# DiagViB-6: A Diagnostic Benchmark Suite for Vision Models in the Presence of Shortcut and Generalization Opportunities

Elias Eulig    Piyapat Saranrittichai    Chaithanya Kumar Mummadi
Kilian Rambach    William Beluch    Xiahan Shi    Volker Fischer
Bosch Center for Artificial Intelligence (BCAI)

## A. Supplementary material

### A.1. Image generation

We provide a comprehensive overview of all factor classes and respective factor space regions in Tab. A1. The five textures from which the texture crops $f_{\text{text.}}$ are sampled are shown in Fig. A1.

| $\mathcal{F}_i$ | $\mathcal{S}_i$ | $\mathcal{N}_i$ | $\mathcal{C}_{i,j}$ | $\mathcal{S}_{i,j}$ |
|---|---|---|---|---|
| position | $[0,1]^2$ | 9 | top-left | $[1/7, 2/7] \times [1/7, 2/7]$ |
| | | | top-center | $[1/7, 2/7] \times [3/7, 4/7]$ |
| | | | top-right | $[1/7, 2/7] \times [5/7, 6/7]$ |
| | | | center-left | $[3/7, 4/7] \times [1/7, 2/7]$ |
| | | | center-center | $[3/7, 4/7] \times [3/7, 4/7]$ |
| | | | center-right | $[3/7, 4/7] \times [5/7, 6/7]$ |
| | | | bottom-left | $[5/7, 6/7] \times [1/7, 2/7]$ |
| | | | bottom-center | $[5/7, 6/7] \times [3/7, 4/7]$ |
| | | | bottom-right | $[5/7, 6/7] \times [5/7, 6/7]$ |
| hue | $[0, 2\pi)$ | 6 | red | $[345°, 15°]$ |
| | | | yellow | $[45°, 75°]$ |
| | | | green | $[105°, 135°]$ |
| | | | cyan | $[165°, 195°]$ |
| | | | blue | $[225°, 255°]$ |
| | | | magenta | $[285°, 315°]$ |
| lightness | $[0,1]^2$ | 4 | dark | $[0, 1/11] \times [4/11, 5/11]$ |
| | | | darker | $[2/11, 3/11] \times [6/11, 7/11]$ |
| | | | brighter | $[4/11, 5/11] \times [8/11, 9/11]$ |
| | | | bright | $[6/11, 7/11] \times [10/11, 1.]$ |
| scale | $[1/1.45, 1.45]$ | 5 | small | $[1/1.45, 1/1.35]$ |
| | | | smaller | $[1/1.25, 1/1.15]$ |
| | | | normal | $[1/1.05, 1.05]$ |
| | | | larger | $[1.15, 1.25]$ |
| | | | large | $[1.35, 1.45]$ |
| shape | MNIST | 10 | '0' | digits '0' |
| | | | '1' | digits '1' |
| | | | $\vdots$ | $\vdots$ |
| texture | textures | 5 | tiles | tiles texture crops |
| | | | wood | wood texture crops |
| | | | carpet | carpet texture crops |
| | | | bricks | bricks texture crops |
| | | | lava | lava texture crops |

Table A1: Overview of factors $\mathcal{F}_i$, respective factor spaces $\mathcal{S}_i$, number of classes $\mathcal{N}_i$, factor classes $\mathcal{C}_{i,j}$ and respective factor space regions $\mathcal{S}_{i,j}$ used in this work.
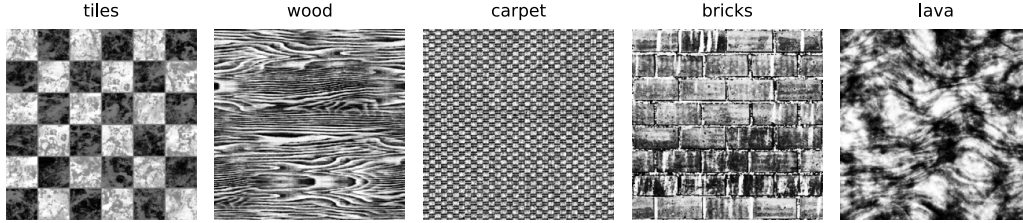
Figure A1: Overview of the textures used in this work. All textures are from CC0Textures.com, licensed under CC0 1.0 Universal.

## A.2. Baselines

We train all baseline architectures using standard cross entropy loss and Adam optimizer [18] using early stopping. All test accuracies reported in this work correspond to the network with lowest validation loss. Batch size and learning rate were optimized on the ZSO study and weight initialization is fixed across all experiments. The output size of all the architectures are modified to predict the factor classes in our studies.

Following ASR, we use a variant of their U-Net architecture for the lens and ResNet18 as feature extractor. Instead of an auxiliary task, we train both the lens and feature extractor using the classification task directly. The lens is trained using least likely adversarial loss using the classification objective and also the reconstruction loss, both with equal weighting. Output of lens is used as inputs for training the feature extractor.

Both VAE methods use architectures similar to [25] with a latent size of 12. We train both VAE with batch size 64 and Adam optimizer [18] with learning rate 0.0001. For the Factor-VAE, we use $\gamma = 20$ and train the discriminator with learning rate 0.00001 throughout our experiments. Some examples of latent traversals from VAE and Factor-VAE are shown in Fig. A2. We experimented with $\beta$-VAEs but observed, that the reconstructions of are poor when we use higher $\beta$ in our settings. One potential reason is the discretely distributed positions in the trained datasets. The position factor of DiagViB-6 has nine possible values with large gap among them violating the Gaussian distribution assumption in the KL-divergence term during training. To empirically demonstrate this claim, we train $\beta$-VAEs with different $\beta$ on two settings: (1) with position factor freely assigned to three different values (2) with position factor fixed to the center of the images. According to the results shown in Fig. A3, the reconstructions from the training with fixed position are significantly better, supporting the aforementioned argument.

## A.3. ZSO and ZGO

The mean accuracies $P_i$ for all factors $\mathcal{F}_i$ on the ZSO together with the mean accuracies $P_{i,j}$ for all factor pairings $(\mathcal{F}_i, \mathcal{F}_j), i \neq j$ on the ZGO study for all baselines are shown in Fig. A5 and A6. Additionally, the aggregated benchmark metrics FAAvg$_i$ and FAMin$_i$ for all factors $\mathcal{F}_i, i \in \{0, 1, ..., 6\}$ together with their respective standard errors for all baselines on the ZSO and ZGO studies are presented in Tab. A2 and A3.

