

A CAUSAL MODEL FOR NR-/FR-IQA

We provide here the causal models for the NR-IQA and FR-IQA settings in Figure 8

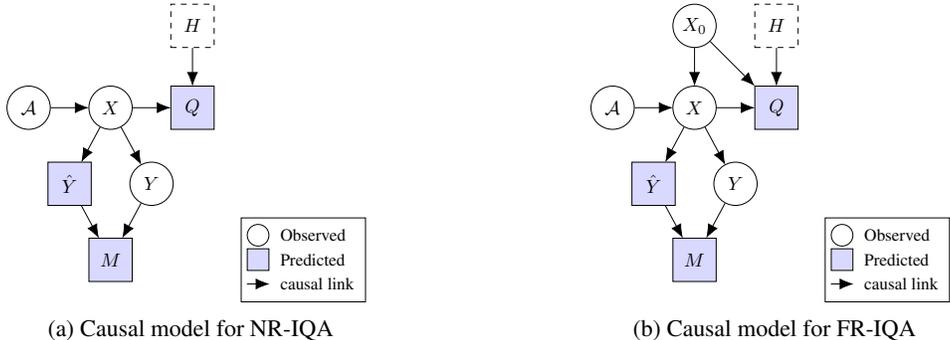


Figure 8: Causal model for the NR-/FR-IQA settings. The FR-IQA setting includes  $X_0$  which is the reference image. In both models,  $H$  indicates human annotator guidance which reflects that IQ metrics are typically calibrated against human perceptual judgements (dashed-line box indicates  $H$  is not used directly in the calculation of  $Q$ ).

While the main manuscript focused on the NR-IQA setting, we can see from the causal models here that the results generalize to the FR-IQA case as well. In particular, the independence of  $Q, M$  given  $X$  is not affected by whether a “clean” reference image ( $X_0$ ) is available for computing  $Q$ . Also, these models also account for the influence of human annotators  $H$  in calibrating the function for computing  $Q$ , but not that this does not change the relationship between  $M, Q$ .

B CAUSAL MODEL FOR IQA WITH LATENT FEATURES

In understanding the difference between the baseline IQA formulation in Figure 1 and the shared features formulation of Figure 2 in Section 3, we provide an expanded version of the baseline DAG in Figure 9. Here we show that the task DNN for computing  $\hat{Y}$  and the function for computing  $Q$  rely on latent features  $Z_{\hat{Y}}$  and  $Z_Q$  respectively.

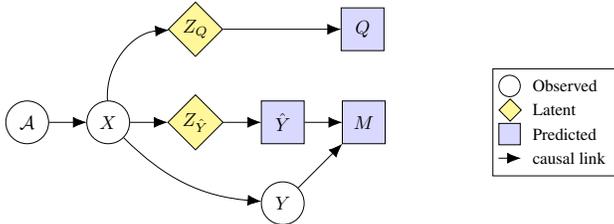


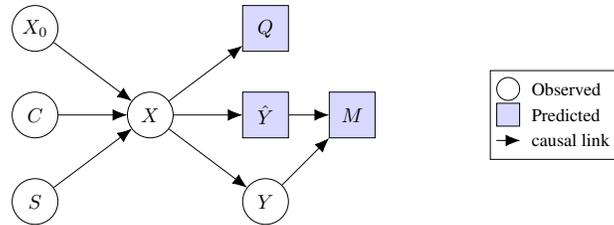
Figure 9: Causal model for IQA that accounts for the use of latent features by the task DNN and IQ metric towards computing  $M$  and  $Q$  respectively.

In this expanded model,  $Z_{\hat{Y}}$  and  $Z_Q$  are independent given  $X$  and not “shared”, and therefore  $Q \perp M \mid X$  as discussed in §3. In contrast, Figure 2 considers the case where  $Z$  represents the features derived from  $X$  that are common between  $Z_{\hat{Y}}$  and  $Z_Q$  shown in the baseline case above. Thus, Figure 2 shows the case where  $X$  does not block all paths between  $Q, M$  since a path exists from  $Q$  to  $M$  through  $Z$ . This ensures that  $Q$  and  $M$  will be correlated given  $X$  unlike in the baseline case above.

C CAUSAL MODEL FOR COMMON CORRUPTIONS ROBUSTNESS EVALUATION

The common corruptions framework (Hendrycks & Dietterich, 2019) is used in our experiments to ensure full control of the image distortion types and severity. Figure 10 shows a version of the

756 baseline IQA causal model customized to account for the corruption process used by this evaluation  
 757 framework. Here, the corrupted image  $X$  is determined by the corruption function (e.g., Gaussian  
 758 noise, defocus blur, fog, contrast, brightness, JPEG compression), the severity ( $S \in \{1, 2, 3, 4, 5\}$ ),  
 759 and the “clean” image  $X_0$ .

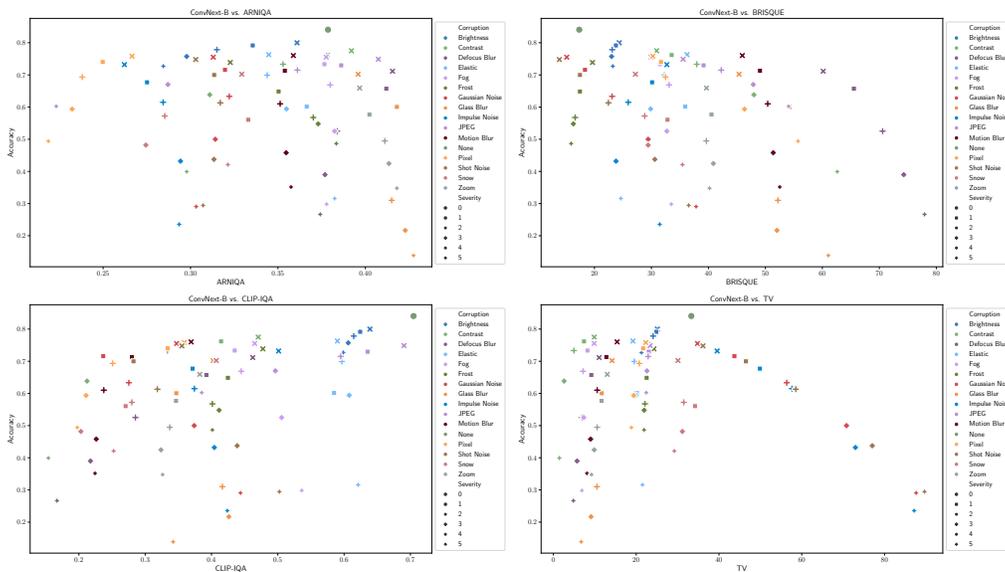


768  
769 Figure 10: Causal model for the common corruptions framework where  $C$  refers to the corruption  
 770 type,  $S$  refers to the corruption severity, and  $X_0$  is the unperturbed, “clean” image.

