# Unified Representation Learning for Efficient Medical Image Analysis

**Ghada Zamzmi**[1], **Sivaramakrishnan Rajaraman**[1], **and Sameer Antani**[1]

[1]National Library of Medicine, National Institutes of Health, Bethesda, MD, USA

## ABSTRACT

Medical image analysis typically includes several tasks such as image enhancement, detection, segmentation, and classification. These tasks are often implemented through separate machine learning methods, or recently through deep learning methods. We propose a novel multitask deep learning-based approach, called unified representation (U-Rep), that can be used to simultaneously perform several medical image analysis tasks. U-Rep is modality-specific and takes into consideration inter-task relationships. The proposed U-Rep can be trained using unlabeled data or limited amounts of labeled data. The trained U-Rep is then shared to simultaneously learn key tasks in medical image analysis, such as segmentation, classification and visual assessment. We also show that pre-processing operations, such as noise reduction and image enhancement, can be learned while constructing U-Rep. Our experimental results, on two medical image datasets, show that U-Rep improves generalization, and decreases resource utilization and training time while preventing unnecessary repetitions of building task-specific models in isolation. We believe that the proposed method (U-Rep) would tread a path toward promising future research in medical image analysis, especially for tasks with unlabeled data or limited amounts of labeled data.

## Introduction

Deep learning has significantly pushed forward the frontiers of automated analysis of several computer vision tasks such as object detection, segmentation, and classification. It has exhibited excellent performance in various fields, including medical image analysis. Several deep learning-based approaches have been developed for various medical imaging modalities to perform region of interest (ROI) detection ([1,2]), segmentation ([2]), registration ([3]), and classification ([2]). Although these approaches achieve promising performance, they typically rely on training individual models for different tasks, which limits the actual impact of deep learning and its ability to transfer knowledge. In addition, using multiple models with high computational complexity and training parameters may inevitably be constrained by limited computational resources, making it difficult for the models to be deployed in clinical settings for potential real-time healthcare applications. These challenges can be alleviated through multi-task learning (MTL), a subfield of machine learning, where multiple tasks could be solved simultaneously, while exploiting the commonalities and differences across the tasks[4].

### Background

MTL has attracted significant attention as it shows remarkable success in improving generalization and performance[4]. MTL allows knowledge transferring by exploiting the commonalities across tasks. Several MTL architectures and loss functions are proposed in the literature. Examples of the state-of-the-art MTL architectures include cross-stitch[5] and progressive networks[6]. Weight uncertainty[7] and GradNorm[8] are other works that propose novel MTL loss functions. MTL has been widely used in various computer vision applications to learn related tasks jointly; e.g., facial landmark detection, head pose estimation, and facial attributes inference[9]; human pose estimation and object recognition[10]; depth prediction, surface normal estimation, and semantic labeling[11]; and scene understanding[12].

Recently, Zamir et al.[13] proposed MTL dictionary-based approach, called Taskonomy, for modeling the structure of space of main visual tasks in computer vision using natural images. Specifically, given that there are N visual tasks, N task-specific networks with the same encoder and decoder architectures are trained from scratch. Then, all feasible transfer functions (first and higher functions) among tasks are trained in the latent space. After transfer modeling, task affinity normalization was generated using Analytic Hierarchy Process (AHP) followed by generating hypergraph or global transfer. This hypergraph is then used to predict the performance of transfers between tasks and optimize for the optimal one. As discussed in the paper, Taskonomy has high computational complexity and it is a completely supervised approach; i.e., it can only be developed using a dataset that has annotations for every task on every image. In case of 3D medical images, Zhou et al.[14] proposed to build a set of models, called Genesis, using self-supervised learning. These methods can be used as the starting point to learn cross-modality tasks. The paper does not address the impact of conflict modalities and tasks on performance; i.e., conflict modalities/task can distract each other.

Although existing works such as Taskonomy[13] and Genesis[14] propose robust MTL frameworks for image analysis and achieve excellent performance, these works have the following main limitations:

- Existing works involve redundancy because they require training of task-specific networks (encoder-decoder) for all tasks[13] or training of multiple networks[14].

- Existing works, such as Taskonomy[13], require having annotations for every task on every image (supervised learning). This scenario is unrealistic in medical image analysis as many medical image datasets do not have labels for every task on every image.

- Existing works, such as Taskonomy[13], use computationally intensive approach to assess inter-tasks relationship. Specifically, the number of higher-order transfers among tasks can lead to a combinatorial explosion. Other works, such as Genesis[14], ignore tasks/modalities relationship and their impact on performance. i.e., conflict modalities/task can distract each other.

- Existing works[13, 14] focus only on solving target tasks (e.g., classification and segmentation). We are not aware of any works that propose to learn pre-processing operations, such as noise reduction and image enhancement, while constructing shared MTL representation.

## Contributions

We demonstrate that a single unified representation (U-Rep) can perform pre-processing operations and then shared to jointly solve key medical image analysis tasks. Our contributions are summarized as follows:

- We propose U-Rep, which can be created ex nihilo (no manual labeling) from a specific image modality using unsupervised learning. U-Rep can also be created using other learning approaches such as supervised learning (manual labeling of a single task), residual learning, and reinforcement learning. U-Rep is then used as a single source model to simultaneously learn the key medical image analysis tasks (target tasks) such as classification and segmentation.

- Our proposed U-Rep is modality-specific and takes into consideration inter-tasks relationship. Our method for modeling task relations is less complex and more suitable for medical image analysis.

- We propose and show that image pre-processing operations, such as noise reduction and image enhancement, can be integrated and learned during the construction of U-Rep. We also propose derivable tasks, which we define as the ones generated from the learned target tasks. Examples of derivable tasks include visual interpretation of machine learning decisions and automated decision recommendation.

- This work is the first to combine unsupervised learning with MTL to jointly perform the key medical image analysis tasks, demonstrated through the use of chest x-ray (CXR) datasets. Specifically, we use a convolutional denoising autoencoder (CDAE) to create U-Rep from the feature space of CXR. U-Rep is then shared to perform the following tasks: segmentation and classification. We focus mainly on these two tasks because they are the major tasks in most medical image applications.

- This work is the first to propose the use of a super-resolution (SR) convolutional neural network (CNN) as a shared representation (U-Rep) to simultaneously solve multiple target tasks. Specifically, we use residual learning to build the encoder from the feature space of CXR. The generated representation is then shared to perform related target tasks.

