

Augmenting Softmax Information for Selective Classification with Out-of-Distribution Data

Supplemental Material

Guoxuan Xia and Christos-Savvas Bouganis

Imperial College London
{g.xia21,christos-savvas.bouganis}@imperial.ac.uk

A Experimental Details

We present detailed information about our experimental setup. Our code is available at <https://github.com/Guoxoug/SIRC>.

A.1 Models and Training

For the main results we train ResNet-50 [4] using the default hyperparameters found in PyTorch’s examples.¹ We train on ImageNet-200 for 90 epochs with a batch size of 256. Stochastic gradient descent is used with a weight decay of 10^{-4} , a momentum of 0.9 and an initial learning rate of 0.1 that steps down by a factor of 10 at epochs 30 and 60. Images are augmented using `RandomResizedCrop` and `RandomHorizontalFlip`. MobileNetV2 [18] uses the same setting, but with an initial learning rate of 0.05. DenseNet-121 is trained with the same settings as ResNet-50 but with Nesterov momentum as per [8]. We perform 5 independent training runs for each architecture, with random seeds $\{1, \dots, 5\}$.

Additionally, we also test on two pre-trained ImageNet-1k models. We use ResNetV2-101 from Google’s Big Transfer² [13], specifically `BiT-S-R101x1`, and DenseNet-121 provided by PyTorch.³ Note that the BiT model takes 480×480 images as input, whereas all other models take standard ImageNet-scale 224×224 images. Note that for evaluating these models we exclude Near-ImageNet-200 and Caltech-45 due to class overlap with ImageNet-1k.

A.2 ImageNet-Scale Datasets

Figure 1 shows a number of random examples from each dataset introduced in the main paper, alongside the number of samples in said dataset. Below we describe the methodology for constructing Colonoscopy and Noise. For the remaining datasets please refer to their original papers for details [7, 10–12, 19]. We note that there is a slight discrepancy between the number of samples reported in [12]

¹ <https://github.com/pytorch/examples/tree/main/imagenet>

² https://github.com/google-research/big_transfer

³ <https://pytorch.org/vision/stable/models.html>

for ImageNet-200 and in the authors’ provided datasets,⁴ but we do not believe this affects the validity of our results.

Noise We randomly generate 10000 square images. All samples are generated independently. Within each image, each value (in space and RGB) is sampled from the same gaussian distribution, with mean 0.5. The standard deviation of said gaussian differs between images. These in turn are generated by sampling from a unit gaussian and squaring the samples. Pixel values are then clipped to be in $[0, 1]$ and mapped to 8-bit integers. The widths of each image are sampled uniformly from $\{2, \dots, 256\}$, and the images are all scaled to 256×256 using the lanczos interpolation method in PIL.⁵ The resulting data thus varies in both scale and contrast (see Fig. 1).

Colonoscopy We separate out frames as individual images from videos provided in [16].⁶ We download the first 10 narrow band imaging (NBI) videos in each class of lesion (hyperplastic, serrated, adenoma) and extract each frame as an individual image. Although the data is not independent in this case, we treat it as such for the purposes of our investigation.

A.3 Confidence Scores

Below we detail all confidence scores S implemented and evaluated in our investigation. There are additional approaches that were omitted from the main paper for the sake of brevity.

- SIRC: for a description of the score see the main paper. We use the whole of the ImageNet-200 *training* set to determine the values of μ_{S_2}, σ_{S_2} . For ImageNet-1k we randomly sample 250,000 images from the training set. Note that for all following methods that require ID data to find parameters, we use the same ID data as for SIRC. We investigate combinations of S_1, S_2 from the cartesian product $\{\text{MSP}, \text{DOCTOR}, \mathcal{H}\} \times \{\|\mathbf{z}\|_1, \text{Residual}\}$.
- Maximum Softmax Probability (MSP)[6]: a baseline score that takes the max value from the softmax $\pi_{\max} = \max_k \pi_k$.
- DOCTOR [3]: the original paper does not directly present it as such, but the confidence score is equivalent to $\|\boldsymbol{\pi}\|_2$.
- Softmax Entropy (\mathcal{H}): measures softmax uncertainty, $\mathcal{H}[\boldsymbol{\pi}] = -\sum_k \pi_k \log \pi_k$. We use $S = -\mathcal{H}[\boldsymbol{\pi}]$ to change it to a measure of confidence.
- l_1 -norm of the features: used in Gradnorm [9], $\|\mathbf{z}\|_1$.
- Residual: used in ViM [9], this score measures the component of the feature vector that is outside of a principal subspace defined using ID data, $\|\mathbf{z}^{P^\perp}\|_2$. We follow [19] in setting the dimensionality of the subspace to 1000 if the

