# Emory CXR
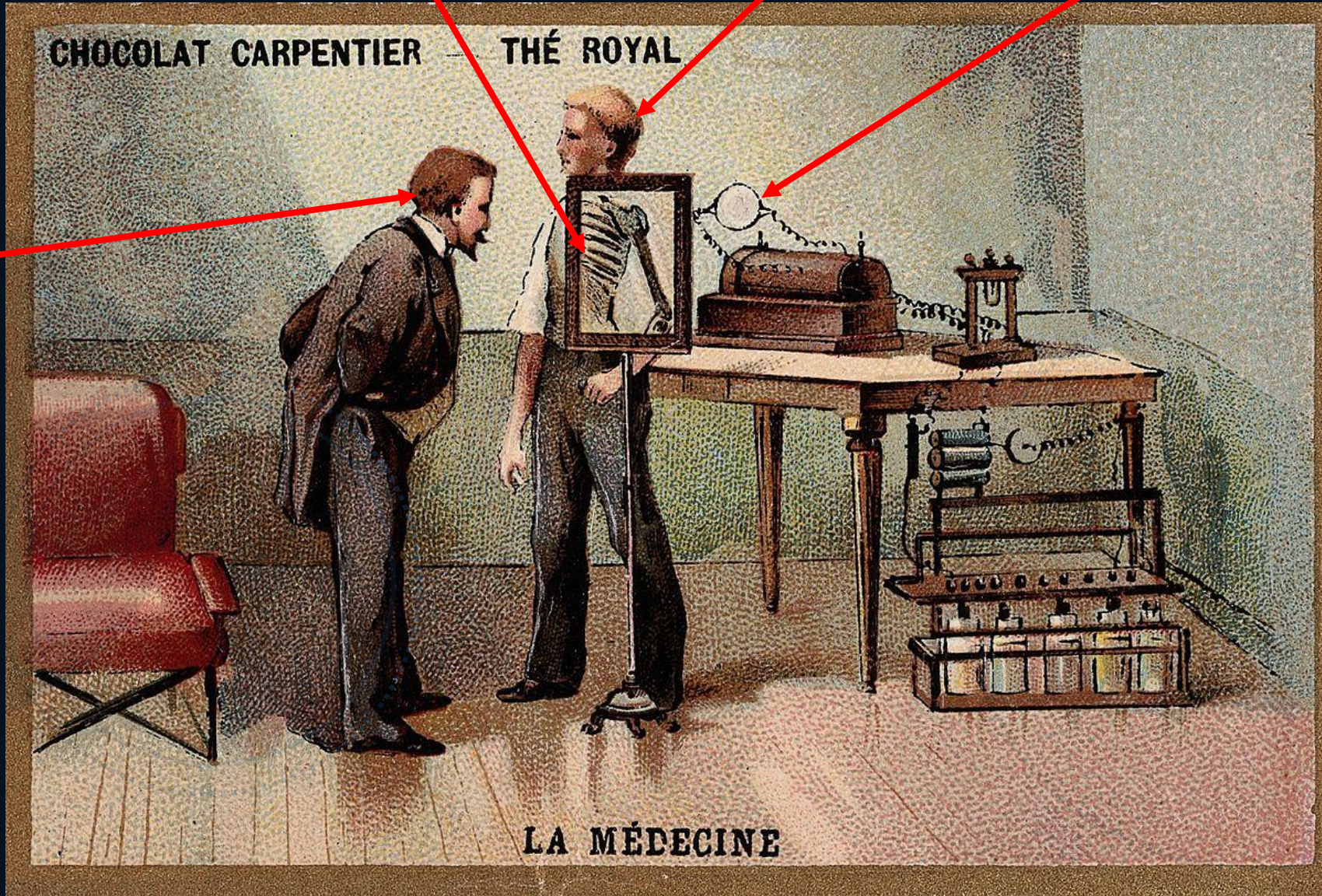## Data Preprocessing

Theo Dapamede, MD, PhD

08/19/2024

X-ray Detector

Patient

X-Ray Generator

Image Processing System

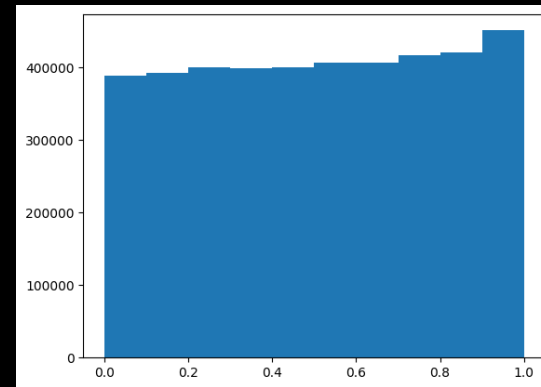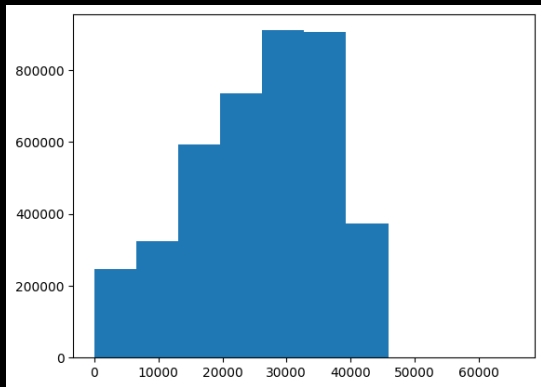CHOCOLAT CARPENTIER — THÉ ROYAL

