

Original Paper

A 3D and Explainable Artificial Intelligence Model for Evaluation of Chronic Otitis Media Based on Temporal Bone Computed Tomography: Model Development, Validation, and Clinical Application

Binjun Chen^{1,2*}, MD; Yike Li^{3*}, MD, PhD; Yu Sun^{4*}, MD, PhD; Haojie Sun^{1,2}, MD; Yanmei Wang^{1,2}, MD, PhD; Jihan Lyu^{1,2}, MD; Jiajie Guo⁵, PhD; Shunxing Bao⁶, PhD; Yushu Cheng⁷, MD; Xun Niu⁴, MD; Lian Yang⁸, MD; Jianghong Xu^{1,2}, MD, PhD; Juanmei Yang^{1,2}, MD, PhD; Yibo Huang^{1,2}, MD, PhD; Fanglu Chi^{1,2}, MD, PhD; Bo Liang^{8*}, MD; Dongdong Ren^{1,2*}, MD, PhD

¹ENT Institute and Department of Otorhinolaryngology, Eye & ENT Hospital, Fudan University, Shanghai, China

²NHC Key Laboratory of Hearing Medicine Research, Eye & ENT Hospital, Fudan University, Shanghai, China

³Department of Otolaryngology—Head and Neck Surgery, Vanderbilt University Medical Center, Nashville, TN, United States

⁴Department of Otorhinolaryngology, Union Hospital, Tongji Medical College, Huazhong University of Science and Technology, Wuhan, China

⁵State Key Laboratory of Digital Manufacturing Equipment and Technology, School of Mechanical Science and Engineering, Huazhong University of Science and Technology, Wuhan, China

⁶Department of Electrical and Computer Engineering, Vanderbilt University, Nashville, TN, United States

⁷Department of Radiology, Eye & ENT Hospital, Fudan University, Shanghai, China

⁸Department of Radiology, Union Hospital, Tongji Medical College, Huazhong University of Science and Technology, Wuhan, China

*these authors contributed equally

Corresponding Author:

Yike Li, MD, PhD

Department of Otolaryngology—Head and Neck Surgery

Vanderbilt University Medical Center

1215 Medical Center Drive

Nashville, TN, 37232

United States

Phone: 1 6153438146

Email: yike.li.1@vumc.org

Abstract

Background: Temporal bone computed tomography (CT) helps diagnose chronic otitis media (COM). However, its interpretation requires training and expertise. Artificial intelligence (AI) can help clinicians evaluate COM through CT scans, but existing models lack transparency and may not fully leverage multidimensional diagnostic information.

Objective: We aimed to develop an explainable AI system based on 3D convolutional neural networks (CNNs) for automatic CT-based evaluation of COM.

Methods: Temporal bone CT scans were retrospectively obtained from patients operated for COM between December 2015 and July 2021 at 2 independent institutes. A region of interest encompassing the middle ear was automatically segmented, and 3D CNNs were subsequently trained to identify pathological ears and cholesteatoma. An ablation study was performed to refine model architecture. Benchmark tests were conducted against a baseline 2D model and 7 clinical experts. Model performance was measured through cross-validation and external validation. Heat maps, generated using Gradient-Weighted Class Activation Mapping, were used to highlight critical decision-making regions. Finally, the AI system was assessed with a prospective cohort to aid clinicians in preoperative COM assessment.

Results: Internal and external data sets contained 1661 and 108 patients (3153 and 211 eligible ears), respectively. The 3D model exhibited decent performance with mean areas under the receiver operating characteristic curves of 0.96 (SD 0.01) and 0.93 (SD 0.01), and mean accuracies of 0.878 (SD 0.017) and 0.843 (SD 0.015), respectively, for detecting pathological ears on the 2 data sets. Similar outcomes were observed for cholesteatoma identification (mean area under the receiver operating

characteristic curve 0.85, SD 0.03 and 0.83, SD 0.05; mean accuracies 0.783, SD 0.04 and 0.813, SD 0.033, respectively). The proposed 3D model achieved a commendable balance between performance and network size relative to alternative models. It significantly outperformed the 2D approach in detecting COM ($P \leq .05$) and exhibited a substantial gain in identifying cholesteatoma ($P < .001$). The model also demonstrated superior diagnostic capabilities over resident fellows and the attending otologist ($P < .05$), rivaling all senior clinicians in both tasks. The generated heat maps properly highlighted the middle ear and mastoid regions, aligning with human knowledge in interpreting temporal bone CT. The resulting AI system achieved an accuracy of 81.8% in generating preoperative diagnoses for 121 patients and contributed to clinical decision-making in 90.1% cases.

