

CS 598 Deep Learning for Healthcare (Spring, 2025) Project Report

Eric Schrock

ejs9@illinois.edu

Abstract

This report details the results of my final project for CS 598: Deep Learning for Healthcare. My final project had four components.

First, I attempted to reproduce a portion of the findings of the paper "Multi-Label Generalized Zero Shot Learning for the Classification of Disease in Chest Radiographs" (Hayat, Lashen, and Shamout 2021), specifically the AUROC scores of the proposed model on each of the fourteen diseases labeled in the dataset, using the code provided by the paper. Running the model with the provided pre-trained weights reproduced the AUROC scores reported by the paper, but re-training the model did not.

Second, I attempted to rapidly rewrite the proposed model from scratch using an LLM. This was part of a research project on whether LLMs accelerate and enhance healthcare-related data science research. It took me just over one day to get to code that looked correct to me, but the resulting AUROC scores were much worse than those from the original model (pre-trained and re-trained). I saved my LLM chat log¹ for review by the researchers.

Third, I performed an extension study. I replaced the visual encoder used in the proposed model, DenseNet-121 (Huang et al. 2018), with the lighter weight EfficientNet-B0 (Tan and Le 2020) to see if I could achieve faster training times with comparable AUROC scores. Training was significantly faster, but the AUROC scores were significantly worse.

Fourth, I submitted a pull request (PR) to PyHealth (Yang et al. 2023), an open source project dedicated to supporting healthcare-related AI research and applications. My PR added support for the ChestX-ray14 dataset (Wang et al. 2017) and for binary classification tasks on the fourteen diseases labeled in the dataset.

- Video presentation ²
- Project GitHub repository ³
- PyHealth pull request ⁴

¹<https://github.com/EricSchrock/cxr-ml-gzsl/blob/main/report/report-llm-chat-log-for-research.txt>

²https://drive.google.com/file/d/1fTt2B8VNEQrtBT_Iooby2_vioUjf59bX/view

³<https://github.com/EricSchrock/cxr-ml-gzsl>

⁴<https://github.com/sunlabuiuc/PyHealth/pull/392>

1 Introduction

1.1 Summary of the Paper in Question

Deep learning models for classifying diseases from chest X-ray (CXR) images have had great success, achieving results comparable to those of human experts. However, collecting and labeling training data for these models is expensive and time consuming. For rarer diseases, it is often not economical and sometimes even impossible to gather the amount of data needed for supervised learning models. When a new disease emerges, the need for massive data collection slows the response. Human radiologists leverage other sources of knowledge to identify diseases they have previously never seen in X-ray form. Could a deep learning model do the same?

Multi-label generalized zero shot learning (ML-GZSL), which uses semantic information to identify classes not present in the set of labeled images used to train the model, has worked well in similar circumstances (Scheirer et al. 2013; Rahman, Khan, and Barnes 2020; Huynh and Elhamifar 2020). However, these prior works have at least two limitations. First, they "extract a fixed visual representation of the image from a pre-trained visual encoder or a detection network" (Hayat, Lashen, and Shamout 2021), which means they cannot be trained end-to-end. Second, "projecting these extracted visual features to the semantic space shrinks the diversity of the visual information, which gives rise to inherent limitations" (Hayat, Lashen, and Shamout 2021), one of those being the hubness problem (Dinu, Lazaridou, and Baroni 2015). Can these limitations be overcome?

"Multi-Label Generalized Zero Shot Learning for the Classification of Disease in Chest Radiographs" (Hayat, Lashen, and Shamout 2021) proposes the CXR-ML-GZSL model, which addresses both limitations. The result is better performance when classifying both seen and unseen diseases in chest X-ray images, compared to two state-of-the-art ML-GZSL models: LESA (Huynh and Elhamifar 2020) and MLZSL (Lee et al. 2018) (see Figure 1).