⁴ <https://github.com/daintlab/unknown-detection-benchmarks>

⁵ https://pillow.readthedocs.io/en/stable/_modules/PIL/Image.html#Image.resize

⁶ http://www.depeca.uah.es/colonoscopy_dataset/

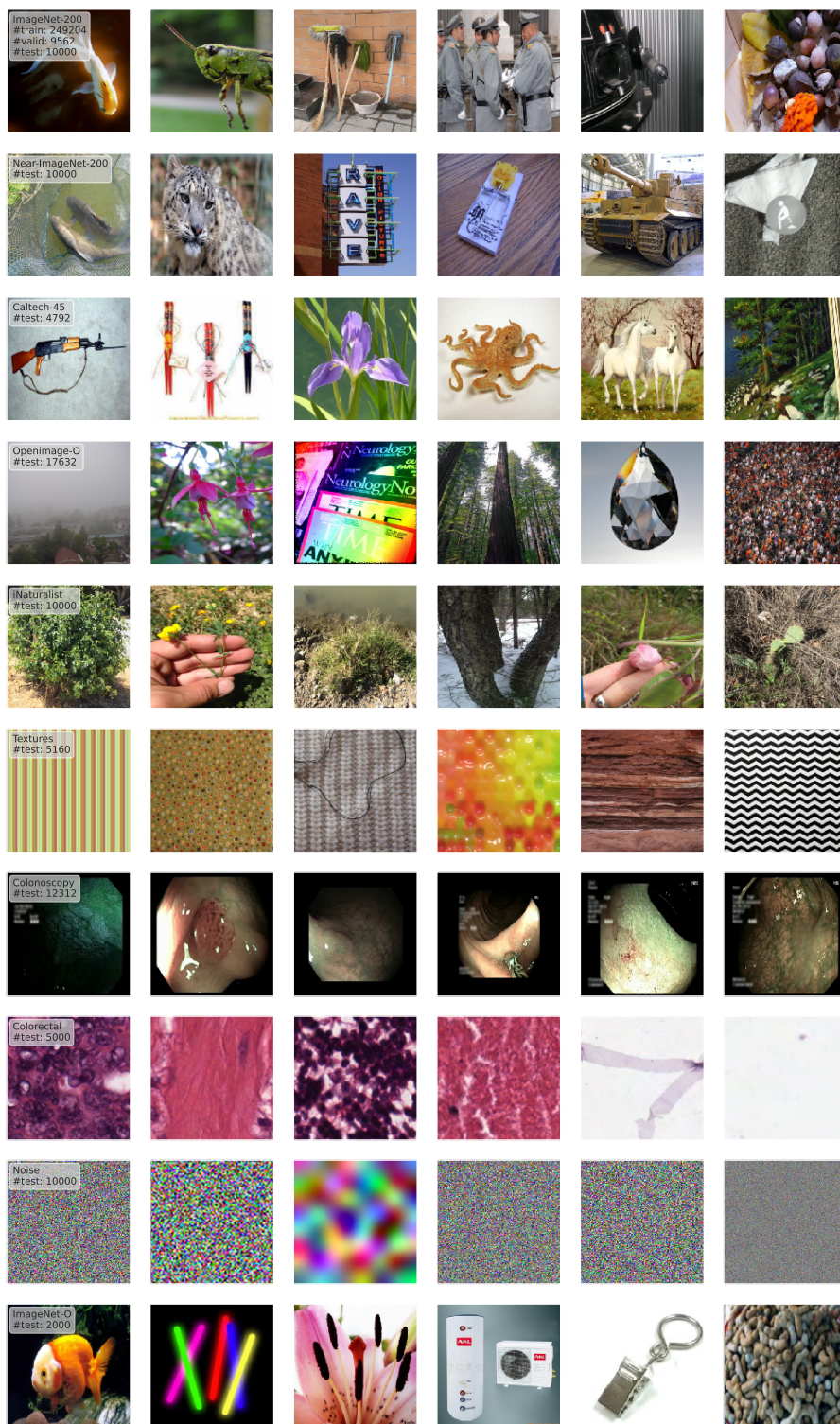


Fig. 1. Random examples from each ImageNet-scale dataset, with the #samples in each.

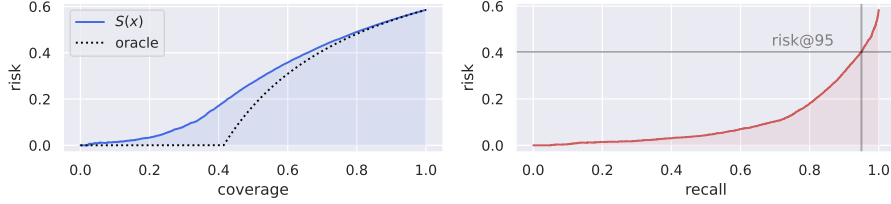


Fig. 2. Visualisations of different evaluation metrics for SCOD. We aim to minimise risk over different selection thresholds t . **Left:** Risk-Coverage curve (coverage is the proportion of all data accepted). We aggregate performance over t by taking the area under the curve. The oracle represents perfect separation of OOD, ID \times | ID \checkmark . **Right:** Risk-Recall curve. We consider both the area under the curve as well as risk@recall=0.95.

dimensionality of \mathbf{z} , $L > 1500$ and 512 otherwise. Like Entropy, we use the negative of the score $S = -\|\mathbf{z}^{P^\perp}\|_2$ as this score is meant to be higher for OOD data. Please refer to Wang et al. [19]’s paper for full details.

- Max Logit [5]: Max Logit is similar to MSP, but the score is taken from the logits before the softmax $v_{\max} = \max_k v_k$.
- Energy [15]: this score aggregates over all logit values as $\log \sum_k \exp v_k$.
- Gradnorm [9]: although this score was originally motivated by gradients, we can view it simply as the combination of two scores, $C = \|\boldsymbol{\pi} - \mathbf{1}/K\|_1 \|\mathbf{z}\|_1$.
- ViM [19]: this linearly combines Energy and Residual, $C = \log \sum_k \exp v_k - c \|\mathbf{z}^{P^\perp}\|_2$. The parameter c is given by the average value of Max Logit divided by the average value of Residual on ID data, which scales the importance of Residual to be similar to that of Energy in the combination.
- Mahalanobis [14]: this score involves building a classwise gaussian mixture model over the features with tied covariance matrix. The confidence is then calculated as $-\min_k (\mathbf{z} - \boldsymbol{\mu}_k)^T \hat{\boldsymbol{\Sigma}} (\mathbf{z} - \boldsymbol{\mu}_k)$. We use the approach in [1, 19] where only the final layer features are considered.

A.4 Evaluation Metrics

Other than the metrics specified in the main paper, we additionally use Area Under the Risk-Coverage Curve (AURC) \downarrow , from [2, 12]. It aggregates risk over all values of *coverage*, which is the proportion of all input data accepted. For AURC there exists an oracle curve, where OOD and ID \times are perfectly disjoint from ID \checkmark . AURC can be reduced either by lowering the oracle curve by reducing the number of ID \times (increasing baseline accuracy of f) or by better separating OOD, ID \times | ID \checkmark (better choice of g) and so bringing the curve closer to the oracle. Thus the metric is suitable for both training based, and post-hoc approaches. Fig. 2 illustrates graphically some of the metrics we use to evaluate SCOD.

B Additional Results

We provide more complete versions of the results presented in the main paper across all architectures and datasets.

B.1 AUROC and FPR@95

We present results across all post-hoc confidence scores in Appendix A.3 for all architectures. We also include $\text{mean} \pm 2\text{std.}$ for experiments with multiple training runs. SIRC performs as expected in all cases – a negligible reduction in $\text{ID}\times$ | $\text{ID}\checkmark$ in exchange for a meaningful uplift in OOD | $\text{ID}\checkmark$ compared to only using S_1 . DOCTOR in general performs somewhere in between MSP and $-\mathcal{H}$, both individually and when used in SIRC, so we relegate it to the appendix. We note that Residual and Mahalanobis perform much better only for ResNetV2-101 (these results are inline with [19]). This may be due to the fact that BiT uses Weight Standardisation and Group Normalisation when training, rather than standard Batch Normalisation. Mukhoti et al. [17] show that limiting the Lipschitz constant of the network during training improves the OOD detection performance of gaussian mixture models, which may be also what is occurring in this example. The Mahalanobis detector performs poorly outside ResNetV2-101 otherwise. There is non-negligible variance between training runs on a number of OOD datasets, highlighting the need to perform multiple training runs. Some datasets (e.g. Noise, Colorectal, SVHN), have especially high variation.

B.2 Varying α and β

We plot performance under varying α and β for all 3 ImageNet-200 architectures (Figs. 3 to 5). We also present the $\text{mean} \pm \text{std.}$ The ability of SIRC to perform consistently better than the baseline generalises across the 3 different CNN architectures. We note that differences in AURC are harder to distinguish, due to the metric considering the proportion of all input data accepted, rather than just the recall of $\text{ID}\checkmark$. The behaviour, however, is similar to AURR in terms of relative performance to the baseline, so we omit AURC from the main results.

B.3 SCOD vs OOD Detection

We plot the change in $\% \text{FPR@95}\downarrow$ relative to the MSP baseline for all architectures and confidence scores (Figs. 6 to 10). The behaviour is as discussed in the main paper, with methods designed for OOD detection achieving gains over the baseline for OOD detection by sacrificing their ability to separate $\text{ID}\times$ | $\text{ID}\checkmark$.