Conclusions: We present a 3D CNN model trained to detect pathological changes and identify cholesteatoma via temporal bone CT scans. In both tasks, this model significantly outperforms the baseline 2D approach, achieving levels comparable with or surpassing those of human experts. The model also exhibits decent generalizability and enhanced comprehensibility. This AI system facilitates automatic COM assessment and shows promising viability in real-world clinical settings. These findings underscore AI's potential as a valuable aid for clinicians in COM evaluation.

Trial Registration: Chinese Clinical Trial Registry ChiCTR2000036300; <https://www.chictr.org.cn/showprojEN.html?proj=58685>

(*J Med Internet Res* 2024;26:e51706) doi: [10.2196/51706](https://doi.org/10.2196/51706)

KEYWORDS

artificial intelligence; cholesteatoma; deep learning; otitis media; tomography, x-ray computed; machine learning; mastoidectomy; convolutional neural networks; temporal bone

Introduction

Chronic otitis media (COM) represents a recurrent inflammatory condition inside the tympanic cavity [1]. COM encompasses various forms, including chronic suppurative otitis media (CSOM) and cholesteatoma, each with unique histological characteristics. CSOM involves the accumulation and discharge of purulent fluid, affecting an estimated 330 million people worldwide, with approximately half experiencing hearing loss [2]. Cholesteatoma is characterized by the buildup of keratinized squamous epithelium, which has the potential to erode auditory structures and exhibits a notable tendency for relapse. Accurate identification and differentiation of COM types are crucial for effective disease management and surgical planning [3]. Mastoidectomy, which involves the removal of part of the temporal bone, is the conventional surgical approach for COM. However, less invasive techniques such as endoscopic tympanoplasty are gaining favor for treating CSOM and other noncholesteatoma conditions due to their potential for reduced structural damage and faster recovery [4-9].

Temporal bone computed tomography (CT) is vital for assessing COM and aiding in surgical planning, especially when initial otoscopic examinations have restricted views and yield inconclusive findings [10]. Offering a cost-effective alternative to magnetic resonance imaging (MRI), CT is instrumental in distinguishing cholesteatoma from CSOM by detecting osseous erosion in the tympanum. Although studies have shown that clinicians are capable of diagnosing COM based on CT alone [11-17], distinguishing between COM subtypes poses greater challenges to the human eye. Moreover, interpreting temporal bone CT scans requires specialized training and experience, which may not be universally available across otolaryngologists.

Artificial intelligence (AI) is making remarkable advancements in health care. Deep learning (DL) models, particularly convolutional neural networks (CNNs), have demonstrated enhanced efficiency and reduced errors in disease diagnoses and prediction of clinical outcomes [18-21]. While a few recent

papers have reported CNN models in evaluating COM with accuracy scores ranging from 0.77 to 0.85, these studies primarily relied on otoscopic or single-layer CT scans [22,23]. These 2D representations may not be optimal for revealing pathological changes in concealed or peripheral anatomical structures, such as the attic space and the mastoid air cells. In addition, the inherent "black box" nature of DL models, where decision-making strategies are challenging to understand, has been a common criticism [24,25]. This lack of comprehensibility hinders the widespread adoption of AI models in clinical practice.

In light of these challenges, this study aimed to create an explainable, 3D CNN model for the automatic interpretation of temporal bone CT scans. The model was designed to pinpoint the region of interest (ROI) and identify pathological and cholesteatomatous conditions in a 3D fashion. Comprehensive benchmarks against baseline methods and human experts on distinct data sets were conducted to demonstrate the robustness and generalizability of this model. In addition, heat map generation was used to highlight potential pathological changes in CT scans and elucidate the model's rationale for making predictions. These features were integrated into an AI system for the automatic, end-to-end evaluation of COM, which was subsequently assessed in clinical settings. The overarching goal of this system is to support clinicians in making informed decisions for common otologic conditions, thereby enhancing efficiency, reliability, and transparency.

