

COVID-19-CT-CXR: A Freely Accessible and Weakly Labeled Chest X-Ray and CT Image Collection on COVID-19 From Biomedical Literature

Yifan Peng¹, Yuxing Tang, Sungwon Lee², Yingying Zhu, Ronald M. Summers³, and Zhiyong Lu

Abstract—The latest threat to global health is the COVID-19 outbreak. Although there exist large datasets of chest X-rays (CXR) and computed tomography (CT) scans, few COVID-19 image collections are currently available due to patient privacy. At the same time, there is a rapid growth of COVID-19-relevant articles in the biomedical literature, including those that report findings on radiographs. Here, we present COVID-19-CT-CXR, a public database of COVID-19 CXR and CT images, which are automatically extracted from COVID-19-relevant articles from the PubMed Central Open Access (PMC-OA) Subset. We extracted figures, associated captions, and relevant figure descriptions in the article and separated compound figures into subfigures. Because a large portion of figures in COVID-19 articles are not CXR or CT, we designed a deep-learning model to distinguish them from other figure types and to classify them accordingly. The final database includes 1,327 CT and 263 CXR images (as of May 9, 2020) with their relevant text. To demonstrate the utility of COVID-19-CT-CXR, we conducted four case studies. (1) We show that COVID-19-CT-CXR, when used as additional training data, is able to contribute to improved deep-learning (DL) performance for the classification of COVID-19 and non-COVID-19 CT. (2) We collected CT images of influenza, another common infectious respiratory illness that may present similarly to COVID-19, and fine-tuned a baseline deep neural network to distinguish a diagnosis of COVID-19, influenza, or normal or other types of diseases on CT. (3) We fine-tuned an unsupervised one-class classifier from non-COVID-19 CXR and performed anomaly detection to detect COVID-19 CXR. (4) From text-mined captions and figure descriptions, we compared 15 clinical symptoms and 20 clinical findings of COVID-19 versus those of influenza to demonstrate the disease differences in the scientific publications. Our database is unique, as the figures are retrieved along with relevant text with fine-grained descriptions, and it can be extended easily in the future. We believe that our work is complementary to existing resources and hope that it will contribute to medical image analysis of the COVID-19 pandemic. The dataset, code, and DL models are publicly available at <https://github.com/ncbi-nlp/COVID-19-CT-CXR>.

Index Terms—COVID-19, chest X-ray, CT

1 INTRODUCTION

THE latest threat to global health is the ongoing outbreak of the COVID-19 caused by SARS-CoV-2 [1]. So far, pneumonia appears to be the most frequent and serious manifestation, and major complications, such as acute

respiratory distress syndrome (ARDS), can present shortly after the onset of symptoms, contributing to the high mortality rate of COVID-19 [2], [3], [4]. Chest X-rays (CXR) and chest computed tomography (CT) scans are playing a major part in the detection and monitoring of these respiratory manifestations. In some cases, CT scans have shown abnormal findings in patients prior to the development of symptoms and even before the detection of the viral RNA [5], [6], [7].

With the shortage of specialists who have been trained to accumulate experiences with COVID-19 diagnosis, there has been a concerted move toward the adoption of artificial intelligence (AI), particularly deep-learning-based methods, in COVID-19 pandemic diagnosis and prognosis, in which well-annotated data always play a critical role [8]. Although there exist large public datasets of CXR [9], [10], [11] and CT [12], there are few collections of COVID-19 images to effectively train a deep neural network [13], [14], [15]. Nevertheless, we have seen a growing number of COVID-19 relevant articles in PubMed [16], [17]. In addition, there is a recent COVID-19 initiative to expand access via PubMed Central Open Access (PMC-OA) Subset to coronavirus-related publications and associated data (<https://www.ncbi.nlm.nih.gov/pmc/about/covid-19-faq/>). As a result,

- Yifan Peng is with the NCBI/NLM/NIH and Department of Population Health Sciences, Weill Cornell Medicine, New York, NY 10065 USA. E-mail: yip4002@med.cornell.edu.
- Yuxing Tang, Sungwon Lee, and Ronald M. Summers are with the Imaging Biomarkers and Computer-Aided Diagnosis Laboratory, Radiology and Imaging Sciences Department, National Institutes of Health (NIH) Clinical Center, Bethesda, MD 20892 USA. E-mail: ytang.cv@hotmail.com, sungwon.lee@nih.gov, rsummers@cc.nih.gov.
- Yingying Zhu is with the Imaging Biomarkers and Computer-Aided Diagnosis Laboratory, Radiology and Imaging Sciences Department, National Institutes of Health (NIH) Clinical Center, Bethesda, MD 20892 USA, and also with the Department of Computer Science and Engineering, University of Texas at Arlington, Arlington, TX 76019 USA. E-mail: yingying.zhu@nih.gov.
- Zhiyong Lu is with the National Center for Biotechnology Information (NCBI), National Library of Medicine (NLM), National Institutes of Health (NIH), Bethesda, MD 20894 USA. E-mail: zhiyong.lu@nih.gov.

Manuscript received 29 June 2020; revised 9 Oct. 2020; accepted 19 Oct. 2020. Date of publication 4 Nov. 2020; date of current version 1 Mar. 2021. (Corresponding author: Zhiyong Lu.) Digital Object Identifier no. 10.1109/TBDATA.2020.3035935

more articles (> 10,000 as of May 9, 2020) relevant to the COVID-19 pandemic or prior coronavirus research were added through PMC-OA with a free-reuse license for secondary analysis.

Non-textual components (e.g., figures and tables) provide key information in many scientific documents and are considered in many tasks, including search engine and knowledge base construction [18], [19]. As such, we have recently seen a growing interest in mining figures within scientific documents [20], [21], [22]. In the medical domain, figures also are a topical interest because they often contain graphical images, such as CXR and CT [23], [24]. Extracting CXR and CT from biomedical publications, however, is neither well studied nor well addressed.

For the above reasons, there is an unmet need to construct the COVID-19 image dataset from PMC-OA to allow researchers to freely access the images along with a description of the text. In this paper, we thus introduce an effective framework to construct a CXR and CT database from PMC-OA and propose a public database, termed COVID-19-CT-CXR. In contrast to previous approaches that relied solely on the manual submission of medical images to the repository, in this work, figures are automatically collected by using the integration of medical imaging and natural-language processing with limited human annotation efforts. In addition, figures in this database are partnered with text that describes these cases with details, a feature not found in other such datasets.

The framework consists of three steps. First, we extracted figures, associated captions, and relevant figure descriptions in the PMC-OA article. Such extraction is non-trivial due to the diverse layout and large volume of articles in the PMC-OA subset. Second, we separated compound figures into subfigures, as medical figures often comprise multiple image panels [21], [24]. Third, we classified subfigures into CXR, CT, or others because a large portion of figures in COVID-19 articles are not CXR or CT. To this end, we designed a deep-learning model to distinguish them from other figure types and to classify them accordingly.

We further demonstrate the utility of COVID-19-CT-CXR through a series of case studies. First, using this database as additional training data, we show that existing deep neural networks can receive benefits in the task of COVID-19/non-COVID-19 classification of CT images. Second, we demonstrate that the database can be used to develop a baseline model to distinguish COVID-19, influenza, and other CT, a less-studied topic. Third, we train an unsupervised one-class classifier from non-COVID-19 CXRs and performed anomaly detection to detect COVID-19 CXRs. Fourth, we extract symptoms and clinical findings from the text, using the natural language-processing methods. The symptoms and clinical findings not only confirm the results that radiologists have found but also potentially identify other findings that may have been overlooked.