LA MÉDECINE

Image source: Wikipedia Commons

# Our experience

- Pixel distribution in **PNG images** show anomalies:
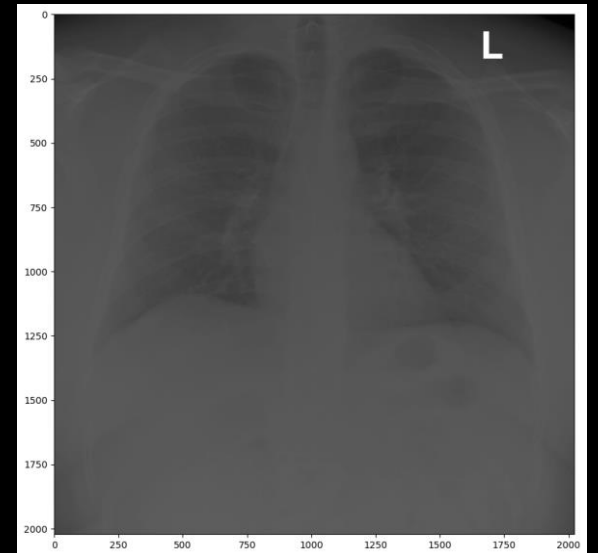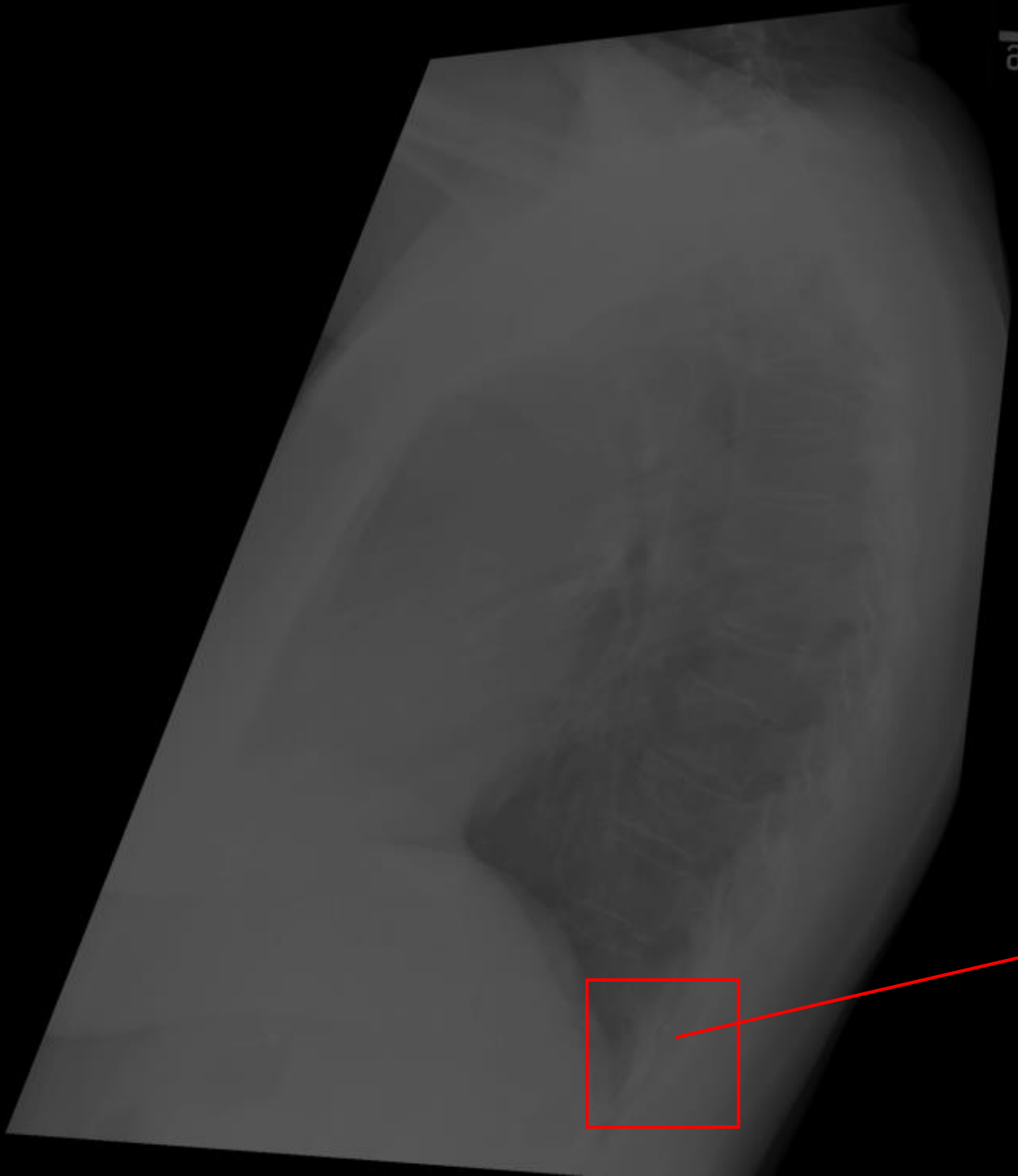