- We use a supervised learning approach to create U-Rep from a representative source task (e.g., task with the largest number of annotations) that is related to other target tasks. The generated representation is then shared to learn target tasks. This is demonstrated with echo Doppler images to jointly learn the following tasks: quality assessment and segmentation. We also demonstrate different derivable tasks.

- We empirically demonstrate the superiority and efficiency of using U-Rep, in terms of performance and computation, as compared to the traditional approach (i.e., individual models for separate tasks) for medical image analysis.

Our proposed U-Rep offers three main advantages for medical image analysis. First, it prevents unnecessary repetitions of learning models in isolation, which leads to a decrease in resource utilization and training time. It also improves generalization across related tasks, and hence performance, since it learns the cross-tasks shared patterns (knowledge transfer). Second, it facilitates cross-institutional model sharing and provides an alternative to data sharing. For example, instead of transferring relatively large datasets, optimized U-Rep backbones built using common imaging modalities can be shared. The shared models can be used to solve different imaging tasks. Third, It allows to integrate different U-Rep models built from different sources,
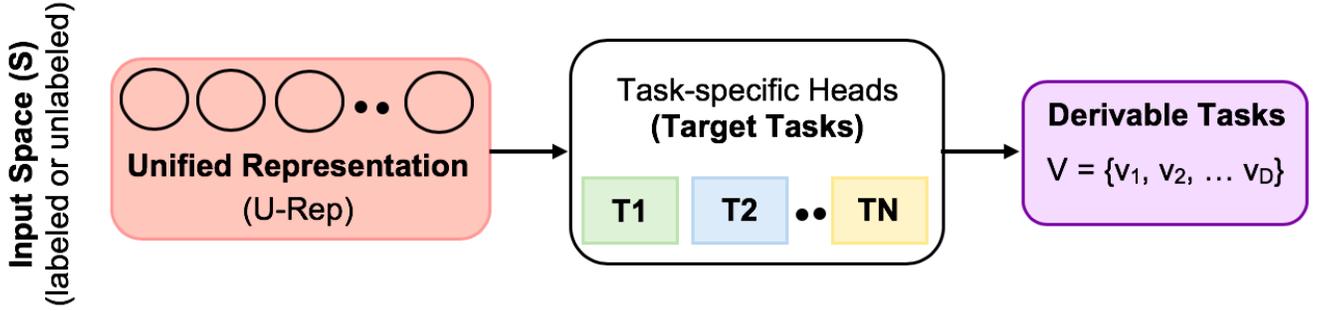
**Figure 1.** Proposed U-Rep for jointly learning the major medical image analysis tasks. $T_{1:N} = \{T_1, .., T_N\}$ and $V_{1:D} = \{V_1, .., V_D\}$ represent target and derivable tasks, respectively.

which can aid in complex decision making. For example, two U-Rep models can be created for 2D cardiac echo and cardiac echo Doppler where each modality is used for different clinical needs with different capability. Then, the outcomes of these U-Rep models could be combined to serve as a potential improvement in augmented decision making over any individual model. We believe U-Rep would tread a path toward promising future research in medical image analysis, particularly when limited labeled data is available, which is a common concern in the field.

## Materials and Methods

### Datasets
Two medical imaging modalities are used to evaluate the proposed U-Rep method: CXR and echo Doppler.

#### CXR collections
We used three publicly available CXR datasets: Radiological Society of North America (RSNA)[15], Montgomery[16], and Shenzhen[16] CXR collections. RSNA dataset includes 26,684 normal and abnormal frontal CXRs provided as DICOM images with $1024 \times 1024$ spatial resolution. The abnormal CXRs have pulmonary opacities indicating pneumonia and other disorders. GT disease bounding boxes are made available for CXRs containing pneumonia-related opacities. Montgomery dataset contains 138 posterior-anterior (PA) CXRs of which 80 CXRs are normal and 58 are abnormal with tuberculosis (TB) manifestations. The images have $4020 \times 4892$ resolution. In addition to normal and abnormal labels, the dataset includes GT binary masks. Shenzhen dataset contains 662 CXR images (abnormal: 326 and normal: 336). Abnormal images include cases with manifestations of TB. The size of the images in this collection varies, but it is approximately $3000 \times 3000$ pixels.

#### Echo Doppler
We used a private echo Doppler dataset containing continuous-wave, pulsed wave, and Tissue Doppler images collected from 100 patients who were referred for echocardiographic examination in the Clinical Center at the National Institutes of Health's (NIH). The use of these de-identified images was approved by the NIH Ethics Review Board (IRB:18-NHLBI-00686). The Doppler traces of the mitral valve flow (MV), mitral annular flow (MA), and tricuspid regurgitation flow (TR) were acquired using different commercial echocardiography systems including Phillips iE33, GE Vivid95, and GE Vivid9. The dataset has a total of 2444 images. All images have a flow type label (TR, MV, or MA) and segmentation mask, provided by an expert technician, which separates the spectral envelope from the background. In addition, the expert technician assessed the quality of a subset of images (814 out of 2444) as low- or good-quality. All GT labels provided by the expert technician were further verified by an expert cardiologist. Prior to training, all images were downsampled to $224 \times 224$ pixel resolution.

### Proposed U-Rep for Medical Image Analysis
#### Notations and Definitions
As shown in Figure 1, U-Rep learns the visual features ($X_S$) of a source input space or imaging modality ($S$). The learned features are then shared to jointly learn $N$ target tasks $T_{1:N} = \{T_1, T_2, ..., T_N\}$. The source modality is defined as $S = \{X_S, P_S(X_S)\}$, where $X_S$ represents the feature space with L dimensions ($\{x_1, x_2, ...., x_L\}$) and $P_S(X_S)$ represents the probability distribution of that modality.

U-Rep aims to improve the learning and generalization of the predictive functions of up to $N$ related tasks ($f_{1:N}(.) = \{f_1(.), f_2(.), .., f_N(.)\}$) by transferring the knowledge from $S$ to $N$ to generate task-specific predictions $Y_{1:N} = \{Y_1, Y_2, ..., Y_N\}$. Since medical image analysis tasks that use a common input or image modality share common features, it follows that the source

and target tasks with different outputs or ground truth (GT) labels ($Y_{GT_S} \neq Y_{GT_{1:N}}$) are related if and only if $P_S(X_S) \sim P_{1:N}(X_{1:N})$. For example, the features extracted by an autoencoder ($Y_{GT_S} = \{\}$) from a specific imaging modality (S) can be used to jointly learn two target tasks $T_1$ and $T_2$ with different label spaces, where $P_S(X_S) \sim P_1(X_1) \sim P_2(X_2)$. A new task-specific head or target task can be added on-the-fly with minimal effort if this constraint is satisfied. An obvious reason for this constraint is that sharing information between unrelated tasks would result in performance degradation (negative transfer). The loss function for jointly learning multiple related tasks can be calculated by adding the losses across these tasks.