The remainder of the paper is organized as follows. Section 2 presents the material and methods to build the dataset. Section 3 contains the details of the statistics of the dataset, results of the image type classification, and the use cases. Finally, Sections 4 and 5 provide the discussion, conclusions, and recommendations for future work.

TABLE 1
An Overview of the COVID-19 Relevant
Articles as of May 9, 2020

| Characteristics | <i>n</i> |
|--------------------------------------|----------|
| COVID-19 relevant articles in PMC-OA | 5,381 |
| Prevention | 2,089 |
| Mechanism | 577 |
| Diagnosis | 546 |
| Case Report | 355 |
| Transmission | 354 |
| General | 238 |
| Epidemic Forecasting | 64 |
| Others | 1,158 |
| Journals | 1,145 |
| Figures | 4,407 |

2 MATERIAL AND METHODS

2.1 COVID-19 Relevant Articles on PMC-OA

Articles in this study were collected from the PMC-OA Subset. PubMed Central[®] (PMC) is a free, full-text archive of biomedical and life sciences journal literature (<https://www.ncbi.nlm.nih.gov/pmc/>). PMC-OA is a well-known portion of the PMC articles under a Creative Commons license (or custom license of the Public Health Emergency COVID-19 Initiative in PMC due to the COVID pandemic) that allows for text mining, secondary analysis, and other types of reuse (<https://www.ncbi.nlm.nih.gov/pmc/about/covid-19-faq/>). In this study, we collected COVID-19 relevant articles using LitCovid [16], a curated literature hub for tracking up-to-date scientific information about the 2019 novel coronavirus. LitCovid screens the search results of the PubMed query: "coronavirus" [All Fields] "ncov" [All Fields] OR "cov" [All Fields] OR "2019-nCoV" [All Fields] OR "COVID-19" [All Fields] OR "SARS-CoV-2" [All Fields]. Relevant articles are identified and curated with assistance from an automated machine-learning and text-classification algorithm. As of May 9, 2020, there were 5,381 PMC-OA articles in the collection (Table 1). The topics of articles ranged from diagnosis to treatment to case reports.

2.2 Overview of the COVID-19-CT-CXR Construction

Fig. 1 shows the overview pipeline of the development. For a given PMC-OA article, we first extract figures, associated captions, and relevant figure descriptions in the PMC-OA article. Then, if figures are compound, we separate them into subfigures. We further classify the individual figures into CT, CXR, or other types of scientific images, using a deep-learning model. The final database includes figures with their types and relevant descriptions in the manuscript.

2.3 Text Extraction

In this step, we identify figure captions and relevant text with the referenced figures. To facilitate the automated processing of full-text articles in PMC-OA, [25] convert PMC articles to BioC format, a data structure in XML for text sharing and processing. Each article in BioC format is encoded in UTF-8, and Unicode characters are converted to strings of ASCII characters. The article also includes section types, figures, tables, and references [26]. In this study, we

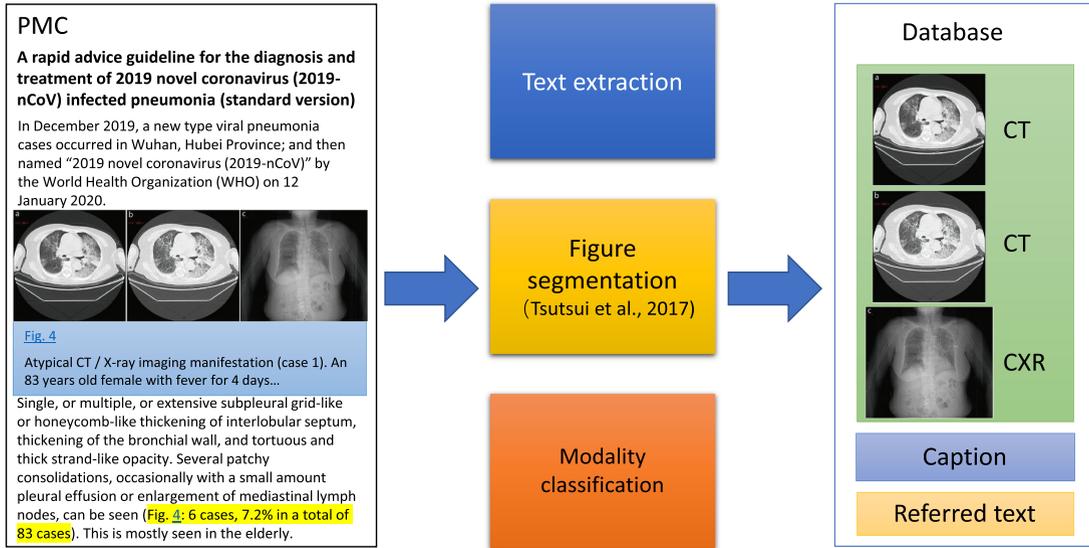


Fig. 1. The overview of the pipeline to collect the images with text.

downloaded the PMC-OA articles through the RESTful web service (<https://www.ncbi.nlm.nih.gov/research/bionlp/APIs/BioC-PubMed/>). We parsed these articles to locate figures with their figure numbers and their captions. We then used the figure number and regular expressions to find where the figure is cross-referenced in the document. Fig. 2 shows an example of a typical biomedical image in the article, “A rapid advice guideline for the diagnosis and treatment of 2019 novel coronavirus (2019-nCoV) infected pneumonia (standard version)” [27]. The examples contain CXR, CT, a figure caption, and text that describes the case with rich information, such as fever, symptoms, and clinical findings.

2.4 Subfigure Separation

Most of the figures in the PMC-OA articles are compound figures. A key challenge here is that one figure may have individual subfigures of the same category (e.g., four CT images) or several categories (e.g., one CXR and one CT

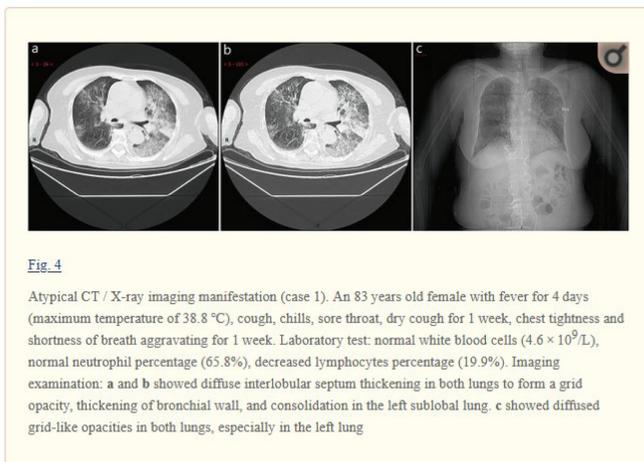


Fig. 2. Examples of CT and CXR that are positive for COVID-19. The figures are from the article, “A rapid advice guideline for the diagnosis and treatment of 2019 novel coronavirus (2019-nCoV) infected pneumonia (standard version)” [27].

image placed side by side). For example, Fig. 2 contains a compound figure with three subfigures [27]. Figs. 2a and 2b are CT images, and Fig. 2c is a CXR. Notably, it is a requirement to decompose compound figures into subfigures before modality classification. In this study, we used a convolutional neural network developed by [24] to separate compound figures. The model was pretrained on the ImageCLEF Medical dataset with an accuracy of 85.9 percent [28].