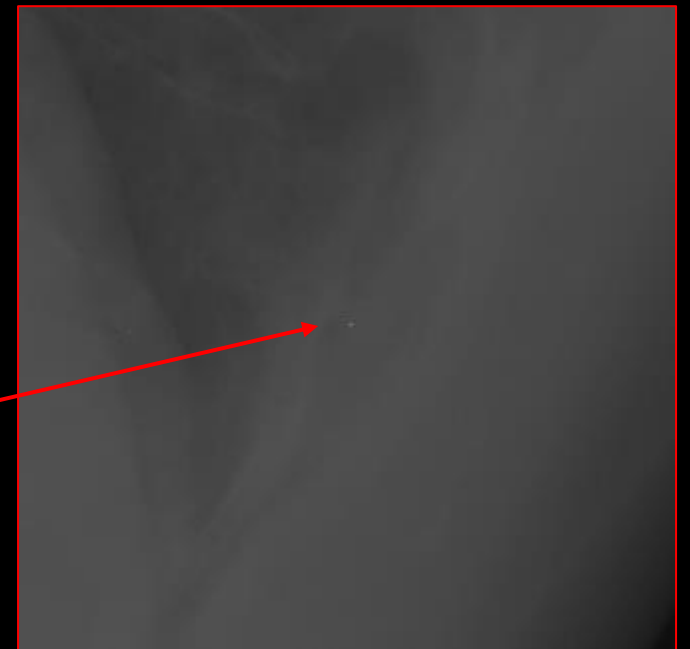  - Some pixels have values >>> 100 SD
- Equalisation, normalisation or standardisation methods performed on the PNG images don't result in the optimal output
