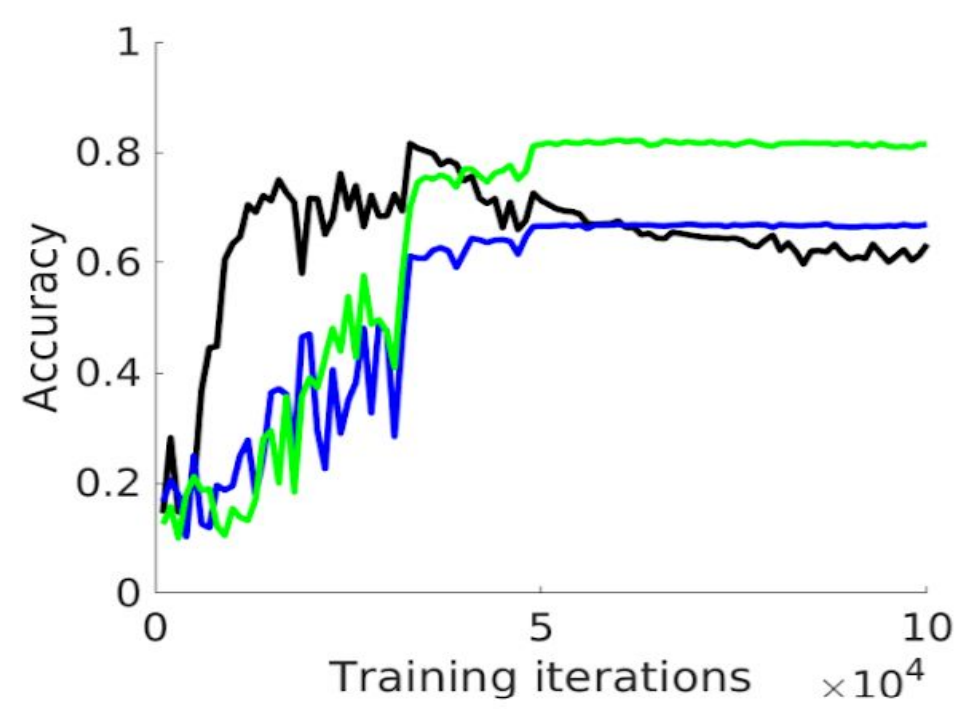
(a) ResNet20: Testing set (higher is better).



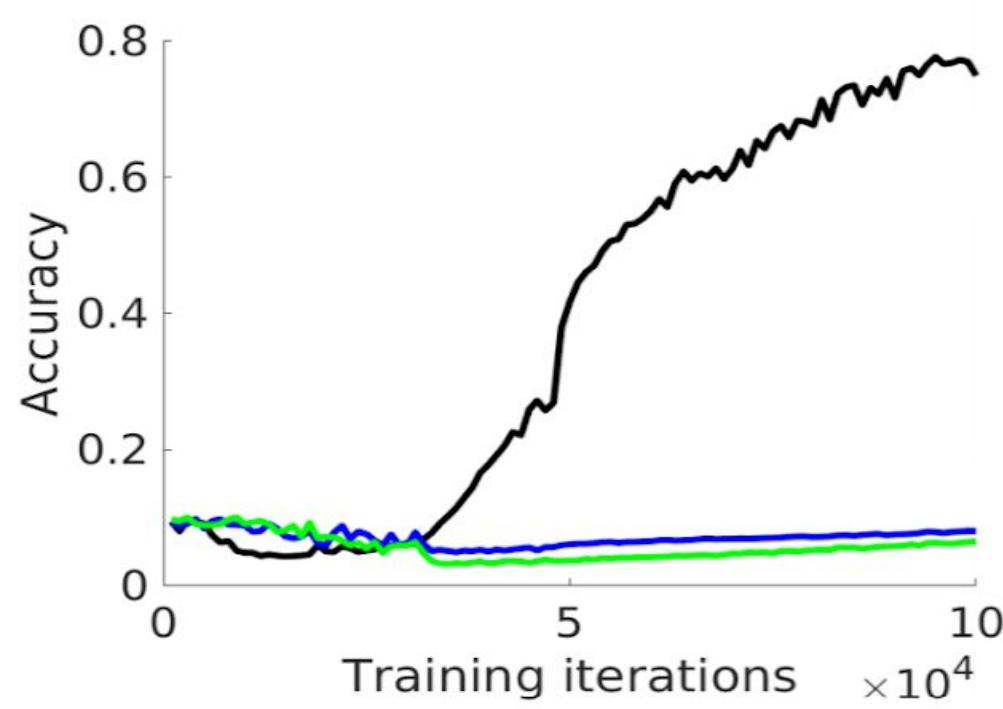
(b) ResNet20: Noisy subset (lower is better).