We applied the model on the figures obtained in previous steps and filtered the subfigures with a size smaller than 224×224 pixels. We consider that subfigures with fewer pixels might be deformed, and most state-of-the-art neural networks in image analysis, such as Inception-v3 [29] and DenseNet [30], require an input size of 224 or larger.

2.5 Image Modality Classification

A large portion of figures in the PMC-OA articles are not CXR or CT images. To distinguish them from other types of scientific figures, we designed a scientific figure classifier that was fine-tuned on a newly created dataset (<https://github.com/ncbi-nlp/COVID-19-CT-CXR>). Table 2 shows the breakdown of the figures by their category in the training and test set. This dataset consists of 2,700 figures in three categories: CXR, CT, and Other scientific figure types. A total of 500 CXRs are randomly picked from the NIH Chest

TABLE 2
Summary of the Dataset for Image Modality Classification

| Modality | Training | Test |
|-----------------------------------|----------|------|
| CXR | | |
| NIH Chest X-ray [11] | 399 | 101 |
| PMC-OA | 38 | 7 |
| CT | | |
| DeepLesion [12] | 415 | 85 |
| PMC-OA | 225 | 21 |
| Other scientific document figures | | |
| DocFigure [31] | 386 | 114 |
| PMC-OA | 737 | 172 |
| Total | 2,200 | 500 |

TABLE 3
Summary of the COVID-19-CT-CXR Dataset

| Characteristics | n |
|------------------------------|--------|
| PMC-OA articles with figures | 1,831 |
| Subfigures | 10,650 |
| CXR | 263 |
| CT | 1,327 |
| Others | 9,060 |

X-ray [11], and 500 CT images are randomly picked from DeepLesion [12]. Other scientific figures are randomly picked from DocFigure [31]. The original DocFigure annotated figures of 28 categories, such as Heat map, Bar plots, and Histogram. Here, we combined these categories into one for simplicity of training the classifier. In addition, we curated 1,200 figures from PMC-OA, using the annotation tool developed by [32].

Our framework uses DenseNet121 to classify image types [33]. The weights (or parameters) were pretrained on ImageNet [34]. We replaced the last classification layer with a fully connected layer with a softmax operation that outputs the approximate probability that an input image is a CXR, CT, or other scientific figure type. All images were resized to 224 x 224 pixels. The hyperparameters include a learning rate of 0.0001, a batch size of 16, and 50 training epochs. All experiments were conducted on a server with an NVIDIA V100 128G GPU from the NIH HPC Biowulf cluster (<http://hpc.nih.gov>). We implemented the framework using the Keras deep-learning library with TensorFlow backend (<https://www.tensorflow.org/guide/keras>).

2.6 Qualification and Statistical Analysis

The performance metrics include the area under the receiver operating characteristic curve (AUC), sensitivity, specificity (recall), precision (positive predictive value), and F1 score. For the classification problem, we chose the label with the highest probability when required in computing the metrics. Each of the models was fine-tuned and tested five times, using the same parameters, training, and testing images each time. The validation set was randomly selected from 10 percent of the training set. Fisher’s exact test was used to determine whether there are nonrandom associations between COVID-19 and influenza’s symptoms and clinical findings [35]. We conduct above statistical analysis using numpy, scipy, matplotlib, and scikit-learn built on Python.

3 RESULTS

3.1 COVID-19-CT-CXR Characteristics

Table 3 shows the breakdown of the figures by modality. We obtained 1,327 CT images and 263 CXR text-mined labeled as positive for COVID-19 from 1,831 PMC-OA articles. These images have different sizes. The minimum, maximum, and average heights are 224, 2,703, and 387.5 pixels, respectively. The minimum, maximum, and average widths are 224, 1,961, and 472.4, respectively. For each article, we also include major elements, such as DOI, title, journal, and publication date for reference. Fig. 3 A shows the cumulative numbers of articles and figures on a weekly basis. We analyzed the proportional

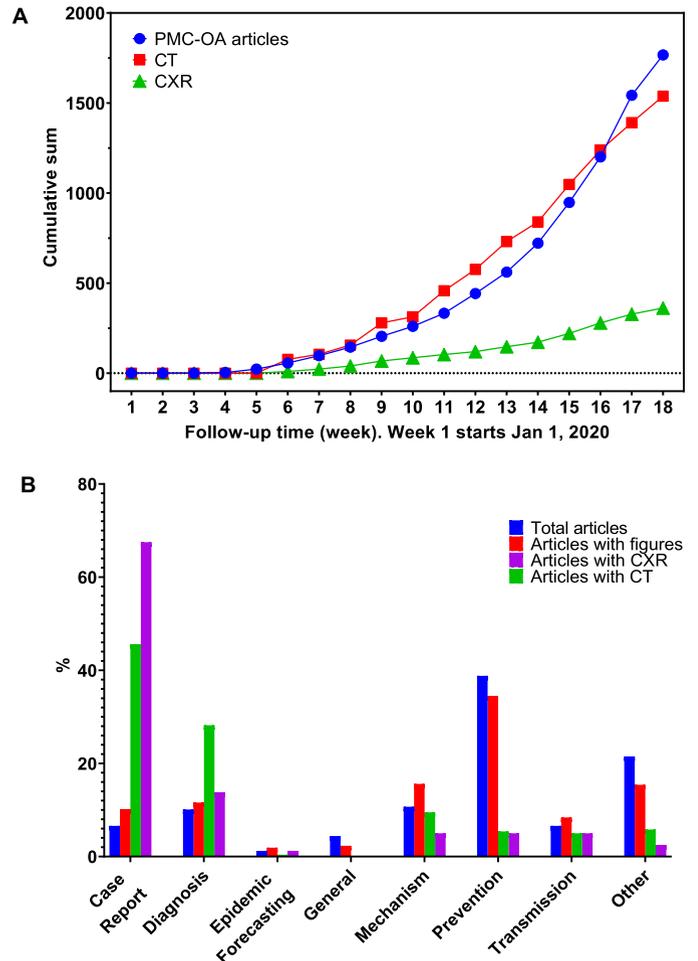


Fig. 3. Characteristics of the COVID-19-CT-CXR. (A) The rapid growth of the number of COVID-19-relevant articles, CT, and CXR in PMC-OA from January 1, 2020 (Week 1). (B) The distribution of categories in COVID-19-relevant PMC-OA articles and articles with figures, CT, and CXR.

distribution of categories in COVID-19 relevant PMC-OA articles, and articles with figures, CT, and CXR. Fig. 3 B shows that the “Case Report” category contains higher proportional articles with CXR/CT.

3.2 Image Modality Classification

Table 4 shows the performance of the model to classify image modality. The macro average F -score is 0.996. The F -score was 0.993 ± 0.004 for CT, 1.000 ± 0.000 for CXR, and 0.998 ± 0.001 for other scientific figure types.