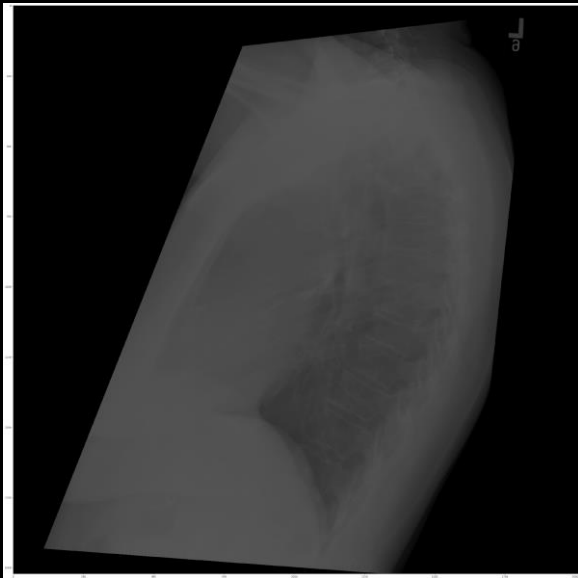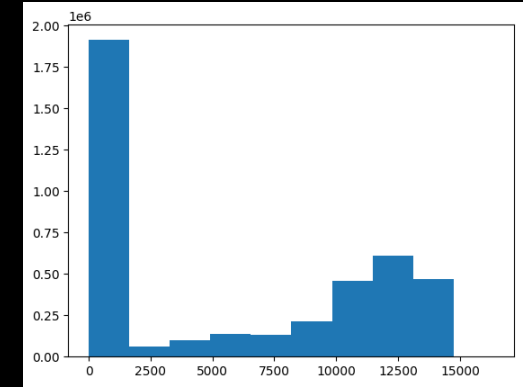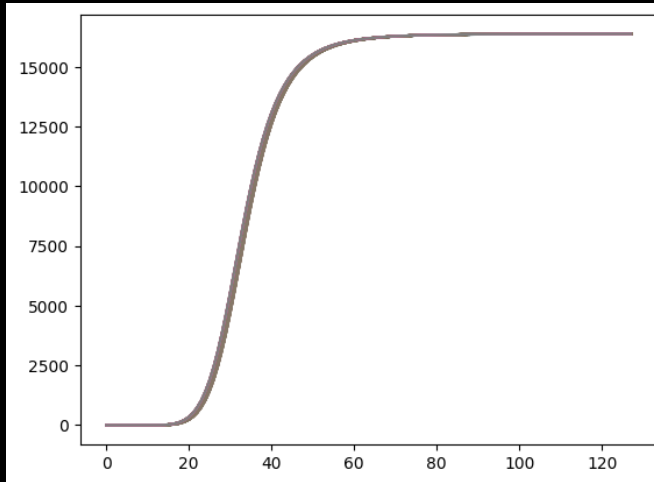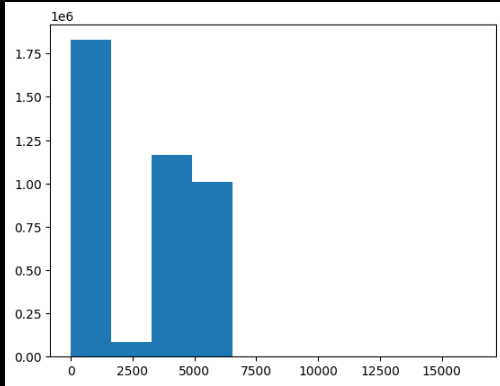  - Hence, datapoints are deleted from the dataset