From the exemplar VAE decoder reconstructions (Fig. A3) and ASR lens outputs (Fig. A4) we find that the high-frequency texture information is not preserved, leading to the observed low classification accuracies on the texture task.
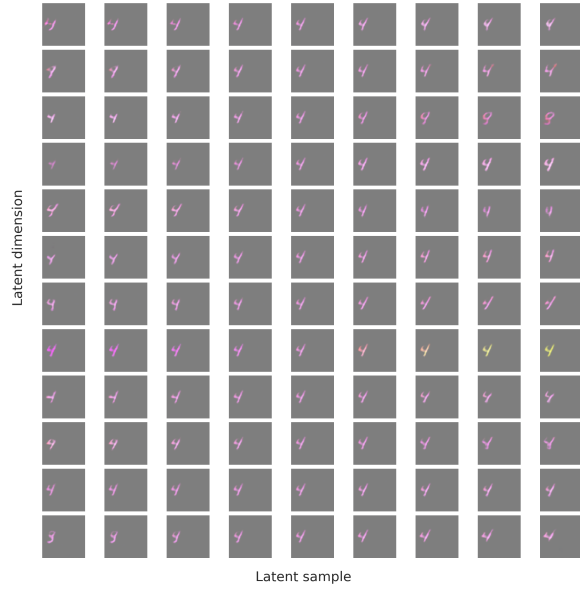
## A.4. Compositional-based generalization opportunities

The mean accuracies $P_i$ for all factors $\mathcal{F}_i$ on the ZSO together with the mean accuracies $P_{i,j}$ for all factor pairings $(\mathcal{F}_i, \mathcal{F}_j), i \neq j$ on the CGO and ZGO studies for all baselines are shown in Fig. A5 and A6. Additionally, the aggregated benchmark metrics FAAvg$_i$ and FAMin$_i$ for all factors $\mathcal{F}_i, i \in \{0, 1, ..., 6\}$ together with their respective standard errors for all baselines on the CGO studies are presented in Tab. A4 to A6.

## A.5. Frequency-based generalization opportunities

The mean accuracies $P_i$ for all factors $\mathcal{F}_i$ on the ZSO together with the mean accuracies $P_{i,j}$ for all factor pairings $(\mathcal{F}_i, \mathcal{F}_j), i \neq j$ on the FGO and ZGO studies for all baselines are shown in Fig. A7 and A8. Additionally, the aggregated benchmark metrics FAAvg$_i$ and FAMin$_i$ for all factors $\mathcal{F}_i, i \in \{0, 1, ..., 6\}$ together with their respective standard errors for all baselines on the FGO studies are presented in Tab. A7 to A9.

As discussed in Sec. 5, we find that ASR fails to exploit the presented GO for hue and lightness. This is likely attributed to the pixel-wise reconstruction loss that is used as a regularizer in the ASR objective. We present qualitative results of this behaviour in Fig. A9.

(a) Latent traversal plot from VAE

(b) Latent traversal plot from Factor-VAE

(c) Latent traversal plot from VAE

(d) Latent traversal plot from Factor-VAE

Figure A2: Comparison of latent traversals from VAE and Factor-VAE

(a) Input image (Sample 1)

(b) Reconstructed images (Sample 1)

(c) Input image (Sample 2)

(d) Reconstructed images (Sample 2)

(e) Input image (Sample 3)

(f) Reconstructed images (Sample 3)

(g) Input image (Sample 4)

(h) Reconstructed images (Sample 4)

Figure A3: Reconstruction images from the $\beta$-VAE networks trained with datasets of free and fixed positions from different $\beta$ values



Figure A4: Examples of Automatic shortcut removal Lens output for ZSO study with factor `texture`.

|            | position    | hue       | lightness | scale     | shape      | texture   |
|------------|-------------|-----------|-----------|-----------|------------|-----------|
| RN18       | $100 \pm 0$ | $100 \pm 0$ | $99 \pm 0$ | $99 \pm 0$ | $100 \pm 0$ | $62 \pm 6$ |
| RN50       | $100 \pm 0$ | $100 \pm 0$ | $89 \pm 9$ | $99 \pm 0$ | $100 \pm 0$ | $77 \pm 3$ |
| RN50-IN    | $99 \pm 0$  | $99 \pm 0$  | $94 \pm 1$ | $74 \pm 2$ | $83 \pm 3$  | $54 \pm 1$ |
| AlexNet    | $100 \pm 0$ | $100 \pm 0$ | $98 \pm 1$ | $97 \pm 0$ | $99 \pm 0$  | $68 \pm 3$ |
| DenseNet   | $100 \pm 0$ | $100 \pm 0$ | $99 \pm 0$ | $99 \pm 0$ | $100 \pm 0$ | $70 \pm 3$ |
| WRN        | $100 \pm 0$ | $100 \pm 0$ | $99 \pm 0$ | $99 \pm 0$ | $100 \pm 0$ | $86 \pm 2$ |
| ASR (RN18) | $100 \pm 0$ | $77 \pm 6$  | $82 \pm 2$ | $92 \pm 1$ | $91 \pm 1$  | $45 \pm 1$ |
| Factor-VAE | $100 \pm 0$ | $97 \pm 1$  | $93 \pm 2$ | $84 \pm 1$ | $91 \pm 2$  | $34 \pm 0$ |
| VAE        | $100 \pm 0$ | $98 \pm 1$  | $94 \pm 1$ | $86 \pm 2$ | $93 \pm 2$  | $34 \pm 0$ |

Table A2: FAAvg and respecive standard error for all baselines on the ZSO study.

|            | position | | hue | | lightness | | scale | | shape | | texture | |
|------------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
|            | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin |
| RN18       | $100 \pm 0$ | $99 \pm 1$  | $74 \pm 2$ | $0 \pm 0$  | $57 \pm 4$ | $31 \pm 5$ | $33 \pm 2$ | $0 \pm 0$ | $41 \pm 2$ | $0 \pm 0$ | $2 \pm 1$  | $0 \pm 0$ |
| RN50       | $100 \pm 0$ | $100 \pm 0$ | $72 \pm 2$ | $0 \pm 0$  | $58 \pm 4$ | $27 \pm 1$ | $32 \pm 2$ | $0 \pm 0$ | $42 \pm 2$ | $0 \pm 0$ | $3 \pm 1$  | $0 \pm 0$ |
| RN50-IN    | $68 \pm 3$  | $18 \pm 2$  | $85 \pm 2$ | $64 \pm 3$ | $72 \pm 4$ | $50 \pm 6$ | $28 \pm 2$ | $4 \pm 1$ | $33 \pm 5$ | $8 \pm 4$ | $10 \pm 1$ | $1 \pm 0$ |
| AlexNet    | $100 \pm 0$ | $100 \pm 0$ | $72 \pm 2$ | $0 \pm 0$  | $62 \pm 6$ | $31 \pm 5$ | $32 \pm 2$ | $0 \pm 0$ | $39 \pm 2$ | $0 \pm 0$ | $3 \pm 1$  | $0 \pm 0$ |
| DenseNet   | $98 \pm 1$  | $88 \pm 6$  | $74 \pm 3$ | $7 \pm 3$  | $61 \pm 6$ | $34 \pm 8$ | $34 \pm 1$ | $0 \pm 0$ | $40 \pm 2$ | $0 \pm 0$ | $2 \pm 0$  | $0 \pm 0$ |
| WRN        | $99 \pm 1$  | $96 \pm 3$  | $73 \pm 3$ | $4 \pm 3$  | $61 \pm 6$ | $37 \pm 8$ | $32 \pm 2$ | $0 \pm 0$ | $41 \pm 2$ | $0 \pm 0$ | $2 \pm 1$  | $0 \pm 0$ |
| ASR (RN18) | $99 \pm 0$  | $96 \pm 1$  | $27 \pm 3$ | $1 \pm 0$  | $43 \pm 2$ | $25 \pm 0$ | $51 \pm 2$ | $0 \pm 0$ | $60 \pm 1$ | $1 \pm 0$ | $8 \pm 0$  | $0 \pm 0$ |
| Factor-VAE | $100 \pm 0$ | $100 \pm 0$ | $65 \pm 4$ | $0 \pm 0$  | $64 \pm 4$ | $27 \pm 0$ | $32 \pm 1$ | $0 \pm 0$ | $35 \pm 2$ | $0 \pm 0$ | $4 \pm 0$  | $0 \pm 0$ |
| VAE        | $100 \pm 0$ | $100 \pm 0$ | $65 \pm 4$ | $0 \pm 0$  | $64 \pm 4$ | $26 \pm 0$ | $33 \pm 2$ | $0 \pm 0$ | $37 \pm 3$ | $0 \pm 0$ | $4 \pm 0$  | $0 \pm 0$ |

Table A3: FAAvg and FAMin together with their respective standard errors for all baselines on the ZGO study.