3.3 Use Cases

To demonstrate the utility of COVID-19-CT-CXR, we conducted four case studies. (1) We combined COVID-19-CT-CXR with previously curated data at <https://github.com/UCSD-AI4H/COVID-CT> [36] and fine-tuned a deep neural network to perform the classification of COVID-19 and non-COVID-19 CT. (2) We collected CT of influenza, using a similar method, and fine-tuned a deep neural network to distinguish among the diagnoses of COVID-19, influenza, and normal or other types of diseases on CT. (3) We fine-tuned an unsupervised one-class learning model, using only non-COVID-19 CXR to perform anomaly detection, to detect

TABLE 4
The Performance of Image Type Classification

| Metrics | CT | CXR | Other scientific figures | Macro Avg |
|--------------------|---------------|---------------|--------------------------|---------------|
| Precision | 0.989 ± 0.004 | 1.000 ± 0.000 | 0.999 ± 0.001 | 0.996 ± 0.002 |
| Recall/Sensitivity | 0.998 ± 0.004 | 1.000 ± 0.000 | 0.996 ± 0.001 | 0.998 ± 0.002 |
| Specificity | 0.997 ± 0.001 | 1.000 ± 0.000 | 0.999 ± 0.002 | 0.999 ± 0.001 |
| F-score | 0.993 ± 0.004 | 1.000 ± 0.000 | 0.998 ± 0.001 | 0.997 ± 0.002 |

The test set is the combination of NIH Chest X-ray, DeepLesion, DocFigure, and PMC-OA.

COVID-19 CXR. (4) We extracted 15 clinical symptoms and 26 clinical findings from the captions and relevant descriptions. We then compared their frequencies to those described in articles on influenza, another common infectious respiratory illness that may present similarly to COVID-19.

3.3.1 Classification of COVID-19 and non-COVID-19 on CT

In the context of the COVID-19 pandemic, it is important to separate patients likely to be infected with COVID-19 from other non-COVID-19 patients. As it is time-consuming for specialists to both accumulate experiences and read a large volume of CT scans to diagnose COVID-19, many studies use machine learning to separate COVID-19 patients from non-COVID-19 patients [14], [37], [38], [39], [40]. In this work, we hypothesize that our creation of additional training data from existing articles can improve the performance of the system and reduce the effort of manual image annotation. To test this hypothesis, we compared the performance of deep neural networks fine-tuned on the existing benchmark [36] and COVID-19-CT-CXR (Table 5). For a fair comparison, we added additional training examples only in the training set and used the same test set as described in [14].

In this experiment, DenseNet121 was pre-trained on ImageNet, fine-tuned, and evaluated on the training and test sets. We then replaced the last classification layer with a single neuron with sigmoid that outputs the approximate probability that an input image is COVID-19 or non-COVID-19. Other experimental settings are the same as that of fine-tuning the image modality classifier. Fig. 4 shows that the model significantly outperforms the baseline when PMC-OA CT figures were added for fine-tuning. Specifically, we achieved the highest performance of 0.891 ± 0.012 in AUC, 0.780 ± 0.074 in recall, 0.816 ± 0.053 in precision, and 0.792 ± 0.015 in F-score (Table 6).

3.3.2 Classification of COVID-19, Influenza, and Other Types of Disease on CT

As the COVID-19 outbreak continues to evolve, there is an increasing number of studies that compare COVID-19 with

other viral pneumonias, such as influenza [41]. Distinguishing patients infected by COVID-19 and influenza is important for public health measures because the current treatment guidelines are different [42]. This task is non-trivial because both viruses have a similar radiological presentation. To assist clinicians at triage, several studies have proposed to use deep learning to distinguish COVID-19 from influenza and no-infection with 3D CT scans [43]. In this paper, we aim to establish a baseline model to distinguish COVID-19 from influenza on single CT figures. To collect CT figures with influenza, we searched the PMC using the query “(Influenza[Title] OR (flu[Title] AND pneumonia[Title])) AND open access[Filter]” and extracted the most recent 10,000 PMC-OA articles. We used the same method to extract CT and its caption and relevant text from the articles (called Influenza-CT). Taken together, we construct a dataset with 983 CT for training and 242 CT for testing (Table 7).

To obtain the baseline model, we use the same model and experimental settings as described in the “Image modality

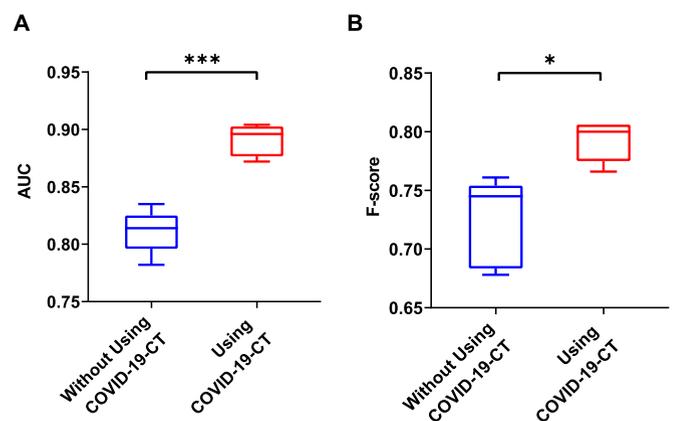


Fig. 4. Comparison of AUC and F-score by models fine-tuned with and without using additional COVID-19 CT extracted from PMC-OA. *: $P \leq 0.05$; ***: $P \leq 0.001$ (t-test).

TABLE 6
Performance Metrics for Classification of COVID-19 and non-COVID-19 CT

| Metrics | Without using COVID-19-CT | Using COVID-19-CT |
|--------------------|---------------------------|-------------------|
| AUC | 0.811 ± 0.017 | 0.891 ± 0.012 |
| Precision | 0.742 ± 0.029 | 0.816 ± 0.053 |
| Recall/Sensitivity | 0.714 ± 0.083 | 0.780 ± 0.074 |
| Specificity | 0.764 ± 0.059 | 0.827 ± 0.073 |
| F-score | 0.724 ± 0.034 | 0.792 ± 0.015 |

TABLE 5
Summary of the Dataset for Classification of COVID-19 and non-COVID-19 CT

| Dataset | COVID-19 | Non-COVID-19 |
|-------------|----------|--------------|
| Training | 251 | 292 |
| COVID-19-CT | 542 | 67 |
| Test | 98 | 105 |

TABLE 7
Summary of the Dataset for Classification of COVID-19, Influenza, and Others in CT

| Dataset | COVID-19 | Influenza | Normal or other diseases |
|----------|----------|-----------|--------------------------|
| Training | 488 | 177 | 318 |
| Test | 118 | 45 | 79 |

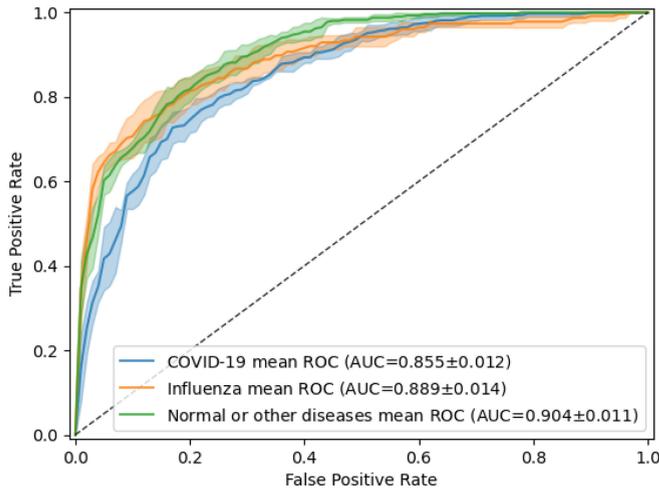


Fig. 5. Receiver operating characteristic (ROC) curves of the classification of COVID-19, influenza, and normal or other types of diseases in CT. The model was fine-tuned and tested 5 times, using the same training and testing images each time. The mean ROC curve is shown together with its standard deviation (shaded area).

classification” section. Fig. 5 shows the performance of the deep-learning model by its receiver operating characteristic (ROC) curves. The AUC was 0.855 ± 0.012 for COVID-19 detection and 0.889 ± 0.014 for influenza detection. Table 8 shows more detail for the results. We achieved the highest precision (0.845 ± 0.026) for COVID-19 detection and high recall (0.711 ± 0.053) for influenza detection.