- Found 1 pixel with extremely high value;
- Not always found in the marker as previously thought;

# From PNGs Back to DICOMs



VOI LUT

Raw Pixel Data
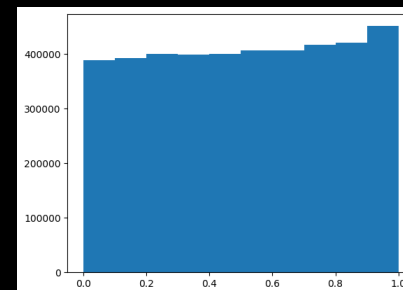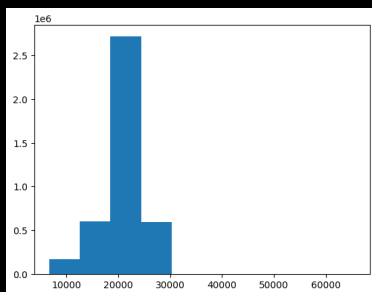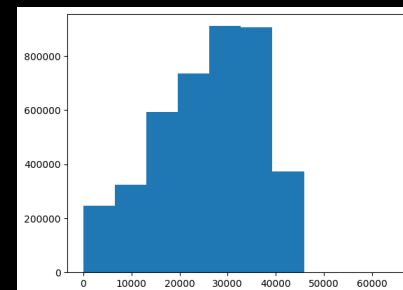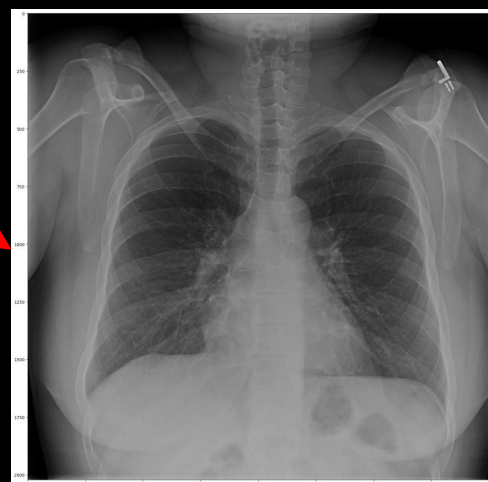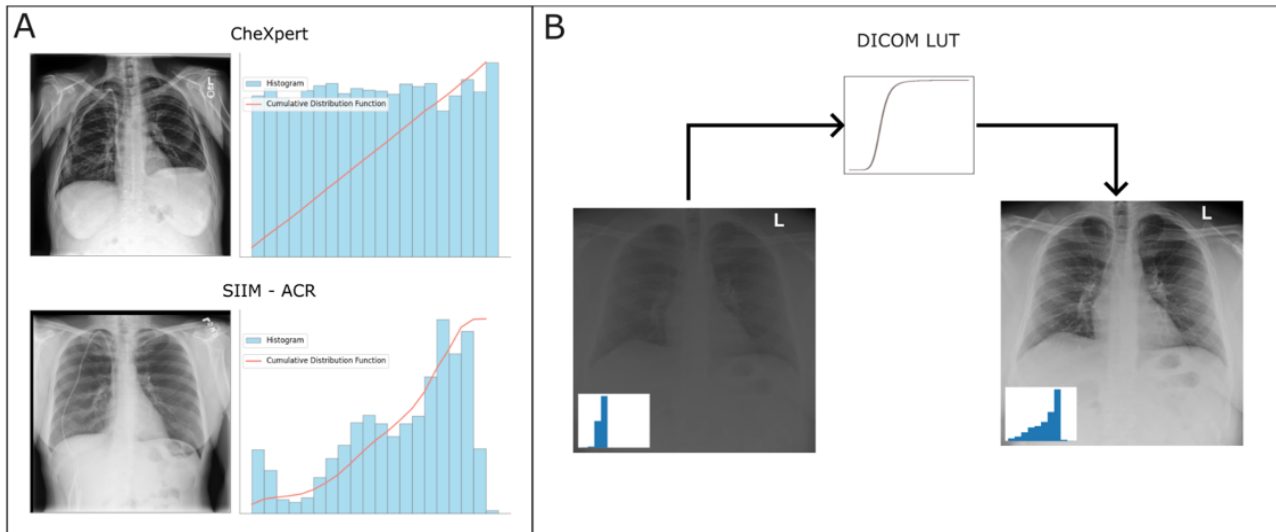
Histogram Equalization

