













## ORIGINAL ARTICLE

# Disorder-Free Data Are All You Need — Inverse Supervised Learning for Broad-Spectrum Head Disorder Detection

Yuwei He , Ph.D.,<sup>1,2</sup> Yuchen Guo , Ph.D.,<sup>1</sup> Jinhao Lyu , M.S.,<sup>3</sup> Liangdi Ma , M.S.,<sup>1,2</sup> Haotian Tan ,<sup>1,2</sup> Wei Zhang ,<sup>4</sup> Guiguang Ding , Ph.D.,<sup>1,2</sup> Hengrui Liang , Ph.D.,<sup>5</sup> Jianxing He , Ph.D.,<sup>5</sup> Xin Lou , Ph.D.,<sup>3</sup> Qionghai Dai , Ph.D.,<sup>1</sup> and Feng Xu , Ph.D.<sup>1,2</sup>

Received: September 18, 2023; Revised: December 18, 2023; Accepted: February 3, 2024; Published: March 28, 2024

## Abstract

**BACKGROUND** The development of artificial intelligence (AI)-based medical systems heavily relies on the collection and annotation of sufficient data containing disorders. However, the preparation of data with complete disorder types and adequate annotations presents a significant challenge, limiting the diagnostic capabilities of existing AI-based medical systems. This study introduces a novel AI-based system that accurately detects a broad spectrum of disorders without requiring any disorder-containing data.

**METHODS** We obtained a training dataset of 21,429 disorder-free head computed tomography (CT) scans and proposed a learning algorithm called inverse supervised learning (ISL). This algorithm learns and understands disorder-free samples instead of disorder-contained ones, enabling the identification of all types of disorders. We also developed a diagnosis and visualization software for clinical usage on the basis of the system's ability to provide visually understandable clues.

**RESULTS** The system achieved area under the receiver operating characteristic curve (AUC) values of 0.883, 0.868, and 0.866 on retrospective (127 disorder types, 9967 scans), prospective (117 disorder types, 3054 scans), and cross-center (46 disorder types, 554 scans) datasets, respectively. These results demonstrate that the system can detect far more disorder types than previous AI-based systems. Furthermore, the ISL-based systems achieved AUC values of 0.893 and 0.895 on pulmonary CT and retinal optical coherence tomography, respectively, demonstrating that ISL can generalize well to nonhead and non-CT images.

**CONCLUSIONS** Our novel AI-based system utilizing ISL can accurately and broadly detect disorders without requiring disorder-containing data. This system not only outperforms previous AI-based systems in terms of disorder detection but also provides visually understandable clues, enhancing its clinical utility. The successful application of ISL to

Dr. Yuwei He, Dr. Yuchen Guo, Mr. Jinhao Lyu, and Ms. Liangdi Ma contributed equally to this article.

The author affiliations are listed at the end of the article.

Prof. Lou can be contacted at [louxin@301hospital.com.cn](mailto:louxin@301hospital.com.cn); Prof. Dai can be contacted at [daiqionghai@tsinghua.edu.cn](mailto:daiqionghai@tsinghua.edu.cn); and Prof. Xu can be contacted at [xufeng2003@gmail.com](mailto:xufeng2003@gmail.com).

This article was updated on May 31, 2024 at [ai.nejm.org](https://ai.nejm.org).

nonhead and non-CT images further demonstrates its potential for broad-spectrum medical applications. (Funded by the National Key R&D Program of China and the National Natural Science Foundation of China.)

## Introduction

Over the past decade, artificial intelligence (AI) has made significant strides and has been applied in various fields. In the medical field, the accumulation of medical image data has enabled many AI diagnostic techniques to achieve radiologist-level performance in recognizing, classifying, and quantifying specific diseases. For example, AI has been used for cerebral hemorrhage recognition<sup>1</sup> and coronavirus disease 2019 (Covid-19) recognition<sup>2</sup> from computed tomography (CT) images. These breakthroughs have led us to envision that AI diagnostic techniques can assist in clinical decision-making from medical images and alleviate the severe shortage of expert radiologists in many areas and hospitals.

Despite the significant progress made in AI techniques, there is still a gap between these techniques and real clinical decision-making. Current AI techniques primarily focus on recognizing specific types of disorders from input medical data. However, for a clinical decision-making workflow, the most basic and essential task is to identify all possible disorder types that could be diagnosed from the medical image. For instance, in the case of brain CT, more than 100 types of disorders could be diagnosed from it. Therefore, a decision-making diagnostic system for brain CT must be capable of detecting a broad spectrum of disorders, because missing the detection of any disorder type is unacceptable. Existing medical AI techniques are developed on the basis of widely used AI paradigms, which involve deciding the disorder types to be handled, collecting sufficient disorder-contained samples, and constructing recognition/localization/segmentation models for the disorders. This paradigm works well when the disorder types are limited and the samples are easily accessible. However, developing a broad-spectrum disorder detection system using this paradigm requires collecting data and constructing models for all types of disorders, which are extremely difficult and inefficient, especially for unusual diseases. Therefore, it is impractical to use the widely used AI paradigms to achieve real clinical decision-making.

To address the challenges mentioned, we propose a novel AI solution called inverse supervised learning (ISL). Instead of using disorder-contained data, which require hundreds of disorder types and a large number of samples for each type, we use disorder-free medical images for supervision. In theory, we use the opposite problem (detecting no-disorder samples) to replace the original problem (detecting all types of disorders). Therefore, instead of training hundreds of models to recognize all possible types of disorders, we train just one model to understand the concept of disorder free fully. Consequently, all disorders can be identified as they differ from the disorder-free samples used in training. With our paradigm, the challenges mentioned are fully resolved because there is no need for samples of all possible disorder types.

To achieve ISL, we utilize the traditional computer vision task of image inpainting in a novel framework. Image inpainting aims to restore the content of a partially missing image on the basis of the context of nonmissing information. Specifically, in this case, an image inpainting network is trained to replenish masked regions in a medical image, where the replenished content always reflects healthy tissue because the training dataset contains only disorder-free medical images. If any disorder exists in the image and the disorder region is masked off, the reconstructed disorder-free image would be significantly different from the original one. Conversely, for an image without disorders, no matter which region is masked, the reconstructed image should always be similar to the original one because they are both healthy and consistent with the rest of the healthy images. By masking, inpainting, and comparing all the image regions, ISL can detect various types of disorders and locate the disorder regions. Our proposed solution, ISL, does not require the deliberate collection of data for any disorder type; ensures that the data used to develop systems are easily accessible; does not require experts to manually annotate the data; enables the developed system to recognize broad-spectrum disorders rather than specific ones; and provides experts with clinical clues, such as disorder locations.

In this study, we utilized ISL to construct a system for broad-spectrum disorder detection on unenhanced brain CT scans.<sup>3,4</sup> CT is a first-line diagnostic modality for assessing brain abnormalities because of its quick acquisition and noninvasive nature. The ISL-based system was developed using only disorder-free head CT images. It achieved expert-level accuracies on a retrospective dataset with 127 disorder types and a prospective dataset with

116 disorder types, surpassing the number of detectable disorder types in previous works. We also applied ISL to build two additional systems: one for pulmonary disorder detection in CT images and another for retinal disorder detection in optical coherence tomography (OCT) images. The results demonstrate that ISL can generalize well to nonbrain- and non-CT-based disorder detection.

## Results

### BUILDING AN ISL-BASED SYSTEM FOR CLINICALLY APPLICABLE BROAD-SPECTRUM HEAD DISORDER DETECTION

Our proposed ISL-based disorder detection system for brain CT comprises a dedisorder network (DeDN), a disorder recognition network (DRN), and a disorder locating module. First, a CT image is processed with specific window width (WW) and window locations (WLs) and then fed into the DeDN to generate its corresponding dedisorder image. Next, we obtain a difference image by subtracting the original and generated images. Finally, the difference image is inputted into the DRN to determine whether any disorder exists in the image. Additionally, the disorder locating module can be used to locate the disorder. Our goal is to provide an effective tool that can assist physicians and researchers in quickly identifying images that may contain disorders from a large volume of CT images for further analysis and diagnosis.

To develop the system, we collected CT scans from the Chinese PLA General Hospital (PLAGH), a leading national hospital that serves patients throughout China. We constructed a training dataset of 21,429 healthy brain CT scans (March 2012 to July 2019) retrieved from the picture archiving and communication systems (PACS). The retrieval process involved matching the fixed description (“no abnormality is observed”) of healthy CT scans with historical diagnosis reports, resulting in a training dataset that was efficiently obtained without requiring expert effort or disorder annotation.

### PERFORMANCE ON THE BROAD-SPECTRUM HEAD DISORDER DETECTION

To evaluate the system, we obtained a retrospective test dataset from the PLAGH (9967 scans, 88.23% with 127 types of disorders, March 2012 to July 2019) and a prospective test dataset from the PLAGH (3054 scans,

88.70% with 116 types of disorders, July 2019 to August 2021). To demonstrate the clinical applicability of our system, we counted all types of disorders described in clinical reports from the PLAGH using a rule-based natural language processing (NLP) algorithm and manual selection by radiologists. We sorted out 127 and 116 types of disorders for testing, respectively. To our knowledge, these test datasets have the broadest coverage of head disorder types. The numbers of scans for each disorder are shown in Tables S1 and S2 in the Supplementary Appendix.

We used a disorder-contained/free classification testing strategy for each type of disorder, with testing carried out at the scan level. This means that the system predicted whether the entire scan contained any disorder or not. Scan-level classification is practical for clinical use because it enables radiologists to quickly identify the presence of disorders, which is particularly useful in emergency treatment.<sup>5</sup> The label of each scan was determined using disorder-related keywords in its associated report and then confirmed by radiologists on the basis of the report and CT images.

For the retrospective and prospective test datasets, the area under the receiver operating characteristic curves (AUCs) with 95% confidence intervals (CIs) for the two datasets, along with the true positive rate, false positive rate, and overall receiver operating characteristic (ROC) curves, are presented in Figure S1 and Tables S1, S2, and S4. Additionally, Tables S1 and S2 also present the sensitivity and specificity with 95% CIs for the disorders. On the retrospective dataset, the system achieved an AUC of greater than 0.95 for 43 disorders and an AUC of greater than 0.90 for 74 disorders. On the prospective dataset, the system achieved an AUC of greater than 0.95 for 30 disorders and an AUC of greater than 0.90 for 50 disorders. These results demonstrate that our system is capable of detecting a broad spectrum of disorders in brain CT.

### ANALYSIS OF LESION DETECTION EFFICACY BY SIZE AND URGENCY OF TREATMENT

In our comprehensive analysis, we delved into the challenges of disorder detection. We identified two primary categories of challenging cases in disorder detection: those that are small and easily missed and those that do not require immediate treatment, which may also be overlooked because of their subtler characteristics. To conduct a thorough analysis, we divided the cases into three groups on the basis of these dimensions.