3.3.3 Anomaly Detection of COVID-19 in CXR Using One-Class Learning

As they lack annotated COVID-19 CXR for training powerful deep-learning classifiers, unsupervised and semi-supervised approaches are highly desired for automated COVID-19 diagnosis. The presence of COVID-19 can be considered a novel anomaly in CXR for the NIH Chest X-ray dataset, in which no COVID-19 cases are available. In this experiment, we performed anomaly detection [44], [45] to detect COVID-19 CXR. We trained a one-class classifier, using only non-COVID-19 CXR, and used this classifier to

TABLE 9
Summary of Dataset Used for Anomaly Detection of COVID-19 in CXR in Unsupervised One-Class Classification

| Dataset | COVID-19 | Non-COVID-19 |
|----------|----------|--------------|
| Training | 0 | 37,829 |
| Test | 184 | 184 |

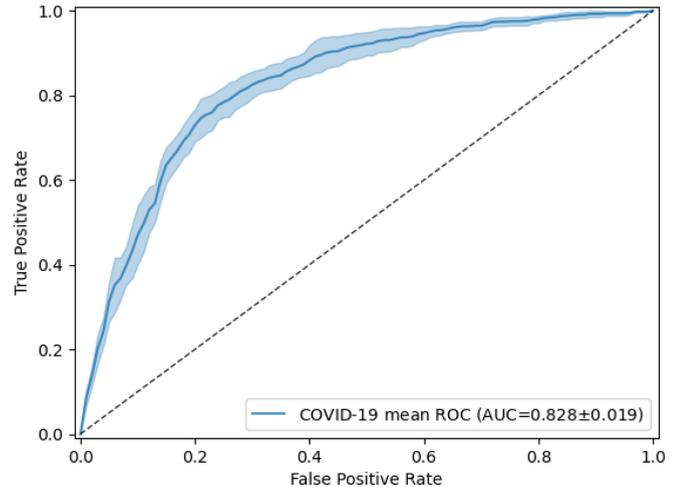


Fig. 6. Receiver operating characteristic (ROC) curves of the classification of COVID-19 anomaly detection in CXR. The model was fine-tuned and tested 5 times, using the same training and testing images each time. The mean ROC curve is shown together with its standard deviation (shaded area).

distinguish COVID-19 CXR from non-COVID-19 CXR. The non-COVID-19 images were a subset extracted from the NIH Chest X-ray dataset by combining 14 abnormalities and a no-finding category. The detailed numbers of training and testing CXR are shown in Table 9. We adopted the generative adversarial one-class learning approach from [46]. Fig. 6 shows the performance of the unsupervised one-class learning by its ROC curves. Table 10 shows more detail for the results. Our model achieved 0.828 ± 0.019 in AUC, 0.767 ± 0.020 in precision, 0.772 ± 0.017 in recall, and 0.769 ± 0.018 in F -score for COVID-19 anomaly detection.

3.3.4 Extraction of Clinical Symptoms and Findings Using Text-Mining

In this case, we extracted clinical symptoms or signs from the figure captions and relevant text that describes the case. A total of 15 symptoms or signs were collected from [3] and the CDC website (<https://www.cdc.gov/coronavirus/2019-ncov/symptoms-testing/symptoms.html>), including chest pain, constipation, cough, diarrhea, dizziness, dyspnea,

TABLE 8
Performance Metrics for Classification of COVID-19, Influenza, and Normal or Other Types of Diseases in CT

| Metrics | COVID-19 | Influenza | Normal or other diseases | Macro Avg |
|--------------------|-------------------|-------------------|--------------------------|-------------------|
| AUC | 0.855 ± 0.012 | 0.889 ± 0.014 | 0.904 ± 0.011 | 0.879 ± 0.010 |
| Precision | 0.845 ± 0.026 | 0.609 ± 0.033 | 0.642 ± 0.021 | 0.699 ± 0.019 |
| Recall/Sensitivity | 0.597 ± 0.030 | 0.711 ± 0.053 | 0.861 ± 0.033 | 0.723 ± 0.022 |
| Specificity | 0.895 ± 0.024 | 0.895 ± 0.013 | 0.767 ± 0.025 | 0.852 ± 0.009 |
| F -score | 0.699 ± 0.018 | 0.655 ± 0.034 | 0.735 ± 0.015 | 0.696 ± 0.018 |

TABLE 10
Anomaly Detection Performance of COVID-19
Versus non-COVID-19 Using Unsupervised
One-Class Learning

| Metrics | COVID-19 vs Non-COVID-19 |
|--------------------|--------------------------|
| AUC | 0.828 ± 0.019 |
| Precision | 0.767 ± 0.020 |
| Recall/Sensitivity | 0.772 ± 0.017 |
| Specificity | 0.765 ± 0.023 |
| F-score | 0.769 ± 0.018 |

fatigue, fever, headache, myalgia, proteinuria, runny nose, sputum production, throat pain, and vomiting.

Extracting these symptoms from text is a challenging task because their mentions in the text can be positive or negative. For example, “fever” is negative in the sentence, “She experienced headache and pharyngalgia but no fever on 29 January.” To discriminate between positive and negative mentions, we applied our previously developed tool, Neg-Bio, on the figure caption and referred text [47]. In short, NegBio utilizes patterns in universal dependencies to identify the scope of triggers that are indicative of negation; thus, it is highly accurate for detecting negative symptom mentions. Fig. 7 A shows the proportion of symptoms for COVID-19 and influenza. The most common symptoms are fever, cough, dyspnea, and myalgia.

We then extracted the radiographic findings from the figure caption and text. The findings (and their synonyms) are based on 20 common thoracic disease types, which are expanded from NIH Chest X-ray 14 labels [11]. Fig. 7 B shows the 20 findings in both COVID-19 and influenza datasets. Both illnesses can result in lung opacity, pneumonia, and consolidation. COVID-19 more likely results in ground-glass opacification (GGO), while influenza more likely results in infiltration than does COVID-19 (Fisher’s exact test, $p < 0.0001$).

4 DISCUSSION

In this abrupt outbreak of SARS-CoV-2, the demand for chest radiographs and CT scans is growing rapidly, but there is a shortage of experienced specialists, radiologists, and researchers. Further, we are still new to this virus and have yet to discover the full radiologic features and prognosis of this disease. The tremendous increase in the number of patients has led to a substantial increase of COVID-19-related PMC-OA articles over the past few months (Figur 3 A), especially in the case report and diagnosis-relevant articles (Fig. 3 B). These articles contain rich chest radiographs and CT images that are helpful for scientists and clinicians in describing COVID-19 cases. Thus, it is important to analyze these images and text to construct a large-scale database. By using the quickly increasing dataset, AI methods can help to find significant features of COVID-19 and speed up the clinical workload. Among others, deep learning is undoubtedly a powerful approach in dealing with a pandemic outbreak of COVID-19.