Methods

Ethical Considerations

This study was conducted in accordance with the principles of the Declaration of Helsinki. Ethical approval was granted by the institutional review boards at Vanderbilt University Medical Center (191804) and the Eye, Ear, Nose and Throat (EENT) Hospital of Fudan University (2019076). Informed consent was waived as all data were de-identified. The observational study, which aimed to assess the model's viability in aiding

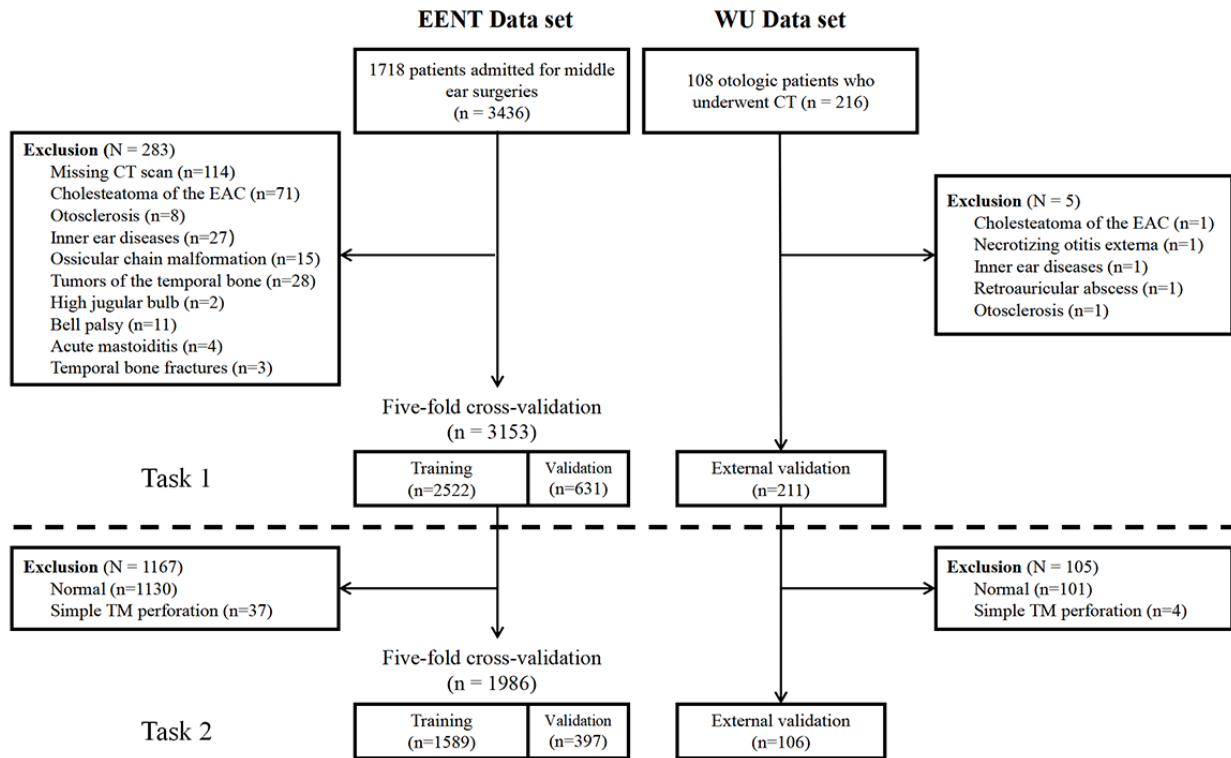
preoperative assessment, was registered with the Chinese Clinical Trial Register (ChiCTR: 2000036300). No compensation was provided to any study participants.

Participants

Data were retrospectively obtained from patients admitted for middle ear surgeries from December 2015 to July 2021 at EENT

Hospital. Patients diagnosed with acute otitis media, any inner or external ear diseases, or those with missing temporal bone CT scan were excluded, resulting in 1661 patients eligible for model development. An extra data set containing 108 patients with COM was collected from Wuhan Union (WU) Hospital for external validation (Figure 1).