We define D derivable tasks $V_{1:D} = \{V_1, V_2, ..., V_D\}$ as the ones that use the information learned by a single or a combination of task-specific heads. An example of a derivable task is the interpretation or visual explanation generated from classification prediction. Such a task is crucial for augmenting medical decision making[17]. Another derivable task is automated decision recommendation, which can be generated by combining the outputs of different task-specific heads. A new derivable task can also be added on-the-fly and with minimal effort as long as it uses information of previously learned task-specific heads to produce the output.

### *Inter-task Relatedness*

Unlike existing methods, our proposed U-Rep is modality-specific and takes into consideration tasks relationship. To assess inter-task relatedness, we used statistical distributions, human intuition/expertise, and empirical evaluations as follows. First, we consider two tasks related if and only if they both operate on the same modality or have similar distributions. For example, a U-Rep created using unsupervised learning from fetal ultrasound images (source) can be used, without relearning image characteristics, to perform fetal ultrasound segmentation or B-mode echocardiography classification since both the source and target tasks operate on ultrasound images or have similar distributions. Mathematically, this can be formulated as $P_S(X_S) \sim P_{1:N}(X_{1:N})$, where $P_S(X_S)$ represents the distribution of the source modality and $P_{1:N}(X_{1:N})$ represents the distributions of N target tasks. Similarly, we consider two target tasks ($T1$ and $T2$) with different outputs or ground truth (GT) labels ($Y_{GT_{T1}} \neq Y_{GT_{T2}}$) related if and only if $P_{T1}(X_{T1}) \sim P_{T2}(X_{T2})$. This means both tasks operate on the same image domain or have similar distributions. We then use human intuition/expertise to group tasks together. For example, we can assume, based on human intuition, that segmentation and classification are related because they both are 2D classification-based tasks; intuitively, semantic segmentation is a pixel-level classification. Finally, we measure the impact of tasks on each other empirically by comparing the performance of different pairing of tasks.

We believe this approach for assessing task relatedness is more suitable for medical images than the approach proposed in Taskonomy[13] because medical images often contain similar anatomy and have a limited number of tasks as compared to natural images.

### *Architecture and Optimization*

Given that U-Rep holds the promise for learning *N* related tasks, a question that arises is: what makes a specific unified backbone better? The answer depends on the input space and the tasks at hand. For example, given a dataset with different label spaces for different tasks, U-Rep can be optimized using the task with the largest number of labeled data (supervised learning). Instead of using the labels of a specific task to learn and optimize U-Rep, an autoencoder can be optimized for a given source modality (S) and used as a shared backbone for jointly learning multiple target tasks (unsupervised learning). In addition to supervised and unsupervised learning, other learning techniques, such as reinforcement, can be used to construct U-Rep.

The architecture of U-Rep can either be a customized or a state-of-the-art convolutional neural network (e.g., VGG). This shared backbone can be optimized using optimization techniques, such as grid, genetic, and probabilistic optimization. As depicted in Figure 1, the optimized U-Rep can be used as a unified representation to jointly learn *N* related tasks. In addition, several pre-processing operations, such as image enhancement, can be learned while constructing and optimizing U-Rep. Using the proposed method (U-Rep) for learning pre-processing operations and simultaneously solving related tasks prevents unnecessary repetitions, reduces computation, and enhances generalizability, especially for tasks with limited amounts of labeled data.

Algorithm 1 and Algorithm 2 provide the steps for constructing and optimizing U-Rep and task-specific heads, respectively. Note that task-specific layers, which provide the task-specific outputs, need to be added prior to training the target tasks.

## Experiments and Results

We evaluate the proposed method (U-Rep) using two medical imaging modalities: CXR and echo Doppler. For CXR, we create an optimized model (U-Rep 1) that capture the representation of CXR data using CDAE. The optimized model is then used as a shared backbone to simultaneously learn two tasks: lung segmentation and abnormality classification. For echo Doppler, we create a customized optimized backbone (U-Rep 2) that capture the representation of the task that is related to other tasks and has the largest amount of labeled data (Doppler flow classification). The optimized model is then used as a shared backbone to simultaneously learn two related tasks: envelope segmentation and quality assessment. For all classification tasks, the models

**Algorithm 1** U-Rep optimization using source $S$
___
   **Input:** $S = \{X_S, Y_{GT_S}\}$, $Y_{GT_S} = \{\}$ for unsupervised
   **Output:** Optimized U-Rep ($W_S$ and $\Theta_S$)
   Initialize weights $W_S$ and $\Theta_S$
   **while** *epochs* **do**
      Search for optimal $W_S$ and $\Theta_S$ that minimizes Loss (validation)
      **if** $Loss_{itr+1} < Loss_{itr}$ **then**
         Update and save $W_S$ and $\Theta_S$
      **else**
         Keep iterating
      **end if**
   **end while**
   Save $W_S$ and $\Theta_S$
___

**Algorithm 2** Task-specific heads training
___
   **Input:** Optimized U-Rep ($W_S$, $\Theta_S$),
   $T_i = \{X_i, Y_{GT_i}\}$, where $i = 1 : N$
   **Output:** $Y_i$
   **Parameters:** $W_S$, $\Theta_S$, $W_i$, $\Theta_i$
   Initialize task-specific weights $W_i$ and $\Theta_i$
   **while** *patience* **do**
      Calculate task-specific Loss (validation)
      **if** $Loss_{itr+1} < Loss_{itr}$ **then**
         Update $W_S$, $W_i$ and $\Theta_i$
      **else**
         Keep iterating
      **end if**
   **end while**
   Save task-specific model
___

are optimized to minimize the categorical cross-entropy (CCE) loss. For segmentation tasks, we propose to optimize the models through minimizing a combination of binary cross-entropy (BCE) and Dice losses as follows:

$$Loss_{Segmentation} = Loss_{BCE} + Loss_{Dice} \tag{1}$$

where,

$$L_{BCE} = -\frac{1}{j}\Sigma_{i=1}^{j} Y_i \times log(P(Y_i)) + (1 - Y_i) \times log(1 - P(Y_i)) \tag{2}$$

and,

$$L_{Dice} = 1 - \frac{2|Y_{GT} \cap Y_i|}{|Y_{GT}| + |Y_i|} \tag{3}$$

where $Y_i$ and $Y_{GT}$ in equation 3 denote pixel-level predictions and ground truth (GT) annotations. We empirically found that the combination of BCE and Dice losses (equation 1) improved model optimization due to the interplay between the global and local feature extraction capabilities of these loss functions. The loss functions of all tasks are then combined.