Although deep learning has shown promise in diagnosing/screening COVID-19, using CT, it remains difficult to collect large-scale labeled imaging data, especially in the public domain. In this work, we present a set of repeatable

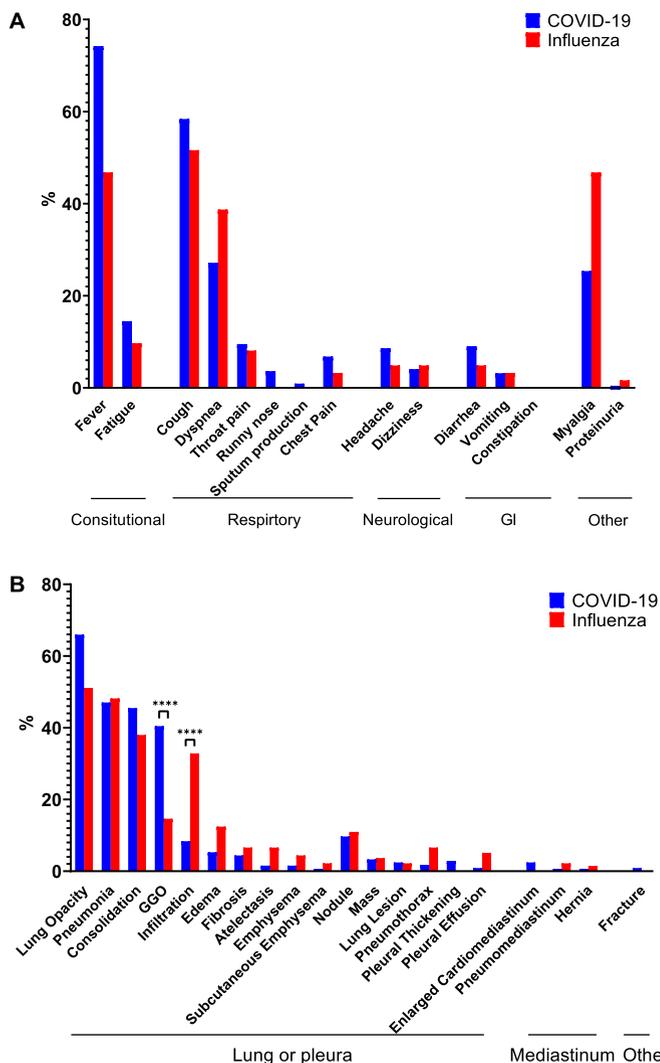


Fig. 7. The frequencies of (A) 15 symptoms and (B) 20 clinical findings text mined from the figure captions and relevant text from the collection of COVID-19- and influenza-relevant articles. ****: $p \leq 0.0001$ (Fisher exact test).

techniques to rapidly build a CT and CXR dataset of COVID-19 from PMC-OA COVID-19-relevant articles. The strength of the study lies in its multidisciplinary integration of medical imaging and natural-language processing. It provides a new way to annotate large-scale medical images required by deep-learning models.

An additional strength includes a highly accurate model for image type classification. As a large portion of figures in the PMC-OA articles are not CXR or CT images, we provided a model to classify these two types from other scientific figure types. Our model achieved both high precision and high recall (Table 4).

To assess the hypothesis that deep neural network fine-tuning on this additional dataset enables us to diagnose COVID-19 with almost no hand-labeled data, we conducted several experiments. First, we showed that this additional data enable significant performance gains to classify COVID-19 versus non-COVID-19 lung infection on CT (Fig. 4 and Supplementary Table 6, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TBDATA.2020.3035935>). For our own system,

we show that our baseline performance compares favorably to the results in [14]. Then, we added more automatically labeled training data and achieved the highest performance of 0.891 ± 0.012 in AUC. The comparison shows that, with additional data, both precision and recall substantially improve (7.4 and 6.6 percent, respectively). This observation indicates that additional COVID-19 CT helps to not only find more but also to restrict the positive predictions to those with the highest certainty in the model.

In a more challenging scenario, we built a baseline system to distinguish COVID-19, influenza, and no-infection CT, which is a more clinically interesting but also more challenging task. We observed that we could achieve high AUCs for both COVID-19 and influenza detection. The recall of COVID-19 detection and the precision of influenza, however, are low (0.597 ± 0.030 and 0.609 ± 0.033 , respectively). Although several studies have tackled this problem [43], to the best of our knowledge, there is no publicly available benchmarking. The differentiation between COVID-19 and influenza on CXR/CT without associated context is challenging. In the experiment on classification of COVID-19, influenza, and other types of disease on CT, we found that although many of the CT findings had overlapping findings, “mixed GGO (Ground glass opacity)” were mostly found in the COVID-19 dataset and “pleural thickening” and “linear opacities” were mostly found in the influenza dataset. It is also worthy to note that the images from PMC-OA may not represent the typical pool of influenza pneumonia real-world images, since researchers may report extreme cases instead of typical cases. While our work only scratches the surface of the classification of COVID-19, influenza, and normal or other types of diseases, we hope that it sheds light on the development of generalizable deep-learning models that can assist frontline radiologists.

In addition, we presented a one-class learning model for anomaly detection of COVID-19 in CXR by learning only from non-COVID-19 radiographs. Compared to the CT-based method, the one-class model achieves comparable performance, showing great potential in discriminating COVID-19 from CXR. The performance of our model, however, is worse than that of [45], suggesting that this weakly labeled dataset should be used as additional training data obtained without additional annotation cost from existing entries in curated databases.

The unique characteristic of our database is that figures are retrieved along with relevant text that describes these cases in detail. Thus, text mining can be applied to extract additional information that confirms the existing results and potentially identifies other findings that may have been overlooked. As proof of this concept, we extracted clinical symptoms and findings from the text. We found that the most common symptoms of COVID-19 were fever and cough (Fig. 7 A), which are consistent with the clinical characteristics in [15]. Other common symptoms include dyspnea (shortness of breath), fatigue, and throat pain. These symptoms are consistent with those reported by the CDC. When comparing the frequencies of these 20 clinical findings to those described in articles on influenza, Fig. 7 shows that both conditions cause lung opacity, pneumonia, and consolidation. Further, GGO appears more frequently for COVID-19, whereas “infiltration” appears more frequently for influenza. This is because radiologists use the term GGO

to describe most COVID-19 findings. In addition, the influenza articles are older than are the COVID articles, and, according to Fleischner Society recommendations, the use of the term *infiltrate* remains controversial, and it is recommended that it no longer be used in reports [48].

In terms of limitations, first, the subfigure segmentation model needs to be improved. In this study, we applied a deep-learning model that was pretrained on an ImageCLEF Medical dataset to this task [24]. Although this model is robust to variations in background color and spaces between subfigures, it sometimes fails to recognize similar subfigures that are aligned very closely. Unfortunately, these cases appear more frequently in our study than in others (e.g., several CT images are placed in a grid). Other errors occur when the model incorrectly treated the spine as spaces in the anteroposterior (AP) chest X-ray and split the large figure into two subfigures. In the future, the figure synthesis approach should be applied to augment the training datasets. Another limitation is that this work extracted only the passage that contains the referred figure. Sometimes, the case is not described in this passage. In the future, we plan to text mine the associated case description in the full text. Finally, while a figure is typically copyrighted with the original article and using previously published figures is not a common practice in scholarly publications, it is possible that one image is reused in different papers or reused in one paper for different purposes. In the future, we plan to develop a model to remove duplicated images in the collection.