Figure 1. Flowchart of data retrieval. CT: computed tomography; EAC: external auditory canal; EENT: Eye, Ear, Nose, and Throat Hospital of Fudan University; TM: tympanic membrane; WU: Wuhan Union Hospital.



Temporal Bone CT Scans

As part of the routine preoperative assessment, each patient underwent at least 1 temporal bone CT, conducted from the lower margin of the external auditory meatus to the top margin of the petrous bone using a SOMATOM Sensation 10 CT scanner (Siemens Inc) at the EENT Hospital. The scanning parameters were as follows: matrix (512 × 512), field of view (220 mm × 220 mm), tube voltage (140 kV), tube current (100 mAs), section thickness (0.6-0.75 mm), window width (4000 HU), and window level (700 HU). CT scans from the WU Hospital were obtained using a SOMATOM Plus 4 model (Siemens Inc) with different settings for field of view (100 mm), voltage (120 kV), and thickness (0.75 mm). All images were saved in the DICOM format.

Label Assignment

All eligible ears were treated as independent cases and assigned ground truth labels based on their diagnoses (Table 1). Each

label was verified according to intraoperative findings and pathology reports for operated ears and using a combination of history, ear examination, audiogram results, and imaging findings for unoperated ears. In cases of unoperated ears, a “normal” label was assigned when there was an absence of ear symptoms, hearing loss, or signs of inflammation. A diagnosis of CSOM was assigned when chronic purulent discharge, conductive hearing loss, and the presence of a perforated tympanic membrane or soft tissue shadow in the tympanic cavity were observed. Cholesteatoma was considered if keratin debris was identified, or if there were signs of osseous damage along with retraction or perforation of the pars flaccida [22]. Two otolaryngology residents with full access to patients’ medical records independently reviewed these labels as unblinded annotators. Any discrepancies were addressed with senior specialists until a consensus was reached. All data were deidentified and stored on password-protected computers.

Table 1. Summary of patient characteristics and label assignment.

Characteristics	EENT ^a data set (N=1661; number of ears=3153)	WU ^b data set (N=108; N=211)
Patient age (years), mean (SD)	41.1 (16.6)	39.8 (14.0)
Patient sex, n (%)		
Male	832 (50.1)	49 (45.4)
Female	829 (49.9)	59 (54.6)
Diagnosis per ear, n (%)		
Normal	1130 (35.8)	101 (47.9)
Cholesteatoma	728 (23.1)	30 (14.2)
CSOM ^c	1011 (32.1)	69 (32.7)
Tympanosclerosis	142 (4.5)	2 (0.1)
Cholesterol granuloma	72 (2.3)	1 (0.05)
OME ^d	41 (1.3)	7 (3.3)
Adhesive otitis media	29 (0.1)	1 (0.05)
Task 1 labels, n (%)		
Normal	1130 (35.8)	101 (47.9)
Pathological	2023 (64.2)	110 (52.1)
Task 2 labels, n (%)		
Cholesteatoma	728 (36.7)	28 (26.4)
Noncholesteatoma	1258 (63.3)	78 (73.6)

^aEENT: Eye, Ear, Nose, and Throat Hospital of Fudan University.

^bWU: Wuhan Union Hospital.

^cCSOM: chronic suppurative otitis media.

^dOME: otitis media with effusion.