Table 1. Performance for Disorder Types across Different Lesion Sizes.*			
Lesion Size	AUC (95% CI)	Sensitivity (95% CI)	Specificity (95% CI)
Large	0.941 (0.941–0.942)	0.883 (0.882–0.885)	0.865 (0.864–0.867)
Medium	0.943 (0.943–0.944)	0.885 (0.883–0.886)	0.875 (0.873–0.876)
Small	0.887 (0.887–0.888)	0.885 (0.884–0.887)	0.771 (0.770–0.773)

\* AUC denotes area under the receiver operating characteristic curve; and CI, confidence interval.

In terms of lesion size, we classified the cases as large, medium, or small. The classification outcomes are detailed in Tables S7 and S8. We computed the AUC values with 95% CIs (Table 1) and plotted ROC curves (Fig. 1) for each size category. The AUC results for different lesion sizes demonstrate AUC accuracies of 0.941, 0.943, and 0.887 for large, medium, and small lesions, respectively. These figures underscore our model’s high accuracy in detecting even smaller lesions, maintaining a commendable level of recognition precision.

When classifying on the basis of the urgency of treatment, we sorted the cases into high, medium, and low urgency levels, calculating the corresponding AUC values for each. The categories were defined as follows. The emergency intervention group encompasses severe disorders necessitating immediate medical attention, such as certain cancers and other conditions that could be life threatening. Disorders in the selective intervention category may not require urgent treatment but could necessitate medical intervention as they evolve. The nonintervention group includes disorders that generally do not require treatment and have minimal impact on patient quality of life. Detailed classification results are shown in Table S8.

The AUC accuracies and ROC curves for lesions of high, medium, and low urgency are shown in Figure 2 and Table 2. The AUC results were 0.946, 0.859, and 0.861, respectively. These results indicate that our model is proficient in identifying lesions with varying degrees of urgency, effectively recognizing even those with less pronounced features.

### EVALUATION OF SYSTEM GENERALIZABILITY

To be practical, an AI-based system should be able to generalize to new data from different centers and hospitals. To evaluate the generalizability of the ISL-based system, we constructed a cross-center test dataset from the Brain Hospital of Hunan Province, which served as an independent test cohort from PLAGH. This dataset consisted of 554 scans, of which 59.01% had 46 different types of disorders. It is worth noting that in the cross-center experiment, we made efforts to collect as much available data as possible to ensure the comprehensiveness of the tested disorders. However, this approach resulted in a smaller number of samples for certain disorders (e.g., the total sample size for basal ganglia cerebral infarction was five). As a result, the performance of these specific disorder types may deviate when compared with the retrospective test set.

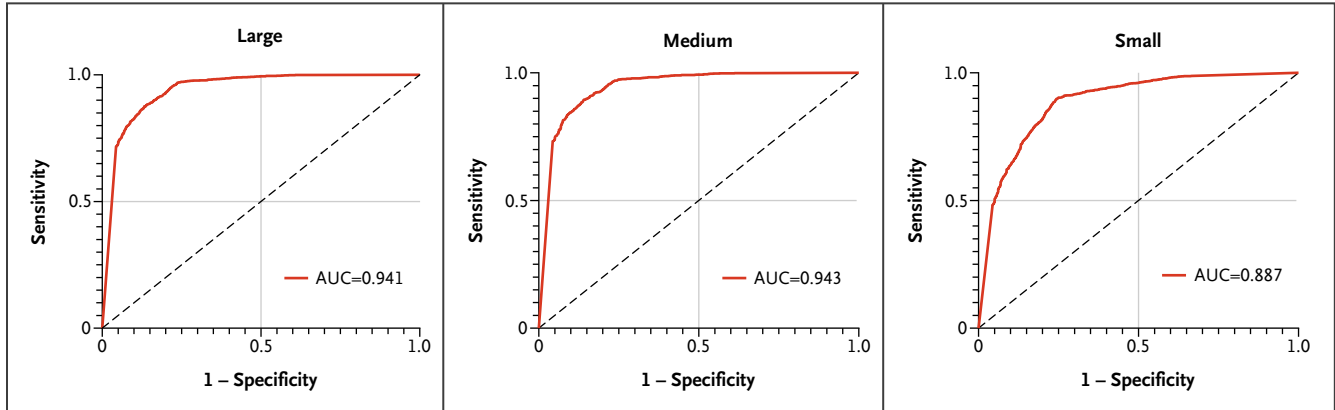


Figure 1. Receiver Operating Characteristic Curves for Disorder Types across Different Lesion Sizes. AUC denotes area under the receiver operating characteristic curve.

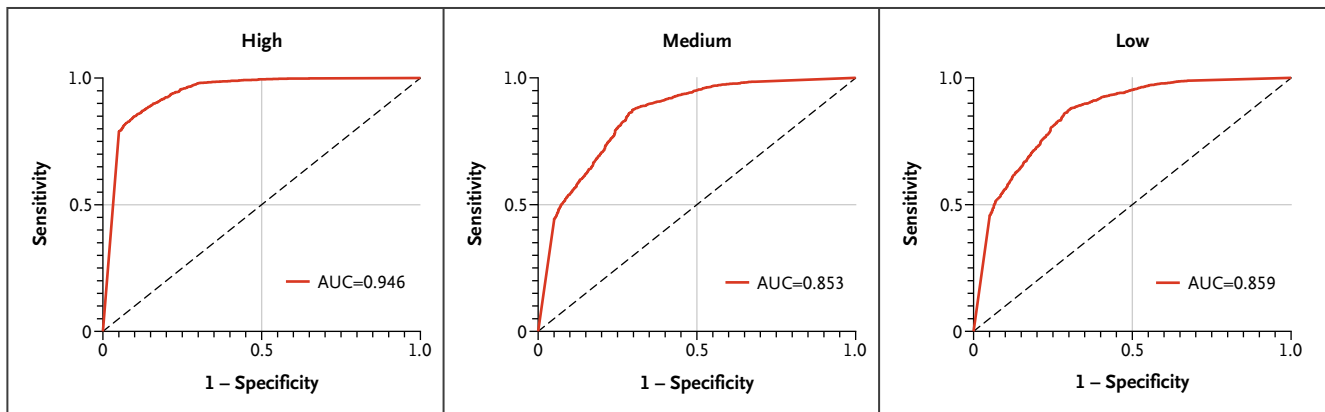


Figure 2. Receiver Operating Characteristic Curves for Disorder Types on the Basis of Urgency of Treatment.

AUC denotes area under the receiver operating characteristic curve.

The AUCs with 95% CIs for the 46 types of disorders along with the overall ROC curve are presented in Figure S1 and Table S3. The average AUC was 0.866, which was only 0.017 lower than that of the retrospective intracenter test. These results demonstrate the generalizability of the system across different centers.

### EVALUATION OF IMPROVING EXPERT PERFORMANCE

In clinical practice, a computer-aided diagnosis (CAD) system should provide understandable clues to support prediction results. Our model can quickly and intuitively locate the disorder region on the basis of the generated dedisorder image, as shown in [Figure 3A](#). Additionally, we developed a CAD software for clinical use, as shown in [Figure 3B](#). The software takes a CT scan as input and outputs possible disorder regions, improving the diagnosis performance of radiologists.

To quantitatively evaluate the improvement, we conducted an experiment involving four radiologists with diverse levels of experience, ranging from 5 to 14 years.

Each radiologist was tasked with independently reviewing a set of 300 randomly chosen samples from our cross-center test dataset, which comprised 100 cases with identified disorders and 200 cases deemed healthy. Initially, the radiologists performed their assessments without the support of our software, relying solely on their expertise. Subsequently, we introduced the diagnostic suggestions provided by our software to examine its influence on the radiologists' ability to diagnose accurately.

The incorporation of the software's insights led to a notable enhancement in diagnostic precision. The average sensitivity across the four radiologists increased by 0.035, whereas the specificity saw a marginal improvement of 0.006. The advancements are visually represented in [Figure 4](#).

The observed improvements underscore the potential of our software to serve as a valuable tool for radiologists, particularly in the accurate detection of disorders. The integration of our software into the diagnostic workflow promises to refine disorder screening processes and support radiologists in delivering more precise and reliable diagnoses. The radiologists reported that the system

Urgency of Treatment	AUC (95% CI)	Sensitivity (95% CI)	Specificity (95% CI)
High	0.942 (0.941–0.942)	0.859 (0.858–0.861)	0.897 (0.895–0.898)
Medium	0.853 (0.853–0.854)	0.849 (0.848–0.851)	0.727 (0.726–0.729)
Low	0.859 (0.859–0.860)	0.838 (0.836–0.840)	0.737 (0.736–0.739)

\* AUC denotes area under the receiver operating characteristic curve; and CI, confidence interval.