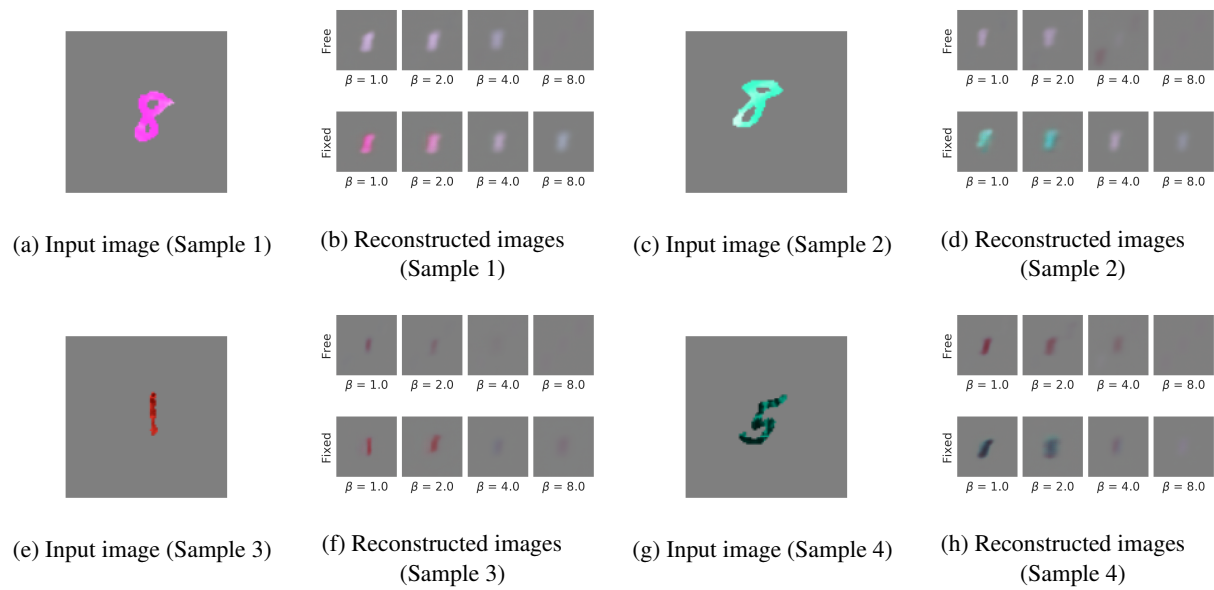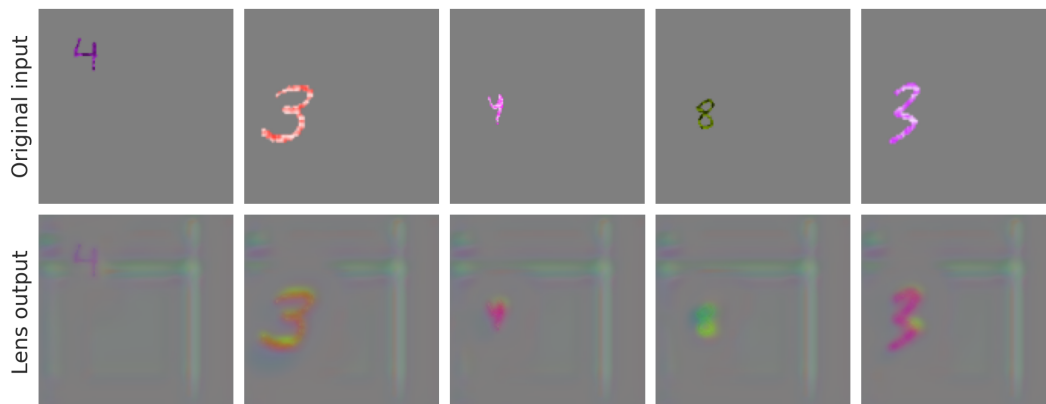Figure A5: Mean accuracies $P_i$ for all factors $\mathcal{F}_i$ on the ZSO study and mean accuracies $P_{i,j}$ for all factor pairings $(\mathcal{F}_i, \mathcal{F}_j), i \neq j$ for the RN18, RN50, RN50-IN, AlexNet and DenseNet on the CGO and ZGO studies.
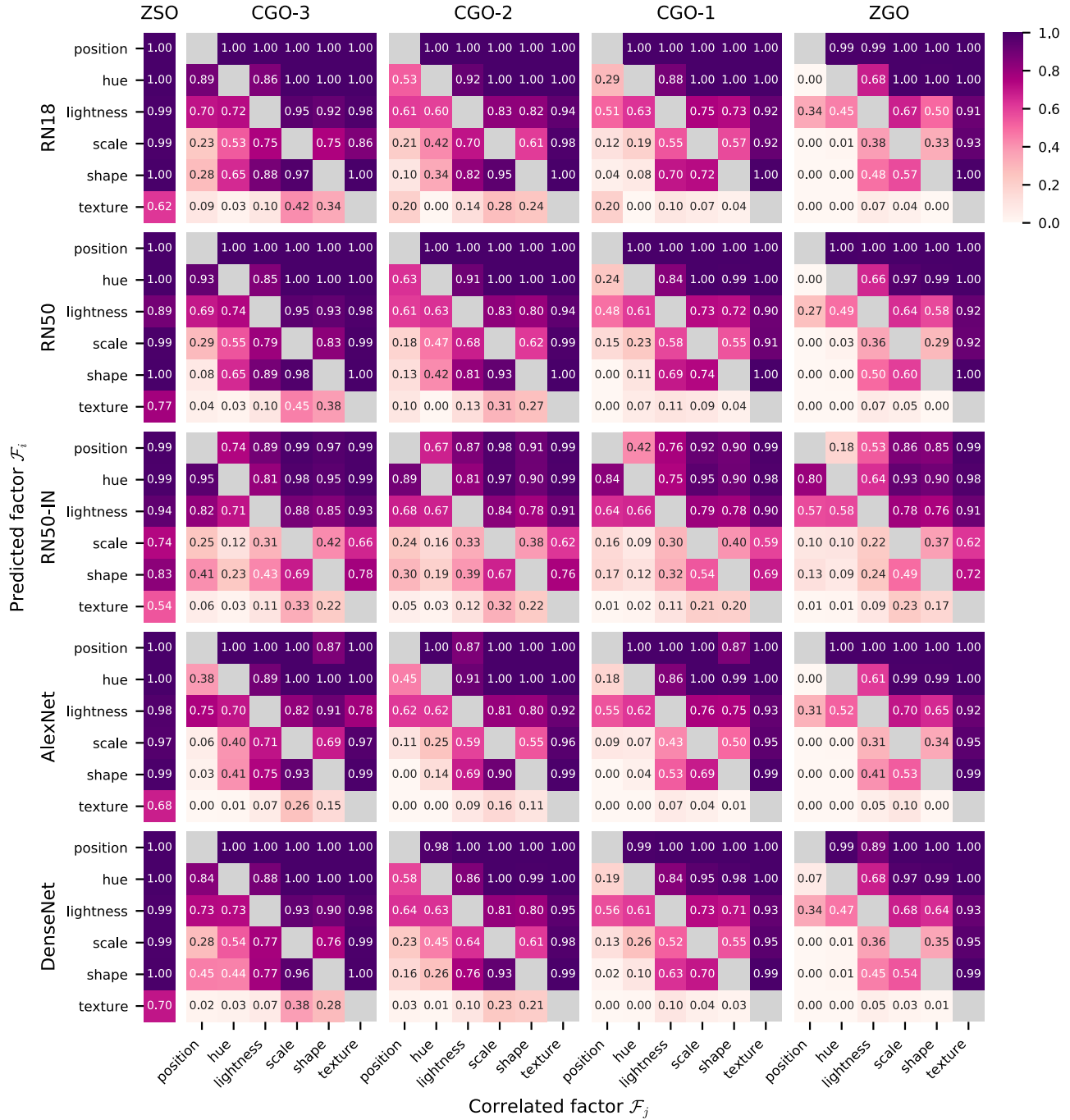
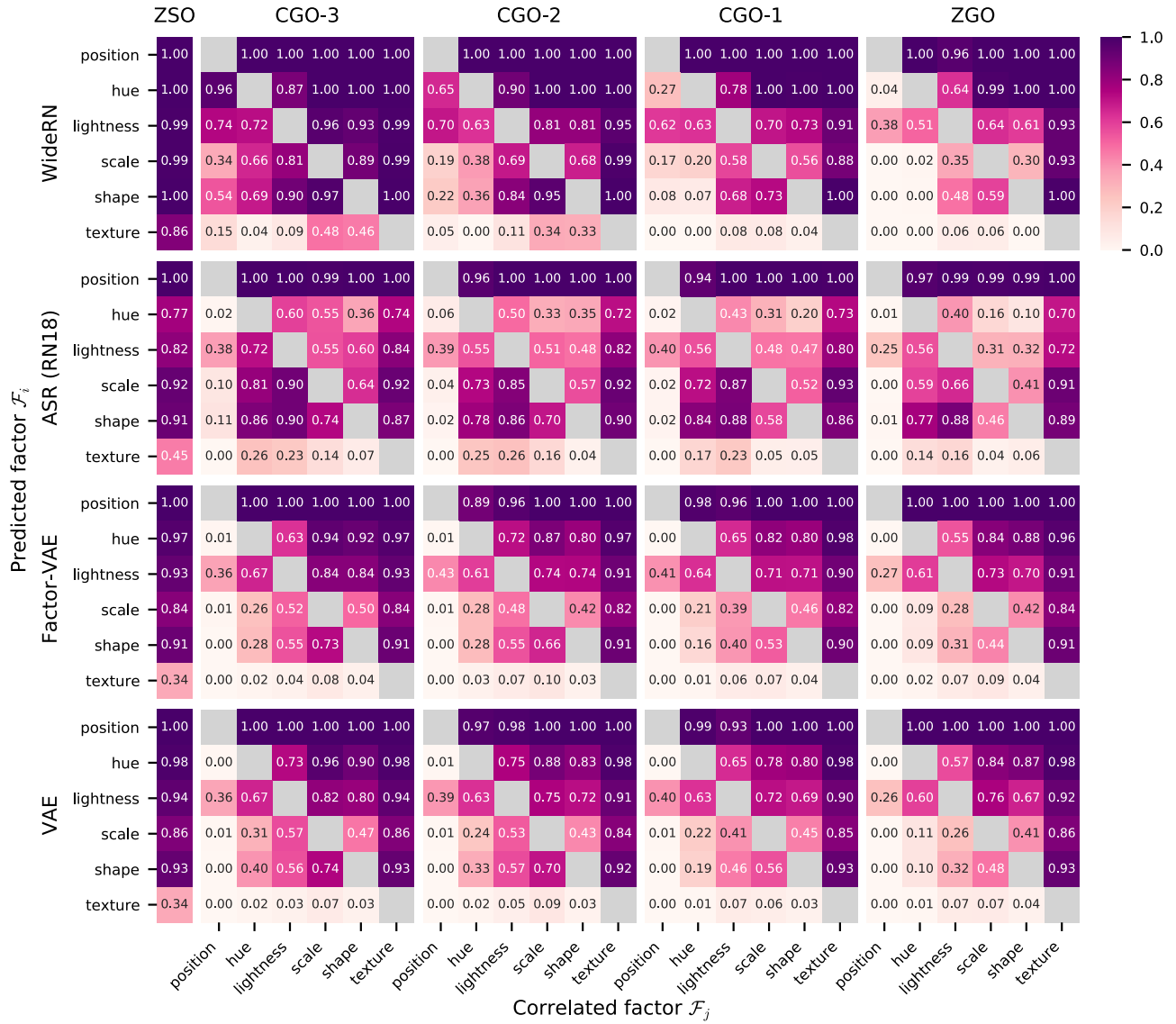Figure A6: Mean accuracies $P_i$ for all factors $\mathcal{F}_i$ on the ZSO study and mean accuracies $P_{i,j}$ for all factor pairings $(\mathcal{F}_i, \mathcal{F}_j), i \neq j$ for the WideRN, ASR (RN18), Factor-VAE and VAE on the CGO and ZGO studies.