In all experiments, we investigate the following questions:

- How does training/learning different tasks based on the optimized U-Rep compare to that of using individual models separately (traditional approach)?

- How does the performance of U-Rep compare to the traditional approach?

- How does U-Rep impact the performance of different tasks, especially those with a limited amount of labeled data?
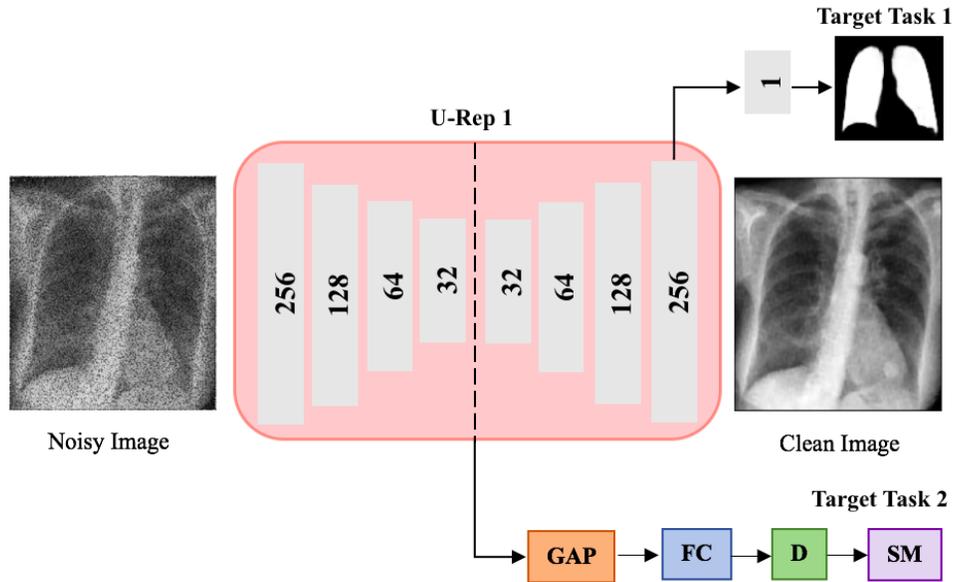
**Figure 2.** U-Rep 1 (pink box) for jointly learning CXR tasks. GAP, FC, D, and SM indicate global average pooling, fully connected, dropout, and Softmax layers. Target task 1 is lung segmentation and target task 2 is abnormality classification.

We describe next the structure as well as the steps for training and optimizing the two U-Rep models (U-Rep 1 and U-Rep 2). We also present their experimental setups and provide empirical answers to the above questions. We conclude the section by comparing our proposed method with the state of the art. Before proceeding, we would like to note that we assessed the relatedness between tasks in all experiments as described in the previous section. The complete implementation of this work and the training parameters will be made publicly available upon acceptance.

## U-Rep 1: Model Construction via Unsupervised Learning

We trained and optimized U-Rep 1 on RSNA CXR dataset using unsupervised learning. While training and constructing U-Rep 1, noise reduction pre-processing operation was learned. We then used the trained single representation (U-Rep 1) to preform lung segmentation and abnormality classification.

### Noise Reduction

We used CDAE to reduce the noise present in CXR images. Noise reduction is an essential pre-processing operation that improves the quality and diagnostic value of medical images while maintaining spatial resolution. In our experiments, we added Gaussian noise (mean = 0 and standard deviation = 0.03) to the images and used CDAE to reconstruct clean images. We quantitatively measured the performance of reconstruction using peak signal-to-noise ratio (PSNR). The typical range of PSNR for images with 8-bit depth is between 30 to 50 dB, where higher is better[18]. We observed that the PSNR value computed between the noisy and reconstructed images was 39.74dB for the test set. This indicated that the optimized CDAE faithfully reconstructed the test images. As the output images were clean, CDAE learned to extract useful CXR features and ignore noise information.

### U-Rep 1 Architecture and Training

We trained CDAE in an unsupervised manner using RSNA CXR dataset. CDAE consists of two parts: (1) an encoder that encodes the input information to its latent space representation and (2) a decoder that reconstructs the image from the latent space representation with the minimal reconstruction error. As depicted in Figure 2, we constructed CDAE with four convolutional layers in the encoder to compress the input to its latent space representation. We used strided convolutions (2, 2) instead of max-pooling layers to increase the expressive capacity of the network, which would improve the overall performance without increasing the number of parameters[19]. We used batch normalization layers to improve generalization and ReLU to speed up model training, resulting in faster convergence. We used upsampling layers in a symmetrical decoder to reconstruct the input from the latent space representation. The model was optimized using Talos optimization tool[20] to minimize the mean squared error (MSE) and reconstruct the input with minimal reconstruction error. We optimized the kernel size [3, 5, 7] and optimizer [SGD, Adam, RMSprop].

| Model | Dropout | Optimizer | Kernel Size | Stride | Validation Loss | Training Time |
|---|---|---|---|---|---|---|
| *Shared optimized U-Rep 1 (Proposed)* | | | | | | |
| Source Modality (auto-encoder) | - | RMSprop* | 3* | 2 | 0.0001 | 18132 seconds |
| Target Task 1 (Segmentation) | - | RMSprop | ↑ | ↑ | 0.0649 | 886.92 seconds |
| Target Task 2 (Classification) | 0.3 | SGD | ↑ | ↑ | 0.4170 | 210.39 seconds |
| *Two Models Optimized Individually (Traditional)* | | | | | | |
| Optimized Task 1 (Segmentation) | - | RMSprop* | 5* | 2 | 0.1016 | 4238.98 seconds |
| Optimized Task 2 (Classification) | 0.1* | SGD* | 5* | 2 | 0.5185 | 305.72 seconds |

**Table 1.** Training parameters and validation loss for U-Rep 1 and traditional approach. CDAE was optimized on RSNA CXR dataset (source modality) and used with two target tasks. ∗ indicates optimized parameters and ↑ indicates inherited parameters. All models have an initial learning rate of 1-e3 that is reduced when validation loss plateau.

| Approach | Tasks | Accuracy | AUC | Sensitivity | Precision | F-score | IoU |
|---|---|---|---|---|---|---|---|
| U-Rep 1 | Target Task 1: Segmentation | 0.9871 | - | - | - | - | 0.9771 |
| | Target Task 2: Classification | 0.8292 | 0.8800 | 0.8292 | 0.8308 | 0.8288 | – |
| Traditional | Optimized Task 1: Segmentation | 0.9702 | - | - | - | - | 0.9512 |
| | Optimized Task 2: Classification | 0.800 | 0.8500 | 0.7990 | 0.8043 | 0.7989 | – |

**Table 2.** Performance of lung segmentation and abnormality classification tasks using U-Rep 1 (denoising AE) and the traditional approach.