1.2 Scope of Reproducibility

I was able to access the dataset used by the paper in question (details in Section 2.2) and reproduce the reported AUROC scores for the fourteen diseases against which the CXR-ML-GZSL model was tested using the provided code and pre-

Method	k=2			k=3			AUROC		
	r@k	p@k	f1@k	r@k	p@k	f1@k	S	U	H
LESA (2020)	0.14	0.10	0.03	0.21	0.11	0.05	0.51	0.50	0.50
MLZSL (2018)	0.20	0.19	0.16	0.30	0.17	0.20	0.72	0.54	0.62
<i>OUR_{e2e}</i>	0.36	0.33	0.32	0.47	0.28	0.34	0.79	0.66	0.72

Figure 1: Performance of LESEA, MLZSL, and CXR-ML-GZSL (“*OUR_{e2e}*”) on the ChestX-ray14 dataset. Metrics are precision, recall, f1, and AUROC (split by seen, unseen, and the harmonic mean of the two) (Hayat, Lashen, and Shamout 2021).

trained weights (details in Section 2.3). However, when I re-trained the model using the code provided, I was not able to reproduce the reported AUROC scores. Furthermore, when I rewrote the data processing, model, training, and evaluation code from scratch with the help of an LLM, I was unable to reproduce the reported AUROC scores.

2 Methodology

2.1 Environment

Table 1 shows the Python version and libraries used to implement the CXR-ML-GZSL model. The “original” versions are based on the environment defined by the provided code⁵.

	Original	Reproduction
Python	3.6.12	3.11.12
matplotlib	3.3.3	3.10.0
numpy	1.19.2	2.0.2
pandas	1.1.3	2.2.2
pillow	8.0.0	11.1.0
sklearn	0.23.2	1.6.1
torch	1.4.0	2.6.0+cu124
torchvision	0.5.0	0.21.0+cu124
tqdm	4.55.1	4.67.1

Table 1: Environment

2.2 Data

The ChestX-ray14 dataset (Wang et al. 2017) contains 112,120 chest X-ray images labeled for the presence or absence of fourteen different diseases. The average number of diseases present in each X-ray is 0.72. Table 2 shows how frequently each disease appears in the dataset. Figure 2 shows an example chest X-ray from the dataset.

The dataset is hosted by the National Institutes of Health (NIH)⁶ on Box⁷. It includes a Python script⁸ to automatically download the images.

⁵<https://github.com/nyuad-cai/CXR-ML-GZSL/blob/master/environment.yml>

⁶<https://www.nih.gov/>

⁷<https://nihcc.app.box.com/v/ChestXray-NIHCC/folder/36938765345>

⁸<https://nihcc.app.box.com/v/ChestXray-NIHCC/file/371647823217>

Disease	Count	Percentage
None	60,361	53.8%
Atelectasis	11,559	10.3%
Cardiomegaly	2,776	2.5%
Consolidation	4,667	4.2%
Edema	2,303	2.1%
Effusion	13,317	11.9%
Emphysema	2,516	2.2%
Fibrosis	1,686	1.5%
Hernia	227	0.2%
Infiltration	19,894	17.7%
Mass	5,782	5.2%
Nodule	6,331	5.6%
Pleural Thickening	3,385	3.0%
Pneumonia	1,431	1.3%
Pneumothorax	5,302	4.7%

Table 2: Dataset Disease Frequency



Figure 2: Chest X-ray Displaying Cardiomegaly

2.3 Model

“Multi-Label Generalized Zero Shot Learning for the Classification of Disease in Chest Radiographs” (Hayat, Lashen, and Shamout 2021) provides a GitHub repository⁹ that implements the CXR-ML-GZSL model and training, as well as pre-trained weights¹⁰ that reproduce the reported results. My project GitHub repository¹¹ provides links to pre-trained weights that reproduce my results along with instructions.

CXR-ML-GZSL, shown in Figure 3, is comprised of a pre-trained text encoder to convert class labels to a semantic embedding space, a trainable visual encoder to convert X-ray images to a visual embedding space, and mapping mod-

⁹<https://github.com/nyuad-cai/CXR-ML-GZSL/>

¹⁰<https://drive.google.com/file/d/17ioJMW3qNx1KtMr-hXn-eqP431cm49Rm/view>

¹¹<https://github.com/EricSchrock/cxr-ml-gzsl>

els into a shared latent embedding space. The output is a score for each possible class, representing how relevant it is to the input image.