771  
772 In this setting, we see that  $C, S$  replace the original set of imaging factors  $\mathcal{A}$  in the graph in Figure 1.  
 773 As such, the analysis from §3 holds in the common corruptions framework and allows us to study the  
 774 relationship between  $Q, M$  in a setting where we can precisely control the imaging conditions.

## 775 D RELATIONSHIP OF NR-IQA AND DNN PERFORMANCE METRICS

776  
777 In Figures 11, 12, and 13 we show the relationship of additional NR-IQA metrics with DNN  
 778 performance for additional architectures and metrics. In general, we see weak trends in accuracy vs.  
 779 average IQ suggesting that these metrics are most consistent with the causal model in Figure 1.



783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799 Figure 11: Comparison of ConvNext-B accuracy with (clockwise) ARNIQA, BRISQUE, CLIP-IQA,  
 800 and TV. Little correlation is observed between group-wise accuracy and each NR-IQA metric.

## 801 E RELATIONSHIP OF STRONG TASK-GUIDED IQA AND DNN PERFORMANCE METRICS

802  
803  
804  
805 In Figures 14, 15, 16, we examine the relationship between DNN performance and the strong task-  
 806 guided metrics ( $Q_p, Q_h, Q_l$ ) described in §5. Each figure pairs the task DNN under consideration  
 807 with a pre-trained task model used to compute the quality metric.  
 808  
809

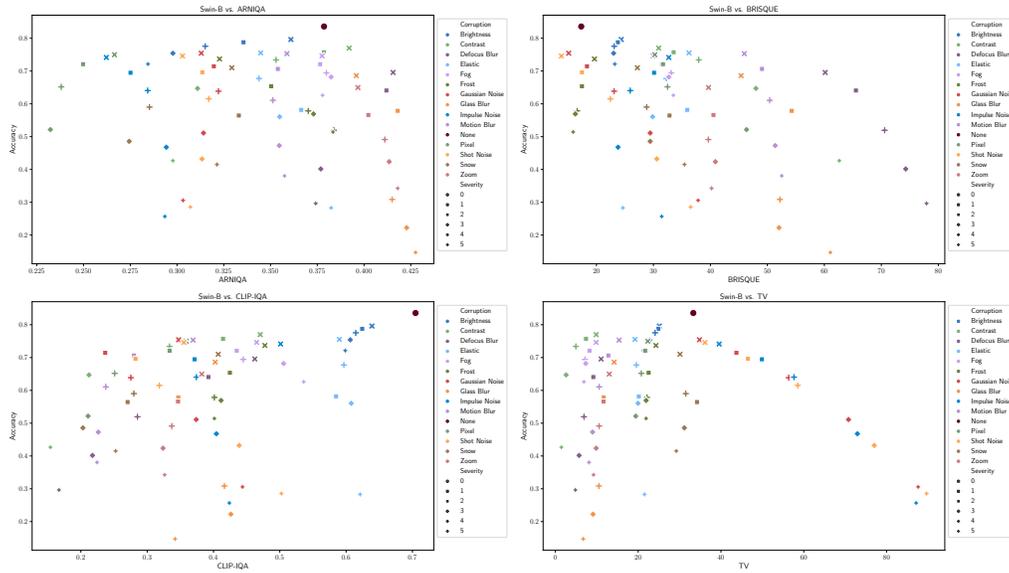


Figure 12: Comparison of Swin-B accuracy with (clockwise) ARNIQA, BRISQUE, CLIP-IQA, and TV. Little correlation is observed between group-wise accuracy and each NR-IQA metric.

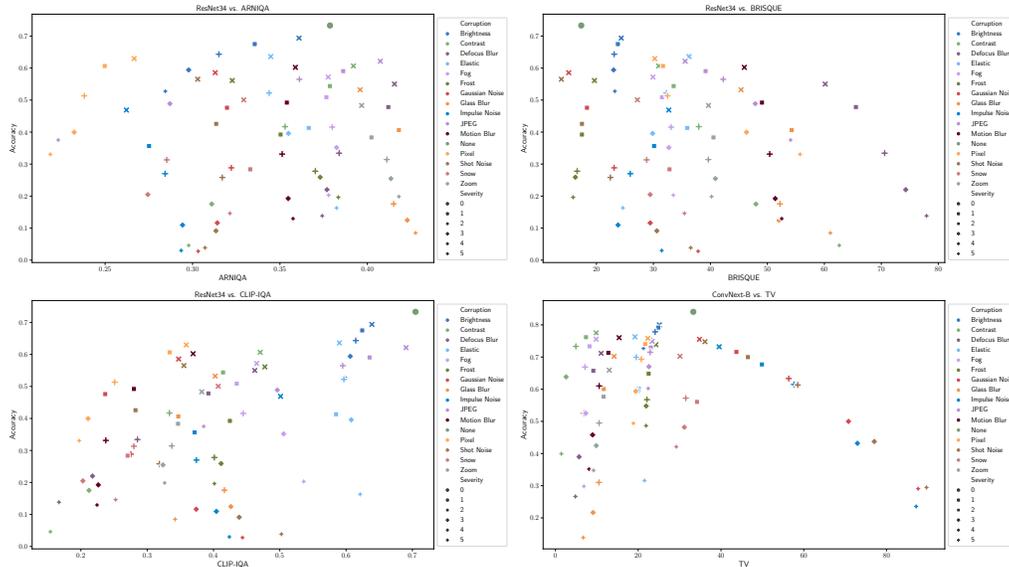


Figure 13: Comparison of ResNet34 accuracy with (clockwise) ARNIQA, BRISQUE, CLIP-IQA, and TV. Little correlation is observed between group-wise accuracy and each NR-IQA metric.

Table 4 also shows the point-wise predictability results for the strong task-guided IQA case. This table extends Table 2 for additional task DNNs. Results here show that strong task-guided IQA metrics are highly correlated with DNN performance and that predictability remains high regardless of whether the pre-trained DNN used to compute  $Q$  is the same DNN used to obtain  $M$ .

## F RELATIONSHIP OF WEAK TASK-GUIDED IQA AND DNN PERFORMANCE METRICS

We provide Figures 17, 18, 19 showing the relationship between DNN performance the weak task-guided ZSCLIP-IQA metric from §6.