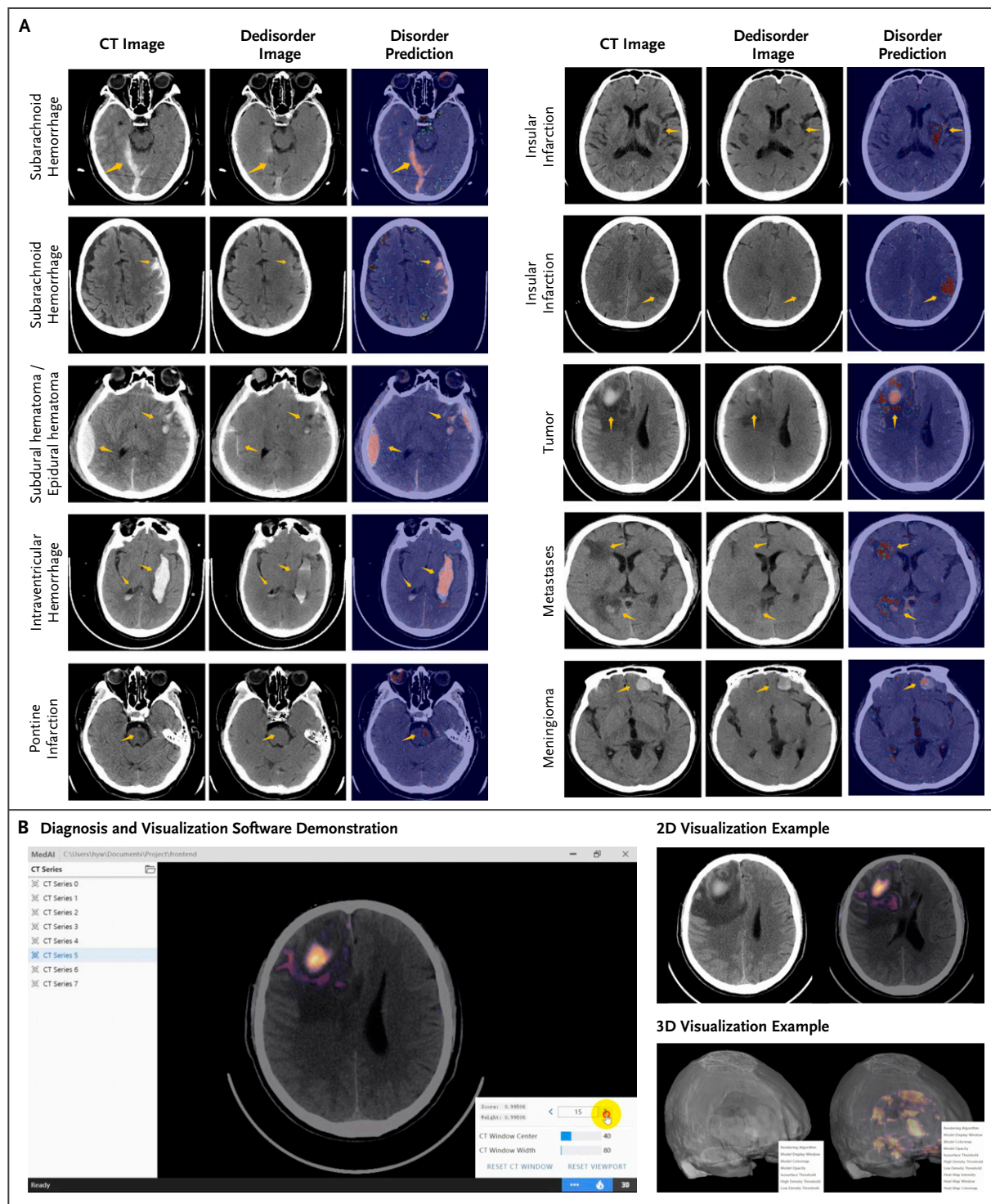
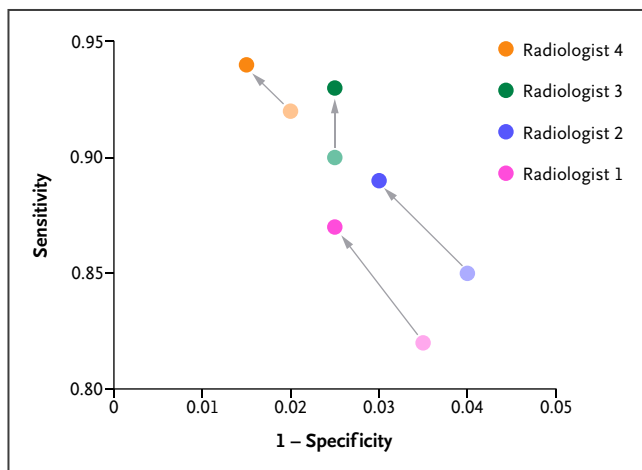


Figure 3. Visualization Examples for Head Disorder Detection.

Typical examples of our system's performance are shown, including the original computed tomography (CT) image, the corresponding dedisorder image generated by our system, and the heat map indicating the probability of containing disorders (Panel A). Warmer colors in the heat map indicate a higher probability of disorders. The heat maps provide visual clues to the system's decision. We also developed a diagnosis and visualization software that takes a CT scan as input and outputs possible disorder locations in the form of a heat map (Panel B). The heat map can be displayed on a two-dimensional (2D) slice or on a three-dimensional (3D) reconstruction scan.



**Figure 4. The Performance of Four Radiologists before and after Considering the System Recommendation.**

The radiologists have 5 (pink), 7 (blue), 10 (green), and 14 (orange) years of working experience.

effectively reduced their workload by accurately identifying a broad spectrum of disorders and contributed to lowering the rate of missed diagnoses. They appreciated the system's ability to provide visually understandable clues, which greatly assisted in their diagnosis process.

### ANALYSIS OF SYSTEM EXPLAINABILITY

Our system not only detects the disorder location in a slice but also provides the disorder distribution in a scan. [Figure 5B](#) shows two example scans with different disorder distributions. On the basis of the distribution curves supplied by our system, we can observe that the disorders in the two scans have centralized and dispersive distributions, respectively. [Figure 5C](#) shows the average disorder distributions for some typical disorders in the retrospective dataset. The average distributions are close to the occurrence frequency at different brain tissues in reality, demonstrating the effectiveness of the system's explainability.

### PERFORMANCE CONTRIBUTIONS FROM DIFFERENT MODULES

This section elaborates on the reasons for adopting each module and demonstrates their contributions to the final performance. The results are presented in [Figure 6](#).

#### *Performance Contribution from the DeDN*

In ISL, the evaluation of the probability of disorder containment depends on the difference image  $\mathbf{x}_{dif}$ , where  $\mathbf{x}_{dd}$

is a dedisorder image generated by a DeDN (see Methods for details). Original medical images are too complex for a system to learn disorder-related information solely on the basis of them. Therefore, we do not directly apply original images for evaluation. Using difference images for evaluation is more intuitive, because the greater the difference between the original and dedisorder images, the higher the probability of disorder containment is.

To numerically demonstrate the effectiveness of the difference image generated by the DeDN, we also applied the original image-based method for evaluation. As shown in [Figure 6](#), on the prospective test dataset, the average AUCs with 95% CIs are 0.657 and 0.752, respectively, where the result from the original image-based method (original CT input) is significantly lower than that from the difference image-based method (D-score evaluation). The improved result highlights the value of the DeDN.

#### *Performance Contribution from the DRN*

After obtaining the difference image  $\mathbf{x}_{dif}$  with a DeDN, we used a DRN to evaluate the probability of disorder containment. Although we could determine the probability directly on the basis of the pixel value sum of the difference image, we did not adopt this strategy. This is because a DeDN cannot produce perfectly healthy tissue, meaning that even for a healthy area, the pixel value sum of that area may still be positive. As a result, the accumulated pixel value sum of all healthy areas would negatively influence the probability evaluation.

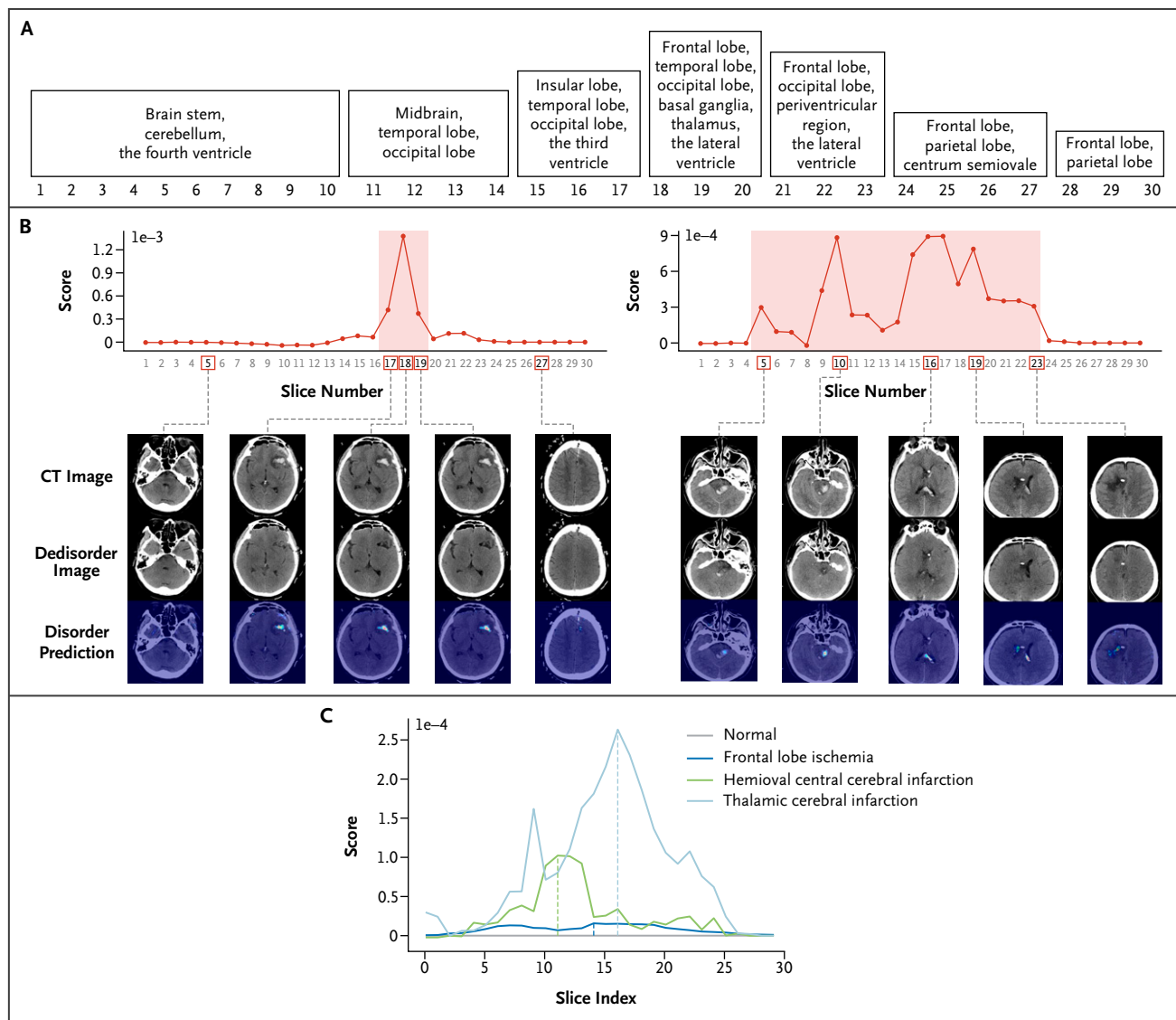
Instead, we used the pixel sum-based evaluation (D-score evaluation) and DRN-based evaluation (network evaluation) on the basis of the difference image, as shown in [Figure 6](#). The average AUCs of the two methods on the prospective dataset were 0.752 and 0.868, respectively, demonstrating the superiority of the DRN.

### EVALUATION OF ISL GENERALIZABILITY

To assess the generalizability of ISL across different body parts and medical image types, we used it to develop two additional systems. The first system is designed for detecting pulmonary disorders in CT images, whereas the second system is designed for detecting retinal disorders in OCT images.

#### *Performance of Pulmonary Disorder Detection*

In addition to brain CT, we developed an ISL-based system for detecting disorders in pulmonary CT scans.