| | position | | hue | | lightness | | scale | | shape | | texture | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin |
| RN18 | $100\pm0$ | $100\pm0$ | $83\pm1$ | $29\pm5$ | $71\pm8$ | $49\pm9$ | $47\pm3$ | $6\pm2$ | $51\pm3$ | $2\pm1$ | $8\pm2$ | $0\pm0$ |
| RN50 | $100\pm0$ | $100\pm0$ | $82\pm1$ | $24\pm6$ | $69\pm8$ | $46\pm8$ | $48\pm4$ | $9\pm3$ | $51\pm3$ | $0\pm0$ | $6\pm2$ | $0\pm0$ |
| RN50-IN | $80\pm2$ | $42\pm4$ | $88\pm2$ | $73\pm3$ | $75\pm5$ | $59\pm8$ | $31\pm3$ | $8\pm2$ | $37\pm6$ | $12\pm5$ | $11\pm1$ | $1\pm0$ |
| AlexNet | $97\pm2$ | $87\pm12$ | $81\pm2$ | $18\pm7$ | $72\pm8$ | $52\pm10$ | $41\pm3$ | $2\pm1$ | $45\pm3$ | $0\pm0$ | $3\pm0$ | $0\pm0$ |
| DenseNet | $100\pm0$ | $99\pm1$ | $79\pm1$ | $19\pm1$ | $71\pm8$ | $50\pm10$ | $48\pm4$ | $9\pm3$ | $49\pm2$ | $2\pm1$ | $3\pm0$ | $0\pm0$ |
| WRN | $100\pm0$ | $100\pm0$ | $81\pm2$ | $27\pm4$ | $72\pm9$ | $51\pm10$ | $48\pm3$ | $10\pm3$ | $51\pm3$ | $4\pm2$ | $4\pm1$ | $0\pm0$ |
| ASR (RN18) | $99\pm1$ | $94\pm3$ | $34\pm3$ | $2\pm1$ | $54\pm4$ | $40\pm5$ | $61\pm2$ | $2\pm1$ | $64\pm3$ | $2\pm1$ | $10\pm1$ | $0\pm0$ |
| Factor-VAE | $99\pm0$ | $93\pm2$ | $65\pm3$ | $0\pm0$ | $67\pm6$ | $41\pm5$ | $38\pm3$ | $0\pm0$ | $40\pm3$ | $0\pm0$ | $4\pm0$ | $0\pm0$ |
| VAE | $98\pm1$ | $92\pm3$ | $64\pm3$ | $0\pm0$ | $67\pm6$ | $40\pm6$ | $39\pm3$ | $1\pm0$ | $43\pm4$ | $0\pm0$ | $3\pm0$ | $0\pm0$ |