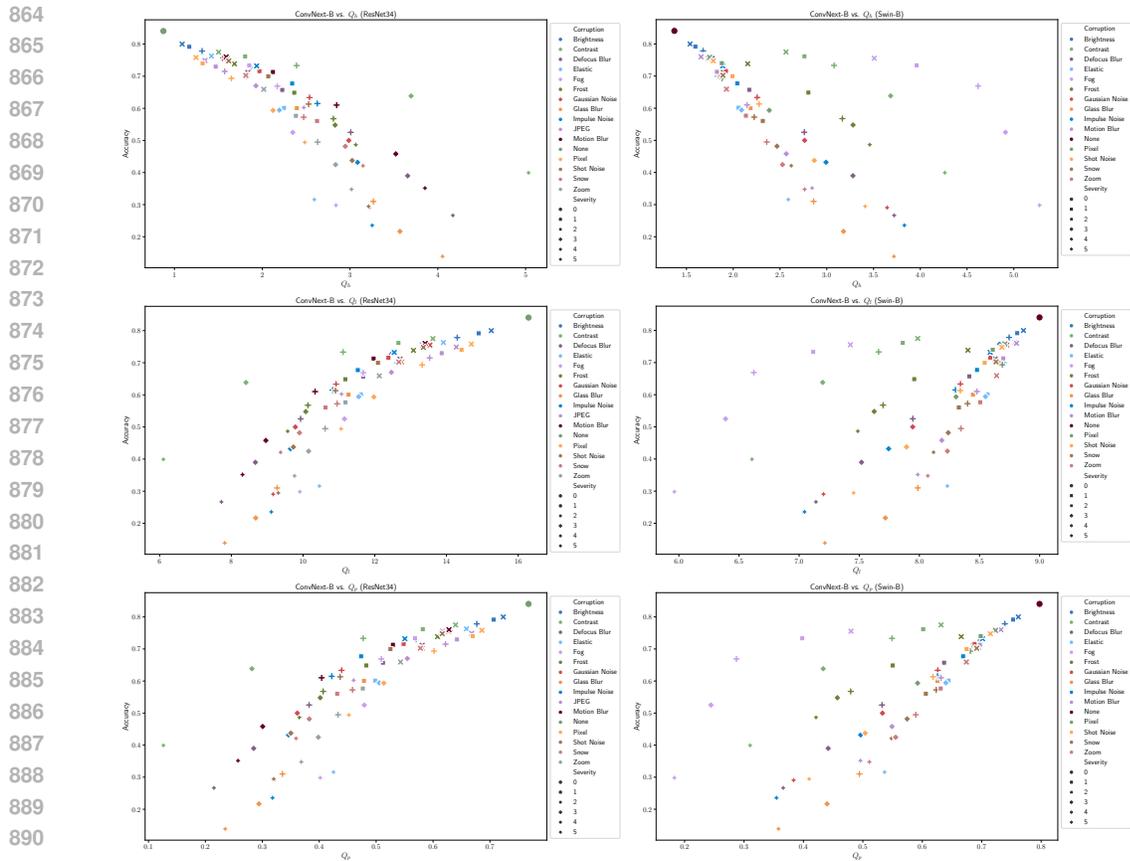


Figure 14: Comparison of ConvNext-B accuracy with (row)  $Q_h, Q_l, Q_p$  computed using (col) ResNet34, Swin-B. High correlation is observed between each IQA metric and accuracy.

## G CONTROLLING FOR IMAGE CONTENT WHEN EVALUATING PREDICTABILITY

The experiments of §4.6 examined predictability by modeling  $P(M|Q, X)$ . Since the content of  $X$  may be a confounder for both  $M, Q$ , we attempt to control for it in two ways.

In the first case, we take advantage of the synthetic nature of the IN-C dataset by looking at the predictability of  $M$  from  $Q$  for each image in the dataset separately which allows us to control the content precisely and only change the quality characteristics. For this experiment, we train a logistic regression classifier to predict  $P(M|Q)$  for individual image IDs trained using only  $M, Q$  computed from the set of distorted variants of each specific image ID. Given the original clean image and 15 corruptions with 5 severity levels each, we run 5-fold cross-validation with an 80/20 train/test split of the 76 total images. We repeat this for all 50k image IDs in the ImageNet validation set. The results shown in Figure 20a are an average of the AUC over all image IDs and folds.

We see that even when controlling for the image content the weak task-guided IQA generally achieves the highest  $mAUC$  with the lowest variance. Overall, this supports our hypothesis and causal analysis that weak task-guidance provides a means to associate  $M, Q$  even when conditioning on the image directly.

In the second case, we adjust for image content by controlling for the image label  $Y$ . Here, we train a separate classifier to model  $P(M|Q)$  for each of the 1000 labels in the ImageNet dataset. Each classifier is trained on the aggregate of 50 images per label along with all 15 corruptions at 5 severity levels (3751 total per label). We again use an 80/20 train/test split and perform 5-fold cross-validation. The results shown in Figure 20b are an average of the  $AUC$  over all labels and folds.