B.4 Plotting S_2 against S_1

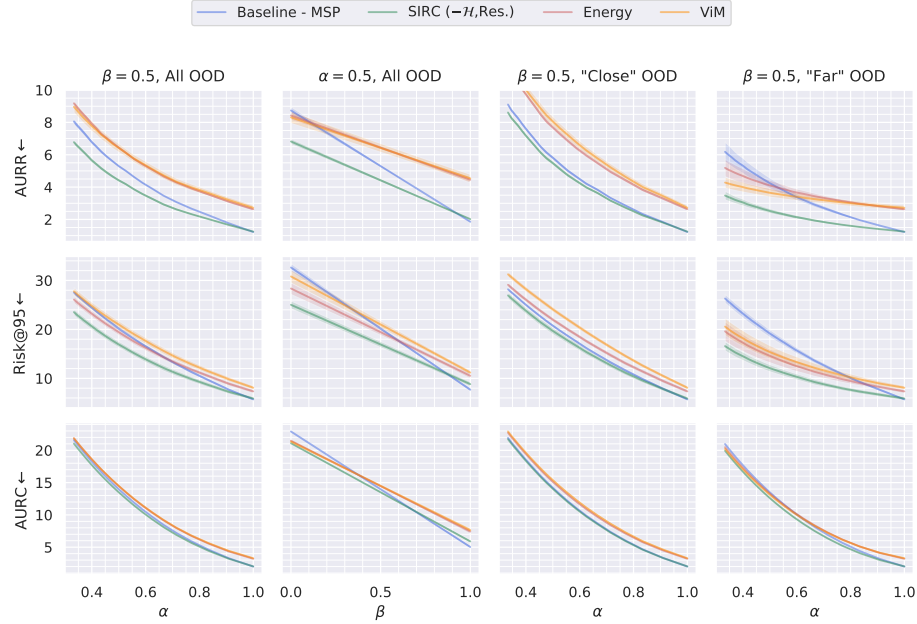
We plot different SIRC combinations on the S_1, S_2 -plane for different experimental configurations (Figs. 11 to 14). If there are multiple training runs, we plot

Table 1. Full %AUROC and %FPR@95 results for all models trained on ImageNet-200. We show the mean \pm 2std. over 5 independent training runs. **Bold** indicates best performance, underline 2nd or 3rd best.