Model Architecture

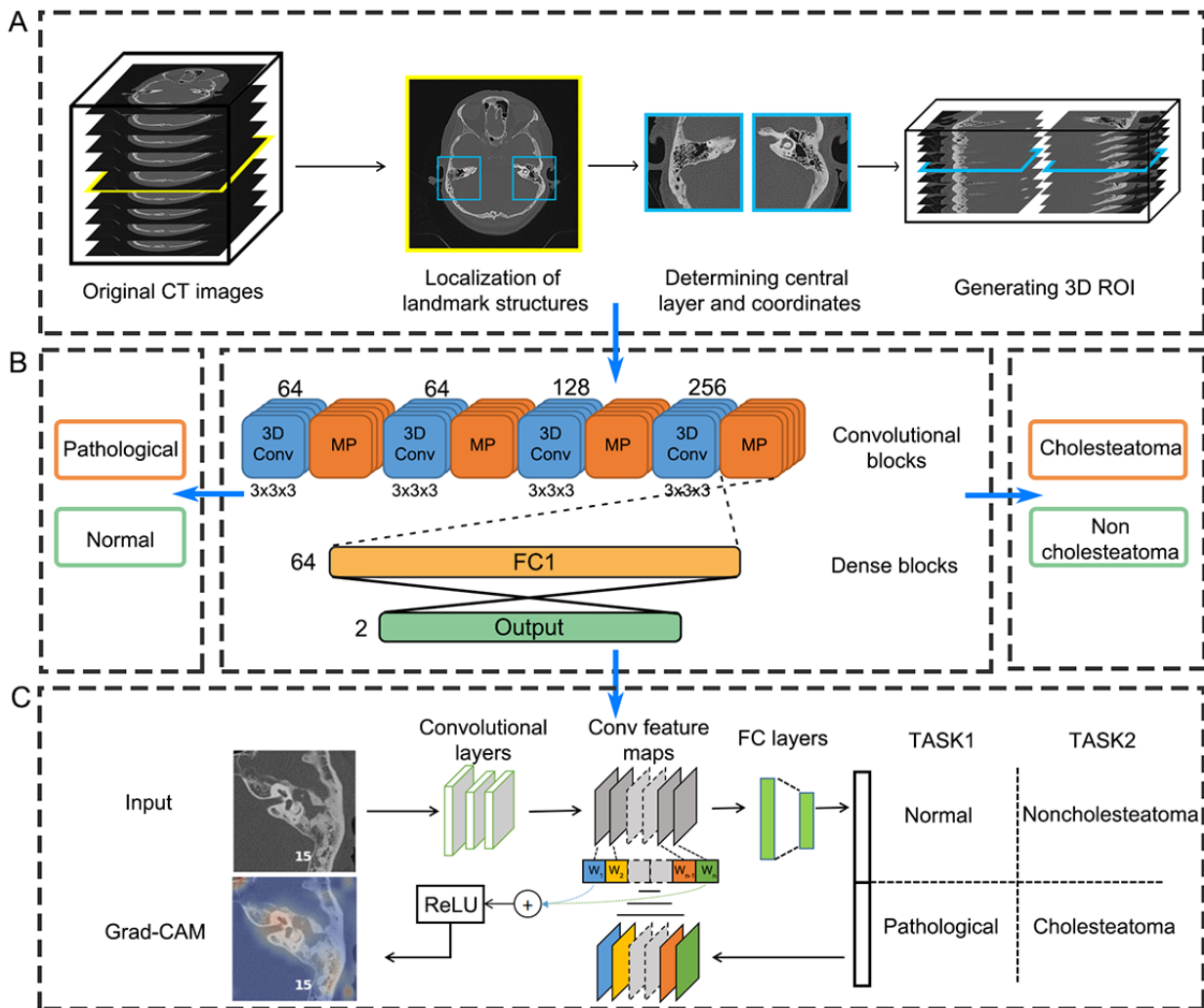
The framework consists of 2 functionally distinct units: a region proposal network for 3D segmentation of ROI, and a classification network for generating predictions. Both networks are established based on CNN models.

Region Proposal Network

This network is designed to extract the middle ear on each side from a full set of temporal bone CT scan (Figure 2A). It contains a YOLO (You Only Look Once; v5) model that is trained to

detect and locate 2 auditory structures, including the internal auditory canal and the horizontal semicircular canal, in a series of 2D axial CT scans [26]. These landmarks, positioned at or around the central level of the middle ear, possess unique graphical appearances recognizable by the object detection model. In our recent study, this model demonstrated a 100% success rate in identifying the middle ear region from temporal bone CT scans [22]. Subsequently, a 3D data matrix (150 × 150 × 32) of the ROI is extracted based on the center coordinates of these 2 structures on each side.

Figure 2. An overview of the AI framework. (A) The region proposal network used to locate landmark structures and segment the 3D ROI from the original CT scans. (B) The classification network based on a 3D convolutional neural network architecture and trained to perform 2 classification tasks. (C) The gradient heatmaps generated to highlight the critical regions for decision-making. Conv: convolution; CT: computed tomographic; FC: fully connected; MP: max pooling; ReLU: rectified linear unit; ROI: region of interest.