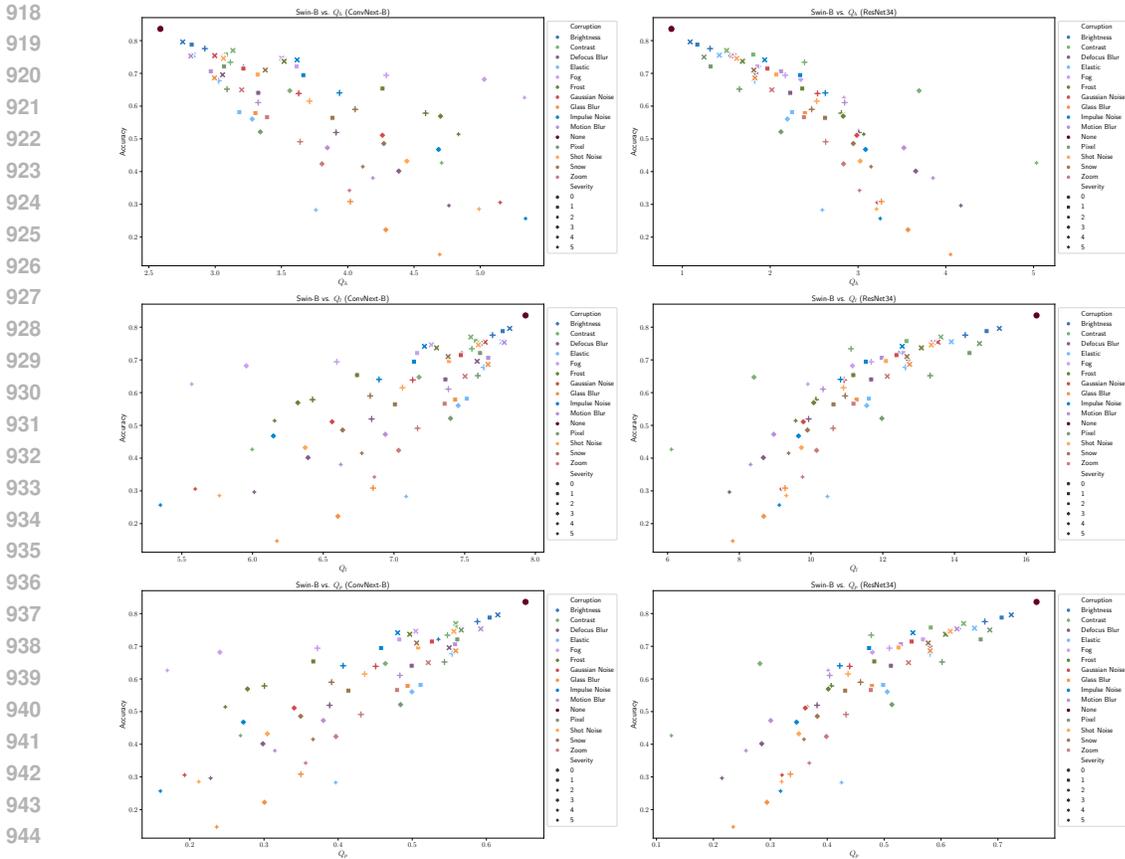


Figure 15: Comparison of Swin-B accuracy with (row)  $Q_h, Q_l, Q_p$  computed using (col) ConvNext-B, ResNet34. High correlation is observed between each IQA metric and accuracy.

We find here that across all IQ metrics evaluated,  $mAUC$  is barely above chance. For the traditional NR-IQA metrics, this supports our analysis and main experiments which show little correlation between  $M, Q$ . For the ZSCLIP-IQA metric, we refer again to the causal model (Fig. 6) and see that while the task selection variable ensures the association between  $M, Q$ , it is the conditioning on  $Y$  (and  $X$ ), as done here, that blocks all paths between  $M, Q$  and once again removes the association.

## H PREDICTABILITY OF DNN PERFORMANCE FOR MILDLY CORRUPTED DATASETS

To show that  $D1$  is satisfied even in the case of mildly corrupted data, we plot the distributions of  $Q$  in Figure 21. Across all variants, even in cases where the likelihood is low, each IQA metric exhibits sensitivity to corruption (D1).

While some IQA metrics are more sensitive to the overall image corruption, this does not necessarily translate to higher predictability. In fact, ZSCLIP-IQA appears to have smaller differences in IQA distribution across variants compared to other IQA metrics, yet the highest predictability of  $M$ . Figure 22 shows the predictability of  $M$  from  $Q$  for variants of the IN-C benchmark created as described in §6.3

Results show that weak task-guided IQA metrics are able to achieve higher  $AUC$  even when the number of corrupted images in the dataset is low. In comparison, conventional NR-IQA metrics achieve lower  $AUC$  and are more sensitive to the total level of corruption.

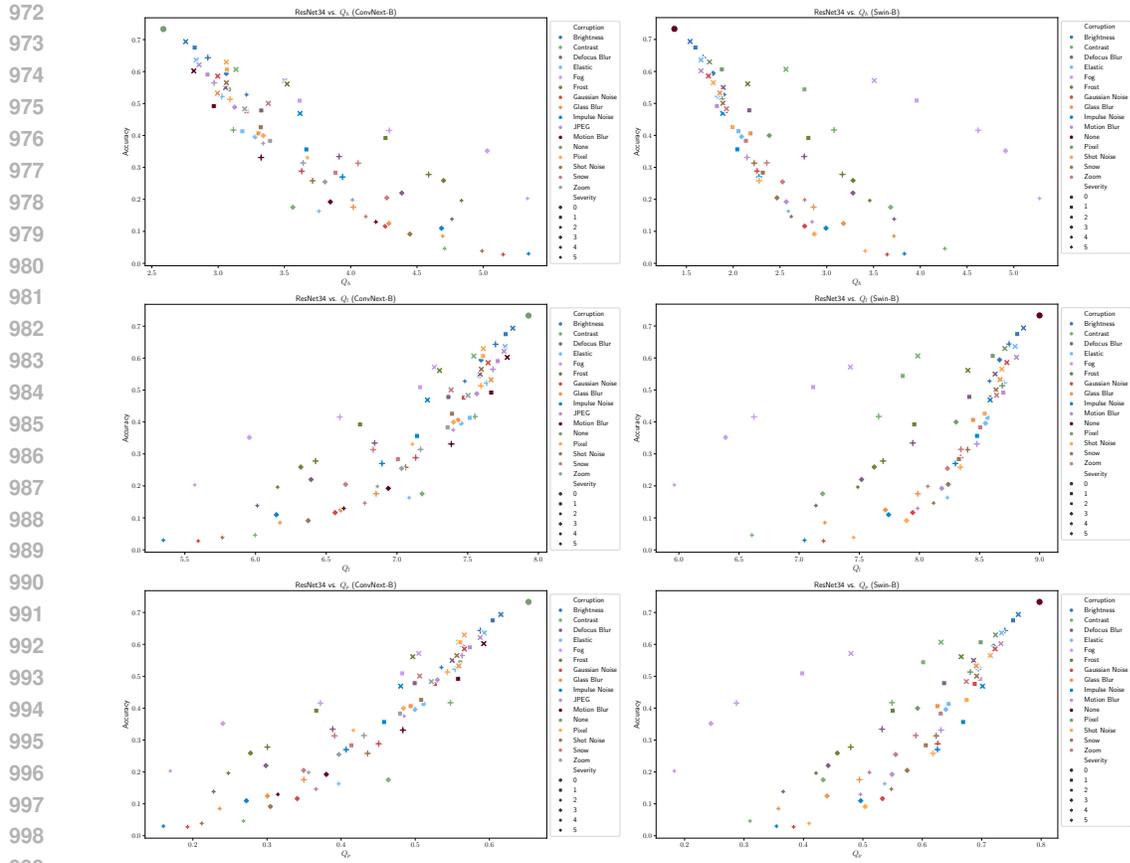


Figure 16: Comparison of ResNet34 accuracy with (row)  $Q_h, Q_l, Q_p$  computed using (col) ConvNext-B, Swin-B. High correlation is observed between each IQA metric and accuracy.

## I COMPUTE RESOURCES

All experiments were run using a single NVIDIA A40 GPU with 48GB of memory. Predictability analysis can be conducted on CPU-only machine with at least 8 cores.

