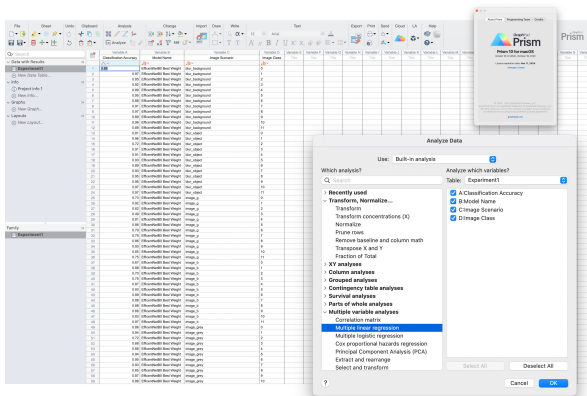# XIMAGENET-12: An Explainable Visual Benchmark Dataset for Model Robustness Evaluation

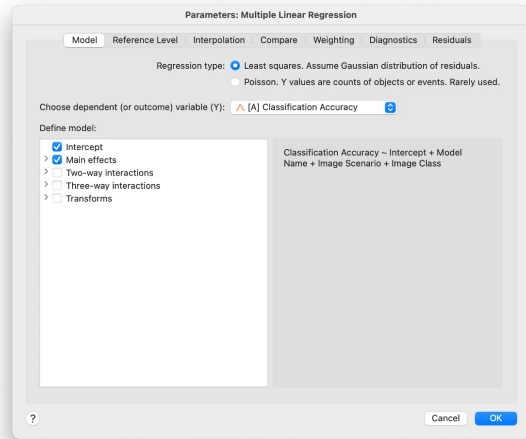## Supplementary Material

## 7. Appendix

In this section we provide the supplementary compiled together with the main paper includes:
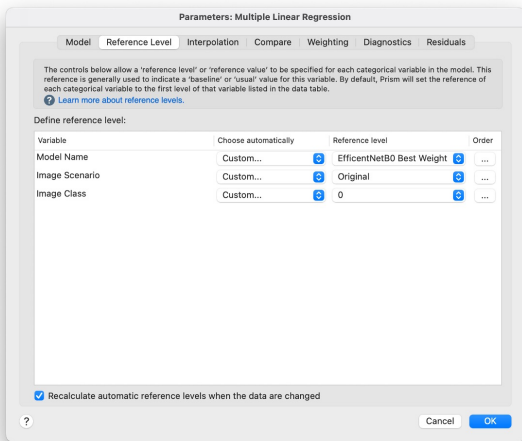
- The illustration of how we use Multiple Linear Regression to verify our hypothesis: from raw data input, for example, in GraphPad Prism, to interpreting examples and residual plots, etc;
- The training details (Accuracy and Loss Plots) and hyperparameters within scenarios and cross-scenario experiments, including diffusion metrics for evaluations, density maps of State-of-the-Art accuracy drop (e.g., referring to our particular experiment, EX1, EX2);
- The ablation study addresses the industry pain points, illustrating robust model selection for challenging scenarios, particularly due to factors such as background variations, camera shifts, color changes, and lighting conditions;
- Sample image of our XIMAGENEt-12 AI-generated image, comprising 12,248 images for AI-generated scenarios using the latest Stable Diffusion XL model and involving uniform promotion and manual selection and filtering.
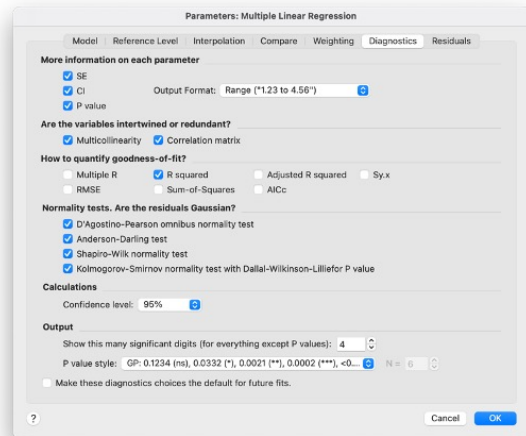
Define the raw data type and variable into statistic software (GraphPad Prism)
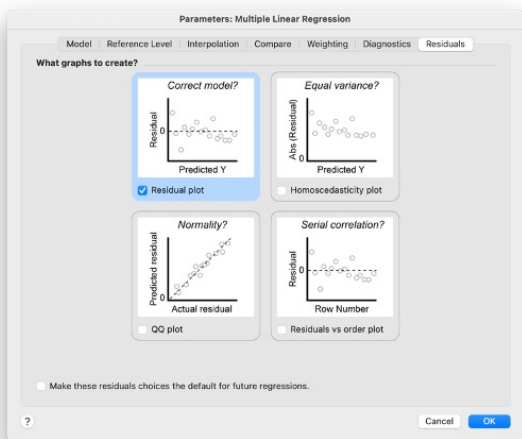


