| | Method | Dataset | | | | | |
|---|---|---|---|---|---|---|---|
| | | ASCADv1 (fixed) | ASCADv1 (random) | DPAv4 (Zaid version) | AES-HD | OTiAiT | OTP |
| | Random | 111.6 ± 0.3 | 108 ± 5 | 13 ± 2 | 127 ± 1 | 1.21 ± 0.04 | 1.05 ± 0.02 |
| First-order parametric methods | SNR | 117.2 ± 0.6 | 116.7 ± 0.7 | 11.4 ± 0.2 | 126 ± 2 | 1.10 ± 0.02 | 1.0125 ± 0.0007 |
| | SOSD | 114.9 ± 0.5 | 105 ± 2 | 8.0 ± 0.8 | 126 ± 2 | 1.14 ± 0.03 | 1.027 ± 0.002 |
| | CPA | 111.5 ± 0.4 | 114 ± 1 | 11.5 ± 0.3 | 126 ± 2 | 1.49 ± 0.04 | 1.0125 ± 0.0007 |
| Neural net attribution | GradVis | 107.0 ± 0.5 | 95 ± 2 | 12.1 ± 0.3 | 127 ± 1 | 1.4 ± 0.2 | 1.0142 ± 0.0008 |
| | Saliency | 107.1 ± 0.5 | 95 ± 2 | 11.8 ± 0.3 | 127 ± 1 | 1.39 ± 0.04 | 1.014 ± 0.001 |
| | Input $*$ Grad | 107.2 ± 0.5 | 95 ± 2 | 11.8 ± 0.4 | 127 ± 1 | 1.36 ± 0.04 | 1.0141 ± 0.0009 |
| | LRP | 107.2 ± 0.5 | 95 ± 2 | 11.8 ± 0.4 | 127 ± 1 | 1.36 ± 0.04 | 1.0141 ± 0.0009 |
| | 1-Occlusion | 107.1 ± 0.5 | 95 ± 2 | 10.1 ± 0.2 | 127 ± 1 | 1.36 ± 0.04 | 1.0141 ± 0.0009 |
| | 5-Occlusion | 107.4 ± 0.4 | 94 ± 2 | 9.6 ± 0.2 | 127 ± 1 | 1.43 ± 0.03 | 1.013 ± 0.002 |
| | 17-Occlusion | 108.7 ± 0.4 | 96 ± 2 | 9.5 ± 0.2 | 127 ± 1 | 1.51 ± 0.02 | 1.021 ± 0.002 |
| | 65-Occlusion | 111.6 ± 0.7 | 99 ± 2 | 10.1 ± 0.2 | 127 ± 1 | 1.60 ± 0.01 | 1.026 ± 0.007 |
| | 257-Occlusion | 118.0 ± 0.8 | 104 ± 1 | 10.1 ± 0.2 | 127 ± 1 | 1.7 ± 0.2 | 1.031 ± 0.006 |
| | $2^{nd}$-order 1-Occlusion | 107.0 ± 0.4 | 95 ± 2 | 10.0 ± 0.2 | 127 ± 1 | 1.34 ± 0.04 | 1.0138 ± 0.0008 |
| | OccPOI | TODO | TODO | TODO | TODO | TODO | TODO |
| | GradVis (ZaidNet) | 108.8 ± 0.8 | n/a | 9.3 ± 0.2 | 126 ± 2 | n/a | n/a |
| | Saliency (ZaidNet) | 108.8 ± 0.8 | n/a | 9.3 ± 0.2 | 126 ± 2 | n/a | n/a |
| | Input $*$ Grad (ZaidNet) | 109.0 ± 0.6 | n/a | 9.2 ± 0.2 | 126 ± 2 | n/a | n/a |
| | 1-Occlusion (ZaidNet) | 109.3 ± 0.6 | n/a | 9.2 ± 0.2 | 126 ± 2 | n/a | n/a |
| | 5-Occlusion (ZaidNet) | 109 ± 1 | n/a | 10.2 ± 0.3 | 126 ± 2 | n/a | n/a |
| | 17-Occlusion (ZaidNet) | 111 ± 1 | n/a | 9.9 ± 0.3 | 126 ± 2 | n/a | n/a |
| | 65-Occlusion (ZaidNet) | 113 ± 1 | n/a | 9.9 ± 0.4 | 126 ± 2 | n/a | n/a |
| | 257-Occlusion (ZaidNet) | 120 ± 2 | n/a | 11 ± 1 | 126 ± 2 | n/a | n/a |
| | $2^{nd}$-order 1-Occlusion (ZaidNet) | 108.8 ± 0.4 | n/a | 9.2 ± 0.2 | 126 ± 2 | n/a | n/a |
| | OccPOI (ZaidNet) | TODO | n/a | TODO | TODO | n/a | n/a |
| | GradVis (WoutersNet) | 109.9 ± 0.5 | n/a | 9.6 ± 0.3 | 126 ± 2 | n/a | n/a |
| | Saliency (WoutersNet) | 109.8 ± 0.5 | n/a | 9.6 ± 0.3 | 126 ± 2 | n/a | n/a |
| | Input $*$ Grad (WoutersNet) | 109.7 ± 0.5 | n/a | 9.4 ± 0.3 | 126 ± 2 | n/a | n/a |
| | 1-Occlusion (WoutersNet) | 109.7 ± 0.4 | n/a | 9.4 ± 0.3 | 126 ± 2 | n/a | n/a |
| | 5-Occlusion (WoutersNet) | 110.1 ± 0.7 | n/a | 10.4 ± 0.4 | 126 ± 2 | n/a | n/a |
| | 17-Occlusion (WoutersNet) | 111.6 ± 0.5 | n/a | 10.0 ± 0.5 | 126 ± 2 | n/a | n/a |
| | 65-Occlusion (WoutersNet) | 114.1 ± 0.7 | n/a | 10.0 ± 0.3 | 126 ± 2 | n/a | n/a |
| | 257-Occlusion (WoutersNet) | 118 ± 1 | n/a | 11.0 ± 0.2 | 126 ± 2 | n/a | n/a |
| | $2^{nd}$-order 1-Occlusion (WoutersNet) | 109.2 ± 0.2 | n/a | 9.4 ± 0.3 | 126 ± 2 | n/a | n/a |
| | OccPOI (WoutersNet) | TODO | n/a | TODO | TODO | n/a | n/a |
| | ALL (ours) | 107.5 ± 0.3 | 101 ± 2 | 12.2 ± 0.4 | 126 ± 2 | 1.23 ± 0.03 | 1.0161 ± 0.0009 |

Table 1: Performance of leakage localization algorithms according to the Fwd-DNNO (forward DNN occlusion) test (smaller is better). To compute this metric, we first train a supervised DNN classifier to map emission traces to the sensitive variable. We then occlude all its inputs and incrementally un-occlude them from most- to least-leaky as estimated by the leakiness assessment under test, and at each step compute its performance (quantified by rank, lower is better) on the test dataset. The Fwd-DNNO metric is given by the average value of these performance assessments (lower is better, because it indicates that claimed leaky features indeed had utility to the classifier). Of the two DNN occlusion metrics, this is more sensitive to *true/false positive* leakiness measurements because the performance of the classifier tends to jump and stay up as soon as it sees leaky measurements. Best result is boxed and best deep learning result is underlined. Results are reported as mean ± std. dev. over 5 random seeds. This metric appears to have high variance and little discriminative power compared to the oSNR and Rev-DNNO metrics (as indicated by the large number of tied 'best' methods), and there is no clear best method according to Fwd-DNNO.