1 **To All Reviewers:** We thank all reviews for your insightful feedback and your appreciation of our $MCR^2$ formulation.

2 We will incorporate suggestions on minor corrections, references, footnotes, and presentations in the final version.

3 **Why diverse intra-class representations?** This work aims to introduce a new objective (i.e., $MCR^2$) for learning

4 representations not only discriminative between classes as with cross-entropy loss, but also *diverse* within class. We

5 believe identifying more discriminative features lead to more reliable classification since the most discriminative

6 feature may not be present in all samples. We rigorously prove that this can be achieved with the proposed $MCR^2$ loss

7 function. Furthermore, we empirically demonstrate this objective can be used to train deep networks that have good

8 properties in handling label noise (in supervised setting) and achieve SOTA for clustering (in unsupervised setting).

9 **Robustness to label noise:** The initial motivation of $MCR^2$ is to

10 promote learning rich discriminative features. It is a nice surprise

11 that so learned deep features are more robust than existing learning

12 objectives including cross entropy and many others shown in Ta-

13 bles 1, 2. Unlike cross entropy that fits labels of individual samples,

14 $MCR^2$ compresses samples of each class *collectively*. As mentioned

15 in Section 4, given the compelling empirical evidence, a rigorous

16 justification of the robustness is an exciting problem for future work.

Table 1: Comparison with OLE and Large Margin [EKM+18] on learning from noisy labels.

| RESNET18 | RATIO=0.1 | RATIO=0.2 | RATIO=0.3 | RATIO=0.4 | RATIO=0.5 |
|---|---|---|---|---|---|
| OLE | 91.04% | 86.01% | 80.69% | 71.79% | 61.06% |
| [EKM+18] | 90.10% | 87.42% | 83.77% | 78.51% | 72.48% |
| $MCR^2$ | **91.16%** | **89.70%** | **88.18%** | **86.66%** | **84.30%** |

Table 2: Comparison with Trimmed Loss [SS19] on learning from noisy labels.

| WRN16 | RATIO=0.1 | RATIO=0.3 | RATIO=0.5 | RATIO=0.7 |
|---|---|---|---|---|
| [SS19] | 90.33% | 88.23% | 82.51% | 64.74% |
| $MCR^2$ | **91.55%** | **88.81%** | **84.25%** | **67.09%** |

17 **To Reviewer #1:** Please refer to the top of the rebuttal for the motivations of larger intra-class subspace in $MCR^2$.

18 **Q1:** *Compare with OLE: " 1). It is not clear why a larger ... these connections are not crystal clear in the paper. 2).*

19 *Does the OLE type loss have the same property as Theorem 1? 3). The authors should show more comparison ... OLE."*

20 **A1:** 1). We will make these connections more clear in the final version; 2). As mentioned in the paper (line 209-213),

21 OLE loss does *not* have the diversity property of $MCR^2$ given in Theorem 1; 3). In Table 1, we compare $MCR^2$ with

22 OLE on the corrupted label task using the same network architecture. $MCR^2$ achieves significantly better performance.

23 **Q2:** *Gaussian assumption of data: "I have a concern whether the rate distortion function ... to be self-contained."*

24 **A2:** Thank you for your suggestion. As shown in [MDHW07], the rate distortion function can serve as a tight and

25 accurate approximation for a wide range of subspace-like distributions. We will give more details in the final version.

26 **Q3:** *The paper can be considered as applying existing objective/criterion ... into learning of deep features.*

27 **A3:** We disagree. Our $MCR^2$ objective is new and is different from those in previous works such as OLE. To our best

28 knowledge, $MCR^2$ is *the first* objective theoretically shown to guarantee both diverse and discriminative properties.

29 **Q4:** *In fact, a core problem in understanding deep learning ... on it.*

30 **A4:** Thanks for your comment. Precisely, we believe identifying a diverse and discriminative representation from the

31 data is an important step to gaining better understandings of the generalizability and robustness of deep learning.

32 **To Reviewer #2:**

33 **Q1:** *Relationship with information bottleneck (IB) framework: " In Section 2, the authors seek to ... Gaussians."*

34 **A1:** Both $MCR^2$ and IB are information-theoretic objectives. However, the goal of IB is to find a *minimal* set of most

35 informative representations while $MCR^2$ aims to capture both diverse and discriminative representations, which is very

36 different. We will better clarify relationships with mutual information-based approaches in our final version.

37 **Q2:** *Related work on label noise: "The label noise robustness experiments ... iterative trimmed loss minimization [1]."*

38 **A2:** In Table 2, we compare $MCR^2$ with [SS19] using the same network. $MCR^2$ achieves better performance *without*

39 any noise ratio dependent hyperparameters as required by [SS19]. We will add the comparison in the final version.

40 **To Reviewer #3:** Please refer to the top of the rebuttal for clarifying the objectives and motivations of our work.

41 **Q1:** *"My main concern is that, I don't see the benefits ... lie on a union of subspaces)."*

42 **A1:** First of all, we do *not* model the original data by subspaces. $MCR^2$ can guide a deep network to map real data on

43 complicated nonlinear submanifolds to a union of orthogonal subspaces. Secondly, once the subspaces are learned, the

44 nearest subspace classification is computationally *efficient*. Finally, compared with *hidden* representations learned by

45 cross-entropy, the union of discriminative subspaces learned by $MCR^2$ is geometrically and statistically meaningful.

46 **Q2:** *"While the theoretical analysis reveals interesting properties ... the loss, see e.g. [ZF14]."*

47 **A2:** Our theoretical analysis reveals that the proposed $MCR^2$ is optimized only when features are the most diverse and

48 discriminative. Our experiments have clearly shown that using $MCR^2$, deep features learned from real data such as

49 CIFAR10 have the same nice properties that are predicted by our theoretical results. We plan to rigorously justify this

50 phenomenon by studying the interplay of the $MCR^2$ objective and the choice of network architectures in future work.

51 **Q3:** *Related work on clustering: "While there is no dedicated related work ... be included and compared against."*

52 **A3:** $MCR^2$ outperforms [HMT+17, JHV19] on both CIFAR10 and CIFAR100 by a large margin. For STL10, [HMT+17]

53 applied pretrained ImageNet models and $MCR^2$ outperforms [JHV19] when using the same amount of training data.

54 **Answers to minor comments:** We will add the above comparison and references, and also compared $MCR^2$ to

55 [EKM+18] on CIAFR10 with label noise (see Table 1). We did not encounter any computation issue when dealing with

56 $\log \det$ and the optimization is stable. The $\Pi$ is defined by the labels and satisfies the simplex constraint (footnote 15).

57 **To Reviewer #4:** Please refer to the top of the rebuttal for the question regarding the robustness of $MCR^2$.

58 **Q1:** *Applying the MCR reduction to the large-scale dataset seems computationally very hard.*

59 **A1:** The computation only increases *linearly* in the number of classes for the supervised learning setting.