Choose Regression type and define the base independent variables (Model, Image Scenario, Image Class)
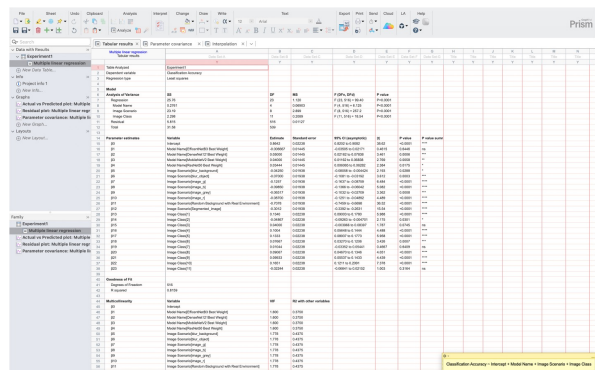


Select reference level for each independent variables (Model, Image Scenario, Image Class)



Set Parameters for Multiple Liner Regression, such as Confidence Level etc



Create target residual plot graph for simulating the regression results



Generate the analyse and interpretation report includes Estimates and P Value for each variables

Figure 8. Multiple Linear Regression Workflow and Example of Interpretations.

| Parameter estimates | Variable | Estimate | Standard error | 95% CI (asymptotic) | |t| | P value | P value summ |
|---|---|---|---|---|---|---|---|
| β0 | Intercept | 0.8642 | 0.02238 | 0.8202 to 0.9082 | 38.62 | <0.0001 | **** |
| β1 | Model Name[EfficentNetB3 Best Weight] | -0.006667 | 0.01445 | -0.03505 to 0.02171 | 0.4615 | 0.6446 | ns |
| β2 | Model Name[DenseNet121Best Weight] | 0.05000 | 0.01445 | 0.02162 to 0.07838 | 3.461 | 0.0006 | *** |
| β3 | Model Name[MobileNetV2 Best Weight] | 0.04000 | 0.01445 | 0.01162 to 0.06838 | 2.769 | 0.0058 | ** |
| β4 | Model Name[ResNet50 Best Weight] | 0.03444 | 0.01445 | 0.006065 to 0.06282 | 2.384 | 0.0175 | * |
| β5 | Image Scenario[blur_background] | -0.04250 | 0.01938 | -0.08058 to -0.004424 | 2.193 | 0.0288 | * |
| β6 | Image Scenario[blur_object] | -0.07000 | 0.01938 | -0.1081 to -0.03192 | 3.612 | 0.0003 | *** |
| β7 | Image Scenario[image_g] | -0.1257 | 0.01938 | -0.1637 to -0.08759 | 6.484 | <0.0001 | **** |
| β8 | Image Scenario[image_b] | -0.09850 | 0.01938 | -0.1366 to -0.06042 | 5.082 | <0.0001 | **** |
| β9 | Image Scenario[image_grey] | -0.06517 | 0.01938 | -0.1032 to -0.02709 | 3.362 | 0.0008 | *** |
| β10 | Image Scenario[image_r] | -0.08700 | 0.01938 | -0.1251 to -0.04892 | 4.489 | <0.0001 | **** |
| β11 | Image Scenario[Random Background with Real Environment] | -0.7078 | 0.01938 | -0.7459 to -0.6698 | 36.52 | <0.0001 | **** |
| β12 | Image Scenario[Segmented_image] | -0.3012 | 0.01938 | -0.3392 to -0.2631 | 15.54 | <0.0001 | **** |
| β13 | Image Class[1] | 0.1340 | 0.02238 | 0.09003 to 0.1780 | 5.988 | <0.0001 | **** |
| β14 | Image Class[2] | -0.04867 | 0.02238 | -0.09263 to -0.004701 | 2.175 | 0.0301 | * |
| β15 | Image Class[3] | 0.04000 | 0.02238 | -0.003966 to 0.08397 | 1.787 | 0.0745 | ns |
| β16 | Image Class[4] | 0.1004 | 0.02238 | 0.05648 to 0.1444 | 4.488 | <0.0001 | **** |
| β17 | Image Class[5] | 0.1333 | 0.02238 | 0.08937 to 0.1773 | 5.958 | <0.0001 | **** |
| β18 | Image Class[6] | 0.07667 | 0.02238 | 0.03270 to 0.1206 | 3.426 | 0.0007 | *** |
| β19 | Image Class[7] | 0.01044 | 0.02238 | -0.03352 to 0.05441 | 0.4667 | 0.6409 | ns |
| β20 | Image Class[8] | 0.09067 | 0.02238 | 0.04670 to 0.1346 | 4.051 | <0.0001 | **** |
| β21 | Image Class[9] | 0.09933 | 0.02238 | 0.05537 to 0.1433 | 4.439 | <0.0001 | **** |
| β22 | Image Class[10] | 0.1651 | 0.02238 | 0.1211 to 0.2091 | 7.378 | <0.0001 | **** |
| β23 | Image Class[11] | -0.02244 | 0.02238 | -0.06641 to 0.02152 | 1.003 | 0.3164 | ns |

