| Method | Total FLOPs | A6000 minutes per trial | | | | | |
|---|---|---|---|---|---|---|---|
| | | ASCADv1 (fixed) | ASCADv1 (random) | DPAv4 (Zaid version) | AES-HD | OTiAiT | OTP |
| Supervised training[†]: $t_{\text{sup}}$ | $C_{\text{sup}} := \Theta(Nn_{\text{sup}}(C_F+C_B))$ | $0.64 \pm 0.03$ | $1.27 \pm 0.02$ | $0.67 \pm 0.03$ | $0.52 \pm 0.02$ | $0.062 \pm 0.005$ | $0.071 \pm 0.004$ |
| GradVis | $C_{\text{sup}} + \Theta(N(C_F+C_B))$ | $t_{\text{sup}} + 0.0655 \pm 0.0005$ | $t_{\text{sup}} + 0.2582 \pm 0.0005$ | $t_{\text{sup}} + 0.0106 \pm 0.0001$ | $t_{\text{sup}} + 0.0379 \pm 0.0002$ | $t_{\text{sup}} + 0.0080 \pm 0.0001$ | $t_{\text{sup}} + 0.0639 \pm 0.0003$ |
| Saliency | $C_{\text{sup}} + \Theta(N(C_F+C_B))$ | $t_{\text{sup}} + 0.075 \pm 0.003$ | $t_{\text{sup}} + 0.2930 \pm 0.0007$ | $t_{\text{sup}} + 0.0111 \pm 0.0002$ | $t_{\text{sup}} + 0.048 \pm 0.003$ | $t_{\text{sup}} + 0.009 \pm 0.001$ | $t_{\text{sup}} + 0.081 \pm 0.004$ |
| Input * Grad | $C_{\text{sup}} + \Theta(N(C_F+C_B))$ | $t_{\text{sup}} + 0.074 \pm 0.002$ | $t_{\text{sup}} + 0.2937 \pm 0.0006$ | $t_{\text{sup}} + 0.01130 \pm 0.00007$ | $t_{\text{sup}} + 0.047 \pm 0.002$ | $t_{\text{sup}} + 0.00870 \pm 0.00009$ | $t_{\text{sup}} + 0.081 \pm 0.005$ |
| LRP | $C_{\text{sup}} + \Theta(N(C_F+C_B))$ | $t_{\text{sup}} + 0.075 \pm 0.002$ | $t_{\text{sup}} + 0.2936 \pm 0.0007$ | $t_{\text{sup}} + 0.014 \pm 0.002$ | $t_{\text{sup}} + 0.047 \pm 0.002$ | $t_{\text{sup}} + 0.011 \pm 0.003$ | $t_{\text{sup}} + 0.079 \pm 0.004$ |
| $m$-Occlusion | $C_{\text{sup}} + \Theta(NC_FT)$ | $t_{\text{sup}} + 0.1200 \pm 0.0006$ | $t_{\text{sup}} + 0.941 \pm 0.001$ | $t_{\text{sup}} + 0.0929 \pm 0.0002$ | $t_{\text{sup}} + 0.1134 \pm 0.0005$ | $t_{\text{sup}} + 0.0174 \pm 0.0001$ | $t_{\text{sup}} + 0.2456 \pm 0.0006$ |
| $2^{\text{nd}}$-order $m$-Occlusion[†,*] | $C_{\text{sup}} + \Theta(NC_FT^2)$ | $t_{\text{sup}} + 18.03 \pm 0.09$ | $t_{\text{sup}} + 374 \pm 1$ | $t_{\text{sup}} + 138.6 \pm 0.2$ | $t_{\text{sup}} + 34.28 \pm 0.02$ | $t_{\text{sup}} + 4.903 \pm 0.005$ | $t_{\text{sup}} + 79.9 \pm 0.5$ |
| OccPOI* | $C_{\text{sup}} + O(NC_FT^2)$ | TODO | TODO | TODO | TODO | TODO | TODO |
| ALL (Ours)[†] | $\Theta(2Nn_{\text{all}}(C_F+C_B))$ | $6.06 \pm 0.36$ | $8.4 \pm 0.6$ | $2.6 \pm 0.5$ | $4.2 \pm 0.4$ | $2.5 \pm 0.5$ | $2.0 \pm 0.1$ |

Table 1: A comparison of the computational cost of the considered methods. We denote by $C_F$ and $C_B$ the cost of a forward and backward pass through a neural net respectively, $N$ the dataset size, $n_{\text{sup}}$ and $n_{\text{all}}$ the number of epochs for supervised learning and adversarial leakage localization respectively, and $T$ the data dimensionality. All neural net attribution methods require first doing supervised training; we report their runtime as $t_{\text{sup}} + t_{\text{resid}}$ where $t_{\text{sup}}$ denotes the time to do supervised training (listed in top row) and $t_{\text{resid}}$ denotes the time to run the method given the trained neural net. Parametric statistical methods (omitted) are done on the CPU and take negligible time compared to the deep learning methods. Runtimes are reported as mean $\pm$ std. dev. of 5 total runtime measurements, with metrics, logging and validation disabled. [†]Estimated by linearly extrapolating the runtime of 100 minibatches. *While both methods require $O(T^2)$ passes through the dataset, all passes must be done sequentially for OccPOI, whereas all may be done in parallel for $2^{\text{nd}}$-order occlusion.