Model	Method	ID \times		OOD mean		Near-IN-200		Caltech-45		OpenImage-0		INaturalist	
		%AUROC	%FPR95	%AUROC	%FPR95	%AUROC	%FPR95	%AUROC	%FPR95	%AUROC	%FPR95	%AUROC	%FPR95
ResNet-50 ID %Error: 1.0/1	SIRC	(MSP, ℓ_1)	90.34 \pm 0.2	52.70 \pm 0.2	91.51	40.27	85.56 \pm 0.6	59.76 \pm 0.9	91.36 \pm 0.6	41.44 \pm 0.0	92.28 \pm 0.5	41.36 \pm 0.2	94.80 \pm 0.3
		(MSP,Res)	90.43 \pm0.3	<u>52.10 \pm0.0</u>	92.56	34.98	85.52 \pm 0.6	60.03 \pm 0.4	91.19 \pm 0.6	42.27 \pm 0.2	92.57 \pm 0.6	39.95 \pm 0.3	94.10 \pm 0.3
		(DR, ℓ_1)	90.29 \pm 0.3	52.54 \pm 0.4	91.83	37.08	85.68 \pm 0.6	<u>58.05 \pm0.9</u>	91.67 \pm 0.5	37.81 \pm 0.2	92.59 \pm 0.4	<u>37.82 \pm0.6</u>	<u>95.18 \pm0.3</u>
		(DR,Res)	90.40 \pm 0.4	51.81 \pm0.9	92.83	<u>31.76</u>	85.62 \pm 0.6	58.19 \pm 0.2	91.44 \pm 0.6	38.92 \pm 0.2	92.87 \pm 0.6	36.36 \pm0.6	94.32 \pm 0.4
		(-H, ℓ_1)	90.00 \pm 0.4	54.26 \pm 0.2	92.24	35.85	85.88 \pm 0.6	58.50 \pm 0.3	92.19 \pm 0.5	36.08 \pm0.2	92.87 \pm 0.4	37.83 \pm 0.3	95.38 \pm0.2
		(-H,Res)	90.13 \pm 0.4	54.01 \pm 0.2	93.36	30.05	<u>85.85 \pm0.6</u>	58.93 \pm 0.3	92.11 \pm 0.5	<u>36.76 \pm0.0</u>	93.25 \pm0.6	36.36 \pm 0.4	94.82 \pm 0.3
	DOCTOR	MSP	<u>90.41 \pm0.3</u>	52.13 \pm 0.0	91.00	43.25	85.59 \pm 0.6	59.74 \pm 0.0	91.13 \pm 0.6	42.72 \pm 0.8	91.95 \pm 0.5	43.55 \pm 0.4	94.23 \pm 0.3
		DOCTOR	90.39 \pm 0.3	<u>51.87 \pm0.6</u>	91.26	40.22	85.73 \pm 0.6	57.89 \pm0.8	91.41 \pm 0.5	39.22 \pm 0.2	92.20 \pm 0.5	40.22 \pm 0.2	94.51 \pm 0.3
		-H	90.07 \pm 0.4	54.05 \pm 0.9	91.81	38.24	85.91 \pm0.6	<u>58.47 \pm0.3</u>	92.01 \pm 0.5	37.20 \pm 0.6	92.59 \pm 0.5	40.10 \pm 0.9	<u>94.90 \pm0.3</u>
		ℓ_1	48.06 \pm 1.1	94.70 \pm 1.4	78.22	58.70	82.27 \pm 0.7	94.58 \pm 0.5	70.26 \pm 1.6	77.53 \pm 1.8	72.23 \pm 0.4	71.31 \pm 2.6	83.65 \pm 2.7
		Residual	47.59 \pm 1.8	96.45 \pm 1.1	58.45	78.97	44.30 \pm 1.1	96.79 \pm 0.4	47.76 \pm 1.4	94.83 \pm 0.9	50.65 \pm 4.0	86.85 \pm 2.2	40.07 \pm 0.3
	Max Logit	Energy	83.21 \pm 0.6	65.16 \pm 0.4	92.33	<u>34.15</u>	82.68 \pm 0.7	65.37 \pm 0.6	92.48 \pm0.6	36.50 \pm 0.1	91.49 \pm 0.4	43.27 \pm 0.1	94.57 \pm 0.3
		Energy	82.05 \pm 0.6	69.79 \pm 0.9	92.06	35.32	81.90 \pm 0.7	68.70 \pm 0.4	<u>92.15 \pm0.6</u>	38.62 \pm 0.9	90.92 \pm 0.4	46.28 \pm 0.3	94.13 \pm 0.4
		GradNorm	60.17 \pm 1.5	87.88 \pm 2.5	85.22	44.41	62.90 \pm 0.7	86.89 \pm 0.8	81.11 \pm 1.7	59.23 \pm 0.3	81.09 \pm 1.8	57.82 \pm 0.7	91.00 \pm 1.4
		VIM	80.62 \pm 0.7	78.13 \pm 0.3	92.34	38.14	78.90 \pm 0.8	80.30 \pm 0.2	90.54 \pm 0.7	54.70 \pm 0.5	91.87 \pm 1.2	43.84 \pm 0.6	90.13 \pm 1.8
		Mahal	49.96 \pm 2.0	96.36 \pm 0.9	61.66	78.92	46.57 \pm 1.5	96.83 \pm 0.5	50.34 \pm 1.6	95.01 \pm 0.6	63.66 \pm 0.7	86.30 \pm 2.0	47.42 \pm 0.5
ResNet-50 ID %Error: 1.0/1	SIRC	(MSP, ℓ_1)	90.34 \pm 0.2	52.70 \pm 0.2	93.64 \pm 0.7	32.02 \pm 0.3	95.93 \pm 1.0	25.33 \pm 0.4	95.84 \pm 0.3	24.39 \pm 0.7	90.72 \pm 0.0	49.63 \pm 0.8	83.44 \pm 0.9
		(MSP,Res)	90.43 \pm0.3	<u>52.10 \pm0.0</u>	96.00 \pm 0.5	19.81 \pm 0.1	95.52 \pm 0.7	27.31 \pm 0.3	95.32 \pm 0.0	26.97 \pm 0.5	98.21 \pm 0.2	10.97 \pm 0.7	84.62 \pm 0.9
		(DR, ℓ_1)	90.29 \pm 0.3	52.54 \pm 0.4	94.01 \pm 0.7	28.62 \pm 0.6	96.34 \pm 1.0	20.94 \pm 0.6	96.28 \pm 0.2	20.30 \pm 0.3	91.08 \pm 0.8	47.75 \pm 0.6	83.64 \pm 1.0
		(DR,Res)	90.40 \pm 0.4	51.81 \pm0.9	<u>96.28 \pm0.5</u>	<u>17.29 \pm0.0</u>	95.82 \pm 0.6	23.07 \pm 0.7	95.62 \pm 1.1	23.40 \pm 0.8	<u>98.63 \pm0.9</u>	<u>7.23 \pm0.0</u>	<u>84.90 \pm0.9</u>
		(-H, ℓ_1)	90.00 \pm 0.4	54.26 \pm 0.2	91.38 \pm 0.7	27.38 \pm 0.7	99.07 \pm 0.8	16.87 \pm 0.4	96.71 \pm 1.2	18.71 \pm 0.3	91.74 \pm 0.2	45.84 \pm 0.7	84.01 \pm 0.9
		(-H,Res)	90.13 \pm 0.4	54.01 \pm 0.2	<u>96.68 \pm0.5</u>	<u>15.70 \pm0.1</u>	96.72 \pm 0.6	18.10 \pm 0.3	96.41 \pm 0.6	20.42 \pm 0.6	<u>99.02 \pm1.5</u>	<u>1.89 \pm0.5</u>	<u>85.33 \pm0.9</u>
	DOCTOR	MSP	<u>90.41 \pm0.3</u>	52.13 \pm 0.0	92.88 \pm 0.8	36.61 \pm 0.1	95.75 \pm 0.8	26.52 \pm 0.2	94.86 \pm 0.5	30.28 \pm 0.6	89.33 \pm 0.5	56.83 \pm 0.2	83.29 \pm 0.9
		DOCTOR	90.39 \pm 0.3	<u>51.87 \pm0.6</u>	93.16 \pm 0.8	33.46 \pm 0.6	96.14 \pm 0.8	22.07 \pm 0.3	95.16 \pm 0.5	27.21 \pm 0.4	89.51 \pm 0.5	54.83 \pm 0.4	83.47 \pm 0.9
		-H	90.07 \pm 0.4	54.05 \pm 0.9	93.77 \pm 0.8	30.79 \pm 0.3	96.87 \pm 0.7	17.55 \pm 0.4	95.93 \pm 0.2	23.43 \pm 0.6	90.47 \pm 0.4	51.63 \pm 0.7	83.89 \pm 0.9