1026  
 1027  
 1028  
 1029  
 1030  
 1031  
 1032  
 1033  
 1034  
 1035  
 1036  
 1037  
 1038  
 1039  
 1040  
 1041  
 1042  
 1043  
 1044  
 1045  
 1046  
 1047  
 1048  
 1049  
 1050  
 1051  
 1052  
 1053  
 1054  
 1055  
 1056  
 1057  
 1058  
 1059  
 1060  
 1061  
 1062  
 1063  
 1064  
 1065  
 1066  
 1067  
 1068  
 1069  
 1070  
 1071  
 1072  
 1073  
 1074  
 1075  
 1076  
 1077  
 1078  
 1079

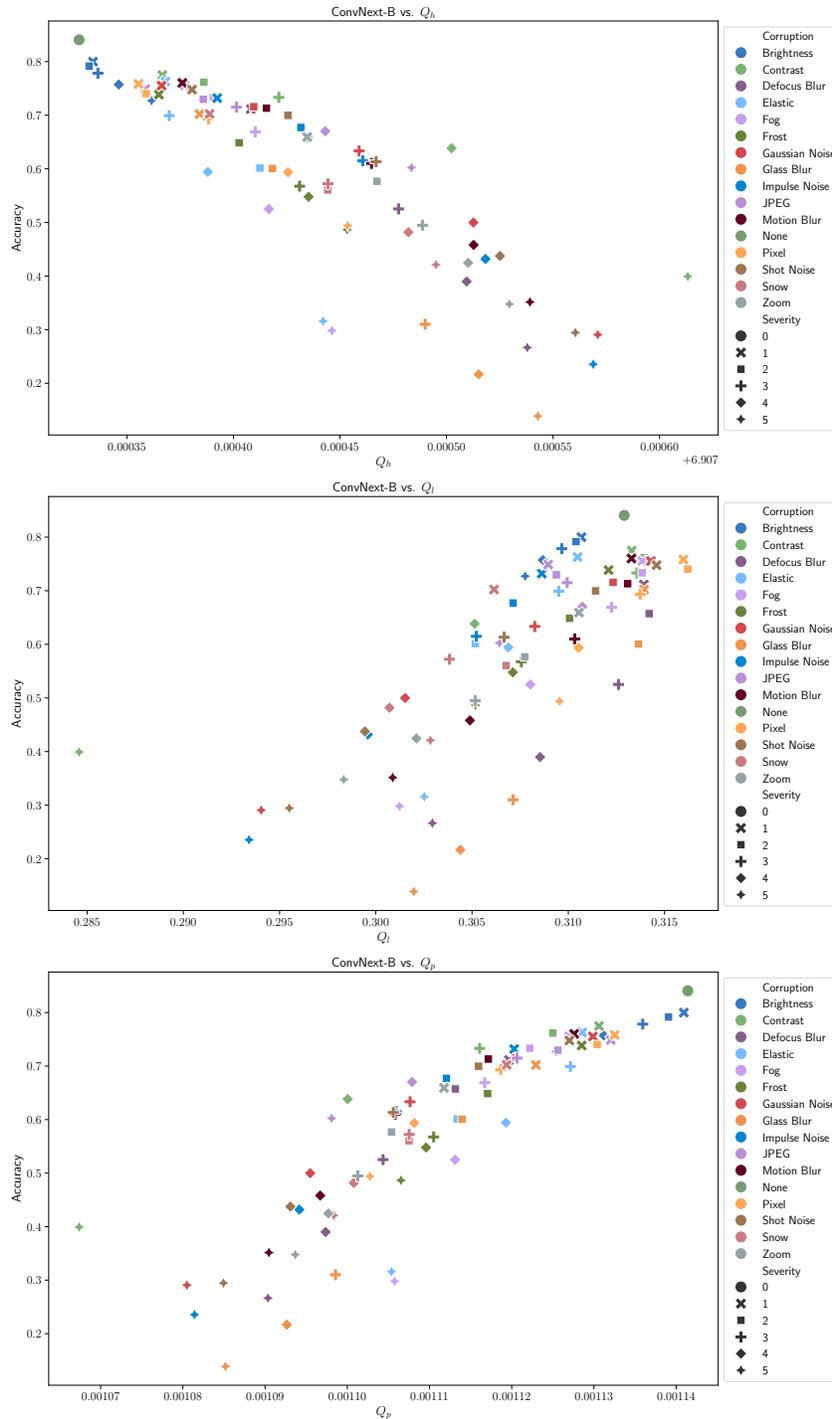


Figure 17: Comparison of ConvNext-B accuracy with (top to bottom)  $Q_h$ ,  $Q_l$ ,  $Q_p$  based on ZSCLIP-IQA. High correlation is observed between each ZSCLIP-IQA variant and accuracy.

1080  
 1081  
 1082  
 1083  
 1084  
 1085  
 1086  
 1087  
 1088  
 1089  
 1090  
 1091  
 1092  
 1093  
 1094  
 1095  
 1096  
 1097  
 1098  
 1099  
 1100  
 1101  
 1102  
 1103  
 1104  
 1105  
 1106  
 1107  
 1108  
 1109  
 1110  
 1111  
 1112  
 1113  
 1114  
 1115  
 1116  
 1117  
 1118  
 1119  
 1120  
 1121  
 1122  
 1123  
 1124  
 1125  
 1126  
 1127  
 1128  
 1129  
 1130  
 1131  
 1132  
 1133