**Figure 5. Patterns of Disorder Distribution in CT Scans.**

The x axes of the graphs represent the slice index of a computed tomography (CT) scan, whereas the y axes represent a slice's abnormal score. A higher score indicates a greater likelihood of lesions in the slice. The correspondence between slice indexes and brain tissues is shown (Panel A). Two example scans with dispersive and centralized distributions are presented (Panel B). The average disorder distributions of some typical disorders in the retrospective dataset are displayed (Panel C).

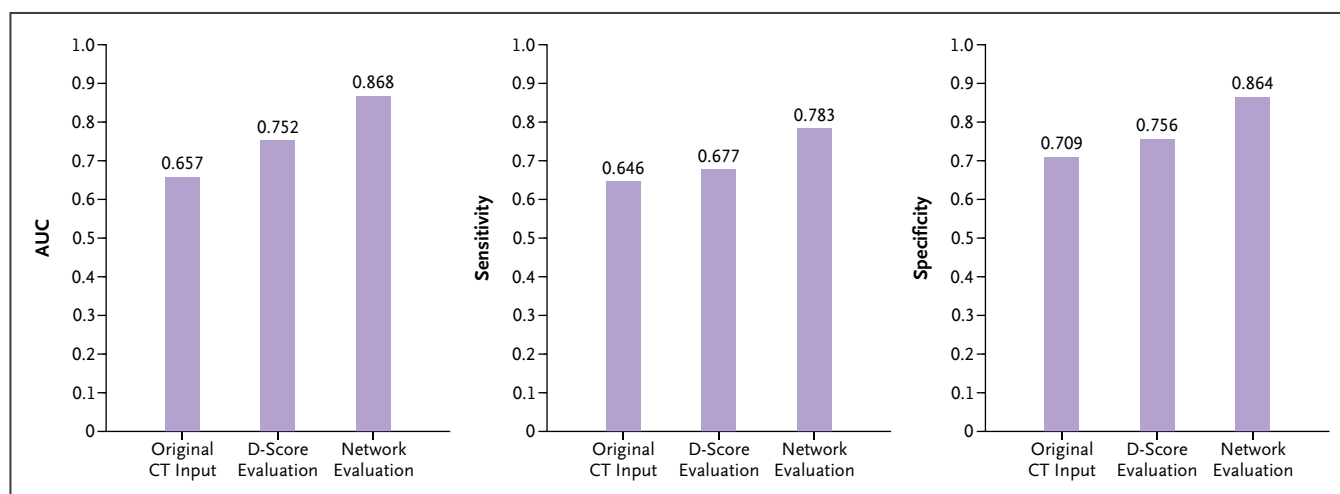
The data used for system development were collected from the First Affiliate Hospital of Guangzhou Medical University, another leading national hospital that serves patients from across China. We constructed a training dataset consisting of 3410 healthy pulmonary CT scans and a test dataset that included six types of pulmonary disorders (82 pneumothorax, 86 pneumonia, 96 bronchiectasis, 88 bullae, 82 atelectasis, and 46 effusion) as well as 600 healthy scans. The AUCs and detection examples for each type of disorder are presented in [Figure 7A](#) and

[Table 3](#). The average AUC was 0.893, indicating that ISL can generalize well across different body parts.

### Performance of Retinal Disorder Detection

To demonstrate the ability of ISL to generalize across different medical image types, we developed a retinal disorder detection system on the basis of OCT images. We used the dataset collected by Kermany et al.,<sup>6</sup> which includes 108,312 images (37,206 with choroidal





**Figure 6. Iterative Performance Improvement by DeDN and DRN.**

Original CT input indicates original computed tomography (CT) slices plus DRN. D-score evaluation indicates DeDN plus pixel value sum. Network evaluation indicates DeDN plus DRN. The comparison was performed on the prospective test dataset. AUC denotes area under the receiver operating characteristic curve; DeDN, dedisorder network; and DRN, disorder recognition network.

neovascularization, 11,349 with diabetic macular edema, 8617 with drusen, and 51,140 normal). Following the development procedure of ISL, we used only the normal OCT images as the training dataset. The model was tested with 1000 images (250 from each category) from 633 patients, as in Kermany et al.<sup>6</sup> The AUCs with 95% CIs on the scan level are summarized in [Table 4](#). The AUCs for choroidal neovascularization, diabetic macular edema, and drusen were 0.939, 0.913, and 0.827, respectively. Despite being developed using only normal OCT images, our system achieved clinically acceptable results, indicating that ISL is applicable to different medical image types. Detection examples for each type of disorder are shown in [Figure 7B](#).

## Discussion

We introduced a learning strategy called ISL and utilized it to develop a head disorder detection system that requires no disorder data or annotation during the development process. The system's detectable disorder coverage is comparable with that of a human expert. Additionally, the system's excellent generalizability and explainability enhance its clinical applicability.

### ANNOTATED AND DISORDER-CONTAINED DATA

Most existing deep-learning medical systems rely on supervised learning, which requires a substantial amount

of annotated data to achieve generalizability, accuracy, and recognition gratuity. However, obtaining sufficient annotated data in medical image research is challenging because of the time-consuming and expert knowledge-intensive nature of the annotating process. For instance, even for an experienced expert, it may take several minutes to annotate a medical image at the region level, which provides strong supervision for disorder detection by indicating the exact lesion region. Consequently, research works that rely on region-level annotation, such as Nikolov et al.<sup>7</sup> and Monteiro et al.,<sup>8</sup> are limited by the amount of annotated data, which hinders the generalizability and accuracy of the system.

To reduce the dependence on annotated data, researchers have explored alternative learning strategies for medical image research. For instance, weakly supervised learning<sup>9-11</sup> allows each training sample to lack a label or have an incorrect label, significantly reducing the annotation cost by experts. Unsupervised learning, on the other hand, uses unannotated training data to enhance the feature representation capacity of a deep-learning network, thereby reducing the number of required annotated samples. Self-supervised learning is a recent representative unsupervised learning method<sup>12,13</sup> that annotates each sample by itself instead of relying on human experts. However, these learning strategies require a substantial amount of disorder-contained data to ensure accuracy. Collecting enough disorder-contained data is challenging for general researchers because of ethical and legal considerations,

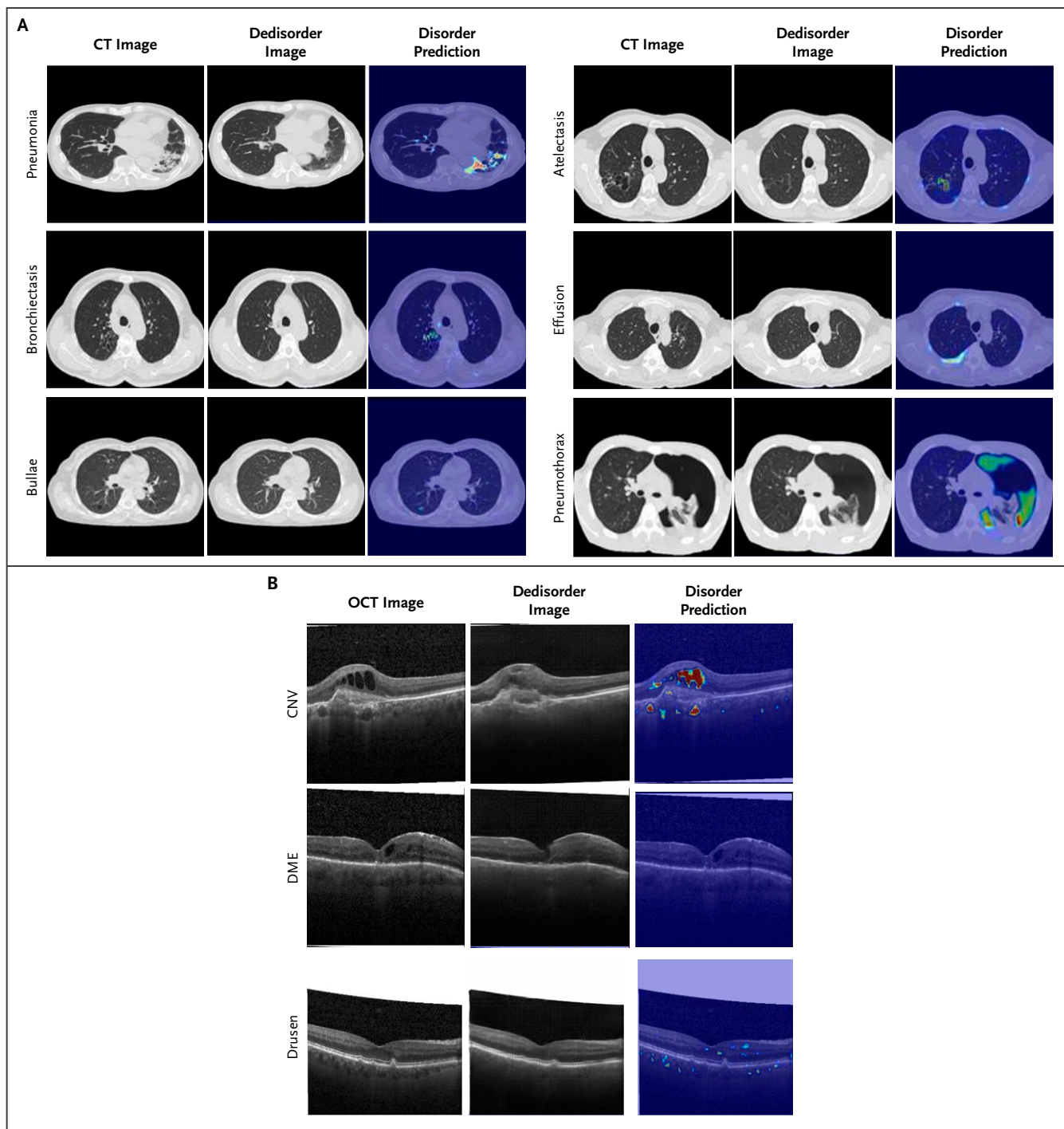


Figure 7. Examples of Our System on CT-Based Pulmonary Disorder Detection (Panel A) and OCT-Based Retinal Disorder Detection (Panel B).

CNV denotes choroidal neovascularization; CT, computed tomography; DME, diabetic macular edema; and OCT, optical coherence tomography.