DICOM VOI LUT
(Values Of Interest)

Figure(s)

**Figure 1.** A) Comparison between pixel distributions of two publicly available datasets: CheXpert which has been HE-preprocessed by default and SIIM-ACR without HE-preprocessing. B) Transforming raw pixel values to clinical standard pixel values using the corresponding DICOM LUT.
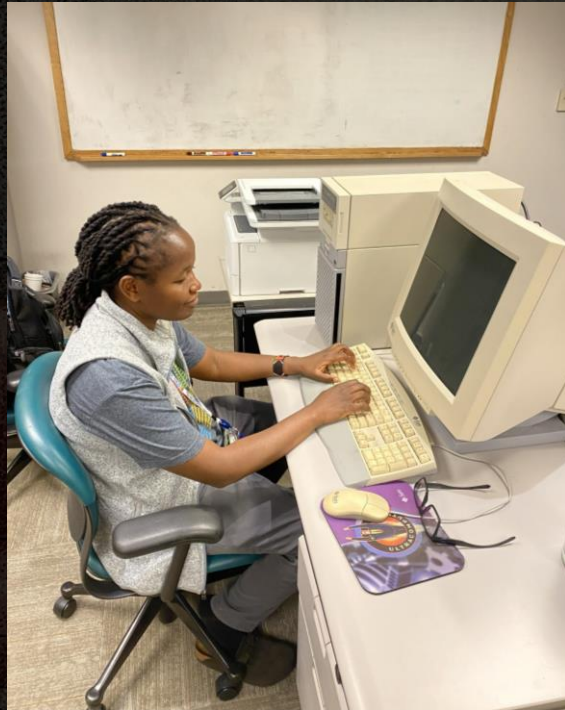
| Training Group | SIIM-ACR Pneumothorax | CheXpert |
|---|---|---|
| **Group 1: LUT (-), HE (-)** | **0.86 [0.85 – 0.87] *** | 0.69 [0.67 – 0.70] |
| **Group 2: LUT (+), HE (-)** | **0.84 [0.83 – 0.85] *** | **0.73 [0.71 – 0.74] *** |
| **Group 3: LUT (-), HE (+)** | 0.79 [0.77 – 0.80] | 0.69 [0.67 – 0.70] |
| **Group 4: LUT (+), HE (+)** | 0.80 [0.79 – 0.81] | 0.67 [0.66 – 0.69] |

*) significantly different to the non-asterisk groups
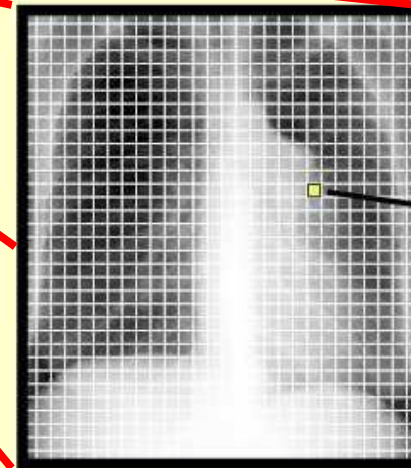
**Table 1**. AUC [CI] of the models trained on different training groups and evaluated on the SIIM-ACR Pneumothorax and CheXpert Datasets

DICOM VOI LUT
1. Increases model performance
2. Increases model generalizability

Metadata

- Patient Information
    - ID, Name, DOB, ...
- Machine Information
- Acquisition Information
    - Date, distances, kVp, ...
- Image Information
    - Encoding algorithm, pixel size, BIT depth, VOI LUT, s...
- ...

Pixel Data

DICOM

CHOCOLAT CARPENTIER --- THÉ ROYAL

LA MÉDECINE

# Hands On

1. Basics of working with DICOM files
2. DICOM Image Preprocessing
3. Standard Normalization Techniques