Figure 9. Examples of Multiple Linear Regression Interpretations: (1) $\beta 0$ (Intercept) estimate equal to 0.8642 means that the base classification accuracy when all predictors are at their reference levels is 86.42%. (2) $\beta 2$ (Model Name [DenseNet121 Best Weight]) estimate equal to 0.05000, $P$ value 0.0006 means that the model [DenseNet121 Best Weight] increases the classification accuracy by 5.000% when compared to the reference level the model [EfficentNetB0 Best Weight]. This effect is also statically significant ($P$ value <0.05). (3) $\beta 11$ (Image Scenario [Random Background with Real Environment]) estimate equal to -0.7078; $P$ value <0.0001 means that this image scenario decreases the classification accuracy by 70.78% when compared to the reference level the Image Scenario [Original] with a significant confidence.
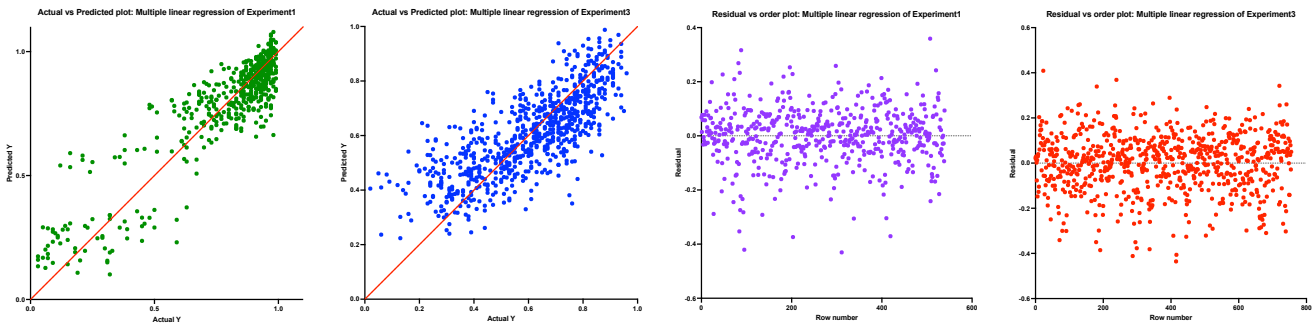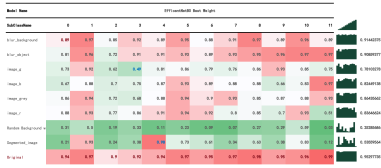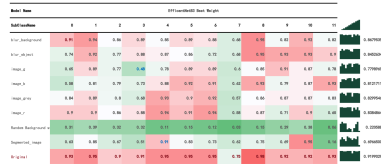