Table 3. Performance on the Pulmonary Computed Tomography Test Dataset for Pulmonary Disorder Detection.*				
Pulmonary CT	Number	AUC (95% CI)	Sensitivity (95% CI)	Specificity (95% CI)
Pneumothorax	82	0.992 (0.992–0.992)	0.996 (0.995–0.996)	0.937 (0.935–0.939)
Pneumonia	86	0.911 (0.909–0.912)	0.874 (0.869–0.878)	0.830 (0.828–0.831)
Bronchiectasis	96	0.811 (0.807–0.811)	0.697 (0.689–0.704)	0.816 (0.815–0.817)
Bullae	88	0.786 (0.782–0.787)	0.589 (0.582–0.595)	0.827 (0.826–0.828)
Atelectasis	82	0.952 (0.951–0.952)	0.984 (0.982–0.986)	0.853 (0.852–0.855)
Effusion	46	0.958 (0.956–0.958)	0.999 (0.999–1.000)	0.883 (0.880–0.885)
Average	6/6	0.893 (0.891–0.894)	0.838 (0.834–0.842)	0.854 (0.853–0.855)

\* AUC denotes area under the receiver operating characteristic curve; CI, confidence interval; and CT, computed tomography.

limiting related research to large medical institutions. For example, Chilamkurthy et al.<sup>5</sup> collected over 300,000 brain CT scans from more than 20 medical centers, which is beyond the reach of most researchers.

Compared with previously adopted learning strategies in medical image research, the proposed ISL tackles a challenging task where no annotated or disorder-contained data are available. The only available data are disorder-free data, which can be easily obtained by any medical institution capable of medical imaging scans.

## DISORDER COVERAGE

The clinical application of medical image research is an important goal. However, most existing works focus on only one or two common disorder types,<sup>5,14</sup> even for systems developed by institutions with abundant medical resources. For instance, the system in the work of Chilamkurthy et al.<sup>5</sup> is derived from over 300,000 scans, yet it can only recognize four types of disorders. This challenge arises from two aspects. First, it is impractical for researchers to construct models for each disorder type because of the difficulty of collecting and annotating medical images. Second, developing models for rare disorders with previous learning strategies is challenging when only a few samples are available. With ISL, researchers do not need to collect data or construct models for specific

disorders, enabling the built system to achieve broad-spectrum disorder detection.

## ANOMALY DETECTION

Distinguishing disorder-contained images from disorder-free ones can be viewed as an anomaly detection problem, which is a popular research field in machine learning. An intuitive assumption is that anomalies lie outside the distribution of normal samples. Therefore, it is natural to train a classifier to differentiate abnormal samples from normal ones.<sup>15,16</sup>

Recent works have utilized generation networks for anomaly detection using two primary approaches: utilizing the latent feature<sup>17–19</sup> and utilizing the reconstructed image.<sup>20–22</sup> In the first approach, a generation network produces a latent feature and a reconstructed image when an image is inputted. The latent features can be used to determine whether the image is abnormal. In the second approach, a generation network produces a corresponding normal image for a given image. If the original image contains abnormal characteristics, it can be recognized on the basis of the difference between the original and generated images. In the field of medical image analysis, two types of methods have achieved certain results in specific diseases.<sup>23–25</sup> For instance, Yao et al.<sup>23</sup> used the second approach to generate healthy pulmonary CT images,

Table 4. Performance on the Retinal Optical Coherence Tomography Test Dataset for Retinal Disorder Detection.*				
Retinal OCT	Number	AUC (95% CI)	Sensitivity (95% CI)	Specificity (95% CI)
CNV	1220	0.939 (0.938–0.940)	0.847 (0.844–0.851)	0.904 (0.900–0.907)
DME	1600	0.913 (0.912–0.914)	0.838 (0.835–0.841)	0.890 (0.888–0.893)
Drusen	1220	0.827 (0.826–0.829)	0.760 (0.757–0.764)	0.762 (0.758–0.765)
Average	3/3	0.895 (0.894–0.896)	0.817 (0.814–0.820)	0.855 (0.852–0.858)

\* AUC denotes area under the receiver operating characteristic curve; CI, confidence interval; CNV, choroidal neovascularization; DME, diabetic macular edema; and OCT, optical coherence tomography.

which were used to determine whether the lungs exhibited Covid-19.

However, both approaches have limitations when applied to broad-spectrum disorder detection in medical images. In the first approach, medical images are complex, which results in complex latent features. Therefore, recognizing disorder-contained images solely on the basis of latent features is challenging. To demonstrate this, we compared our method with a baseline method that directly fed original medical images into the DRN. The baseline method achieved an average AUC of 0.653 on the prospective dataset, which is inferior to the results (AUC of 0.868) obtained by our method. In the second approach, existing generation-based methods reconstruct the original disorder tissues of a medical image because of the strong feature representation capability of generative networks, which fails to achieve abnormal recognition. With the generation strategy in ISL, only context images and global structure information are provided, allowing the generation network to eliminate the interference of the original disorder tissue and conceive healthy tissues like a radiologist. To showcase the efficacy of our system in comparison with existing techniques, we conducted a comparative analysis with other representative reconstruction-based anomaly detection methods, specifically Auto-Encoder,<sup>26</sup> AnoGAN,<sup>17</sup> GANomaly,<sup>27</sup> pix2pix,<sup>28</sup> and Cycle-GAN.<sup>29</sup> The experiment was carried out on the task of detecting pulmonary disorders. The results of the analysis are presented in Table S12. Our system outperformed the baselines, with the highest AUC of 0.846 achieved by GANomaly, which is 0.047 lower than our method. This significant improvement underscores the ability of our ISL-based system to successfully accomplish disorder recognition tasks.

## Methods

### CT SCAN COLLECTION

Initially, we retrieved 954,508 scans from the PACS of the PLAGH between March 2012 and July 2019. These scans contained CT images stored in the digital imaging and communications in medicine (DICOM) format, and all DICOMs were deidentified before data analysis. We then screened the scans by excluding reconstructed scans (processed with algorithms in CT machines), nonaxial section scans (coronal section and sagittal section scans), nonhead scans (scans of breast and full body, etc.), and nonorigin scans (scans of CT angiography and CT perfusion, etc.).

The inclusion and exclusion criteria of the screening process are detailed in [Figure 8](#). After screening, a total of 62,239 CT scans with an average slice number of 28 were selected.

### DISORDER TYPES STATISTICS

Each retrieved scan includes a clinical report written by an interpreting radiologist during the examination. To determine the disorder types to be interpreted by our system, we first applied a rule-based NLP algorithm to the clinical reports. This algorithm counted the occurrence frequencies of different word phrases. We then invited three radiologists to analyze the frequency statistics results and select the disorder types to be evaluated. Ultimately, 127 types of disorders were selected.

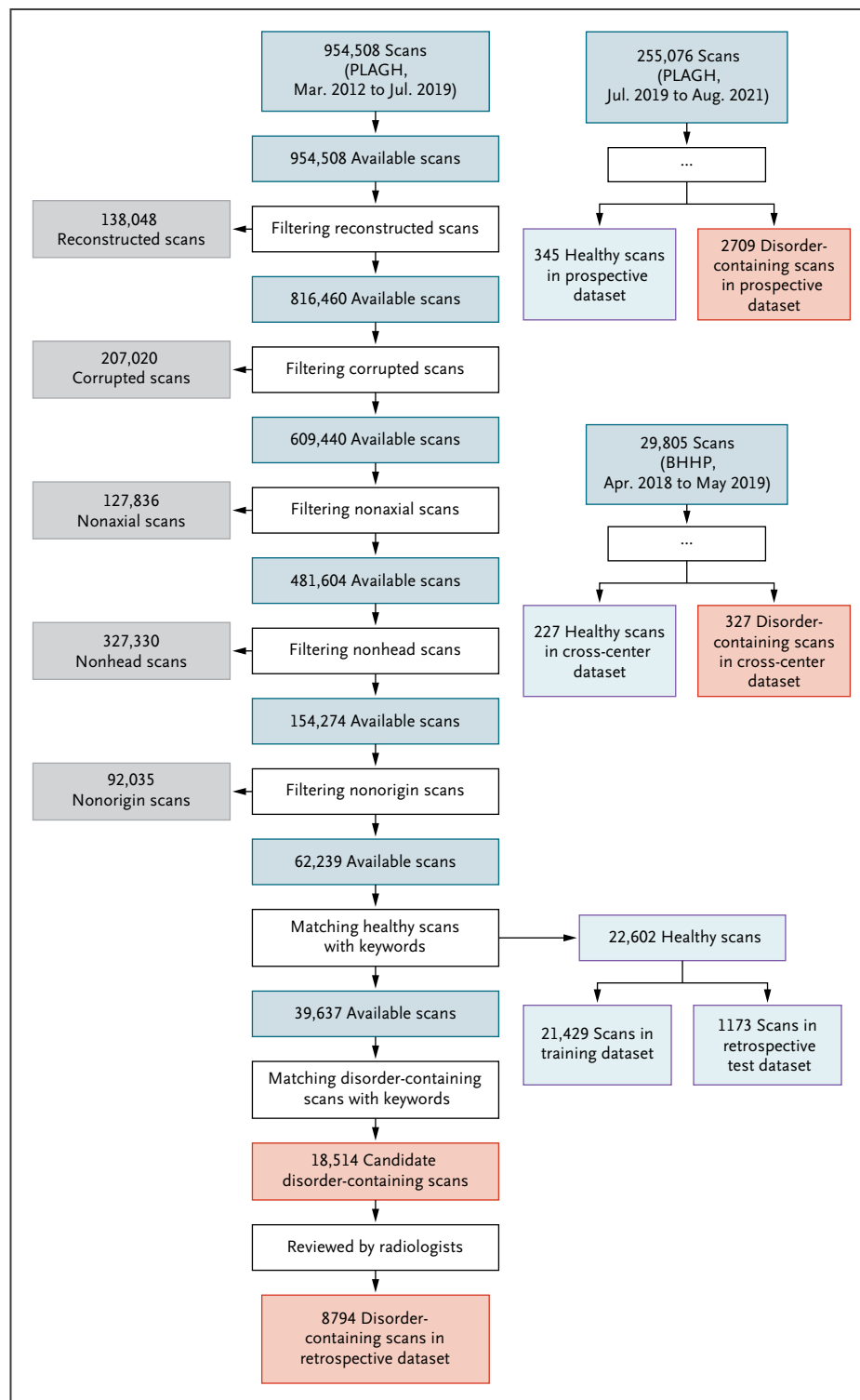
Our method selection was primarily driven by the desire to create a system capable of handling the most common types of disorders encountered in actual clinical practice while also ensuring a comprehensive coverage of various disorder types. Our dataset, which spans nearly 7 years (2012 to 2019) and includes data from 301 hospitals, is believed to encompass most types of disorders. By focusing on the most frequently occurring disorder types, we aimed to enhance the practicality of our system, ensuring that it is well equipped to manage common disorders while maintaining a broad scope of disorder types.

### TRAINING DATASET AND TEST DATASET

The construction of the development and test datasets relied on the clinical reports, which were considered the gold standard. The training dataset only included disorder-free CT scans, and their reports uniformly described them as “no abnormality is observed.” Therefore, we could efficiently obtain disorder-free CT scans. Ultimately, we selected 22,602 disorder-free scans, which were divided into two parts: the training dataset (21,429 scans) and the negative samples in the retrospective test dataset (1173 scans); detailed data statistics are shown in Table S5.