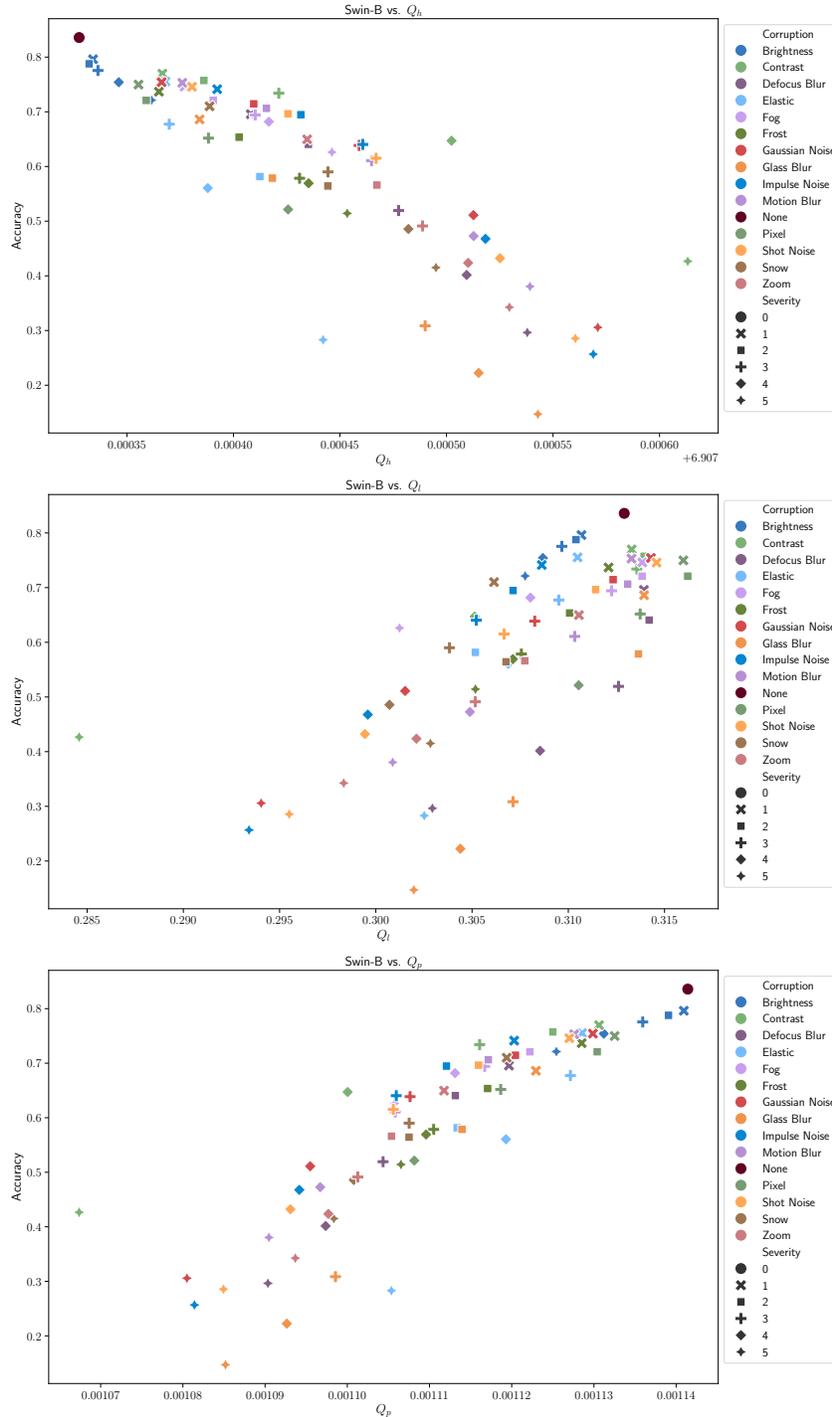


Figure 18: Comparison of Swin-B accuracy with (top to bottom)  $Q_h$ ,  $Q_l$ ,  $Q_p$  based on ZSCLIP-IQA. Some correlation is observed between each ZSCLIP-IQA variant and accuracy.

1134  
 1135  
 1136  
 1137  
 1138  
 1139  
 1140  
 1141  
 1142  
 1143  
 1144  
 1145  
 1146  
 1147  
 1148  
 1149  
 1150  
 1151  
 1152  
 1153  
 1154  
 1155  
 1156  
 1157  
 1158  
 1159  
 1160  
 1161  
 1162  
 1163  
 1164  
 1165  
 1166  
 1167  
 1168  
 1169  
 1170  
 1171  
 1172  
 1173  
 1174  
 1175  
 1176  
 1177  
 1178  
 1179  
 1180  
 1181  
 1182  
 1183  
 1184  
 1185  
 1186  
 1187

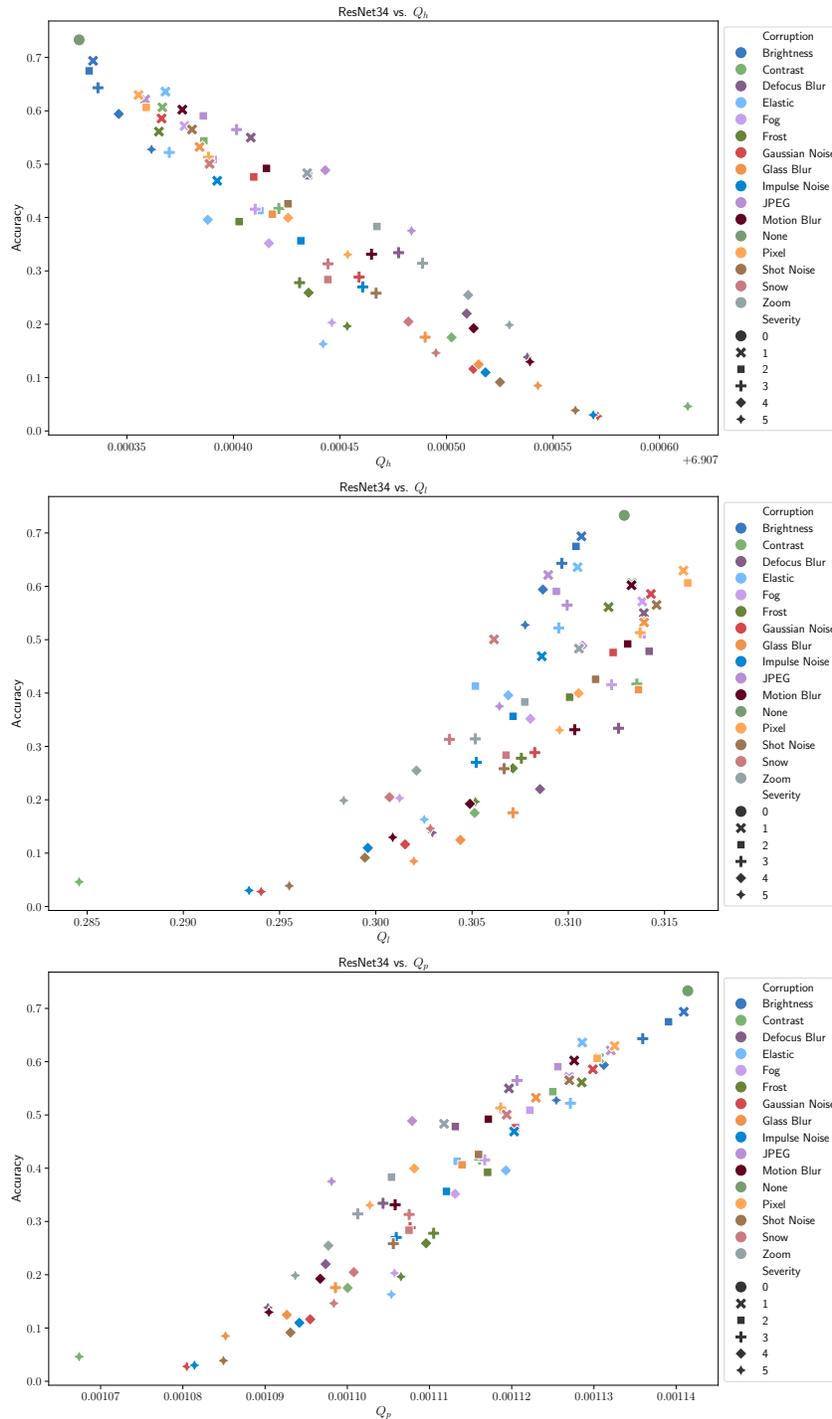


Figure 19: Comparison of ResNet34 accuracy with (top to bottom)  $Q_h$ ,  $Q_l$ ,  $Q_p$  based on ZSCLIP-IQA. High correlation is observed between each ZSCLIP-IQA variant and accuracy.