Table 2 shows the training parameters of the optimized CDAE backbone and the validation loss. Once the optimized weights of CDAE is learned, they are used to simultaneously learn two related tasks: lung segmentation and abnormality classification.

### Task-specific Heads Training

To create a task-specific head for abnormality classification, the optimized CDAE was used as a shared backbone for feature extraction (Figure 2). Specifically, the encoder part of CDAE was instantiated and appended with the following layers: global average pooling (GAP), fully connected (FC), dropout (D), and Softmax (SM) layers. The abnormality classification model was trained and validated to minimize the CCE loss using Montgomery and Shenzen CXR datasets with an aim to improve adaptation, performance, and generalization.

To perform lung segmentation, we instantiated the CDAE with its optimized weights and replaced the final convolutional layer with the one that has a single neuron to generate the binary lung masks. The segmentation model was optimized to minimize the combination of BCE and Dice losses (equation 1) to obtain the highest segmentation accuracy. Table 2 presents the training parameters for the segmentation and abnormality classification tasks. As shown in Table 2, the optimized weights were used to jointly learn segmentation (Target Task 1) and abnormality classification (Target Task 2) tasks. These tasks inherited U-Rep 1 (CDAE) optimized weights and parameters. RMSprop and Stochastic Gradient Descent (SGD) optimizers were used to optimize the segmentation and abnormality classification heads.

As a baseline, we optimized two individual models, one for each task (traditional approach), to separately learn lung segmentation and abnormality classification tasks. The individual segmentation model has the same architecture as CDAE, but with randomly initialized weights and optimized for the search parameters. The individual classification model has the same architecture as that of the encoder of the optimized CDAE and appended with GAP, FC, D, and SM layers. This model was randomly initialized and then optimized for the search parameters. The training parameters of these baseline models are shown in Table 2.

As can be observed from Table 2, U-Rep 1 remarkably decreased the validation loss and training time for both target tasks as compared to the traditional approach. The total training time using our approach is 19,229.31 seconds or 320 minutes (U-Rep 1 + target task 1 + target task 2). Given that a separate model for noise reduction was used, then the total training time using the traditional approach would be 22,676.7 seconds or 378 minutes (noise reduction model + optimized task 1 + optimized task 2).
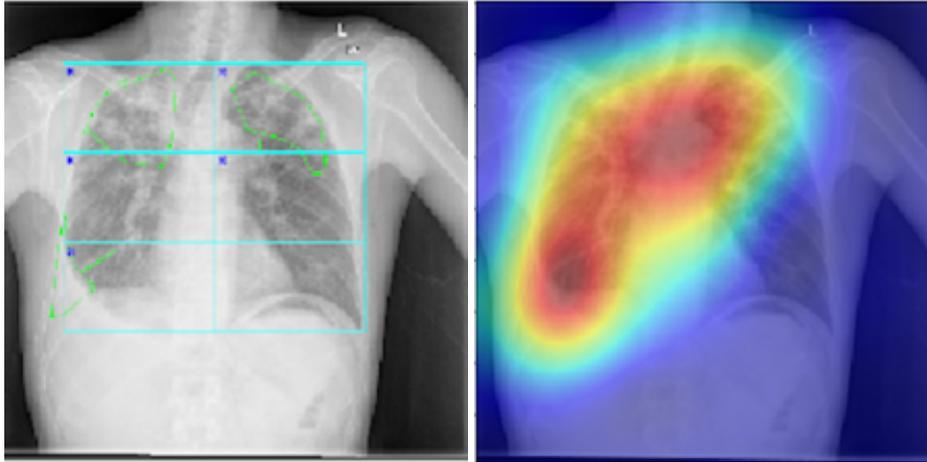
**Figure 3.** Visual explanation task. Left: original image with GT annotation (green outline). Right: Grad-CAM visualization. Note that high activations (red pixels) are observed in the affected ROIs.

### U-Rep 1 Performance

We present here the performance results of using the optimized U-Rep 1 or CDAE for jointly learning two CXR tasks: lung segmentation and abnormality classification. For optimizing CDAE, we used the RSNA CXR dataset with 70%, 20%, 10% patient-level splits for training, validation, and testing, respectively. Montgomery and Shenzhen CXR collections (small datasets) were used for training and validating task-specific heads.

Table 3 presents the performance of using the optimized CDAE (U-Rep 1) as a shared backbone for lung segmentation and abnormality classification tasks. It also presents the performance using individual optimized models (traditional approach). We measured the performance of the classification task using the average accuracy, area under the curve (AUC), sensitivity, precision, and F-score. We used the accuracy (pixel-level) and intersection over union (IoU) metrics to report the performance of the segmentation task. As can be seen, U-Rep 1 remarkably improved the performance of both tasks as compared to the traditional approach. These results provide an evidence that the proposed method, by transferring the knowledge from a representative source to related target tasks, can improve the generalizability across the tasks, and result in better performance as compared to the traditional approach. This is especially true for tasks with a limited number of labeled examples.

After learning the classification task, we used gradient-weighted class activation mapping (Grad-CAM)[21] to localize discriminative features or areas the model looks at when classifying the CXRs into normal and abnormal classes (derivable task). Grad-CAM generates heat maps showing high activated regions denoted by red pixels. Note that Grad-CAM is one approach to visualize classification predictions; several other approaches can be used instead of Grad-CAM for visual explanation. Figure 3 presents the visualization of the classification task using Grad-CAM (task derivable from the classification head).

## U-Rep 2: Model Construction via Supervised Learning

We trained and optimized U-Rep 2 on echo Doppler dataset using supervised learning. We then used this single trained representation (U-Rep 2) to preform Doppler quality assessment and envelope segmentation.