Regarding the positive samples (disorder-containing scans) in the test dataset, we initially retrieved 18,514 scans using stated disorder-related keywords. We then invited 30 board-certified radiologists with 6 to 15 years of experience to label each scan on the basis of the images and its associated report. The radiologists assigned a binary label (i.e., zero, one) to each scan, where one indicated that the scan contained the expected disorder type. Ultimately, 8794 scans were labeled as one and selected





**Figure 8. The Process of Constructing the Datasets.**

The left side shows that 954,508 scans (March 2012 to July 2019) were collected from the Chinese PLA General Hospital (PLAGH) by retrieving head computed tomography–related keywords. After a series of filtering steps, a training dataset and a retrospective dataset were constructed. The training dataset consisted of 21,429 healthy scans, and the retrospective dataset consisted of 1173 healthy scans and 8794 disorder-containing scans. The right side shows that we collected 255,076 scans (July 2019 to August 2021) from the PLAGH and 29,805 scans (April 2018 to May 2019) from the Brain Hospital of Hunan Province (BHHP). We used these data to construct the prospective and cross-center test datasets using the same process.

as the positive samples in the retrospective test dataset. The prospective and cross-center test datasets were constructed similarly. However, because of the smaller data amounts compared with the retrospective dataset, they also contained a smaller number of disorder types, specifically 116 and 46, respectively. Tables S9 to S11 show the detailed data statistics.

## DEVELOPING A SYSTEM WITH ISL

ISL allows for the training of a deep-learning network without accessing disorder-contained samples, enabling researchers with only general and healthy images to build a broad-spectrum disorder detection system. ISL is built on two technologies: missing information completion and data distribution estimation. Missing information completion enables a system to reconstruct healthy tissues of masked parts of a medical image using a DeDN derived from general and healthy images. The scanning medical images of the human body are relatively standardized. Therefore, for healthy images, the reconstructed version should be very close to the original version. Additionally, for a medical image containing any disorders, the reconstructed image will differ significantly from the original. Data distribution estimation requires the estimation of the distribution of healthy difference images, which are calculated using healthy images and their reconstructed images generated by a DeDN. If an image contains any disorders, its difference image will fall outside the distribution and be detected. Notably, unlike many existing disorder detection algorithms, ISL can detect a significantly increased number of disorder types.

## CT SLICE CONVERSION

In our dataset, the pixel values in a CT scan are represented by 14-bit numbers, which exceed the range of human perception. To address this, we converted each CT slice into a three-channel, eight-bit image, conforming to the standard image format and suitable for display. Radiologists typically use specific WLs and WWs to observe various types of disorders. Building on this, we applied specific WL and WW settings for the image conversion, with the specific settings outlined in Table S16.

## PROBLEM FORMULATION

ISL is designed to address binary-classification problems by predicting the probability of the presence of any disorder type in a medical image. For example, in brain CT scans, the input of an ISL-based system is a slice  $\mathbf{x}_i$  from a brain CT scan for a CT scan  $\{\mathbf{x}_i\}_{i=1}^N$ , where  $N$  is the slice

number of this scan. The system output  $p_i$  indicates the probability of any disorder type in the slice  $\mathbf{x}_i$ . During deployment, the slice-level outputs  $\{p_i\}_{i=1}^N$  are aggregated to a scan-level output  $p'$  by averaging the probabilities of all the slices in the scan, where  $p' = \frac{1}{N} \sum_{i=1}^N p_i$ . We adopt the slice-wise processing method because we believe that for the initial assessment of disorders in medical image analysis, the information provided by a single image is already adequate. Slice-wise processing offers a more efficient strategy, where sequential information is utilized to confirm the precise categories of disorders. As the ISL-based system processes individual slices, we have omitted the subscript number of slices in a scan for the sake of conciseness in the following method introduction.

An ISL-based system comprises two networks: a DeDN and a DRN. Given a medical image  $\mathbf{x}$ , we first use a DeDN to generate a dedisorder image  $\mathbf{x}_{dd}$  of  $\mathbf{x}$ . If  $\mathbf{x}$  contains a disorder, the disorder tissues in the area are converted into healthy ones. Then, the difference image  $\mathbf{x}_{diff} = \|\mathbf{x} - \mathbf{x}_{dd}\|$  is inputted into the DRN network, which predicts the probability  $p$  of disorder containment in the image. The overview of ISL is shown in [Figure 9](#).

## DEDN

Given a masked image,  $\bar{\mathbf{x}} = \mathbf{x} \cdot \mathbf{m}$ , where  $\mathbf{m}$  is an image mask of  $\mathbf{x}$ , a deep encoder-decoder network (DeDN) can predict the masked region and generate a reconstructed image  $\hat{\mathbf{x}}$ . The detailed architecture of the DeDN is shown in Table S15, and it has been proven to be effective in many image generation tasks. In our architecture comparison experiment (Table S13), we found that the adopted architecture has already captured the most salient features necessary for generating high-quality medical images. In this study, we utilized the DeDN to generate dedisordered medical images. Specifically, we divided a medical image  $\mathbf{x}$  into  $K \times K$  grids of uniform size. For each grid with coordinates  $(i, j)$ , where  $1 \leq i \leq K$  and  $1 \leq j \leq K$ , we applied a mask  $\mathbf{m}^{(i,j)}$  to erase it and obtain a masked image. The DeDN was then used to reconstruct the masked image and generate the reconstructed image  $\hat{\mathbf{x}}^{(i,j)}$ . Finally, we combined the  $K \times K$  generated images into a reconstructed dedisordered medical image using the following equation:

$$\mathbf{x}_{dd} = \sum_{i=1}^K \sum_{j=1}^K \hat{\mathbf{x}}^{(i,j)} \cdot (1 - \mathbf{m}^{(i,j)}). \quad (1)$$

We will use  $\mathbf{x}_{dd}$  for comparative analysis with the original image  $\mathbf{x}$ . A deep encoder-decoder network (DeDN) takes as input a masked image  $\bar{\mathbf{x}}^{(i,j)}$  and multiple image edge

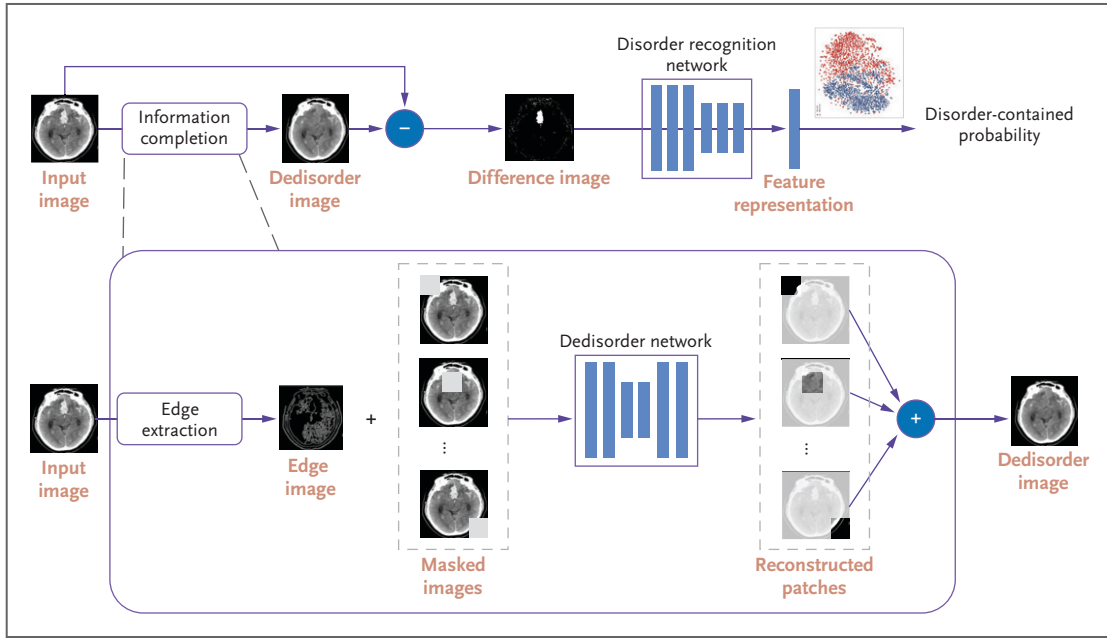


Figure 9. The Overview of ISL, a Learning Algorithm for Developing Broad-Spectrum Disorder Detection Systems.

The training dataset consists only of healthy scans, and a dedisorder network is learned to generate dedisorder images. A disorder recognition network is then used to predict the probability of disorder containment on the basis of the difference image obtained by subtracting the input and dedisorder image. This approach enables the developed model to achieve broad-spectrum disorder detection even without any disorder-contained data. ISL denotes inverse supervised learning.

maps  $\{\mathbf{e}_k\}_{k=1}^{n_e}$  of  $\mathbf{x}$ . Edge maps retain structural information of the masked region, which can improve the quality of the reconstructed image. Edge maps can be constructed using mature image processing schemes, such as the Canny Edge Detector.<sup>30</sup> The DeDN  $G$  generates the reconstructed image  $\hat{\mathbf{x}}^{(i,j)}$  using the following equation:

$$\hat{\mathbf{x}}^{(i,j)} = G(\bar{\mathbf{x}}^{(i,j)}, \{\mathbf{e}_k\}_{k=1}^{n_e}). \quad (2)$$

We trained the network using a joint loss:

$$\mathcal{L}_{de} = \lambda_{\ell_1} \mathcal{L}_{\ell_1} + \lambda_{adv} \mathcal{L}_{adv} + \lambda_p \mathcal{L}_{perc} + \lambda_s \mathcal{L}_{style}, \quad (3)$$

which includes an  $\ell_1$  loss, adversarial loss, perceptual loss, and style loss. The  $\ell_1$  loss minimizes the reconstruction error between  $\hat{\mathbf{x}}^{(i,j)}$  and  $\mathbf{x}$ . The adversarial loss  $\mathcal{L}_{adv}$  ensures the reality of the generated image and is defined as

$$\mathcal{L}_{adv} = \mathbb{E}_{(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{(\mathbf{x}, \mathbf{m}^{(i,j)})} \log [1 - D(G(\bar{\mathbf{x}}^{(i,j)}, \{\mathbf{e}_k\}_{k=1}^{n_e}))], \quad (4)$$

where  $D$  is the discriminator network. We also included perceptual loss  $\mathcal{L}_{perc}$  and style loss  $\mathcal{L}_{style}$ , following Nazari

et al.<sup>31</sup> The perceptual loss  $\mathcal{L}_{perc}$  penalizes reconstructed images that are not perceptually similar to the original ones and is defined as a distance measure between activation maps of a pretrained network:

$$\mathcal{L}_{perc} = \mathbb{E}_{(\mathbf{x})} \left[ \sum_i \frac{1}{n_a} \|\phi_i(\mathbf{x}) - \phi_i(\hat{\mathbf{x}})\|_1 \right], \quad (5)$$

where  $\phi_i$  is the activation map of the  $i$ th layer of the pretrained VGG-19 network and  $n_a$  is the number of layers. We chose the output of the first rectified linear units (ReLU) activation layer in each of the five blocks<sup>31</sup> of VGG-19 pretrained on the ImageNet dataset.<sup>32</sup> This choice was on the basis of its proven effectiveness in capturing image features. The comparison experiment on the cross-center test dataset showed similar results to ResNet34 (Table S14), indicating the robustness of our model to the choice of architecture for perceptual loss calculation.

