



# Benchmarking Differentially Private Residual Networks for Medical Imagery



Sahib Singh<sup>\*1,2,3</sup> Harshvardhan D. Sikka<sup>\*1,3,4</sup> Sasikanth Kotti<sup>1,5</sup> Andrew Trask<sup>1,6</sup>

<sup>\*</sup>Equal Contribution <sup>1</sup>OpenMined <sup>2</sup>Ford R&A <sup>3</sup>Manifold Computing Group <sup>4</sup>Harvard University <sup>5</sup>Indian Institute of Technology Jodhpur <sup>6</sup>University of Oxford.

## Focus

We measure the effectiveness of Differential Privacy (DP) when applied to medical imaging. We compare two robust differential privacy mechanisms: Local-DP and DP-SGD and benchmark their performance, analyzing the trade-off between the accuracy and the level of privacy the model guarantees. We also examine how useful these privacy guarantees prove to be in a real world medical setting.

### Key Preliminaries

**Differential Privacy:** the introduction of randomized noise to ensure privacy through plausible deniability.

**Local DP:** an algorithm  $\pi$  satisfies  $\epsilon$ -LDP where  $\epsilon > 0$  if and only if for any input  $v$  and  $v'$

$$\forall y \in \text{Range}(\pi) : P[\pi(v) = y] \leq e^\epsilon P[\pi(v') = y]$$

**DP - SGD:** a modification of SGD that bounds the sensitivity of each gradient and uses a moments accountant algorithm to amplify and track the privacy loss across weight updates.

**Laplace Distribution:** a symmetric version of the exponential distribution. The distribution centered at 0 (i.e.  $\mu = 0$ ) with scale  $\beta$  has the following probability density function.

$$\text{Lap}(x|b) = \frac{1}{2\beta} \exp\left(\frac{-|x|}{\beta}\right)$$

**Laplace Mechanism:** Laplace Mechanism independently perturbs each coordinate of the output with Laplace noise (from the Laplace distribution having mean zero) scaled to the sensitivity of the function.

Correspondence: [sahibsin@alumni.cmu.edu](mailto:sahibsin@alumni.cmu.edu)

Full Paper: <https://arxiv.org/abs/2005.13099>

## Experimental Setup

**Architecture:** Pretrained Resnet-18 trained in all experiments over 50 epochs with a 0.01 learning rate, and batch size of 128.

**LDP Experiment:** 3 other versions of the dataset were generated by the addition of different perturbations to the images. The tradeoff in accuracy with varying scales of perturbations ( $\beta=1$ ;  $\beta=2$ ;  $\beta=4$ ) were examined.

**DP - SGD Experiment:** We clip the gradient in the  $l_2$  norm, add random noise to it and then multiply it by the learning rate before updating the model parameters. Perturbations of ( $\beta=1$ ;  $\beta=2$ ;  $\beta=4$ ) were examined.

### Results: the Accuracy-Privacy Tradeoff

In some cases the training accuracy is lower than the test accuracy (See Figure 4 below). One possible explanation for this is that Local-DP adds noise to the training data, making the latent features harder to learn. Later when we run the model on the test data it performs better because the latent features are now relatively easy to capture since the model has already learned representations in a noisy scenario.

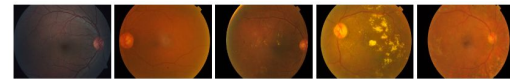
While Local-DP maintains the theoretical guarantee of Differential Privacy, it does not always provide the visual privacy we expect. In some cases the image was completely blurred out while in a few others there was hardly any visual change to the image (See Chest X-Rays Dataset above). DP-SGD is the mechanism of choice for ensuring more robust privacy.

Model Accuracy		Original	LDP			DP-SGD		
			$\beta=1$	$\beta=2$	$\beta=4$	$\beta=1$	$\beta=2$	$\beta=4$
Pneumonia	Train	99.76%	89.24%	83.66%	78.60%	85.87%	83.55%	82.53%
	Test	90.22%	91.03%	89.58%	91.67%	79.81%	77.72%	79.65%
APTOS	Train	98.39%	77.19%	65.52%	58.55%	43.77%	29.39%	17.82%
	Test	80.21%	76.26%	73.26%	72.17%	50.89%	49.25%	43.38%

Figure 4. Model Accuracy (%) — DP mechanism vs Dataset

## Datasets

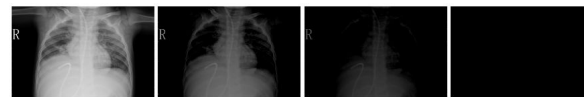
**APTOS: Retinal Scans for Diabetic Retinopathy.**



(a) No DR (b) Mild DR (c) Moderate DR (d) Severe DR (e) Proliferative DR

Figure 1. Examples from the APTOS Blindness Detection Dataset. Samples progress in severity of Diabetic Retinopathy.

**Chest X-Rays Dataset: Chest Radiography for Pneumonia**



(a) Original Image (b)  $\beta=1$  (c)  $\beta=2$  (d)  $\beta=4$

Figure 2. Comparison of sample image from dataset before and after Local DP based Obfuscation.