### U-Rep 2 Architecture and Training

We used a customized CNN as the shared backbone and optimized it for the flow classification task since it is related to all tasks and has the largest amount of labeled data. As can be seen in Figure 4, the shared backbone has six convolutional layers with the same padding. Dilated kernels were used in the $4^{th}$, $5^{th}$, and $6^{th}$ convolutional layers to capture wider context at a reduced computational cost that would help with the related segmentation task[22]. We used ReLU after each convolutional layer to speed-up model training and convergence. The output of the deepest convolutional layer from the optimized CNN was fed to the GAP and FC layers. The output of the FC layer was fed to a dropout layer to reduce overfitting. The last FC layer has three neurons corresponding to three classes: TR, MV, and MA. The model was optimized to minimize the CCE loss and increase the prediction probabilities. We used Talos optimization[20] to optimize the following parameters: kernel size [3, 5, 7], dilation rate [2, 3], dropout ratio [0.1, 0.3, 0.5], and optimizer [SGD, Adam, RMSprop]. Table 3 shows the training parameters and validation loss of the optimized U-Rep 2 backbone.
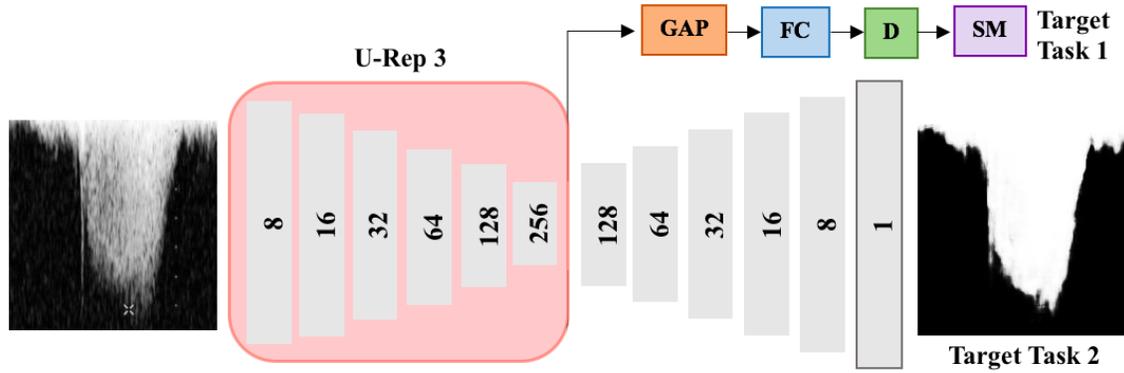
**Figure 4.** U-Rep 2 (pink box) for jointly learning echo Doppler tasks. Target task 1 is quality assessment and target task 2 is envelope segmentation.

| Tasks | Dropout | Optimizer | Kernel Size | Dilation | Loss | Accuracy | Training Time |
|---|---|---|---|---|---|---|---|
| *Shared Optimized U-Rep 2 (Proposed)* | | | | | | | |
| Optimized U-Rep 2 (Flow Classification) | 0.5* | Adam* | 3* | 2* | 0.0035 | 1.0 | 1871.78 seconds |
| Target Task 1 (Quality Assessment) | 0.5 | SGD | ↑ | ↑ | 0.2943 | 0.9318 | 1039.14 seconds |
| Target Task 2 (Envelope Segmentation) | - | Adam | ↑ | ↑ | 0.2219 | 0.9768 | 3107.04 seconds |
| *Three Models Optimized Individually (Traditional)* | | | | | | | |
| Optimized Task 1 (Flow Classification) | 0.5* | Adam* | 3* | 2* | 0.0035 | 1.0 | 1871.78 seconds |
| Optimized Task 2 (Quality Assessment) | 0.3* | Adam* | 3* | 2* | 0.3458 | 0.9035 | 1978.95 seconds |
| Optimized Task 3 (Envelope Segmentation) | - | RMSprop* | 5* | 3* | 0.2252 | 0.9528 | 6035.92 seconds |

**Table 3.** Training parameters and validation loss/accuracy for the proposed (U-Rep 2) and traditional approaches. A customized CNN was optimized using the labeled data of the flow classification task and shared to jointly learn two tasks: quality assessment and envelop segmentation. ∗ indicates the optimized parameters and ↑ indicates inherited parameters. All models have an initial learning rate of 1-e3 that is reduced when validation loss plateau.

| Approach | Tasks | Accuracy | AUC | Sensitivity | Precision | F-score | IoU |
|---|---|---|---|---|---|---|---|
| U-Rep 2 | Source Task: Classification | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | - |
|  | Target Task 1: Quality Assessment | 0.9139 | 0.9800 | 1.0 | 0.8125 | 0.9166 | – |
|  | Target Task 2: Segmentation | 0.9568 | – | – | – | – | 0.972 |
| Traditional | Optimized Task 2: Quality Assessment | 0.8931 | 0.8800 | 0.9139 | 0.9257 | 0.9027 | – |
|  | Optimized Task 3: Segmentation | 0.9357 | – | – | – | – | 0.954 |

**Table 4.** Performance of quality assessment and envelope segmentation tasks using U-Rep 2 and traditional approaches.

### Task-specific Heads Training

The weights of the optimized U-Rep backbone, which was generated using the labeled data of the flow classification task, were used to jointly learn quality assessment (Target Task 1) and envelope segmentation (Target Task 2) tasks. These tasks inherited the optimized weights, kernel size, and dilation rate from U-Rep 2.

To create a task-specific head for quality assessment, the optimized U-Rep 2 was truncated at the deepest convolutional
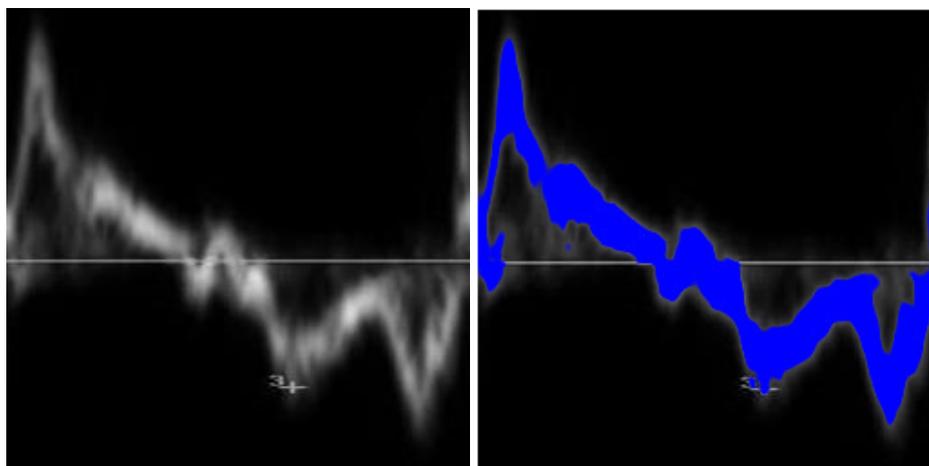
**Figure 5.** Example of the mask, which was generated using our proposed U-Rep, superimposed on MA flow envelope.

layer and appended with the GAP, FC, D, and SM layers. This task-specific head was trained using SGD optimizer to minimize the CCE loss. To create the task specific-head for envelope segmentation, the optimized U-Rep 2 was truncated at the deepest convolutional layer; a symmetrical decoder was constructed and appended. This task-specific head was trained using Adam optimizer to optimize a combination of BCE and Dice losses (equation 1). Table 3 shows the training parameters and validation loss/accuracy of both task-specific heads.