The style loss is calculated on the basis of these activation maps and is an effective tool to alleviate the “checkerboard” artifacts caused by transposed convolution layers. The loss

measures the differences between covariances of the activation maps and is defined as

$$\mathcal{L}_{\text{style}} = \mathbb{E}_{(\mathbf{x})} \left[ \sum_i \frac{1}{n_a} \|G_i^{\phi_i}(\mathbf{x}) - G_i^{\phi_i}(\hat{\mathbf{x}})\|_1 \right], \quad (6)$$

where  $G_i^{\phi_i} = \phi_i \phi_i^T$  is a Gram matrix constructed from the activation map  $\phi_i$ . The Gram matrix of an activation map captures the correlation between different channels and the texture structure of its corresponding image. For a real medical image, the Gram matrix resembles the identity matrix, with larger diagonal values indicating strong correlations within the same feature and smaller off-diagonal values reflecting feature independence. Conversely, a blurry and texture-lacking generated image results in a constant Gram matrix with similar values for each element, indicating a lack of feature differentiation. To optimize the model, we minimize the difference between the Gram matrices of the real and generated images.

### DRN

To minimize the impact of reconstructed noise on disorder detection and improve the performance further, we developed a DRN that takes the difference image  $\mathbf{x}_{\text{dif}}$  as input and extracts an embedded representation in the latent space. We designated the embedding distribution of difference images from disorder-free data as the reference distribution. The DRN should ensure that the embeddings of disorder-free data are centralized and compact, whereas the embeddings of disorder-contained data are random and as far as possible from the reference distribution. In this case, the distance between an embedded representation and the center of the reference distribution can effectively indicate the possibility of disorder presence.

Inspired by the support vector data description algorithm<sup>33</sup> and the contrastive learning method,<sup>34</sup> we developed the DRN on the basis of augmentation views. In addition to a given healthy medical image  $\mathbf{x}$ , the DRN uses two augmented views,  $\mathbf{x}^-$  and  $\mathbf{x}^+$ , generated from  $\mathbf{x}$  for network training.  $\mathbf{x}^-$  is produced with rotation and flipping transformations, yielding an appearance akin to  $\mathbf{x}$ .  $\mathbf{x}^+$ , on the other hand, is generated with cutout transformation, which can damage healthy tissues in  $\mathbf{x}$ . As a result,  $\mathbf{x}^+$  disrupts the inherent distribution of the healthy image  $\mathbf{x}$  and is thus considered a disorder-contained view. For DRN, the main basis for judgment is the size of the pixel difference and the range of difference. Therefore, by applying cutout transformation to the original healthy image, we

can obtain images with large pixel differences and a wide range of differences. This difference or “anomaly” is what DRN is trained to detect.

After training the DeDN, the input of the DRN consists of three parts: the reference difference image  $\mathbf{x}_{\text{dif}}$ , the negative difference image  $\mathbf{x}_{\text{dif}}^- = \|\mathbf{x}^- - \hat{\mathbf{x}}_{\text{dif}}^-\|$ , and the positive difference image  $\mathbf{x}_{\text{dif}}^+ = \|\mathbf{x}^+ - \hat{\mathbf{x}}_{\text{dif}}^+\|$ , where  $\hat{\mathbf{x}}^-$  and  $\hat{\mathbf{x}}^+$  denote the reconstructed images of  $\mathbf{x}^-$  and  $\mathbf{x}^+$ , respectively. The DRN extracts embedded representations of the difference images, denoted as  $\mathbf{h}$ ,  $\mathbf{h}^-$ , and  $\mathbf{h}^+$ . The DRN learns reasonable embedded representations of disorder-free input by maximizing the similarity between  $\mathbf{h}^-$  and  $\mathbf{h}$  while distinguishing  $\mathbf{h}^+$  from  $\mathbf{h}$ . In this study, we used the Euclidean distance as the metric to measure the similarity of the embeddings. We pretrained the network using MoCo<sup>12</sup> and averaged the embeddings of the wide-ranging training dataset. The averaged embedding  $\mathbf{c}$  is considered the center of the reference distribution. Setting the center as an anchor, we designed a compactness loss to maximize the similarity between the negative embeddings:

$$\mathcal{L}_{\text{com}} = \mathbb{E}_{(\mathbf{h})} [\|\mathbf{h} - \mathbf{c}\|_2] + \mathbb{E}_{(\mathbf{h}^-)} [\|\mathbf{h}^- - \mathbf{c}\|_2], \quad (7)$$

where  $\mathcal{L}_{\text{com}}$  minimizes the distances between the embeddings  $\mathbf{h}$ ,  $\mathbf{h}^-$ , and the reference center, which ensures that the DRN can extract consistent features for disorder-free difference images. To further improve the discriminative ability of the network, we used a discrimination loss  $\mathcal{L}_{\text{dis}}$ , which forces the network to maximize the distance between the reference center and the embedded representation of  $\mathbf{x}^+$ :

$$\mathcal{L}_{\text{dis}} = -\mathbb{E}_{(\mathbf{h}^+)} [\|\mathbf{h}^+ - \mathbf{c}\|_2]. \quad (8)$$

The overall loss function utilized to train the DRN is defined as

$$\mathcal{L} = \lambda_c \mathcal{L}_{\text{com}} + \lambda_d \mathcal{L}_{\text{dis}}, \quad (9)$$

where  $\lambda_c$  and  $\lambda_d$  are the weights of the loss functions  $\mathcal{L}_{\text{com}}$  and  $\mathcal{L}_{\text{dis}}$ , respectively.

### DISORDER VISUALIZATION

The ISL-based system is capable of identifying the locations of disorders, which is crucial for clinical applications.<sup>1,35,36</sup> Higher pixel values in the image regions of  $\mathbf{x}_{\text{dif}}$  indicate a higher likelihood of the presence of a disorder. To enhance the visual appeal of the results, we conducted several postprocessing steps on  $\mathbf{x}_{\text{dif}}$ : eliminating the bias caused by the normal range reconstruction (pixels with values below a certain threshold  $t$  were set to zero), reducing reconstruction noise (after normalizing the pixel



values to the range of [0, 1], we added the values of the  $s \times s$  region surrounding each pixel to itself; this smoothing technique reduced the noise in the image), and enhancing the disorder area (we utilized an exponential function to manipulate the pixel values, resulting in an amplification of the differences in values among pixels). This process serves to accentuate the presence of disorder within the region of interest.

With processed  $\mathbf{x}_{dif}$ , which is denoted as  $\mathbf{x}_{dif}^*$ , we used the average pixel difference score (APDS) as a metric to quantify the discrepancy between the original and reconstructed images. The APDS is computed by averaging the pixel values of the processed  $\mathbf{x}_{dif}$  within the effective pixel area. This area encompasses human body structures and is differentiated by nonzero pixel values. Formally, given an image  $\mathbf{x}$ , its APDS is calculated as the ratio of the sum of pixel values in the original image  $\mathbf{x}$  to the number of nonzero pixels in the corresponding processed difference image  $\mathbf{x}_{dif}^*$ , denoted as  $sum(\mathbf{x}_{dif}^*)/count(\mathbf{x} > 0)$ . Our experimental results revealed that the APDS for normal images was approximately  $5 \times 10^{-5}$ . In contrast, for images with lesions, this metric typically escalated to an order of  $5 \times 10^{-4}$ . The observed difference in these metrics is substantial enough to effectively distinguish between normal images and those with lesions.

## MODEL SELECTION AND STATISTICAL ANALYSIS

Because we were unable to access data containing disorders for model evaluation during training, we selected the model when the training loss did not decrease for five consecutive epochs. The primary parameters of our system, including the network architectures, hyperparameter values, and optimization strategies, are presented in Tables S15 and S16. To ensure statistical significance, we applied 95% CI. Specifically, for each iteration, we randomly sampled 30% of CT scans from the test dataset for evaluation. We repeated this procedure 1000 times and calculated the 95% CI of the evaluation metrics for the model. To determine the optimal classification threshold, we used a derivative-based method, specifically by maximizing the harmonic mean of sensitivity and specificity. This is expressed in the following optimization criterion: [maximize  $(2 \times \text{sensitivity} \times \text{specificity}) / (\text{sensitivity} + \text{specificity})$ ]. This criterion is known as a variant of the F1 score, which balances sensitivity and specificity to achieve the best trade-off between the two.

## Disclosures

Supported by the National Key R&D Program of China (grant numbers 2020AAA0105500 to Dr. Yuchen Guo, Prof. Guiguang Ding, and Prof. Qionghai Dai and 2018YFA0704000 to Prof. Feng Xu) and the National Natural Science Foundation of China (grant numbers U21B2013 to Dr. Yuchen Guo, 62021002 to Prof. Feng Xu, 81825012, 82327803, and 82151309 to Prof. Xin Lou, and 82271952 to Mr. Jinhao Lyu).

Author disclosures and other supplementary materials are available at [ai.nejm.org](https://ai.nejm.org).

We thank the radiologists from the Chinese PLA General Hospital for their efforts in labeling the test data.