5 CONCLUSION

We have developed a framework for rapidly constructing a CXR/CT database from PMC full-text articles. Our database is unique, as figures are retrieved along with relevant text that describes these cases in detail, and it can be extended easily in the future. Hence, the work is complementary to existing resources. Applications of this database show that our creation of additional training data from existing articles improves the system performance on COVID-19 versus non-COVID-19 classification in CT and CXR. We hope that the public dataset can facilitate deep-learning model development, educate medical students and residents, help to evaluate findings reported by radiologists, and provide additional insights for COVID-19 diagnosis. With an ongoing commitment to data sharing, we anticipate increasingly adding CXR and CT images to be made available as well in the coming months. The code that extracts the text from PMC, segments subfigures, and classifies image modality is openly available at <https://github.com/ncbi-nlp/COVID-19-CT-CXR>.

ACKNOWLEDGMENTS

This work was supported in part by the Intramural Research Programs of the National Library of Medicine (NLM) and National Institutes of Health (NIH) Clinical Center. This work also was supported by NLM under Grant 4R00LM013001. This work utilized the computational resources of the NIH HPC Biowulf cluster (<http://hpc.nih.gov>). This material is also based upon the work supported by Google Cloud.

REFERENCES

- [1] A. S. Fauci, H. C. Lane, and R. R. Redfield, "COVID-19 - Navigating the uncharted," *New England J. Med.*, vol. 382, pp. 1268–1269, Mar. 2020.
- [2] N. Chen *et al.*, "Epidemiological and clinical characteristics of 99 cases of 2019 novel coronavirus pneumonia in Wuhan, China: A descriptive study," *Lancet*, vol. 395, pp. 507–513, Feb. 2020.
- [3] W.-J. Guan *et al.*, "Clinical characteristics of coronavirus disease 2019 in China," *New England J. Med.*, vol. 382, pp. 1708–1720, Apr. 2020.
- [4] D. Wang *et al.*, "Clinical characteristics of 138 hospitalized patients with 2019 novel coronavirus-infected pneumonia in Wuhan, China," *J. Amer. Med. Assoc.*, vol. 323, pp. 1061–1069, Feb. 2020.
- [5] H. Shi *et al.*, "Radiological findings from 81 patients with COVID-19 pneumonia in Wuhan, China: A descriptive study," *Lancet Infect. Dis.*, vol. 20, pp. 425–434, Apr. 2020.
- [6] X. Xie, Z. Zhong, W. Zhao, C. Zheng, F. Wang, and J. Liu, "Chest CT for typical 2019-nCoV pneumonia: Relationship to negative RT-PCR testing," *Radiology*, vol. 12, Feb. 2020, Art. no. 200343.
- [7] X. Mei *et al.*, "Artificial intelligence-enabled rapid diagnosis of patients with COVID-19," *Nat. Med.*, vol. 26, no. 8, pp. 1224–1228, Aug. 2020.
- [8] F. Shi *et al.*, "Review of artificial intelligence techniques in imaging data acquisition, segmentation and diagnosis for COVID-19," *IEEE Rev. Biomed. Eng.*, to be published, doi: [10.1109/RBME.2020.2987975](https://doi.org/10.1109/RBME.2020.2987975).
- [9] J. Irvin *et al.*, "CheXpert: A large chest radiograph dataset with uncertainty labels and expert comparison," in *Proc. AAAI Conf. Artif. Intell.*, 2019, pp. 590–597.
- [10] A. E. W. Johnson *et al.*, "MMIC-CXR-JPG, A large publicly available database of labeled chest radiographs," 2019, *arXiv:1901.07042*.
- [11] X. Wang, Y. Peng, L. Lu, Z. Lu, M. Bagheri, and R. M. Summers, "ChestX-Ray8: Hospital-scale chest X-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 3462–3471.
- [12] K. Yan, X. Wang, L. Lu, and R. M. Summers, "DeepLesion: Automated mining of large-scale lesion annotations and universal lesion detection with deep learning," *J. Med. Imag.*, vol. 5, Jul. 2018, Art. no. 036501.
- [13] J. P. Cohen, P. Morrison, and L. Dao, "COVID-19 image data collection," 2020, *arXiv:2006.11988*.
- [14] X. He *et al.*, "Sample-efficient deep learning for COVID-19 diagnosis based on CT scans," to be published, doi: [10.1101/2020.04.13.20063941](https://doi.org/10.1101/2020.04.13.20063941).
- [15] K. Zhang *et al.*, "Clinically applicable AI system for accurate diagnosis, quantitative measurements, and prognosis of COVID-19 pneumonia using computed tomography," *Cell*, vol. 181, pp. 1423–1433, 2020.
- [16] Q. Chen, A. Allot, and Z. Lu, "Keep up with the latest coronavirus research," *Nature*, vol. 579, Mar. 2020, Art. no. 193.
- [17] L. L. Wang *et al.*, "CORD-19: The COVID-19 open research dataset," 2020, *arXiv:2004.10706*. PMID: 32510522.
- [18] S. R. Choudhury *et al.*, "A figure search engine architecture for a chemistry digital library," in *Proc. 13th ACM/IEEE-CS Joint Conf. Digit. Libraries*, 2013, pp. 369–370.
- [19] C. L. Smith, J. A. Blake, J. A. Kadin, J. E. Richardson, C. J. Bult, and M. G. D. Group, "Mouse Genome database (MGD)-2018: Knowledgebase for the laboratory mouse," *Nucleic Acids Res.*, vol. 46, pp. D836–D842, Jan. 2018.
- [20] Z. Ahmed, S. Zeeshan, and T. Dandekar, "Mining biomedical images towards valuable information retrieval in biomedical and life sciences," *Database*, vol. 2016, 2016, Art. no. baw118.
- [21] P. Li, X. Jiang, and H. Shatkay, "Figure and caption extraction from biomedical documents," *Bioinformatics*, vol. 35, pp. 4381–4388, Nov. 2019.
- [22] N. Siegel, N. Lourie, R. Power, and W. Ammar, "Extracting scientific figures with distantly supervised neural networks," in *Proc. 18th ACM/IEEE Joint Conf. Digit. Libraries*, 2018, pp. 223–232.
- [23] L. D. Lopez *et al.*, "A framework for biomedical figure segmentation towards image-based document retrieval," *BMC Syst. Biol.*, vol. 7, 2013, Art. no. S8.
- [24] S. Tsutsui and D. Crandall, "A data driven approach for compound figure separation using convolutional neural networks," in *Proc. IAPR Int. Conf. Document Anal. Recognit.*, 2017.
- [25] D. C. Comeau, C.-H. Wei, R. I. Doğan, and Z. Lu, "PMC text mining subset in BioC: About three million full-text articles and growing," *Bioinformatics*, vol. 35, pp. 3533–3535, Sep. 2019.
- [26] Ş. Kafkas, X. Pi, N. Marinos, F. Talo', A. Morrison, and J. R. McEntyre, "Section level search functionality in europe PMC," *J. Biomed. Semantics*, vol. 6, 2015, Art. no. 7.
- [27] Y.-H. E. A. Jin, "A rapid advice guideline for the diagnosis and treatment of 2019 novel coronavirus (2019-nCoV) infected pneumonia (standard version)," *Mil. Med. Res.*, vol. 7, Feb. 2020, Art. no. 4.
- [28] A. G. S. De Herrera, S. Bromuri, R. Schaer, and H. Müller, "Overview of the medical tasks in ImageCLEF 2016," *CLEF Work. Notes. Evora, Portugal*, 2016.
- [29] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 2818–2826.
- [30] F. Iandola, M. Moskewicz, S. Karayev, R. Girshick, T. Darrell, and K. Keutzer, "DenseNet: Implementing efficient convnet descriptor pyramids," 2014, *arXiv:1404.1869*.
- [31] K. V. Jobin, A. Mondal, and C. V. Jawahar, "DocFigure: A dataset for scientific document figure classification," in *Proc. Int. Conf. Document Anal. Recognit. Workshops*, 2019, pp. 74–79.
- [32] Y.-X. Tang *et al.*, "Automated abnormality classification of chest radiographs using deep convolutional neural networks," *NPJ Digit. Med.*, vol. 3, 2020, Art. no. 70.
- [33] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 4700–4708.
- [34] O. Russakovsky *et al.*, "ImageNet large scale visual recognition challenge," *Int. J. Comput. Vis.*, vol. 115, no. 3, pp. 211–252, 2015.
- [35] R. A. Fisher, "On the interpretation of χ^2 from contingency tables, and the calculation of P," *J. Roy. Statist. Soc.*, vol. 85, no. 1, Jan. 1922, Art. no. 87.
- [36] J. Zhao, Y. Zhang, X. He, and P. Xie, "COVID-CT-Dataset: A CT scan dataset about COVID-19," 2020, *arXiv:2003.13865*.
- [37] J. Chen *et al.*, "Deep learning-based model for detecting 2019 novel coronavirus pneumonia on high-resolution computed tomography," *Scientific Rep.*, vol. 10, no. 1, pp. 1–11, 2020.
- [38] C. Jin *et al.*, "Development and evaluation of an AI system for COVID-19 diagnosis," *MedRxiv Reprint*, to be published, doi: [10.1101/2020.03.20.20039834](https://doi.org/10.1101/2020.03.20.20039834).
- [39] S. Wang *et al.*, "A deep learning algorithm using CT images to screen for corona virus disease (COVID-19)," *MedRxiv Preprint*, to be published, doi: [10.1101/2020.02.14.20023028](https://doi.org/10.1101/2020.02.14.20023028).
- [40] C. Zheng *et al.*, "Deep learning-based detection for COVID-19 from chest CT using weak label," 2020.
- [41] Y. Luo *et al.*, "Using a diagnostic model based on routine laboratory tests to distinguish patients infected with SARS-CoV-2 from those infected with influenza virus," *Int. J. Infect. Dis.*, vol. 95, pp. 436–440, May 2020.
- [42] D. Kimberlin, *Red Book 2018–2021: Report of the Committee on Infectious Diseases*. Elk Grove Village, IL, USA: Amer. Acad. Pediatrics, 2018.
- [43] X. Xu *et al.*, "Deep learning system to screen coronavirus disease 2019 pneumonia," *Engineering*, to be published, doi: [10.1016/j.eng.2020.04.010](https://doi.org/10.1016/j.eng.2020.04.010).
- [44] V. Chandola, A. Banerjee, and V. Kumar, "Anomaly detection," *ACM Comput. Surv.*, vol. 41, no. 3, pp. 1–58, Jul. 2009.
- [45] J. Zhang, Y. Xie, Y. Li, C. Shen, and Y. Xia, "COVID-19 screening on chest X-ray images using deep learning based anomaly detection," 2020, *arXiv:2003.12338*.
- [46] Y.-X. Tang, Y.-B. Tang, M. Han, J. Xiao, and R. M. Summers, "Abnormal chest X-Ray identification with generative adversarial one-class classifier," in *Proc. IEEE 16th Int. Symp. Biomed. Imag.*, 2019, pp. 1358–1361.
- [47] Y. Peng, X. Wang, L. Lu, M. Bagheri, R. Summers, and Z. Lu, "NegBio: A high-performance tool for negation and uncertainty detection in radiology reports," *AMIA Joint Summits Translational Sci. Proc.*, vol. 2017, pp. 188–196, 2018. [Online]. Available: <https://arxiv.org/abs/1712.05898>
- [48] J. Bueno, L. Landeras, and J. H. Chung, "Updated fleischner society guidelines for managing incidental pulmonary nodules: Common questions and challenging scenarios," *Radiographics*, vol. 38, pp. 1337–1350, 2018.