Figure 10. Multiple linear Regression Residual Plot and QQ Plot of Ex1 (SOTA Model Classfication) and Ex3 (SOTA Segmentation) on Hypothesis Verification Process. These shows how two distributions (accuray point)' quantiles line up, with our theoretical distribution (e.g., the normal distribution) as the x variable (Scenarios, Image Object, Model Name) and regression model residuals as the y variable. If the points lie on or close to a 45-degree line, it means that the data follow the reference distribution closely, boosting confidence in the regression results.
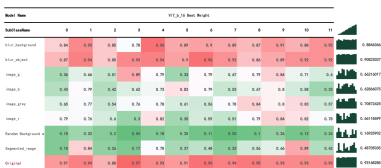
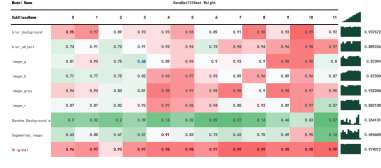The EfficientNetB0 [34] accuracy for each class



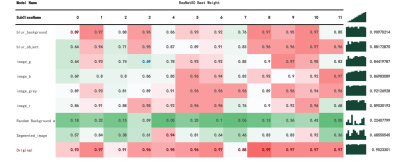The EfficientNetB3 [34] accuracy for each class



The MobileNetV2 [29] accuracy for each class
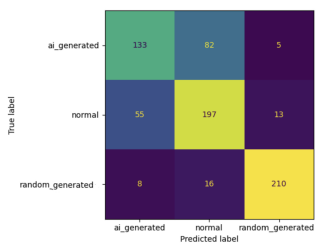


The ViT [8] accuracy for each class



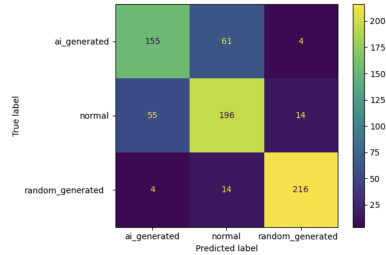The DenseNet121 [13] accuracy for each class
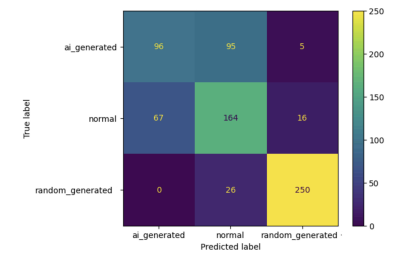


The ResNet50 [11] accuracy for each class

Figure 11. The SOTA models accuracy density map for each class on Experiment 1. Testing images with different background indeed is a challenging scenario for vision models.
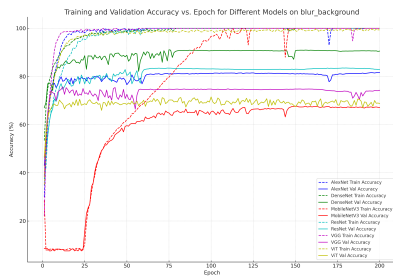