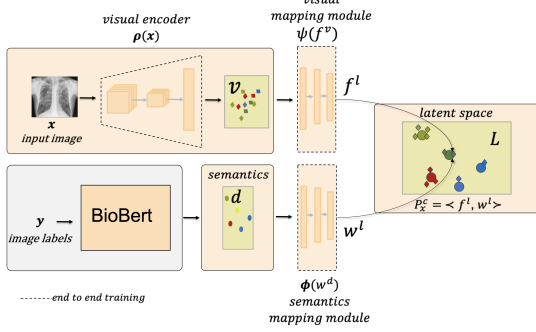


Figure 3: CXR-ML-GZSL (Hayat, Lashen, and Shamout 2021)

Since the visual encoder is trainable, the whole pipeline between the X-ray image and the semantic embedding of the class label is trainable from end-to-end, better tuning the overall model to the task at hand. Additionally, the added latent embedding space is meant to address the limitations of mapping the visual embedding space directly onto the semantic embedding space.

For the pre-trained text encoder, the model uses BioBERT (Lee et al. 2019), which was trained specifically on a biomedical corpora. For the trainable visual encoder, the model uses Densenet-121 (Rajpurkar et al. 2017), as at that time it was the best chest X-ray classifier. For both BioBERT and Densenet-121, the model skips the final classification step in order to get semantic and visual embeddings, respectively, as output instead. The two mapping models are three-layer feedforward neural networks.

CXR-ML-GZSL is the first use of ML-GZSL to classify diseases from chest X-ray images. It is also unique from other ML-GZSL models in the following three ways.

- It can be trained end-to-end, thanks to the trainable visual encoder.
- It maps the semantic and visual embedding spaces into a shared latent embedding space, instead of mapping the visual space onto the semantic space, causing less visual information to be lost.
- It uses BioBERT, which was trained on a biomedical corpora, resulting in semantic embeddings that are tuned for healthcare use cases.

2.4 Training

Hyperparameters According to the paper, the CXR-ML-GZSL model was trained on a range of learning rate (LR) and γ values, but it does not state which values produced the best results. I picked a value from each range for my reproduction attempts.

The paper does not state the batch sizes used for the training, validation, and test datasets, so I used the hard coded values when retraining the provided code and a set batch size

of 64 when training the LLM generated model. Where the hyperparameter values in the paper do not match the code provided (δ and the learning rate decay patience and factor), I used the values listed in the paper.

The CXR-ML-GZSL model was trained with the Adam optimizer (Kingma and Ba 2017), which I also used in my reproduction attempts. Table 3 lists the hyperparameter values used for training.

	Original	Repro	LLM
LR	[1e-4, 5e-5, 1e-5]	1e-4	1e-4
Patience	10 epochs	10 epochs	10 epochs
Decay factor	0.01	0.01	0.01
Batch sizes	16/160/48	16/160/48	64/64/64
γ_1	[0.1, 0.05, 0.01]	0.1	0.1
γ_2	[0.1, 0.05, 0.01]	0.1	0.1
δ	0.5	0.5	0.5
Optimizer	Adam	Adam	Adam

Table 3: Hyperparameters Used (Original, Reproduction Using Original Code, and Reproduction Using LLM Generated Code)

Computational Requirements Table 4 shows the computational requirements for training.