As a baseline, we optimized two individual models, one for each task, to separately learn envelope segmentation and quality assessment tasks (traditional approach). The individual quality assessment and segmentation models have the same architecture as that of the optimized CNN backbone; the decoder of the segmentation model was created symmetrically from the optimized CNN. The training parameters and validation loss/accuracy of the baseline models are shown in Table 3.

### *U-Rep 2 Performance*

We used U-Rep 2 to learn three echo Doppler tasks: flow classification, quality assessment, and envelope segmentation. In all experiments, we used 70%, 20%, 10% patient-level splits for training, validation, and testing, respectively. Table 4 presents the performance of target tasks based on optimized U-Rep 2 and the traditional approach. As observed, U-Rep 2 improved the segmentation performance, and the sensitivity and AUC of the quality assessment task. Figure 5 shows an instance of the mask, which was generated using our proposed method, superimposed on MA flow envelope.

After learning the target tasks, we derived two tasks. The first task is the visual interpretation of the classification predictions. Visual interpretation is important in medical applications because it provides an understanding of why a specific outcome was chosen and assists in model optimization. The second derivable task is a recommendation generated based on combining the outputs of the flow classification and quality assessment tasks. This recommendation is used to decide if the image is suitable for further analysis. For example, the echocardiographer manually excludes low-quality TR flow images with unclear envelopes from further analysis. Figure 6 presents visualizations (derivable task 1) of MA and MV flows predictions using Grad-CAM[21]. The recommendation (derivable task 2) generated based on merging the outputs of the flow classification and quality assessment heads is shown in Figure 7. As can be seen, merging the output of flow classification and quality assessment heads allowed to decide if the image is usable and should be used for further analysis; cardiologists manually exclude non-usable images since they increase the subjectivity and decrease the accuracy of measurements, which can impact diagnosis.

Two main conclusions can be drawn from the empirical evaluation and results of three U-Reps models. First, the proposed U-Rep approach improved the ability to generalize across tasks by transferring the knowledge from a representative source image modality or task (as applicable) to related target tasks, and led to superior overall performance. This is especially true for the tasks with limited amount of labeled data. Second, it reduced the model and computational complexity and prevented unnecessary repetitions of training individual models. The experimental results are promising and provide empirical evidence for the superiority of the proposed method as compared to the traditional approach for medical image analysis.

## Comparison with the State of the Art

We compare our work with two existing recent works: Tasknomy[13] and Genesis[14]. We compared the performance of U-Rep 2 for the echo Doppler dataset with the performance of Tasknomy[13] and Genesis[14]. As shown in Table 5, our approach achieved
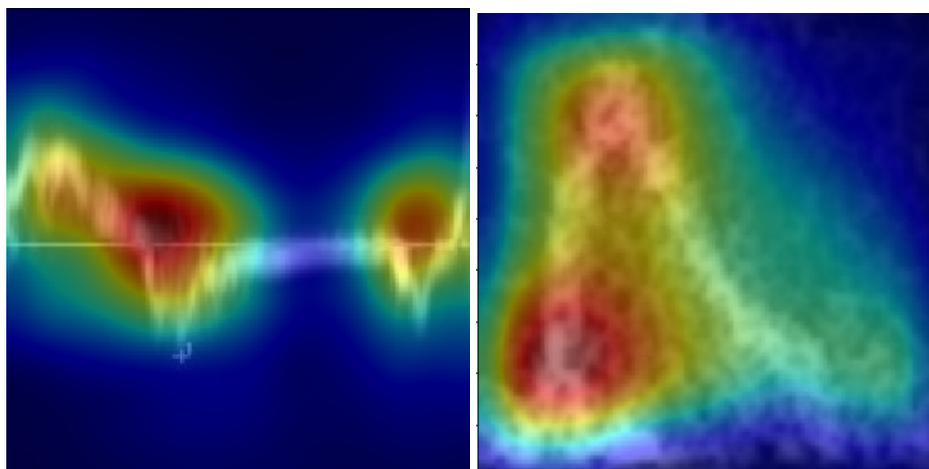
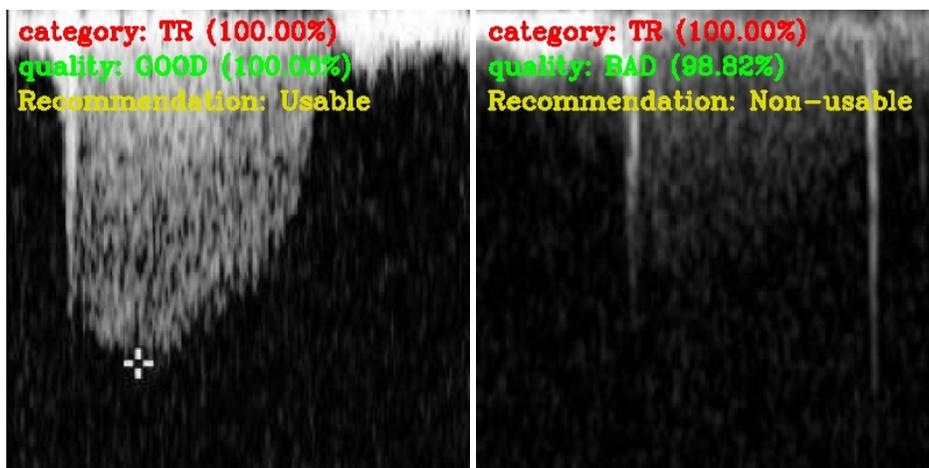**Figure 6.** Derivable task 1. Visual explanation for MA (left) and MV (right) flows.