ResNet50 [11] 74.9%
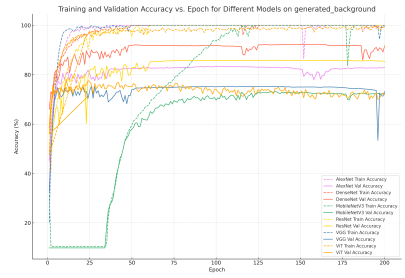


DenseNet [13] 78.8%



MobileNet [29] 70.9%

Figure 12. Model accuracy for classifying normal/AI-generated/random-generated images.
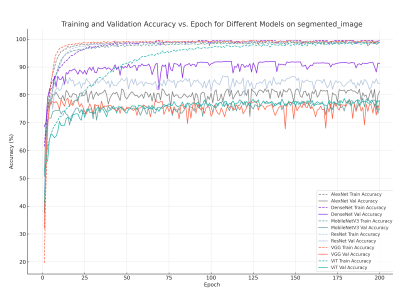
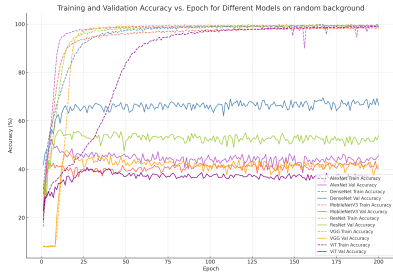Top-1 training and validation accuracy on blurred background scenario

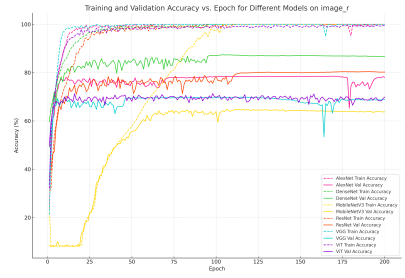Top-1 training and validation accuracy on blurred object scenario

Top-1 training and validation accuracy on AI-generated background scenario

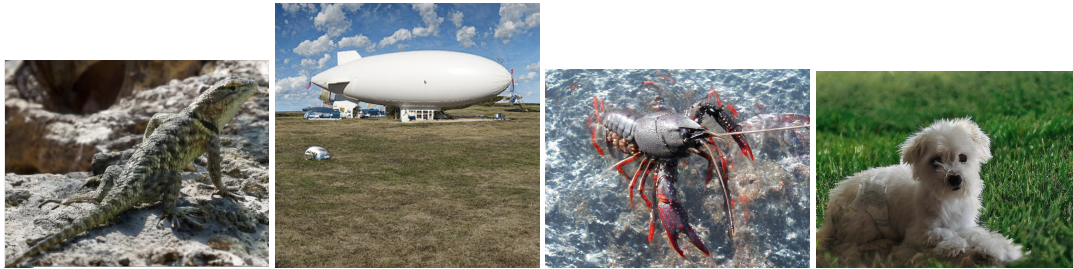Top-1 training and validation accuracy on the segmented scenario

Top-1 training and validation accuracy on random background-generated scenario

Top-1 training and validation accuracy on only one red channel generated scenario

Figure 13. Top-1 training and validation accuracy for SOTA models on Experiment 2: within the same Scenario.

(a) **Prompt**: Generate high-definition pictures like those in the National Geographic magazine, keep the background unchanged.

(b) **Prompt**: Generate a realistic blue sky, and clouds background and please do not change the foreground airship object.

(c) **Prompt**: Generate high resolution images in sea water

(d) **Prompt**: Generate a picture with a foreground and the green grass in the background, similar to the official HD picture released by the state.

(e) **Prompt**: Generate high-resolution pictures like fox in the lawn, National Geographic, keep the background and foreground more simple and real.

(f) **Prompt**: Generate a simple image more realistic in the style of ocean magazine.

(g) **Prompt**: Generate an image with the car in the background and similar to the HD image published by the state.

(h) **Prompt**: Generate high-resolution pictures like those in National Financial Magazine, and keep the background and foreground consistent and the environment more real!

(i) **Prompt**: Please generate high-resolution pictures like those in the National Music Magazine and keep the background unchanged.

(j) **Prompt**: Generate high-resolution pictures in the style of those in National Food Magazine, and keep the background and foreground consistent and the environment more real!

(k) **Prompt**: Generate high-definition pictures like those in the National Geographic magazine, keep the background unchanged.

(l) **Prompt**: Generate high-resolution pictures, such as National Marine Magazine's oceans and whales, to keep the background real.

Figure 14. AI generated images with prompts within XIMAGENET-12 Dataset.