	Original	Repro	LLM
GPU	Quadro RTX 6000	A100	A100
GPU RAM	Unknown	38.8	19.1
System RAM	Unknown	24.6	22.5
Epochs	100	100	100
Training (hrs)	~8	~7.5	~8.5

Table 4: Computational Requirements (Original, Reproduction Using Original Code, and Reproduction Using LLM Generated Code)

Loss Function The increased complexity of CXR-ML-GZSL, compared to other ML-GZSL models, requires a multifaceted training objective, captured in a three-part loss function, where γ_1 and γ_2 are the regularization parameters for the second and third components of the loss function.

$$\min_{\phi, \rho, \psi} \mathcal{L} = \mathcal{L}_{rank} + \gamma_1 \mathcal{L}_{align} + \gamma_2 \mathcal{L}_{con}, \quad (1)$$

\mathcal{L}_{rank} adds penalties for any positive ground-truth relevance scores not larger than all negative ground-truth relevance scores by at least a margin of δ . \mathcal{L}_{align} adds penalties if input images and their labels do not map near each other in the latent embedding space. \mathcal{L}_{con} add penalties if labels do not have similar relationships in both the semantic and latent embedding spaces, as those semantic relationships are key to zero shot learning.

2.5 Evaluation

For my reproductions, I generated AUROC scores for each of the fourteen diseases present in the dataset. Additionally, I generated the mean AUROC scores for the ten diseases seen

by the model during training and the four diseases withheld from training, along with the harmonic mean between the "seen" and "unseen" means.

3 Results

3.1 Reproduction Results

Table 5 shows the individual and mean AUROC scores reported by the original paper vs those from my three reproduction attempts. Using the code AND pre-trained weights provided by the original paper perfectly reproduced the reported AUROC scores. This gives confidence that the provided model code and weights are the same as what was used to generate the paper's findings.

	Orig	Pre	Re	LLM
Means				
"Seen" mean	0.79	0.79	0.77	0.72
"Unseen" mean	0.66	0.66	0.60	0.55
Harmonic mean of means	0.72	0.72	0.67	0.62
"Seen" diseases				
Atelectasis	0.76	0.76	0.72	0.73
Cardiomegaly	0.90	0.90	0.90	0.84
Effusion	0.83	0.83	0.82	0.77
Infiltration	0.70	0.70	0.67	0.63
Mass	0.80	0.80	0.79	0.76
Nodule	0.75	0.75	0.73	0.72
Pneumothorax	0.83	0.83	0.79	0.51
Consolidation	0.69	0.69	0.68	0.79
Pleural thickening	0.72	0.72	0.72	0.61
Hernia	0.90	0.90	0.86	0.83
"Unseen" diseases				
Pneumonia	0.62	0.62	0.55	0.60
Edema	0.67	0.67	0.54	0.49
Emphysema	0.74	0.74	0.66	0.51
Fibrosis	0.60	0.60	0.66	0.58

Table 5: AUROC Scores from the **Original Paper** vs Reproduced with **Pre**-trained Weights vs Reproduced by **Retraining** the Model vs Reproduced with **LLM** Generated Code

Retraining the provided model with the provided training code did not reproduce the paper's findings. This could be for at least two reasons. First, I may have used the wrong learning rate or γ values, as the paper did not specify the exact values used.

Second, the provided training code may not be the same code used to generate the paper's findings. Two things point in this direction. First, the default epoch and δ values in the provided training script do not match those listed in the paper. Second, the provided training code contains at least three bugs.

- The \mathcal{L}_{align} value is hardcoded to 0.0.
- The LR decay patience and factor do not match those listed in the paper.
- There is an off-by-one error that causes training to end one epoch early.

Training the LLM generated model had worse results than retraining the provided model. It is likely the differing batch size had a negative impact. It is also likely I did not catch all the bugs in the LLM generated code, given how rapidly I developed it.

3.2 Extension Study Results

My extension study tested the performance of the CXR-ML-GZSL model with a lighter weight visual encoder called EfficientNet-B0. For the sake of time, I cut both the dataset and the number of training epochs in half. The result was significantly faster training at the cost of significantly worse AUROC scores.

	DenseNet-121	EfficientNet-B0
"Seen" AUROC mean	0.69	0.60
"Unseen" AUROC mean	0.54	0.50
Harmonic mean of means	0.60	0.55
Training time (hours)	1.83	1.48

Table 6: Extension Study Results

4 Discussion

4.1 How Reproducible Is This Paper?

The AUROC scores reported by the paper are reproducible using the provided code and pre-trained weights. I suspect the model training can also be reproduced, but not without significant effort. Section 4.3 and Section 4.4 detail shortcomings in the paper and provided code that make reproducing the model training more difficult than it needs to be.