		ℓ_1	48.06 \pm 1.1	94.70 \pm 1.4	88.90 \pm 1.5	39.67 \pm 0.6	76.97 \pm 0.7	82.24 \pm 0.4	97.28 \pm 0.3	14.64 \pm 0.3	97.36 \pm 0.6	13.51 \pm 0.1	63.00 \pm 1.7
		Residual	47.59 \pm 1.8	96.45 \pm 1.1	82.84 \pm 2.4	46.63 \pm 0.8	38.09 \pm 1.9	90.64 \pm 0.4	53.93 \pm 0.2	88.78 \pm 0.3	91.31 \pm 0.4	20.92 \pm 0.2	68.04 \pm 2.7
	Max Logit	Energy	83.21 \pm 0.6	65.16 \pm 0.4	95.44 \pm 0.8	22.04 \pm 0.2	97.65 \pm0.7	13.56 \pm0.0	<u>98.93 \pm1.0</u>	<u>5.83 \pm0.4</u>	94.73 \pm 0.5	31.53 \pm 0.8	82.98 \pm 0.9
		Energy	82.05 \pm 0.6	69.79 \pm 0.9	95.37 \pm 0.8	22.50 \pm 0.2	<u>97.51 \pm0.8</u>	<u>14.19 \pm0.5</u>	99.07 \pm1.0	<u>5.00 \pm0.5</u>	94.93 \pm 0.4	29.05 \pm 0.8	82.52 \pm 0.9
		GradNorm	60.17 \pm 1.5	87.88 \pm 2.5	93.00 \pm 1.1	26.57 \pm 0.3	90.54 \pm 0.6	42.85 \pm 0.2	<u>98.98 \pm1.1</u>	4.98 \pm0.2	97.59 \pm 0.4	13.05 \pm 0.9	70.78 \pm 1.8
		VIM	80.62 \pm 0.7	78.13 \pm 0.3	98.46 \pm0.4	<u>7.62 \pm0.1</u>	92.45 \pm 1.4	44.55 \pm 0.4	98.04 \pm 1.3	8.81 \pm 0.2	99.82 \pm0.1	0.51 \pm 0.3	88.85 \pm0.9
		Mahal	49.96 \pm 2.0	96.36 \pm 0.9	84.64 \pm 2.1	46.98 \pm 0.2	41.02 \pm 1.2	99.70 \pm 0.3	57.88 \pm 0.2	88.37 \pm 0.1	94.08 \pm 0.4	20.45 \pm 0.3	69.29 \pm 2.5
MobileNet-V2 ID %Error: 2.1/5	SIRC	(MSP, ℓ_1)	89.53 \pm 0.3	55.51 \pm 1.0	92.27	34.82	84.78 \pm 0.3	61.33 \pm 1.1	90.46 \pm 0.3	43.03 \pm 0.9	91.27 \pm 0.4	44.05 \pm 1.1	94.20 \pm 0.8
		(MSP,Res)	89.67 \pm0.3	<u>55.10 \pm1.4</u>	91.78	38.56	84.84 \pm 0.4	61.18 \pm 0.3	90.25 \pm 0.4	44.42 \pm 0.3	91.20 \pm 0.5	44.83 \pm 1.8	93.22 \pm 0.9
		(DR, ℓ_1)	89.40 \pm 0.2	56.49 \pm 0.2	<u>92.66</u>	<u>32.30</u>	84.90 \pm 0.3	61.26 \pm 0.8	90.82 \pm 0.3	<u>40.52 \pm0.2</u>	91.61 \pm 0.4	42.36 \pm1.0	<u>91.63 \pm0.7</u>
		(DR,Res)	<u>89.60 \pm0.3</u>	55.69 \pm 0.2	92.08	36.21	84.98 \pm 0.3	60.92 \pm 1.1	90.58 \pm 0.4	41.98 \pm 0.3	91.51 \pm 0.5	<u>43.22 \pm1.0</u>	93.40 \pm 0.9
		(-H, ℓ_1)	88.90 \pm 0.2	58.64 \pm 0.1	92.92	32.16	84.96 \pm 0.2	62.72 \pm 0.9	<u>93.35 \pm0.3</u>	39.69 \pm0.4	<u>91.82 \pm0.4</u>	43.99 \pm 1.4	94.74 \pm0.7
		(-H,Res)	89.12 \pm 0.3	57.85 \pm 0.1	<u>92.69</u>	<u>34.20</u>	85.08 \pm0.2	62.06 \pm 0.8	91.33 \pm 0.3	39.63 \pm0.2	91.93 \pm0.4	44.01 \pm 0.8	93.74 \pm 0.7
	DOCTOR	MSP	<u>89.64 \pm0.3</u>	55.03 \pm1.3	91.54	37.73	84.84 \pm 0.3	61.03 \pm 0.2	90.17 \pm 0.3	44.77 \pm 1.7	90.91 \pm 0.5	46.34 \pm 1.8	93.57 \pm 0.9
		DOCTOR	89.57 \pm 0.2	55.48 \pm 0.1	91.86	37.43	84.99 \pm 0.3	60.61 \pm0.3	90.52 \pm 0.3	42.46 \pm 0.1	91.20 \pm 0.5	44.96 \pm 1.0	93.91 \pm 0.8
		-H	89.02 \pm 0.2	58.43 \pm 0.2	92.37	36.04	<u>85.03 \pm0.2</u>	62.53 \pm 0.4	91.16 \pm 0.3	41.27 \pm 0.0	91.54 \pm 0.4	46.11 \pm 1.4	92.44 \pm 0.7
		ℓ_1	53.56 \pm 0.7	93.40 \pm 0.4	81.06	53.50	56.05 \pm 0.7	92.65 \pm 0.5	75.15 \pm 1.4	73.17 \pm 2.2	74.05 \pm 1.4	68.93 \pm 1.7	88.03 \pm 1.7
		Residual	41.99 \pm 0.8	97.30 \pm 0.3	41.42	94.11	42.46 \pm 0.7	97.37 \pm 0.3	49.98 \pm 1.2	96.70 \pm 0.9	44.63 \pm 1.1	94.39 \pm 0.6	92.14 \pm 0.7
	Max Logit	Energy	83.14 \pm 0.6	63.85 \pm 1.8	92.08	34.64	81.75 \pm 0.4	67.36 \pm 0.3	91.40 \pm0.2	42.14 \pm 0.1	89.70 \pm 0.8	50.66 \pm 0.3	92.63 \pm 1.0
		Energy	81.87 \pm 0.7	67.98 \pm 0.2	91.68	36.68	80.87 \pm 0.4	70.81 \pm 1.2	90.93 \pm 0.3	45.77 \pm 1.9	88.86 \pm 0.8	54.53 \pm 3.1	91.76 \pm 1.0
		GradNorm	65.27 \pm 1.1	85.73 \pm 1.1	87.25	40.67	66.07 \pm 0.7	85.13 \pm 1.0	83.94 \pm 1.0	56.57 \pm 0.9	81.94 \pm 1.2	58.20 \pm 1.7	90.73 \pm 1.3
		VIM	80.21 \pm 0.4	74.36 \pm 0.1	89.46	51.97	79.15 \pm 0.3	75.78 \pm 1.4	89.17 \pm 0.4	58.45 \pm 0.4	87.66 \pm 1.0	59.54 \pm 1.2	81.93 \pm 0.3
		Mahal	44.44 \pm 1.0	97.14 \pm 0.6	43.65	94.20	44.57 \pm 0.7	97.23 \pm 0.4	42.82 \pm 1.1	96.64 \pm 0.8	48.03 \pm 1.2	94.11 \pm 0.8	
MobileNet-V2 ID %Error: 2.1/5	SIRC	(MSP, ℓ_1)	90.22 \pm 0.8	52.41 \pm 0.7	91.68	38.83	85.28 \pm 0.3	59.35 \pm 1.0	91.33 \pm 0.5	40.80 \pm 0.2	91.88 \pm 0.4	43.34 \pm 1.4	93.52 \pm 0.8
		(MSP,Res)	90.20 \pm 0.8	52.42 \pm 0.7	<u>92.81</u>	<u>32.68</u>	85.12 \pm 0.3	59.65 \pm 1.1	91.29 \pm 0.5	40.51 \pm 0.6	92.55 \pm 0.3	39.24 \pm 1.0	93.42 \pm 0.8
		(DR, ℓ_1)	90.21 \pm 0.8	52.36 \pm 0.9	91.93	38.08	85.28 \pm 0.4	59.49 \pm 1.0	91.63 \pm 0.5	37.15 \pm 1.1	92.14 \pm 0.8	39.89 \pm 1.1	93.85 \pm 0.7
		(DR,Res)	90.18 \pm 0.8	52.46 \pm 0.3	<u>93.01</u>	<u>29.68</u>	85.28 \pm 0.4	59.78 \pm 1.2	91.53 \pm 0				