**Figure 7.** Derivable task 2. Examples of merging the outputs of flow classification and quality assessment heads to select only good quality image with dense envelope.

better classification and segmentation performance than existing methods. These results suggest that U-Rep adds greater value for machine learning applied to medical image analysis.

| | Flow Classifcation | | | Quality Assessment | | | Envelope Segmentation | |
|---|---|---|---|---|---|---|---|---|
| | Accuracy | AUC | F score | Accuracy | AUC | F score | Accuracy | IoU |
| **U-Rep 2** | **1.0** | **1.0** | **1.0** | 0.91 | **0.98** | **0.92** | **0.95** | **0.97** |
| **Tasknomy**[13] | 0.93 | 0.94 | 0.97 | 0.87 | 0.88 | 0.84 | 0.83 | 0.81 |
| **Genesis**[14] | 0.97 | 0.91 | 0.91 | **0.93** | 0.91 | 0.87 | 0.95 | 0.92 |

**Table 5.** Comparison between U-Rep 2 and existing works using echo Doppler dataset; bold values indicates superior performance.

## Possible Extensions and Conclusion

Although we demonstrated the feasibility of using a unified representation (U-Rep) to learn two pre-processing operations (i.e., noise reduction and image enhancement) and simultaneously solved key tasks in medical image analysis, U-Rep is flexible and can be easily extended to integrate other pre-processing operations and learn any number of related tasks. For example,

bounding box regression task can be added to U-Rep by appending a task-specific head that has region pooling layers for extracting region-wise features and FC layers for box classification and regression. Similarly, U-Rep can be easily extended to analyze 3D images using a 3D CNN as the backbone. Finally, several optimization methods (e.g., Bayesian) can be used with U-Rep to learn the optimal set of hyperparameters.

In conclusion, we demonstrated, for the first time, that a single unified representation (U-Rep) can be shared to jointly learn key medical image analysis tasks. We presented empirical evaluations of three different U-Rep models using two datasets: CXR and echo Doppler. In the first case (U-Rep 1), an autoencoder (unsupervised learning) was optimized to learn the feature space of CXR modality and then shared to jointly perform lung segmentation and abnormality classification tasks. In the second case (U-Rep 2), we optimized U-Rep on the task with the largest number of labels (supervised learning) to perform envelope segmentation and quality assessment. We also proposed derivable tasks (e.g., visual explanation and decision recommendation), which are the tasks that generate outputs based on the information learned by a single or a combination of task-specific heads. Our experimental results demonstrated that U-Rep enhanced generalizability and significantly improved performance, especially for the tasks with limited amounts of labeled data. In addition, we observed that U-Rep decreased computations as compared to the traditional approach. These results are promising and prove the efficiency and superiority of the proposed U-Rep for medical image analysis.

## Data Availability

We used three publicly available CXR datasets: RSNA, Montgomery, and Shenzhen. Adequate references of these datasets are provided in CXR dataset section. We also used a private echo Doppler dataset collected in the clinical center at the National Institutes of Health (NIH). The data collection was approved by the NIH Ethics Review Board (IRB18-NHLBI-00686). This dataset is not publicly available as it contains information that could compromise research participant privacy/consent.

## References

1. Zhou, S. K., Greenspan, H. & Shen, D. *Deep learning for medical image analysis* (Academic Press, 2017).

2. Litjens, G. *et al.* A survey on deep learning in medical image analysis. *Med. image analysis* **42**, 60–88 (2017).

3. Balakrishnan, G., Zhao, A., Sabuncu, M. R., Guttag, J. & Dalca, A. V. An unsupervised learning model for deformable medical image registration. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 9252–9260 (2018).

4. Thung, K.-H. & Wee, C.-Y. A brief review on multi-task learning. *Multimed. Tools Appl.* **77**, 29705–29725 (2018).

5. Misra, I., Shrivastava, A., Gupta, A. & Hebert, M. Cross-stitch networks for multi-task learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 3994–4003 (2016).

6. Rusu, A. A. *et al.* Progressive neural networks. *arXiv preprint arXiv:1606.04671* (2016).

7. Kendall, A., Gal, Y. & Cipolla, R. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 7482–7491 (2018).

8. Chen, Z., Badrinarayanan, V., Lee, C.-Y. & Rabinovich, A. Gradnorm: Gradient normalization for adaptive loss balancing in deep multitask networks. In *International Conference on Machine Learning*, 793–802 (2018).

9. Zhang, Z., Luo, P., Loy, C. C. & Tang, X. Facial landmark detection by deep multi-task learning. In *European conference on computer vision*, 94–108 (Springer, 2014).

10. Elhoseiny, M., El-Gaaly, T., Bakry, A. & Elgammal, A. Convolutional models for joint object categorization and pose estimation. *arXiv preprint arXiv:1511.05175* (2015).

11. Eigen, D. & Fergus, R. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *Proceedings of the IEEE international conference on computer vision*, 2650–2658 (2015).

12. Xiao, T., Liu, Y., Zhou, B., Jiang, Y. & Sun, J. Unified perceptual parsing for scene understanding. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 418–434 (2018).

13. Zamir, A. R. *et al.* Taskonomy: Disentangling task transfer learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 3712–3722 (2018).

14. Zhou, Z. *et al.* Models genesis: Generic autodidactic models for 3d medical image analysis. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 384–393 (Springer, 2019).

15. Shih, G. *et al.* Augmenting the national institutes of health chest radiograph dataset with expert annotations of possible pneumonia. *Radiol. Artif. Intell.* **1**, e180041 (2019).

16. Jaeger, S. *et al.* Two public chest x-ray datasets for computer-aided screening of pulmonary diseases. *Quant. imaging medicine surgery* **4**, 475 (2014).

17. Rajaraman, S., Candemir, S., Kim, I., Thoma, G. & Antani, S. Visualization and interpretation of convolutional neural network predictions in detecting pneumonia in pediatric chest radiographs. *Appl. Sci.* **8**, 1715 (2018).

18. Welstead, S. T. *Fractal and wavelet image compression techniques* (SPIE Optical Engineering Press Bellingham, Washington, 1999).

19. Zeiler, M. D., Krishnan, D., Taylor, G. W. & Fergus, R. Deconvolutional networks. In *2010 IEEE Computer Society Conference on computer vision and pattern recognition*, 2528–2535 (IEEE, 2010).

20. Talos. Autonomio talos [computer software]. *Retrieved from http://github.com/autonomio/talos* (2019).

21. Selvaraju, R. R. *et al.* Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision*, 618–626 (2017).

22. Yu, F. & Koltun, V. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122* (2015).

## Acknowledgments

## Author contributions and Correspondence

G.Z., S.R, and S.A. designed the framework, G.Z. and S.R conceived the experiment(s), G.Z. and S.R conducted the experiment(s), G.Z., S.R, and S.A. analysed the results. G.Z. wrote the manuscript. All authors reviewed the manuscript. Correspondence and requests for materials should be addressed to G.Z.