Table A4: Aggregated benchmark metrics $\text{FAAvg}_i$ and $\text{FAMin}_i$ for all factors $\mathcal{F}_i, i \in \{0, 1, ..., 6\}$ together with their respective standard errors for all baselines on the CGO-1 study.

| | position | | hue | | lightness | | scale | | shape | | texture | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin |
| RN18 | $100\pm0$ | $100\pm0$ | $89\pm4$ | $53\pm16$ | $76\pm11$ | $51\pm14$ | $59\pm5$ | $21\pm6$ | $64\pm2$ | $4\pm2$ | $17\pm4$ | $0\pm0$ |
| RN50 | $100\pm0$ | $100\pm0$ | $91\pm4$ | $59\pm17$ | $76\pm10$ | $56\pm14$ | $59\pm4$ | $18\pm7$ | $66\pm3$ | $13\pm5$ | $16\pm4$ | $0\pm0$ |
| RN50-IN | $88\pm2$ | $67\pm6$ | $91\pm3$ | $80\pm4$ | $77\pm8$ | $63\pm12$ | $35\pm5$ | $15\pm4$ | $46\pm6$ | $18\pm5$ | $15\pm2$ | $2\pm1$ |
| AlexNet | $97\pm2$ | $87\pm12$ | $87\pm2$ | $45\pm7$ | $75\pm9$ | $55\pm14$ | $49\pm4$ | $11\pm5$ | $55\pm2$ | $0\pm0$ | $7\pm2$ | $0\pm0$ |
| DenseNet | $100\pm0$ | $98\pm2$ | $89\pm4$ | $56\pm14$ | $77\pm10$ | $58\pm15$ | $58\pm5$ | $23\pm4$ | $62\pm2$ | $8\pm3$ | $12\pm3$ | $0\pm0$ |
| WRN | $100\pm0$ | $100\pm0$ | $91\pm4$ | $62\pm17$ | $78\pm11$ | $58\pm14$ | $59\pm4$ | $19\pm3$ | $67\pm3$ | $11\pm6$ | $17\pm4$ | $0\pm0$ |
| ASR (RN18) | $99\pm1$ | $96\pm3$ | $39\pm6$ | $6\pm6$ | $55\pm7$ | $39\pm9$ | $62\pm2$ | $4\pm2$ | $65\pm3$ | $2\pm1$ | $14\pm2$ | $0\pm0$ |
| Factor-VAE | $97\pm2$ | $89\pm6$ | $67\pm4$ | $1\pm0$ | $68\pm8$ | $43\pm10$ | $40\pm4$ | $1\pm0$ | $48\pm3$ | $0\pm0$ | $5\pm1$ | $0\pm0$ |
| VAE | $99\pm1$ | $97\pm3$ | $69\pm4$ | $1\pm0$ | $68\pm7$ | $39\pm8$ | $41\pm4$ | $1\pm0$ | $51\pm4$ | $0\pm0$ | $4\pm0$ | $0\pm0$ |

Table A5: Aggregated benchmark metrics $\text{FAAvg}_i$ and $\text{FAMin}_i$ for all factors $\mathcal{F}_i, i \in \{0, 1, ..., 6\}$ together with their respective standard errors for all baselines on the CGO-2 study.

| | position | | hue | | lightness | | scale | | shape | | texture | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin |
| RN18 | $100\pm0$ | $100\pm0$ | $95\pm3$ | $81\pm9$ | $86\pm6$ | $65\pm13$ | $62\pm7$ | $22\pm3$ | $76\pm3$ | $28\pm7$ | $19\pm4$ | $1\pm0$ |
| RN50 | $100\pm0$ | $100\pm0$ | $96\pm2$ | $85\pm8$ | $86\pm7$ | $66\pm17$ | $69\pm7$ | $29\pm9$ | $72\pm2$ | $8\pm3$ | $20\pm4$ | $1\pm0$ |
| RN50-IN | $92\pm2$ | $74\pm7$ | $94\pm1$ | $81\pm5$ | $84\pm5$ | $70\pm10$ | $35\pm5$ | $12\pm3$ | $51\pm8$ | $23\pm10$ | $15\pm2$ | $3\pm1$ |
| AlexNet | $97\pm2$ | $87\pm12$ | $86\pm2$ | $38\pm7$ | $79\pm7$ | $54\pm11$ | $56\pm7$ | $6\pm2$ | $62\pm4$ | $3\pm2$ | $10\pm2$ | $0\pm0$ |
| DenseNet | $100\pm0$ | $100\pm0$ | $94\pm2$ | $80\pm6$ | $86\pm7$ | $71\pm13$ | $66\pm7$ | $23\pm4$ | $72\pm4$ | $29\pm11$ | $16\pm3$ | $1\pm0$ |
| WRN | $100\pm0$ | $100\pm0$ | $97\pm2$ | $87\pm6$ | $87\pm6$ | $69\pm13$ | $74\pm6$ | $32\pm7$ | $82\pm3$ | $48\pm9$ | $24\pm6$ | $2\pm0$ |
| ASR (RN18) | $100\pm0$ | $99\pm0$ | $45\pm1$ | $2\pm1$ | $62\pm8$ | $34\pm8$ | $67\pm3$ | $10\pm5$ | $69\pm3$ | $11\pm7$ | $14\pm1$ | $0\pm0$ |
| Factor-VAE | $100\pm0$ | $100\pm0$ | $70\pm3$ | $1\pm1$ | $73\pm6$ | $36\pm7$ | $43\pm4$ | $1\pm1$ | $49\pm3$ | $0\pm0$ | $4\pm1$ | $0\pm0$ |
| VAE | $100\pm0$ | $100\pm0$ | $72\pm3$ | $0\pm0$ | $72\pm7$ | $36\pm8$ | $44\pm5$ | $1\pm0$ | $53\pm4$ | $0\pm0$ | $3\pm1$ | $0\pm0$ |

Table A6: Aggregated benchmark metrics $\text{FAAvg}_i$ and $\text{FAMin}_i$ for all factors $\mathcal{F}_i, i \in \{0, 1, ..., 6\}$ together with their respective standard errors for all baselines on the CGO-3 study.