Table 2. %AUROC and %FPR@95 results for single pre-trained ImageNet-1k models.

Model	Method	ID \times	OOD mean		Openimage-O		INaturalist		Textures		Colonoscopy		Colorectal		Noise		ImageNet-O			
		AUROC \uparrow FPR@95 \downarrow	AUROC \uparrow FPR@95 \downarrow	AUROC \uparrow FPR@95 \downarrow	AUROC \uparrow FPR@95 \downarrow	AUROC \uparrow FPR@95 \downarrow	AUROC \uparrow FPR@95 \downarrow	AUROC \uparrow FPR@95 \downarrow	AUROC \uparrow FPR@95 \downarrow	AUROC \uparrow FPR@95 \downarrow	AUROC \uparrow FPR@95 \downarrow	AUROC \uparrow FPR@95 \downarrow	AUROC \uparrow FPR@95 \downarrow	AUROC \uparrow FPR@95 \downarrow	AUROC \uparrow FPR@95 \downarrow	AUROC \uparrow FPR@95 \downarrow	AUROC \uparrow FPR@95 \downarrow			
ResNetV2-101 ID %Error: 22.63	SIRC	(MSP, $\ z\ _1$)	86.17	63.37	90.08	37.09	90.25	47.09	94.37	29.13	88.74	43.84	97.10	16.69	93.94	30.28	99.26	4.32	66.93	88.30
		(MSP, Res.)	86.34	62.32	92.89	22.23	92.60	33.42	94.61	27.15	96.80	10.04	97.18	14.69	96.88	14.06	99.99	0.00	73.00	74.10
		(DR, $\ z\ _1$)	85.36	66.04	90.44	35.13	90.35	47.22	94.75	26.89	89.19	41.72	97.13	15.22	94.85	23.82	99.48	2.93	67.73	88.10
		(DR, Res.)	85.55	64.66	93.34	22.76	93.25	31.27	94.99	24.30	97.15	8.06	97.36	12.46	97.52	10.10	99.99	0.00	73.12	73.10
		($-H$, $\ z\ _1$)	83.43	69.62	90.98	36.74	90.46	51.57	94.53	33.11	89.06	47.03	98.37	9.66	95.71	23.18	99.48	3.80	69.25	88.80
		($-H$, Res.)	83.50	68.91	93.67	25.61	92.75	40.20	94.63	33.51	97.05	10.19	98.53	9.02	98.12	8.62	100.00	0.00	77.10	77.70
	DOCTOR	MSP	86.35	61.93	89.16	41.81	90.13	47.48	93.70	32.96	87.04	51.90	97.47	14.92	91.55	44.28	97.57	12.53	66.87	88.60
		$-H$	85.67	64.49	89.57	40.48	90.33	47.64	93.95	32.04	87.27	52.64	97.89	12.61	92.34	39.78	97.94	9.74	67.25	88.90
		$\ z\ _1$	83.49	69.09	90.25	41.32	90.23	54.09	93.80	38.97	87.47	54.92	98.52	8.69	94.02	34.70	98.46	7.73	69.49	88.16
		Residual	50.18	94.86	85.59	50.02	80.17	68.38	76.76	80.65	72.67	10.99	67.60	98.55	95.43	25.34	99.95	0.00	81.57	66.20
		Max Logit	77.25	71.07	90.24	41.72	88.11	59.64	91.87	48.90	87.08	55.70	99.04	4.64	96.25	25.98	98.79	6.47	70.57	90.70
		Energy	74.68	77.15	89.41	45.23	85.86	68.88	89.27	59.78	85.85	61.61	99.19	3.17	96.56	23.82	98.83	6.83	70.33	92.55
	ViM	GradNorm	64.64	88.00	84.85	46.15	73.53	76.05	87.99	53.65	85.04	50.85	94.56	29.39	95.94	22.56	99.82	0.67	57.05	89.85
		ViM	70.30	86.87	94.95	25.61	92.08	41.79	91.68	47.40	90.17	3.39	95.59	29.88	99.30	1.26	100.00	0.00	86.80	55.55
		Mahal	56.82	93.95	89.62	46.82	86.43	61.39	85.09	73.14	98.19	9.19	77.36	98.76	96.09	22.40	99.88	0.00	84.28	62.85
DenseNet-121 ID %Error: 25.58	SIRC	(MSP, $\ z\ _1$)	85.99	63.14	89.52	33.01	90.93	39.50	95.36	21.61	89.65	37.34	96.79	17.15	96.06	20.94	99.74	1.10	58.10	93.45
		(MSP, Res.)	85.97	63.33	90.05	31.62	91.17	38.58	94.08	27.50	93.38	22.93	96.18	19.71	95.51	23.60	99.67	0.60	68.45	98.45
		(DR, $\ z\ _1$)	85.77	64.51	90.00	30.09	91.55	36.00	95.48	17.81	90.32	33.12	97.10	14.22	96.88	15.68	99.79	0.83	58.36	93.00
		(DR, Res.)	85.72	65.09	90.43	28.80	91.72	35.28	94.50	24.32	92.85	19.42	96.30	16.79	96.21	17.90	99.62	0.54	60.83	87.35
		($-H$, $\ z\ _1$)	84.90	67.31	90.83	28.43	92.41	34.47	96.52	16.28	91.05	32.85	97.89	9.86	97.79	12.36	99.83	0.68	60.31	92.70
		($-H$, Res.)	84.85	67.87	91.46	26.46	92.64	34.09	95.67	20.33	94.42	19.07	97.45	11.37	97.56	12.30	99.79	0.39	63.77	92.70
	DOCTOR	MSP	86.11	62.67	88.81	36.77	90.26	43.08	94.26	27.56	88.31	43.72	96.90	17.10	94.44	30.72	99.55	1.69	57.97	93.55
		$-H$	85.93	63.43	89.28	34.17	90.82	39.93	94.83	23.95	88.85	41.01	97.33	13.52	95.24	25.94	99.64	1.37	58.23	93.45
		$\ z\ _1$	84.97	66.76	90.39	30.68	91.91	37.18	95.83	19.56	90.08	37.56	97.45	9.42	97.00	17.20	99.76	1.15	60.17	92.70
		Residual	47.53	94.93	78.50	53.82	69.94	70.15	89.06	39.14	84.61	49.73	58.85	89.84	92.89	34.40	99.88	0.49	61.90	92.70
		Max Logit	51.52	94.26	71.96	66.47	69.78	78.27	61.14	93.61	90.21	33.64	37.94	99.37	75.49	70.74	97.32	14.40	71.83	75.25
		Energy	77.97	71.35	91.62	28.87	92.10	38.48	96.07	20.57	91.59	34.32	98.20	8.77	98.62	6.48	99.89	0.49	64.77	92.95
	ViM	GradNorm	76.13	75.77	91.47	30.02	91.54	42.66	95.60	23.50	91.39	35.43	97.87	11.18	98.86	5.02	99.91	0.39	65.12	92.95
		ViM	55.44	92.10	85.31	42.04	78.97	58.55	93.87	25.24	89.62	37.81	81.08	68.36	97.63	13.10	99.96	0.02	56.04	94.01
		Mahal	70.16	88.53	89.58	47.81	88.40	56.49	88.74	66.34	96.64	17.69	82.83	89.17	95.19	31.76	99.57	0.01	75.66	73.20
			57.28	94.10	68.90	81.55	69.02	86.67	49.94	97.79	82.79	55.43	66.51	96.97	58.34	96.34	75.13	68.35	80.53	69.30

**Fig. 3.** Varying α and β for ResNet-50 (ImageNet-200) (values $\times 10^2$).

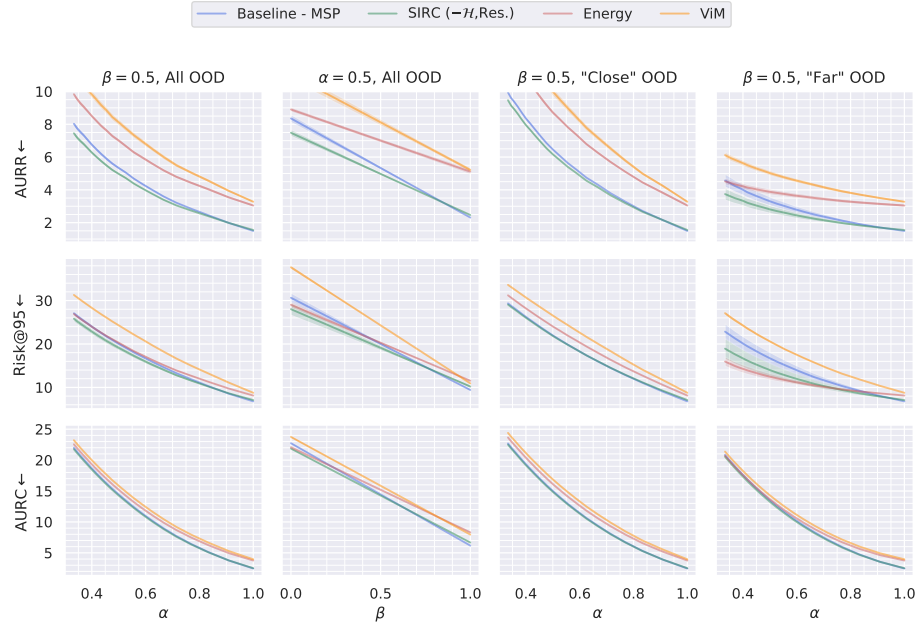


Fig. 4. Varying α and β for MobileNetV2 (ImageNet-200) (values $\times 10^2$).