Yifan Peng received the PhD degree. He is currently an assistant professor with Weill Cornell Medicine. He was a research fellow with the National Center for Biotechnology Information (NCBI), National Library of Medicine (NLM), National Institutes of Health (NIH). His main research interests include biomedical and clinical natural language processing and medical image analysis. He has published many papers in top journals and conferences, including the *Nucleic Acids Research*, *npj Digital Medicine*, *Journal of the American Medical Informatics Association*, *CVPR*, and *MICCAI*. He is also an academic editor of the *PLoS ONE*.



Yuxing Tang received the BS and MS degrees from the Department of Information and Telecommunication Engineering, Beijing Jiaotong University, Beijing, China, in 2009 and 2011, respectively, and the PhD degree in computer science from the Department of Mathematics and Computer Science, École Centrale de Lyon, Écully, France, in 2016. He is a postdoctoral fellow with the Imaging Biomarkers and Computer-Aided Diagnosis (CAD) Laboratory, National Institutes of Health (NIH) Clinical Center. His main research interests include

computer vision and machine learning, in particular, deep learning techniques for visual category recognition, object detection, image segmentation and their application in medical imaging.



Sungwon Lee received the MD and PhD degrees. She is currently a radiologist and research fellow with the National Institutes of Health (NIH). Her research interests include segmentation and classification of medical imaging, especially chest, body, and musculoskeletal images of CT and MRI.



Yingying Zhu received the PhD degree. She is currently a staff scientist with the Department of Radiology, Clinical Center, National Institutes of Health (NIH). Her main research interests include computer vision, medical image analysis, and machine learning. She has published many papers in top journals and conferences, including the *IEEE Transaction on Medical Imaging*, the *Medical Image Analysis*, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *ECCV*, *CVPR*, *IPMI*, and *MICCAI*.



Ronald M. Summers received the MD and PhD degrees. He is currently a senior investigator with the NIH. He joined the Diagnostic Radiology Department, NIH Clinical Center, in 1994. He directs the Imaging Biomarkers and Computer-Aided Diagnosis (CAD) Laboratory. His research interests include virtual colonoscopy, CAD, multi-organ multi-atlas registration, and development of large radiologic image databases. His clinical areas of specialty are thoracic and gastrointestinal radiology and body cross-sectional

imaging. His current research focuses on developing fully-automated interpretation of abdominal CT scans.



Zhiyong Lu received the PhD degree. He is currently a deputy director for Literature Search at the National Center for Biotechnology (NCBI), leading its overall efforts of improving literature search and information access in NCBI's production resources. He is also an NIH senior investigator (early tenure) and directs the Text Mining / Natural Language Processing (NLP) Research Program, NCBI/NLM where they are developing computational methods and software tools for analyzing and making sense of unstructured text

data in biomedical literature and clinical notes towards accelerated discovery and better health.

▷ For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/csdl.