Figure A7: Mean accuracies $P_i$ for all factors $\mathcal{F}_i$ on the ZSO study and mean accuracies $P_{i,j}$ for all factor pairings $(\mathcal{F}_i, \mathcal{F}_j), i \neq j$ for the RN18, RN50, RN50-IN, AlexNet and DenseNet on the FGO and ZGO studies.
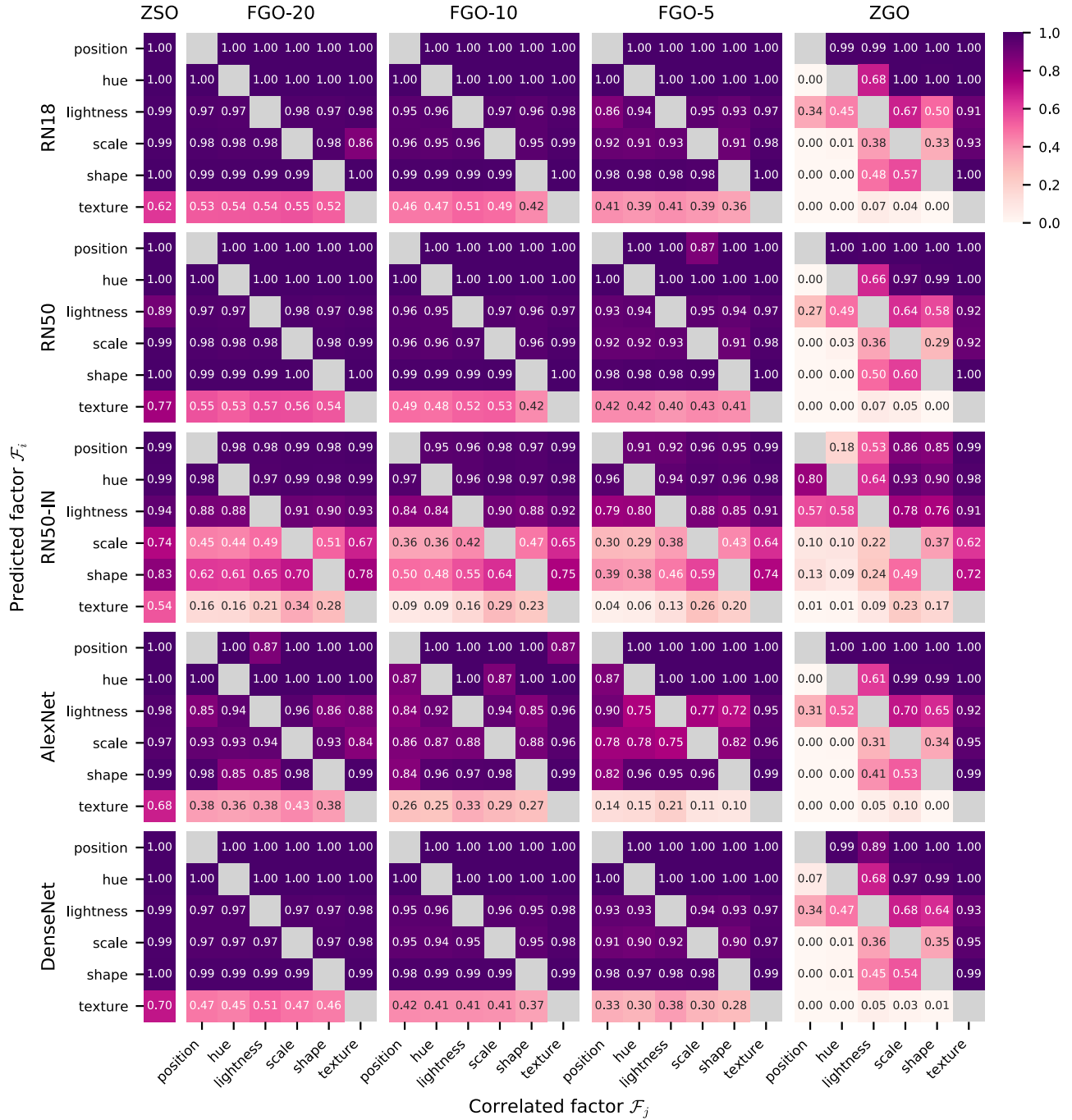
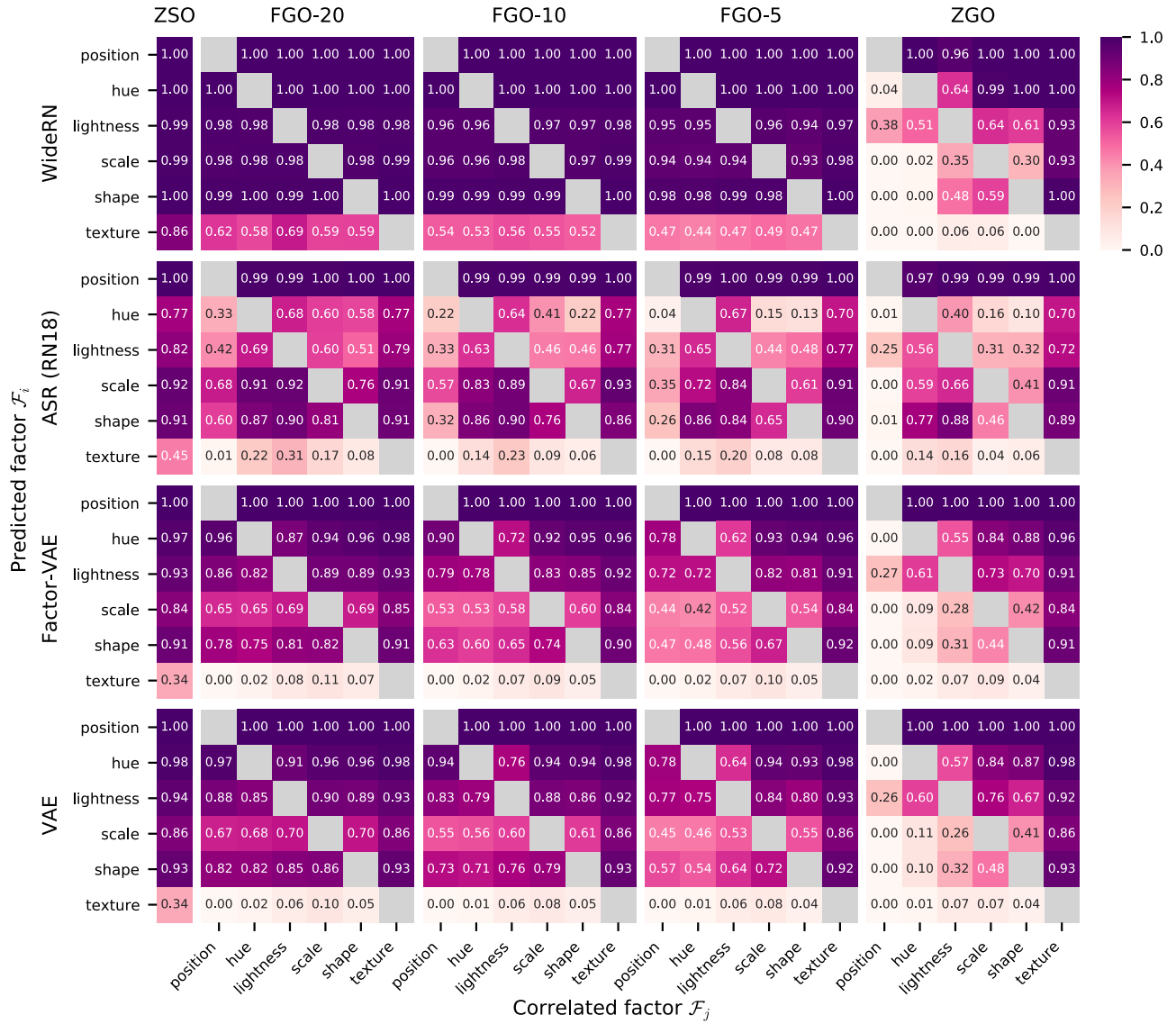Figure A8: Mean accuracies $P_i$ for all factors $\mathcal{F}_i$ on the ZSO study and mean accuracies $P_{i,j}$ for all factor pairings $(\mathcal{F}_i, \mathcal{F}_j), i \neq j$ for the WideRN, ASR (RN18), Factor-VAE and VAE on the FGO and ZGO studies.