4.2 What Was Easy About Reproducing This Paper?

- The material covered in CS 598: Deep Learning for Healthcare was sufficient to understand the paper.
- Updating the provided code for newer Python and library versions was trivial.
- Running the provided model code using the provided weights was straightforward.

4.3 What Was Hard About Reproducing This Paper?

- Getting data to and from the cloud and dealing with cloud computing session timeouts took a lot of trial and error.
- The project deadline, long dataset download and model training times, and expensive computational requirements limited the number of development/test iterations I could perform.
- Missing details (like hyperparameter values) and training code bugs increased the number of iterations required to fully reproduce the reported results.

4.4 How to Improve Reproducibility?

- Use version control from the start to provide a history of changes.
- Check-in the exact code used to generate the results you report. If different results are generated from different commits, tag those commits to make it clear exactly how each result was produced.
- Set default parameter values to match the values used to generate the results you report.
- Train your models on a single epoch to check for off-by-one errors in epoch indexing.
- If you experiment with different hyperparameter values, make it clear which values produced the results you report.

5 Contributions

I completed this project as a team of one. All work is my own unless otherwise stated.

References

- Dinu, G.; Lazaridou, A.; and Baroni, M. 2015. Improving zero-shot learning by mitigating the hubness problem. *arXiv:1412.6568*.
- Hayat, N.; Lashen, H.; and Shamout, F. E. 2021. Multi-Label Generalized Zero Shot Learning for the Classification of Disease in Chest Radiographs. *arXiv:2107.06563*.
- Huang, G.; Liu, Z.; van der Maaten, L.; and Weinberger, K. Q. 2018. Densely Connected Convolutional Networks. *arXiv:1608.06993*.
- Huynh, D.; and Elhamifar, E. 2020. A Shared Multi-Attention Framework for Multi-Label Zero-Shot Learning. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 8773–8783.
- Kingma, D. P.; and Ba, J. 2017. Adam: A Method for Stochastic Optimization. *arXiv:1412.6980*.
- Lee, C.-W.; Fang, W.; Yeh, C.-K.; and Wang, Y.-C. F. 2018. Multi-Label Zero-Shot Learning with Structured Knowledge Graphs. *arXiv:1711.06526*.
- Lee, J.; Yoon, W.; Kim, S.; Kim, D.; Kim, S.; So, C. H.; and Kang, J. 2019. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4): 1234–1240.
- Rahman, S.; Khan, S.; and Barnes, N. 2020. Deep0Tag: Deep Multiple Instance Learning for Zero-Shot Image Tagging. *Trans. Multi.*, 22(1): 242–255.
- Rajpurkar, P.; Irvin, J.; Zhu, K.; Yang, B.; Mehta, H.; Duan, T.; Ding, D.; Bagul, A.; Langlotz, C.; Shpanskaya, K.; Lungren, M. P.; and Ng, A. Y. 2017. CheXNet: Radiologist-Level Pneumonia Detection on Chest X-Rays with Deep Learning. *arXiv:1711.05225*.
- Scheirer, W.; Rocha, A.; Sapkota, A.; and Boulton, T. 2013. Toward Open Set Recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(7): 1757–1772.
- Tan, M.; and Le, Q. V. 2020. EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. *arXiv:1905.11946*.
- Wang, X.; Peng, Y.; Lu, L.; Lu, Z.; Bagheri, M.; and Summers, R. M. 2017. ChestX-Ray8: Hospital-Scale Chest X-Ray Database and Benchmarks on Weakly-Supervised Classification and Localization of Common Thorax Diseases. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 3462–3471. IEEE.
- Yang, C.; Wu, Z.; Jiang, P.; Lin, Z.; Gao, J.; Danek, B.; and Sun, J. 2023. PyHealth: A Deep Learning Toolkit for Healthcare Predictive Modeling. In *Proceedings of the 27th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD) 2023*.