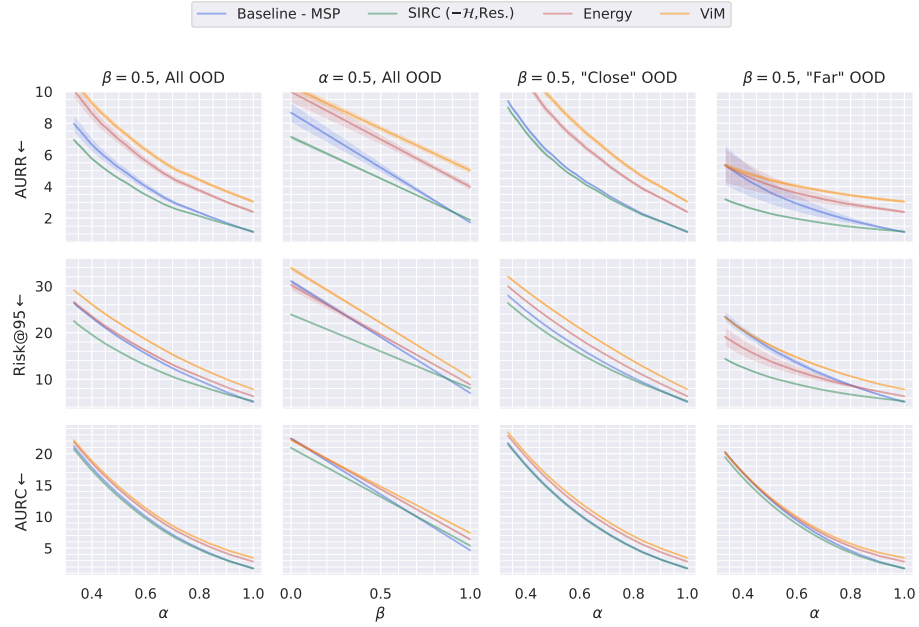


Fig. 5. Varying α and β for DenseNet-121 (ImageNet-200) (values $\times 10^2$).

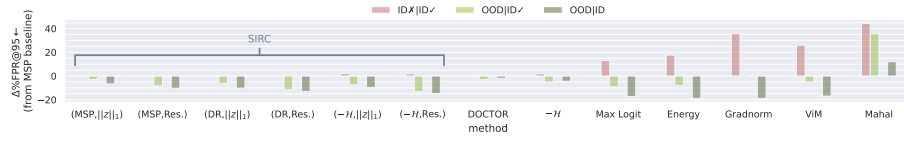


Fig. 6. ResNet-50 (ImageNet-200), comparing the change in %FPR@95 relative to the MSP baseline for different detection methods and data groups.

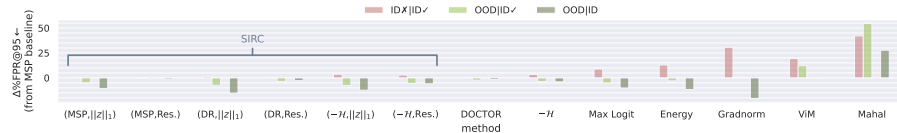


Fig. 7. MobileNetV2 (ImageNet-200), comparing the change in %FPR@95 relative to the MSP baseline for different detection methods and data groups.

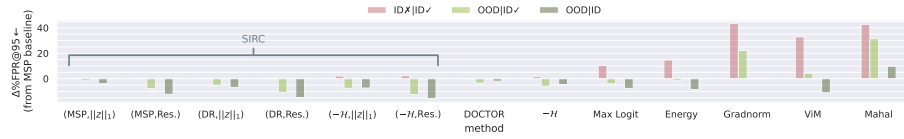


Fig. 8. DenseNet-121 (ImageNet-200), comparing the change in %FPR@95 relative to the MSP baseline for different detection methods and data groups.

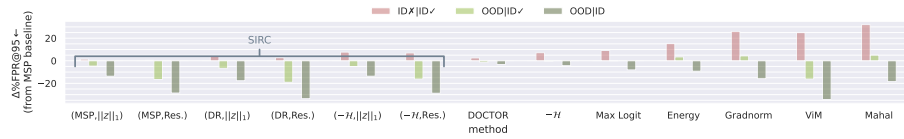


Fig. 9. ResNetV2-101 (ImageNet-1k), comparing the change in %FPR@95 relative to the MSP baseline for different detection methods and data groups.

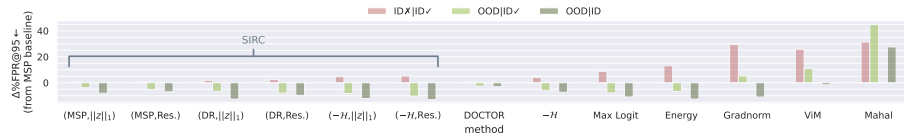


Fig. 10. DenseNet-121 (ImageNet-1k), comparing the change in %FPR@95 relative to the MSP baseline for different detection methods and data groups.

the distributions corresponding to the outputs of the 1st run. Decision contours corresponding to the default parameter setting for SIRC are also overlayed. We note that the inconsistency of Residual can be observed here, where in some cases the OOD distribution is much lower than ID, whilst in others, there is almost complete overlap. In the case of MobileNetV2 on iNaturalist it is in fact higher for OOD than ID, although the nature of SIRC means that it is robust to such S_2 failure (as discussed in the main paper).

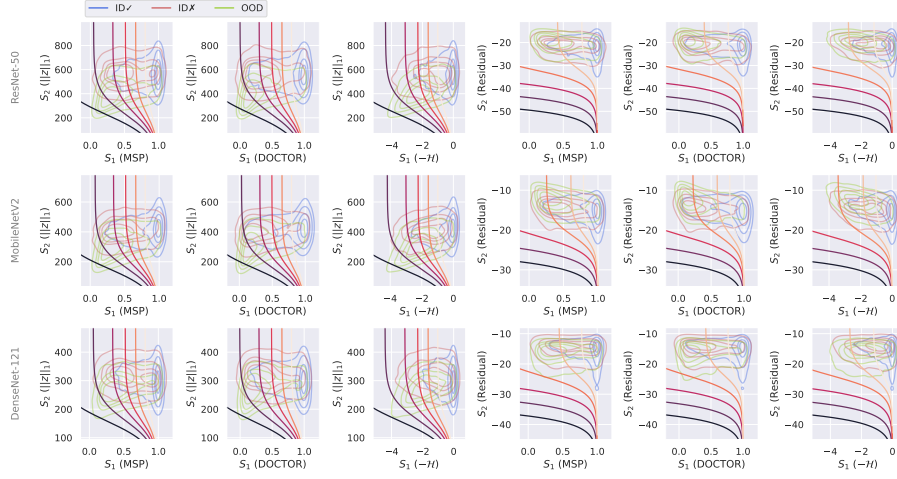


Fig. 11. SIRC combinations on the S_1, S_2 -plane, ID: ImageNet-200, OOD: iNaturalist.