| | position | | hue | | lightness | | scale | | shape | | texture | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin |
| RN18 | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $93 \pm 2$ | $85 \pm 8$ | $93 \pm 1$ | $89 \pm 1$ | $99 \pm 0$ | $98 \pm 0$ | $39 \pm 2$ | $35 \pm 3$ |
| RN50 | $97 \pm 2$ | $87 \pm 12$ | $100 \pm 0$ | $100 \pm 0$ | $95 \pm 1$ | $93 \pm 2$ | $93 \pm 1$ | $90 \pm 1$ | $99 \pm 0$ | $98 \pm 0$ | $41 \pm 2$ | $36 \pm 1$ |
| RN50-IN | $95 \pm 1$ | $91 \pm 2$ | $96 \pm 1$ | $94 \pm 1$ | $85 \pm 3$ | $78 \pm 4$ | $41 \pm 4$ | $28 \pm 4$ | $51 \pm 8$ | $37 \pm 9$ | $14 \pm 1$ | $4 \pm 1$ |
| AlexNet | $100 \pm 0$ | $100 \pm 0$ | $97 \pm 2$ | $87 \pm 12$ | $82 \pm 4$ | $64 \pm 8$ | $82 \pm 3$ | $69 \pm 8$ | $94 \pm 3$ | $82 \pm 11$ | $14 \pm 2$ | $6 \pm 3$ |
| DenseNet | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $94 \pm 1$ | $92 \pm 2$ | $92 \pm 1$ | $88 \pm 1$ | $98 \pm 0$ | $97 \pm 1$ | $32 \pm 2$ | $26 \pm 2$ |
| WRN | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $95 \pm 1$ | $94 \pm 1$ | $95 \pm 1$ | $92 \pm 1$ | $98 \pm 0$ | $98 \pm 0$ | $47 \pm 3$ | $42 \pm 2$ |
| ASR (RN18) | $99 \pm 0$ | $98 \pm 0$ | $34 \pm 3$ | $1 \pm 1$ | $53 \pm 4$ | $28 \pm 1$ | $69 \pm 2$ | $35 \pm 6$ | $70 \pm 2$ | $26 \pm 7$ | $10 \pm 1$ | $0 \pm 0$ |
| Factor-VAE | $100 \pm 0$ | $100 \pm 0$ | $84 \pm 4$ | $61 \pm 9$ | $80 \pm 5$ | $69 \pm 7$ | $55 \pm 4$ | $40 \pm 4$ | $62 \pm 5$ | $44 \pm 7$ | $5 \pm 0$ | $0 \pm 0$ |
| VAE | $100 \pm 0$ | $100 \pm 0$ | $86 \pm 3$ | $63 \pm 8$ | $82 \pm 4$ | $74 \pm 5$ | $57 \pm 4$ | $42 \pm 4$ | $68 \pm 5$ | $52 \pm 6$ | $4 \pm 0$ | $0 \pm 0$ |

Table A7: Aggregated benchmark metrics $\text{FAAvg}_i$ and $\text{FAMin}_i$ for all factors $\mathcal{F}_i, i \in \{0, 1, ..., 6\}$ together with their respective standard errors for all baselines on the FGO-5 study.

| | position | | hue | | lightness | | scale | | shape | | texture | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin |
| RN18 | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $96 \pm 1$ | $95 \pm 1$ | $96 \pm 1$ | $95 \pm 1$ | $99 \pm 0$ | $99 \pm 0$ | $47 \pm 3$ | $41 \pm 4$ |
| RN50 | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $96 \pm 1$ | $95 \pm 1$ | $97 \pm 0$ | $95 \pm 1$ | $99 \pm 0$ | $99 \pm 0$ | $49 \pm 3$ | $41 \pm 2$ |
| RN50-IN | $97 \pm 1$ | $95 \pm 1$ | $97 \pm 1$ | $96 \pm 1$ | $87 \pm 2$ | $83 \pm 3$ | $45 \pm 3$ | $35 \pm 4$ | $59 \pm 7$ | $48 \pm 9$ | $17 \pm 2$ | $8 \pm 1$ |
| AlexNet | $97 \pm 2$ | $87 \pm 12$ | $95 \pm 3$ | $73 \pm 15$ | $90 \pm 3$ | $75 \pm 9$ | $89 \pm 1$ | $85 \pm 2$ | $95 \pm 3$ | $84 \pm 11$ | $28 \pm 2$ | $22 \pm 2$ |
| DenseNet | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $96 \pm 1$ | $95 \pm 1$ | $95 \pm 1$ | $94 \pm 1$ | $99 \pm 0$ | $98 \pm 1$ | $40 \pm 1$ | $36 \pm 1$ |
| WRN | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $97 \pm 1$ | $96 \pm 1$ | $97 \pm 0$ | $95 \pm 0$ | $99 \pm 0$ | $99 \pm 0$ | $54 \pm 5$ | $52 \pm 5$ |
| ASR (RN18) | $99 \pm 0$ | $98 \pm 0$ | $45 \pm 4$ | $11 \pm 4$ | $53 \pm 3$ | $31 \pm 2$ | $78 \pm 1$ | $50 \pm 3$ | $74 \pm 4$ | $32 \pm 7$ | $10 \pm 0$ | $0 \pm 0$ |
| Factor-VAE | $100 \pm 0$ | $100 \pm 0$ | $89 \pm 3$ | $72 \pm 9$ | $83 \pm 4$ | $76 \pm 6$ | $62 \pm 4$ | $51 \pm 5$ | $71 \pm 5$ | $60 \pm 5$ | $5 \pm 0$ | $0 \pm 0$ |
| VAE | $100 \pm 0$ | $100 \pm 0$ | $91 \pm 2$ | $76 \pm 7$ | $86 \pm 3$ | $79 \pm 5$ | $64 \pm 4$ | $53 \pm 4$ | $78 \pm 5$ | $70 \pm 6$ | $4 \pm 0$ | $0 \pm 0$ |