Table 4: Correlation between IQ and accuracy. SRCC, PLCC computed using average accuracy for each (corruption, severity). AUC and CE based on point-wise predictions (95% CI within  $\pm 0.001$ ). SRCC, PLCC values have  $p < 0.05$ .

Model	IQA Metric	AUC $\uparrow$	CE $\downarrow$	$PLCC \uparrow$	$SRCC \uparrow$
ConvNext-B	ConvNext-B $Q_h$	0.772	0.562	$0.822 \pm 0.070$	$0.854 \pm 0.063$
	ConvNext-B $Q_l$	0.778	0.555	$0.826 \pm 0.067$	$0.854 \pm 0.063$
	ConvNext-B $Q_p$	0.826	0.504	$0.888 \pm 0.045$	$0.910 \pm 0.044$
	ResNet34 $Q_h$	0.725	0.601	$0.859 \pm 0.069$	$0.924 \pm 0.045$
	ResNet34 $Q_l$	0.717	0.603	$0.866 \pm 0.051$	$0.925 \pm 0.043$
	ResNet34 $Q_p$	0.719	0.601	$0.875 \pm 0.051$	$0.926 \pm 0.044$
	Swin-B $Q_h$	0.760	0.579	$0.624 \pm 0.165$	$0.724 \pm 0.170$
	Swin-B $Q_l$	0.724	0.601	$0.604 \pm 0.166$	$0.686 \pm 0.176$
	Swin-B $Q_p$	0.791	0.547	$0.742 \pm 0.132$	$0.797 \pm 0.139$
ResNet34	ConvNext-B $Q_h$	0.767	0.557	$0.858 \pm 0.060$	$0.889 \pm 0.055$
	ConvNext-B $Q_l$	0.760	0.563	$0.853 \pm 0.059$	$0.896 \pm 0.055$
	ConvNext-B $Q_p$	0.801	0.522	$0.904 \pm 0.044$	$0.930 \pm 0.041$
	ResNet34 $Q_h$	0.848	0.470	$0.930 \pm 0.028$	$0.969 \pm 0.023$
	ResNet34 $Q_l$	0.827	0.492	$0.951 \pm 0.015$	$0.973 \pm 0.020$
	ResNet34 $Q_p$	0.850	0.461	$0.960 \pm 0.015$	$0.977 \pm 0.021$
	Swin-B $Q_h$	0.751	0.574	$0.643 \pm 0.153$	$0.774 \pm 0.140$
	Swin-B $Q_l$	0.709	0.600	$0.628 \pm 0.146$	$0.754 \pm 0.145$
	Swin-B $Q_p$	0.774	0.551	$0.747 \pm 0.129$	$0.825 \pm 0.112$
Swin-B	ConvNext-B $Q_h$	0.744	0.586	$0.706 \pm 0.129$	$0.768 \pm 0.098$
	ConvNext-B $Q_l$	0.746	0.586	$0.709 \pm 0.127$	$0.768 \pm 0.099$
	ConvNext-B $Q_p$	0.791	0.542	$0.788 \pm 0.102$	$0.834 \pm 0.079$
	ResNet34 $Q_h$	0.722	0.603	$0.828 \pm 0.078$	$0.896 \pm 0.053$
	ResNet34 $Q_l$	0.713	0.604	$0.831 \pm 0.061$	$0.892 \pm 0.053$
	ResNet34 $Q_p$	0.716	0.602	$0.845 \pm 0.062$	$0.897 \pm 0.052$
	Swin-B $Q_h$	0.766	0.578	$0.483 \pm 0.207$	$0.654 \pm 0.174$
	Swin-B $Q_l$	0.732	0.597	$0.458 \pm 0.203$	$0.611 \pm 0.181$
	Swin-B $Q_p$	0.807	0.529	$0.620 \pm 0.184$	$0.732 \pm 0.142$

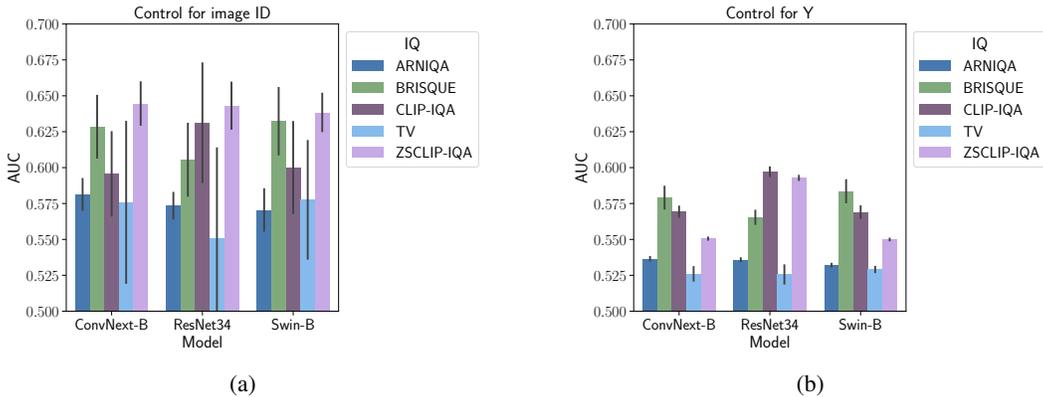


Figure 20: Mean AUC (mAUC) for classifiers trained (a) per image ID ( $X_i$ ) to model  $P(M|Q, X_i)$  and (b) per label  $Y$  to model  $P(M|Q, Y)$ . Averages are taken over all images/labels respectively and cross-validation folds with error bars indicating one standard deviation. Higher variance in (a) is attributed to lower sample sizes for training the logistic regression classifier.

1242  
 1243  
 1244  
 1245  
 1246  
 1247  
 1248  
 1249  
 1250  
 1251  
 1252  
 1253  
 1254  
 1255  
 1256  
 1257  
 1258  
 1259  
 1260  
 1261  
 1262  
 1263  
 1264  
 1265  
 1266  
 1267  
 1268  
 1269  
 1270  
 1271  
 1272  
 1273  
 1274  
 1275  
 1276  
 1277  
 1278  
 1279  
 1280  
 1281  
 1282  
 1283  
 1284  
 1285  
 1286  
 1287  
 1288  
 1289  
 1290  
 1291  
 1292  
 1293  
 1294  
 1295