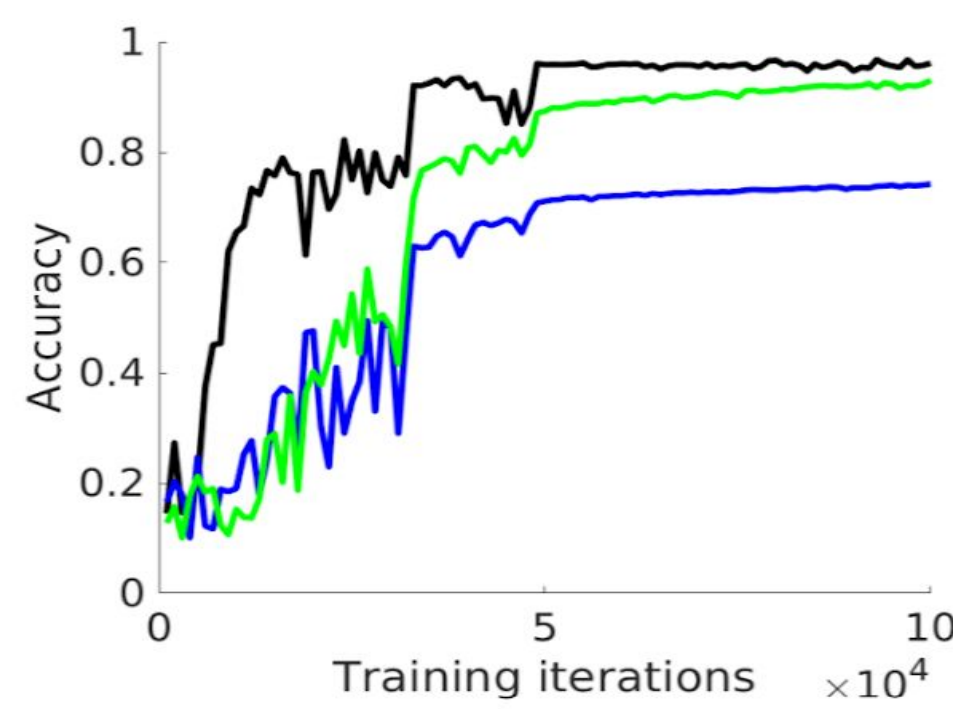
(c) ResNet20: Clean subset (higher is better).



(d) ResNet56: Testing set (higher is better).



(e) ResNet56: Noisy subset (lower is better).



(f) ResNet56: Clean subset (higher is better).

Figure 4: CIFAR-10 with noise rate  $r = 40\%$ . The accuracies on testing set, noisy subset and clean subset of training set along with training iterations. The legend on the top left is shared by all subfigures. *Better viewed in colour.*

Backbone	$r$	Loss	Testing accuracy		Training accuracy: Naive fitting		Hybrid accuracy: Meaningful patterns
			Best	Final	Noisy subset	Clean subset	
ResNet20	0%	CCE	91.5	91.3	—	<b>100</b>	<b>98.5</b>
		MAE	89.3	89.2	—	95.8	94.7
		IMAE	<b>91.7</b>	<b>91.4</b>	—	99.8	98.4
	40%	CCE	81.2	67.0	34.3	93.3	72.6
		MAE	76.2	75.9	6.8	84.6	79.7
		IMAE	<b>84.3</b>	<b>84.0</b>	<b>5.5</b>	<b>94.0</b>	<b>88.2</b>
ResNet56	80%	CCE	43.0	20.3	38.3	57.0	22.0
		MAE	27.7	27.5	<b>9.7</b>	29.4	27.8
		IMAE	<b>52.0</b>	<b>41.0</b>	16.8	<b>64.8</b>	<b>41.5</b>
	0%	CCE	<b>92.4</b>	<b>92.2</b>	—	<b>100</b>	<b>98.7</b>
		MAE	89.0	89.0	—	96.1	94.9
		IMAE	92.2	<b>92.2</b>	—	99.8	98.5
ResNet56	40%	CCE	81.6	63.3	75.0	<b>96.2</b>	63.6
		MAE	67.0	66.9	8.1	74.3	70.2
		IMAE	<b>82.2</b>	<b>81.5</b>	<b>6.5</b>	93.1	<b>86.5</b>
	80%	CCE	<b>38.2</b>	16.4	52.5	<b>62.3</b>	18.7
		MAE	15.2	15.1	<b>9.6</b>	15.6	15.1
		IMAE	37.1	<b>34.0</b>	13.0	44.7	<b>34.8</b>

## 4 Summary

- New Findings:

- MAE emphasises on medium-probability examples, thus being noise-tolerant.
- MAE underfits to meaningful patterns when severe noise exists.

- Solution: IMAE

- An effective and simple solution for improving MAE's learning ability while preserving its noise-robustness.
- It works for different types of abnormal training examples.