Table A8: Aggregated benchmark metrics $\text{FAAvg}_i$ and $\text{FAMin}_i$ for all factors $\mathcal{F}_i, i \in \{0, 1, ..., 6\}$ together with their respective standard errors for all baselines on the FGO-10 study.

| | position | | hue | | lightness | | scale | | shape | | texture | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin | FAAvg | FAMin |
| RN18 | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $98 \pm 1$ | $97 \pm 1$ | $95 \pm 2$ | $84 \pm 11$ | $99 \pm 0$ | $99 \pm 0$ | $54 \pm 5$ | $51 \pm 5$ |
| RN50 | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $98 \pm 1$ | $97 \pm 1$ | $98 \pm 0$ | $98 \pm 0$ | $100 \pm 0$ | $99 \pm 0$ | $55 \pm 5$ | $52 \pm 5$ |
| RN50-IN | $98 \pm 0$ | $98 \pm 1$ | $98 \pm 0$ | $97 \pm 1$ | $90 \pm 2$ | $88 \pm 2$ | $51 \pm 3$ | $44 \pm 3$ | $67 \pm 6$ | $60 \pm 7$ | $23 \pm 2$ | $15 \pm 2$ |
| AlexNet | $97 \pm 2$ | $87 \pm 12$ | $100 \pm 0$ | $100 \pm 0$ | $90 \pm 4$ | $77 \pm 10$ | $91 \pm 3$ | $80 \pm 11$ | $93 \pm 5$ | $85 \pm 11$ | $38 \pm 2$ | $34 \pm 2$ |
| DenseNet | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $97 \pm 1$ | $96 \pm 1$ | $97 \pm 0$ | $97 \pm 0$ | $99 \pm 0$ | $99 \pm 0$ | $47 \pm 3$ | $43 \pm 3$ |
| WRN | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $100 \pm 0$ | $98 \pm 0$ | $97 \pm 1$ | $98 \pm 0$ | $97 \pm 0$ | $99 \pm 0$ | $99 \pm 0$ | $61 \pm 5$ | $57 \pm 6$ |
| ASR (RN18) | $100 \pm 0$ | $99 \pm 0$ | $59 \pm 4$ | $26 \pm 8$ | $60 \pm 3$ | $38 \pm 2$ | $84 \pm 2$ | $67 \pm 5$ | $82 \pm 3$ | $60 \pm 7$ | $16 \pm 1$ | $1 \pm 1$ |
| Factor-VAE | $100 \pm 0$ | $100 \pm 0$ | $94 \pm 2$ | $84 \pm 5$ | $88 \pm 3$ | $82 \pm 5$ | $71 \pm 3$ | $63 \pm 3$ | $81 \pm 4$ | $75 \pm 5$ | $5 \pm 0$ | $0 \pm 0$ |
| VAE | $100 \pm 0$ | $100 \pm 0$ | $96 \pm 1$ | $89 \pm 4$ | $89 \pm 3$ | $84 \pm 4$ | $72 \pm 3$ | $65 \pm 3$ | $85 \pm 3$ | $81 \pm 4$ | $5 \pm 0$ | $0 \pm 0$ |

Table A9: Aggregated benchmark metrics $\text{FAAvg}_i$ and $\text{FAMin}_i$ for all factors $\mathcal{F}_i, i \in \{0, 1, ..., 6\}$ together with their respective standard errors for all baselines on the FGO-20 study.
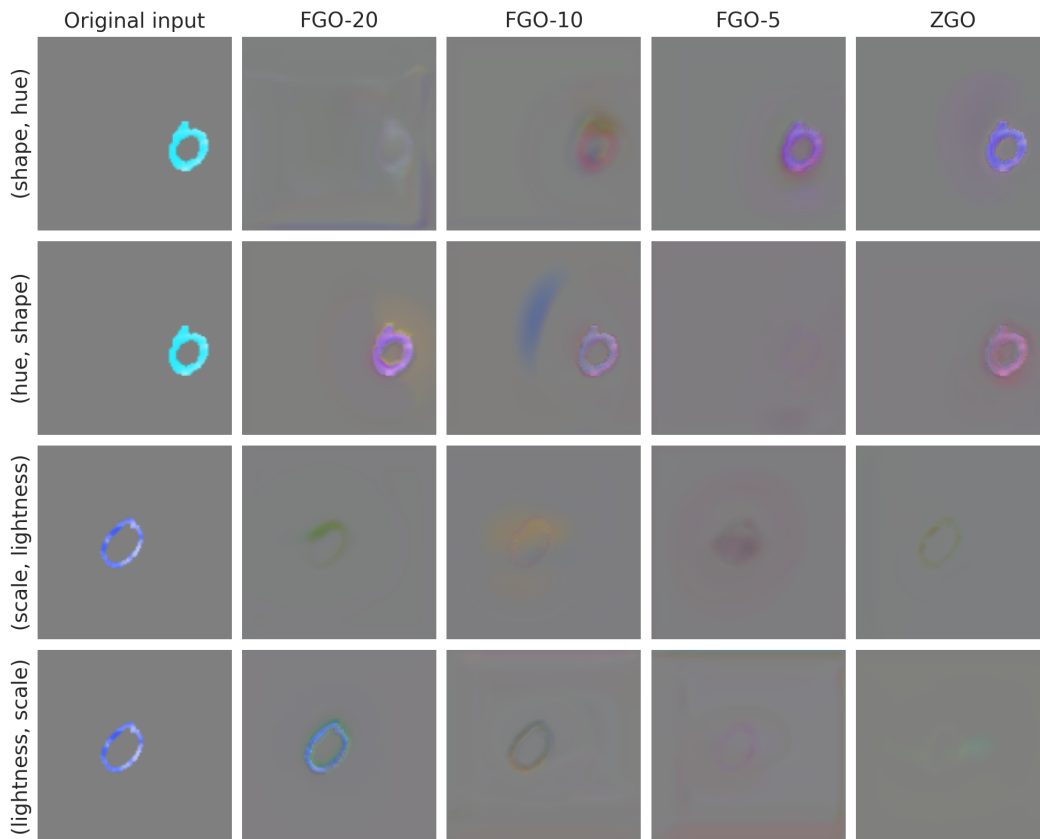
Figure A9: Examples of Automatic shortcut removal Lens output for different studies in our benchmark on datasets with following factor correlations: (shape, hue) (first row), (hue, shape) (second row), (scale, lightness) (third row), (lightness, scale) (fourth row).