Classification Network

A 3D CNN model is built to interpret the extracted ROI and classify different types of conditions (Figure 2B). This model features 4 convolution blocks and 2 dense blocks (Table 2). Each convolution block consists of a 3D convolutional layer to summarize graphical features along all axes of the input image,

followed by a max-pooling layer for downsampling these features and another layer for batch normalization. These high-level features are then pooled and passed to the fully connected layers of the dense blocks, where the diagnosis is predicted based on the calculated probability of each class by a softmax function. A dropout layer is applied to prevent overfitting [27].

Table 2. Architecture of the 3D convolutional neural network model.

Block and kernel inputs	Settings
Convolution 1	
Conv3D ^a	(3,3,3,64)
MaxPooling3D ^b	(2,2,2)
BatchNormalization ^c	
Convolution 2	
Conv3D	(3,3,3,64)
MaxPooling3D	(2,2,2)
BatchNormalization	
Convolution 3	
Conv3D	(3,3,3,128)
MaxPooling3D	(2,2,2)
BatchNormalization	
Convolution 4	
Conv3D	(3,3,3,256)
MaxPooling3D	(2,2,2)
BatchNormalization	
GlobalAveragePooling3D ^d	
Dense 1	
Fully connected	64
Dropout	0.3
Output	
Fully connected	2

^aConv3D: 3D convolutional layer.

^bMaxPooling3D: 3D max pooling layer.

^cBatchNormalization: batch normalization layer.

^dGlobalAveragePooling3D: layer performing global average pooling for 3D data.

Model Training and Testing

Task 1—Detection of COM

The first classification model was trained in a binary task distinguishing between normal and pathological ears in all cases (n=3153). The training and testing procedures involved 5-fold cross-validation on the internal (EENT) data set. Specifically, the data set was evenly partitioned into 5 nonoverlapping subsets in a random, stratified fashion. In each iteration, 1 subset was reserved for testing (n=631), while the remaining 4 were used for training (n=2522). Model performance metrics were averaged over 5 iterations of this process. During each training session, a random 20% of training images (n=504) were allocated for validation. Training was set for 1000 epochs with an initial learning rate of 0.0001, and the Adam optimizer was used to dynamically adjust the algorithm's learning capability and minimize errors [28]. Early termination was implemented if no further decrease in validation loss was observed for a consecutive 10 epochs. These hyperparameters were determined based on the resultant model performance and training efficiency

shown in a preliminary study. The trained model was also evaluated on the external data set (n=211) in each round.

Task 2—Identification of Cholesteatoma

The second classification model was trained to specifically identify cholesteatoma on selected CT scans that displayed signs of inflammation in the middle ears. This task was designed to simulate a common clinical scenario where clinicians need to differentiate cholesteatoma from other types of COM in patients with positive imaging findings. The aim was to provide a preoperative assessment of the risk of cholesteatoma, assisting clinicians in surgical planning [3,29]. For this task, a subset of CT scans with visible soft tissue density or increased opacification in the middle ear or mastoid was selected from both the internal (n=1986) and external sets (n=106). The remaining methods, including extraction of ROI, network architecture, and the training and testing procedures, were consistent with those used in the first task.

Ablation Study

To refine model selection and gain a better understanding of the network's behavior, an ablation study was performed to compare the proposed classification network with 3 alternative models, each incorporating modifications to certain features. Specifically, the number of convolutional blocks was decreased and increased by 1 in alternative model 1 and model 2, respectively, and a different size of filter was applied in model 3 (Tables S1-S3 in [Multimedia Appendix 1](#)). To ensure adequate statistical power for detecting differences across models, experiments were conducted on the main data set using the same methodology as outlined in the preceding sections.

Benchmarking Against the 2D Approach

To investigate whether the use of 3D CT scans may enhance diagnostic performance, a benchmark study was designed to compare the proposed system with a baseline model using 2D images. This baseline model, previously established by our team, uses transfer learning on a pretrained Inception-V3 (Google LLC) model [22]. In this study, the base model of Inception-V3 was retained, and the final classification layer was customized with a binary output. Training and validation were conducted in the same manner as the 3D model, except that only a single CT scan at the central layer of the ROI was used as the input for the 2D model. All image-preprocessing techniques and hyperparameter settings remained consistent with those outlined in the previous study [22].