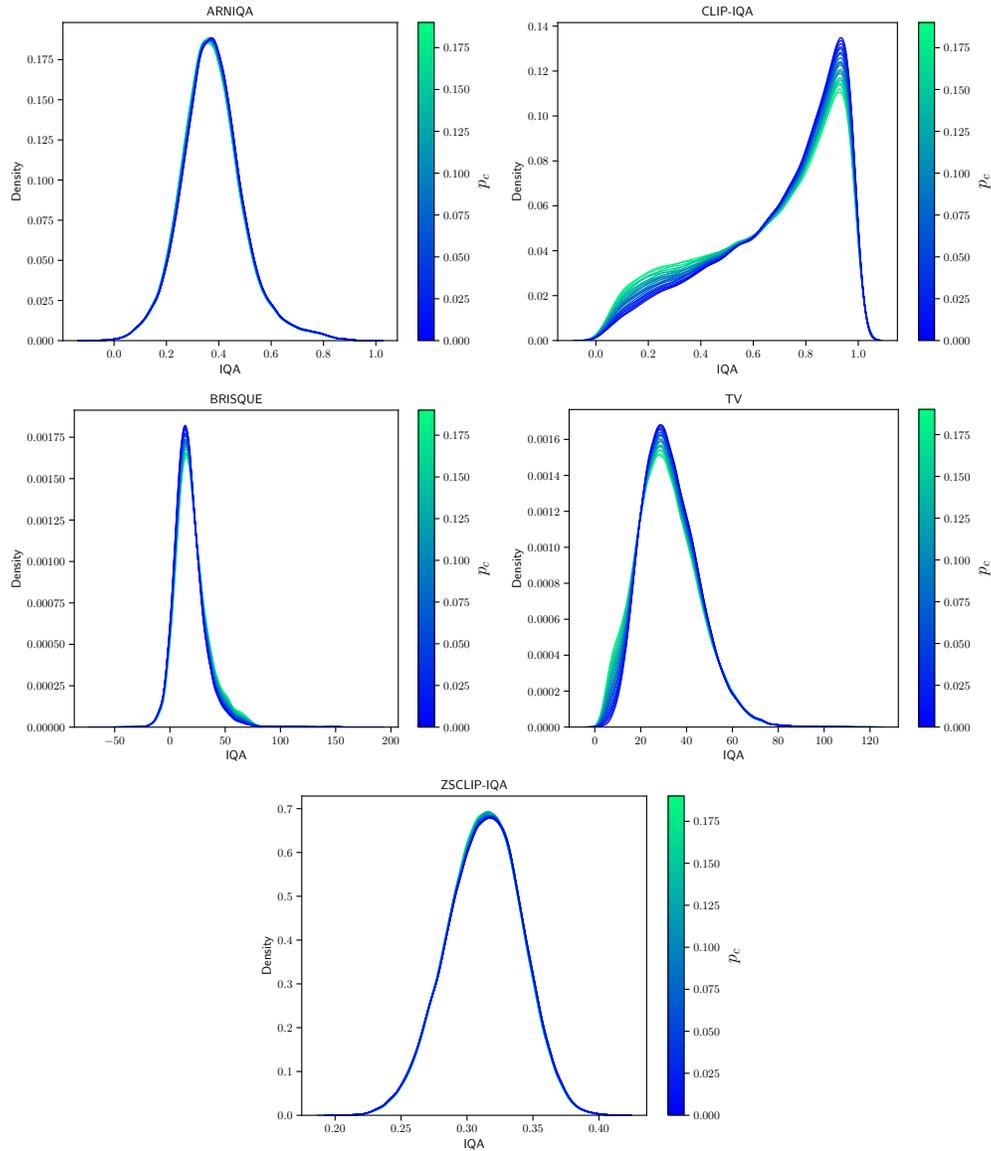


Figure 21: Distribution of IQA for each mildly corrupted variant of IN-val. Line color indicates the likelihood of image corruption  $p_c$  for each variant. Note that the amount of similarity/difference in the IQ distribution across variants does not explain the predictability which is determined by the causal DAG such as in §3, §5, and §6. See Figure 22 for predictability results.

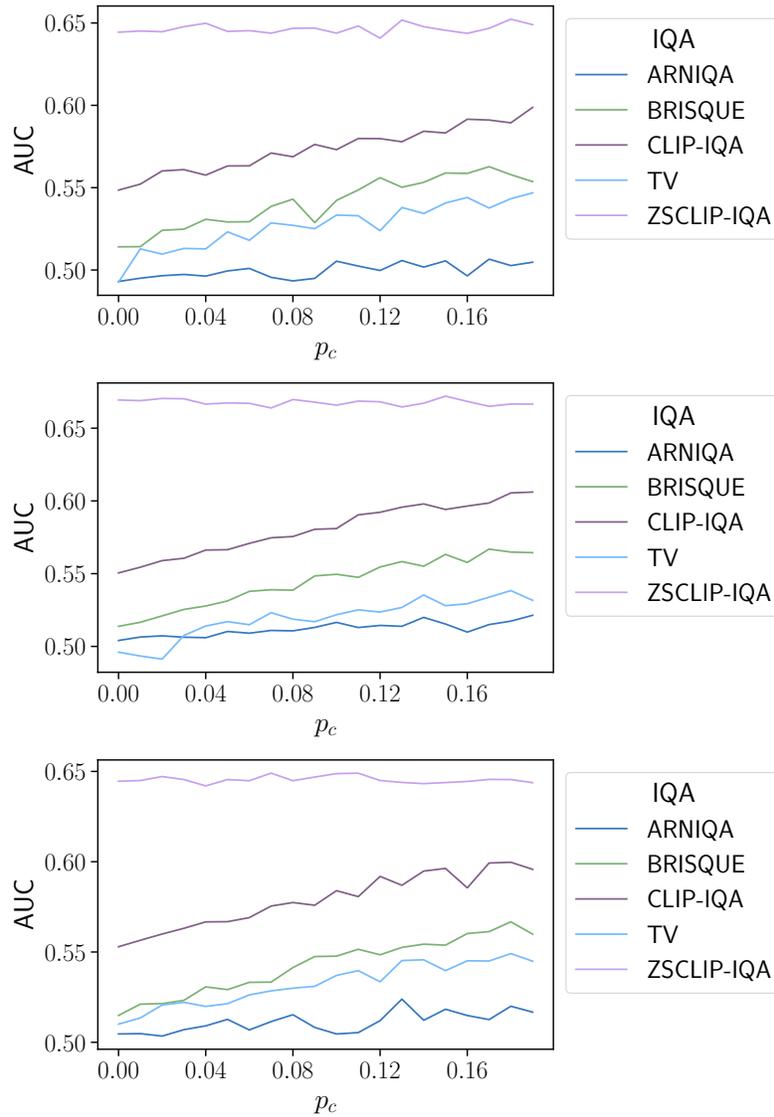


Figure 22: AUC vs.  $p_c$  where  $p_c$  represents the fraction of images in the test set that are mildly corrupted. Results are listed top to bottom: ConvNet-B, ResNet34, Swin-B. Predictability for ZSCLIP-IQA is relatively insensitive to the proportion of corrupted images whereas other metrics only improve as the proportion and diversity of corruptions increases.