## Author Affiliations

- <sup>1</sup> Institute for Brain and Cognitive Sciences, BNRIst, Tsinghua University, Beijing
- <sup>2</sup> School of Software, Tsinghua University, Beijing
- <sup>3</sup> Department of Radiology, Chinese PLA General Hospital, Beijing
- <sup>4</sup> Department of Radiology, Brain Hospital of Hunan Province, Hunan, China
- <sup>5</sup> China State Key Laboratory of Respiratory Disease & National Clinical Research Center for Respiratory Disease, The First Affiliated Hospital of Guangzhou Medical University, Guangzhou, China
- <sup>6</sup> Department of Automation, Tsinghua University, Beijing

## References

1. Lundberg SM, Nair B, Vavilala MS, et al. Explainable machine-learning predictions for the prevention of hypoxaemia during surgery. *Nat Biomed Eng* 2018;2:749-760. DOI: [10.1038/s41551-018-0304-0](https://doi.org/10.1038/s41551-018-0304-0).
2. Liang H, Guo Y, Chen X, et al. Artificial intelligence for stepwise diagnosis and monitoring of COVID-19. *Eur Radiol* 2022;32:2235-2245. DOI: [10.1007/s00330-021-08334-6](https://doi.org/10.1007/s00330-021-08334-6).
3. Motwani M, Dey D, Berman DS, et al. Machine learning for prediction of all-cause mortality in patients with suspected coronary artery disease: a 5-year multicentre prospective registry analysis. *Eur Heart J* 2017;38:500-507.
4. Hadamitzky M, Achenbach S, Al-Mallah M, et al. Optimized prognostic score for coronary computed tomographic angiography: results from the CONFIRM registry (COronary CT Angiography Evaluation For Clinical Outcomes: An International Multicenter Registry). *J Am Coll Cardiol* 2013;62:468-476. DOI: [10.1016/j.jacc.2013.04.064](https://doi.org/10.1016/j.jacc.2013.04.064).
5. Chilamkurthy S, Ghosh R, Tanamala S, et al. Deep learning algorithms for detection of critical findings in head CT scans: a retrospective study. *Lancet* 2018;392:2388-2396. DOI: [10.1016/S0140-6736\(18\)31645-3](https://doi.org/10.1016/S0140-6736(18)31645-3).
6. Kermayn DS, Goldbaum M, Cai W, et al. Identifying medical diagnoses and treatable diseases by image-based deep learning. *Cell* 2018;172:1122-1131.e9. DOI: [10.1016/j.cell.2018.02.010](https://doi.org/10.1016/j.cell.2018.02.010).

7. Nikolov S, Blackwell S, Mendes R, et al. Deep learning to achieve clinically applicable segmentation of head and neck anatomy for radiotherapy. September 12, 2018 (<https://arxiv.org/abs/1809.04430v1>). Preprint.
8. Monteiro M, Newcombe VFJ, Mathieu F, et al. Multiclass semantic segmentation and quantification of traumatic brain injury lesions on head CT using deep learning: an algorithm development and multicentre validation study. *Lancet Digit Health* 2020;2:e314-e322. DOI: [10.1016/S2589-7500\(20\)30085-6](https://doi.org/10.1016/S2589-7500(20)30085-6).
9. Campanella G, Hanna MG, Geneslaw L, et al. Clinical-grade computational pathology using weakly supervised deep learning on whole slide images. *Nat Med* 2019;25:1301-1309. DOI: [10.1038/s41591-019-0508-1](https://doi.org/10.1038/s41591-019-0508-1).
10. Lu MY, Williamson DFK, Chen TY, Chen RJ, Barbieri M, Mahmood F. Data-efficient and weakly supervised computational pathology on whole-slide images. *Nat Biomed Eng* 2021;5:555-570. DOI: [10.1038/s41551-020-00682-w](https://doi.org/10.1038/s41551-020-00682-w).
11. Guo Y, He Y, Lyu J, et al. Deep learning with weak annotation from diagnosis reports for detection of multiple head disorders: a prospective, multicentre study [published correction appears in *Lancet Digit Health* 2022;4:e572]. *Lancet Digit Health* 2022;4:e584-e593. DOI: [10.1016/S2589-7500\(22\)00090-5](https://doi.org/10.1016/S2589-7500(22)00090-5).
12. He K, Fan H, Wu Y, Xie S, Girshick R. Momentum contrast for unsupervised visual representation learning. Paper presented at 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Seattle, WA, USA, June 13-19, 2020. DOI: [10.1109/CVPR42600.2020.00975](https://doi.org/10.1109/CVPR42600.2020.00975).
13. He K, Chen X, Xie S, Li Y, Dollár P, Girshick R. Masked autoencoders are scalable vision learners. Paper presented at IEEE/CVF Conference on Computer Vision and Pattern Recognition, New Orleans, LA, USA, June 18-24, 2022. DOI: [10.1109/CVPR52688.2022.01553](https://doi.org/10.1109/CVPR52688.2022.01553).
14. Lee H, Yune S, Mansouri M, et al. An explainable deep-learning algorithm for the detection of acute intracranial haemorrhage from small datasets. *Nat Biomed Eng* 2019;3:173-182. DOI: [10.1038/s41551-018-0324-9](https://doi.org/10.1038/s41551-018-0324-9).
15. Ruff L, Vandermeulen R, Goernitz N, et al. Deep one-class classification. Paper presented at International Conference on Machine Learning, Stockholm, Sweden, July 10-15, 2018.
16. Schölkopf B, Williamson R, Smola A, Shawe-Taylor J, Platt J. Support vector method for novelty detection. Paper presented at NIPS, Denver, CO, USA, November 29-December 4, 1999.
17. Schlegl T, Seeböck P, Waldstein SM, Schmidt-Erfurth U, Langs G. Unsupervised anomaly detection with generative adversarial networks to guide marker discovery. In: Niethammer M, Styner M, Aylward S, et al., eds. *Information processing in medical imaging*. IPMI 2017. Cham, Switzerland: Springer, 2017:146-157. DOI: [10.1007/978-3-319-59050-9\\_12](https://doi.org/10.1007/978-3-319-59050-9_12).
18. Akcay S, Atapour-Abarghouei A, Breckon TP. Ganomaly: semi-supervised anomaly detection via adversarial training. Paper presented at Asian Conference on Computer Vision, Perth, Australia, December 2-6, 2018.
19. Abati D, Porrello A, Calderara S, Cucchiara R. Latent space autoregression for novelty detection. Paper presented at IEEE/CVF Conference on Computer Vision and Pattern Recognition, Long Beach, CA, USA, June 15-20, 2019.
20. Tian Y, Pang G, Chen Y, Singh R, Verjans JW, Carneiro G. Weakly-supervised video anomaly detection with robust temporal feature magnitude learning. Paper presented at IEEE/CVF International Conference on Computer Vision, Montreal, QC, Canada, October 10-17, 2021. DOI: [10.1109/ICCV48922.2021.00493](https://doi.org/10.1109/ICCV48922.2021.00493).
21. Liznerski P, Ruff L, Vandermeulen RA, Franks BJ, Kloft M, Müller K-R. Explainable deep one-class classification. October 12, 2020 (<https://arxiv.org/abs/2007.01760v2>). Preprint.
22. Fan Y, Wen G, Li D, Qiu S, Levine MD, Xiao F. Video anomaly detection and localization via Gaussian mixture fully convolutional variational autoencoder. *Comput Vis Image Underst* 2020;195:102920. DOI: [10.1016/j.cviu.2020.102920](https://doi.org/10.1016/j.cviu.2020.102920).
23. Yao Q, Xiao L, Liu P, Zhou SK. Label-free segmentation of Covid-19 lesions in lung CT. *IEEE Trans Med Imaging* 2021;40:2808-2819. DOI: [10.1109/TMI.2021.3066161](https://doi.org/10.1109/TMI.2021.3066161).
24. Baur C, Graf R, Wiestler B, Albarqouni S, Navab N. Steganomaly: inhibiting cyclegan steganography for unsupervised anomaly detection in brain MRI. Paper presented at International Conference on Medical Image Computing and Computer-Assisted Intervention, Lima, Peru, October 4-8, 2020. DOI: [10.1007/978-3-030-59713-9\\_69](https://doi.org/10.1007/978-3-030-59713-9_69).
25. Stepec D, Skocaj D. Unsupervised detection of cancerous regions in histology imagery using image-to-image translation. Paper presented at IEEE/CVF Conference on Computer Vision and Pattern Recognition, Nashville, TN, USA, June 19-25, 2021. DOI: [10.1109/CVPRW53098.2021.00419](https://doi.org/10.1109/CVPRW53098.2021.00419).
26. Baur C, Wiestler B, Albarqouni S, Navab N. Deep autoencoding models for unsupervised anomaly segmentation in brain MR images. Paper presented at Brain lesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries: 4th International Workshop, BrainLes 2018, Granada, Spain, September 16, 2018. DOI: [10.1007/978-3-030-11723-8\\_16](https://doi.org/10.1007/978-3-030-11723-8_16).
27. Akcay S, Atapour-Abarghouei A, Breckon TP. Ganomaly: semi-supervised anomaly detection via adversarial training. Paper presented at Computer Vision-ACCV 2018: 14th Asian Conference on Computer Vision, Perth, Australia, December 2-6, 2018.
28. Isola P, Zhu J-Y, Zhou T, Efros AA. Image-to-image translation with conditional adversarial networks. Paper presented at IEEE Conference on Computer Vision and Pattern Recognition, Honolulu, HI, USA, July 21-26, 2017.
29. Zhu J-Y, Park T, Isola P, Efros AA. Unpaired image-to-image translation using cycle-consistent adversarial networks. Paper presented at IEEE International Conference on Computer Vision, Venice, Italy, October 22-29, 2017. DOI: [10.1109/ICCV.2017.244](https://doi.org/10.1109/ICCV.2017.244).
30. Ding L, Goshtasby A. On the canny edge detector. *Pattern Recognit* 2001;34:721-725. DOI: [10.1016/S0031-3203\(00\)00023-6](https://doi.org/10.1016/S0031-3203(00)00023-6).

31. Nazeri K, Ng E, Joseph T, Qureshi FZ, Ebrahimi M. Edgeconnect: generative image inpainting with adversarial edge learning. January 11, 2019 (<https://arxiv.org/abs/1901.00212>). Preprint.
32. Deng J, Dong W, Socher R, Li L-J, Li K, Li F-F. ImageNet: a large-scale hierarchical image database. Paper presented at 22nd IEEE Conference on Computer Vision and Pattern Recognition, Miami, FL, USA, June 20–25, 2009.
33. Tax DM, Duin RP. Support vector data description. *Mach Learn* 2004;54:45–66. DOI: [10.1023/B:MACH.0000008084.60811.49](https://doi.org/10.1023/B:MACH.0000008084.60811.49).
34. He K, Fan H, Wu Y, Xie S, Girshick R. Momentum contrast for unsupervised visual representation learning. November 14, 2019 (<https://arxiv.org/abs/1911.05722v2>). Preprint.
35. Kux L. Clinical and patient decision support software; draft guidance for industry and Food and Drug Administration staff. Washington, DC: U.S. Food and Drug Administration, 2017.
36. Muelly MC, Peng L. Spotting brain bleeding after sparse training. *Nat Biomed Eng* 2019;3:161–162. DOI: [10.1038/s41551-019-0368-5](https://doi.org/10.1038/s41551-019-0368-5).