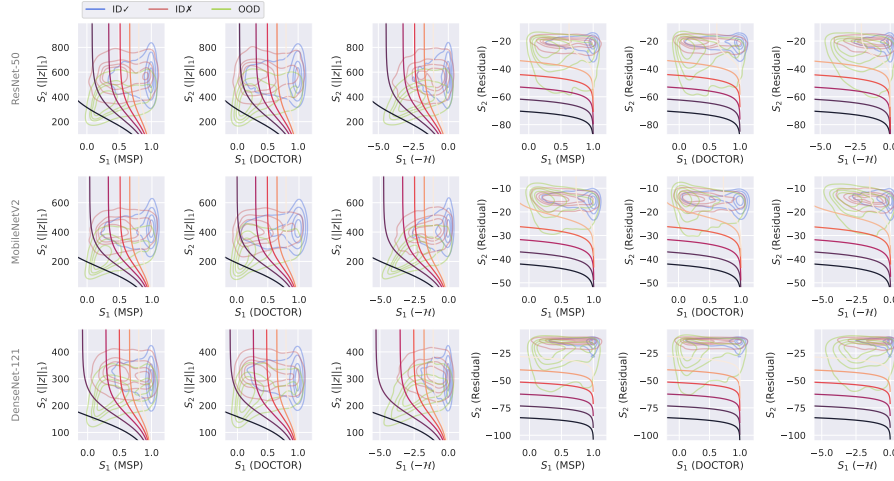


Fig. 12. SIRC combinations on the S_1, S_2 -plane, ID: ImageNet-200, OOD: Textures.

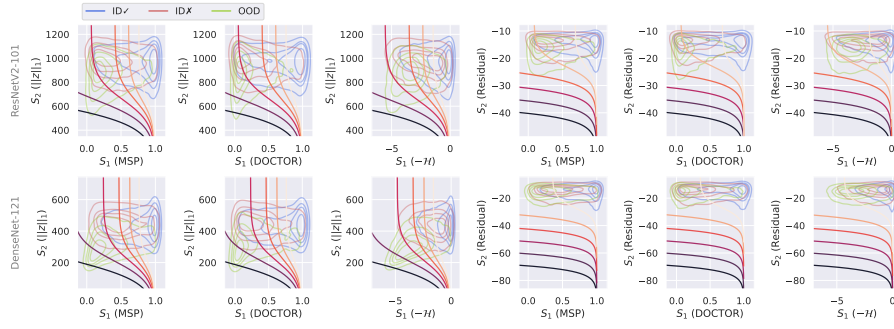


Fig. 13. SIRC combinations on the S_1, S_2 -plane, ID ImageNet-1k, OOD: iNaturalist.

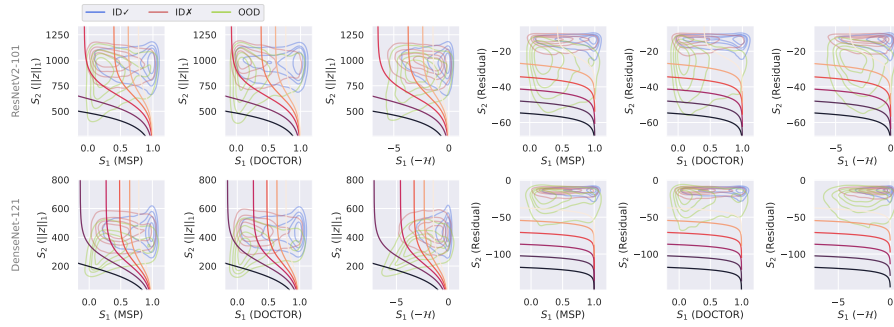


Fig. 14. SIRC combinations on the S_1, S_2 -plane, ID ImageNet-1k, OOD: Textures.

References

- [1] Fort, S., Ren, J., Lakshminarayanan, B.: Exploring the limits of out-of-distribution detection. In: NeurIPS (2021)
- [2] Geifman, Y., El-Yaniv, R.: Selective classification for deep neural networks. In: NIPS (2017)
- [3] Granese, F., Romanelli, M., Gorla, D., Palamidessi, C., Piantanida, P.: Doctor: A simple method for detecting misclassification errors. In: NeurIPS (2021)
- [4] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) pp. 770–778 (2016)
- [5] Hendrycks, D., Basart, S., Mazeika, M., Mostajabi, M., Steinhardt, J., Song, D.X.: Scaling out-of-distribution detection for real-world settings. arXiv: Computer Vision and Pattern Recognition (2020)
- [6] Hendrycks, D., Gimpel, K.: A baseline for detecting misclassified and out-of-distribution examples in neural networks. ArXiv abs/1610.02136 (2017)
- [7] Hendrycks, D., Zhao, K., Basart, S., Steinhardt, J., Song, D.X.: Natural adversarial examples. 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 15257–15266 (2021)
- [8] Huang, G., Liu, Z., Weinberger, K.Q.: Densely connected convolutional networks. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) pp. 2261–2269 (2017)
- [9] Huang, R., Geng, A., Li, Y.: On the importance of gradients for detecting distributional shifts in the wild. In: NeurIPS (2021)
- [10] Huang, R., Li, Y.: Mos: Towards scaling out-of-distribution detection for large semantic space. 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 8706–8715 (2021)
- [11] Kather, J.N., Weis, C.A., Bianconi, F., Melchers, S.M., Schad, L.R., Gaiser, T., Marx, A., Zöllner, F.G.: Multi-class texture analysis in colorectal cancer histology. Scientific Reports 6 (2016)
- [12] Kim, J., Koo, J., Hwang, S.: A unified benchmark for the unknown detection capability of deep neural networks. ArXiv abs/2112.00337 (2021)
- [13] Kolesnikov, A., Beyer, L., Zhai, X., Puigcerver, J., Yung, J., Gelly, S., Houlsby, N.: Big transfer (bit): General visual representation learning. In: ECCV (2020)
- [14] Lee, K., Lee, K., Lee, H., Shin, J.: A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In: NeurIPS (2018)
- [15] Liu, W., Wang, X., Owens, J.D., Li, Y.: Energy-based out-of-distribution detection. ArXiv abs/2010.03759 (2020)
- [16] Mesejo, P., Pizarro, D., Abergel, A., Rouquette, O.Y., Béorchia, S., Poincloux, L., Bartoli, A.: Computer-aided classification of gastrointestinal lesions in regular colonoscopy. IEEE transactions on medical imaging (2016)

- [17] Mukhoti, J., Kirsch, A., van Amersfoort, J.R., Torr, P.H.S., Gal, Y.: Deterministic neural networks with appropriate inductive biases capture epistemic and aleatoric uncertainty. ArXiv abs/2102.11582 (2021)
- [18] Sandler, M., Howard, A.G., Zhu, M., Zhmoginov, A., Chen, L.C.: Mobilenetv2: Inverted residuals and linear bottlenecks. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition pp. 4510–4520 (2018)
- [19] Wang, H., Li, Z., Feng, L., Zhang, W.: Vim: Out-of-distribution with virtual-logit matching. ArXiv abs/2203.10807 (2022)