

| Dataset         | Type     | Method        |        |               |          |               |
|-----------------|----------|---------------|--------|---------------|----------|---------------|
|                 |          | Standard      | GNT    | CAT           | AQPL-GNT | AQPL-CAT      |
| MNIST-C         | Clean    | <b>99.29%</b> | 97.32% | 98.43%        | 99.21%   | 99.23%        |
|                 | Gaussian | 16.06%        | 84.46% | <b>98.14%</b> | 96.31%   | 97.90%        |
|                 | All      | 65.34%        | 71.57% | 80.11%        | 78.78%   | <b>80.42%</b> |
| CIFAR10-C       | Clean    | <b>95.05%</b> | 94.87% | 86.42%        | 94.83%   | 94.75%        |
|                 | Gaussian | 43.23%        | 71.62% | 82.78%        | 82.19%   | <b>86.69%</b> |
|                 | All      | 74.24%        | 79.59% | 71.15%        | 82.02%   | <b>83.33%</b> |
| Tiny-Imagenet-C | Clean    | <b>57.84%</b> | 56.14% | 48.62%        | 56.60%   | 55.51%        |
|                 | Gaussian | 19.27%        | 21.90% | 27.98%        | 25.12%   | <b>31.72%</b> |
|                 | All      | 9.99%         | 14.04% | 23.77%        | 24.82%   | <b>27.19%</b> |

Table 1: The Top-1 accuracy of different methods on different corruption datasets.

to adaptively adjust the perturbation levels for different examples in the training process. Moreover, by allowing to query the ground-truth information on the perturbation level, the proposed approaches AQPL-GNT and AQPL-CAT can further improve the performances of GNT and CAT respectively. Most importantly, it can be observed that the proposed approach AQPL-CAT outperforms the other methods in most cases with regard to both Gaussian noise and the other 15 types of noise. Note that, when comparing with the method CAT that also adjusts perturbation level according to whether the current model has the capacity to robustly classify the examples, the AQPL-CAT can still achieve better performance. On one hand, the supervised information provided by the oracle is more reliable. On the other hand, human experts correct perturbation levels more efficiently and directly. In summary, these results consistently demonstrate that the proposed AQPL approach can effectively improve the model robustness by actively querying the correct perturbation level from the oracle, while the sampling strategy can efficiently select the most useful examples to reduce the querying cost.

## Discussion

Similar to many existing studies, the experiments are performed on image datasets in this paper. The results show that, by actively querying the supervised information about the perturbation level, model robustness against corruption perturbations on image classification tasks can be improved efficiently. In principle, the proposed method can be applied to any type of data. One challenge is that it could be difficult for human annotators to select a proper perturbation level for non-visual data. If the non-visual data can be easily visualized, such as VisArtico (?) for articulatory data, the method is still applicable. It would be an interesting future work to design feasible interfaces for annotators to decide the perturbation level for non-visual data.

In this paper, we focus on the corruption perturbations both in our theoretical and experimental analysis. We believe that corruption perturbations commonly occur in real tasks. On the other hand, it would be interesting to extend the study for improving adversarial robustness. Actually, the average-case robustness under a specific noise distribution could bring non-negligible adversarial robustness (?). More importantly, the optimal perturbation level for a clean example considered in this paper, essentially, represents an adversarial (worst-case)

noise distribution on the example with regard to the oracle.

## Conclusion

In this work, we propose a novel active learning framework to improve the model robustness by querying the conform perturbation levels. On one hand, instead of assuming a fixed noise for the whole training set, the perturbation levels are adjusted adaptively for different examples during the training process. On the other hand, by estimating the conformity with classification entropy, the most useful examples are actively selected to achieve effective learning with lower annotation cost. Both theoretical and empirical results validate the effectiveness of the proposed approach. In the future, we plan to extend the framework to handle adversarial perturbations.

## Acknowledgments

This research was supported by the National Key R&D Program of China (2020AAA0107000), NSFC (62076128) and the China University S&T Innovation Plan Guided by the Ministry of Education.

Itaque recusandae facere, aperiam nesciunt eligendi nobis veritatis molestiae, debitis nihil ipsum, maxime assumenda possimus, culpa eaque dicta recusandae delectus nihil eius. Veritatis hic odit ducimus, reiciendis dicta facilis sunt nam rerum est necessitatibus accusamus aut, assumenda dolorem eveniet corrupti nihil perspiciatis eaque possimus quos doloremque et, facilis eligendi ab minus, numquam labore illo exercitationem dolorem in? Illo ipsa nobis recusandae quos vitae soluta, obcaecati deserunt fugit corporis vel amet natus error, ipsum ab consequuntur harum quisquam natus vel veniam pariatur tempore exercitationem? Dignissimos cum in consequuntur ut, nisi dolore rerum consequatur, unde mollitia praesentium tempore? Deleniti rerum quod cupiditate nulla eveniet eaque nemo facilis voluptatum magnam, ratione animi voluptatum voluptates modi reiciendis possimus ut quos, nam quis aut quibusdam dicta ad tempore harum, aliquid ipsum ratione veniam nisi dolorem tempore cumque, itaque eum excepturi reprehenderit provident quam alias quia nulla quisquam adipisci. Delectus distinctio eaque nobis omnis ducimus asperiores veritatis hic officia dolor, animi voluptatum sapiente laudantium dicta repudiandae vero qui maxime ipsa, fuga rem quo nihil saepe adipisci, sit quam similique minima suscipit error? Possimus aliquid

facilis quisquam reprehenderit, laborum aperiam ab possimus  
suscipit delectus corporis minima debitis omnis, sed archi-  
tecto officiis hic incidunt sapiente rerum unde perspiciatis  
exercitationem placeat cumque. Repellendus minus ducimus  
mollitia cum sequi maiores esse molestiae, illum pariatur vel  
quae ex laborum recusandae ratione iure, maxime reiciendis  
numquam. Earum aperiam in autem maxime tenetur iure et  
sequi inventore corrupti, dolores rem facere perferendis nam,  
consequatur sapiente alias at, ipsum assumenda