Benchmarking Against Human Experts

Another benchmark test was performed against human experts to provide an additional unbiased evaluation of the proposed system. Seven human specialists with a broad range of qualifications were recruited to perform both tasks based on the same image data. The participants included 2 senior otologists, each with 12 years of clinical experience, 1 senior head and neck radiologist with 21 years of experience, 1 attending otologist with 7 years of experience, and 3 otolaryngology residents with 3, 3, and 2 years of experience, respectively. Each expert was provided only with the CT scans and instructed to make a task-specific diagnosis to each ear (task 1: normal or pathological; task 2: cholesteatoma or noncholesteatoma). The test data for clinicians comprised a random selection of 244 ears from the EENT set and all eligible ears from the WU set. To assess intrarater reliability, a random replication of 10% of test cases ($n=48$) was mixed with these data. All test cases ($N=502$) had not been previously seen by any experts. They were anonymized, shuffled, and stored on a password-protected computer along with spreadsheets to record each expert's diagnoses for these cases.

Generation of Heat Maps

Gradient-Weighted Class Activation Mapping was used to visualize model's rationale for decision-making ([Figure 2C](#)).

In essence, this approach leverages the gradients of the target class flowing into the final convolutional layer to produce a coarse localization heat map, highlighting the critical regions in the image [30]. In this study, heat maps were generated in a 3D fashion and rescaled to match the original images using TensorFlow 2.11 in Python 3.91 (Python Core Team) [31].

Clinical Applications

The validated model was integrated into a Python program, enabling the automated assessment of COM from raw CT inputs to the generation of explainable diagnoses in an end-to-end fashion (see the section "Data Availability Statements" and [Multimedia Appendix 2](#)). To evaluate its viability in assisting otologists in clinical settings, this system was used with a prospective cohort of patients undergoing middle ear surgeries at EENT hospital from November 2023 to January 2024 in a single-arm observational study. Preoperative model predictions, along with routine assessments, were provided to 2 senior otologists, who were given autonomy to determine surgical strategies based on their discretion. Surgeons were surveyed regarding the use of model-generated information in their decision-making processes for these cases. Model predictions were used to analyze the selection of surgical approaches and to measure model performance against pathological findings. Hearing gain was assessed by comparing the air conduction threshold at 2 weeks postoperatively with the baseline.

Statistical Analysis

Descriptive statistics were applied as appropriate. The overall predictability of a model was evaluated by the area under the receiver operating characteristic (AUROC) curve. The optimal cutoff threshold on the curve was determined at the point with minimal distance to the upper left corner on the validation set and subsequently applied to the test set. The numbers of correctly and incorrectly classified cases were displayed in a confusion matrix, and these were used to calculate the performance metrics, including accuracy, recall, specificity, precision, and F_1 -score. These metrics offer comprehensive insights into the model's performance, covering overall correctness in identifying both positives and negatives (accuracy), sensitivity in detecting positive cases (recall), capability in ruling in patients (specificity), propensity for preventing false alarms (precision), and effectiveness in identifying positive cases while minimizing false positives and false negatives (F_1 -score). They were derived as shown in [Textbox 1](#). Results are averaged over 5 iterations of cross-validation or external validation and presented as mean (SD). Intrarater consistency was evaluated using Cohen kappa. Significance was determined through pairwise 2-tailed t test for difference in performance between models and via 1-way analysis of variance between the proposed model and human experts. The alpha level was set at .05. Statistical analyses were conducted using Python 3.91 [31].

Textbox 1. The calculation of performance metrics.

$$\text{Accuracy} = (\text{True positive} + \text{True negative}) / \text{Total sample size}$$

$$\text{Recall} = \text{True positive} / (\text{True positive} + \text{False negative})$$

$$\text{Specificity} = \text{True negative} / (\text{True negative} + \text{False positive})$$

$$\text{Precision} = \text{True positive} / (\text{True positive} + \text{False positive})$$

$$F_1\text{-score} = 2 \times \text{True positive} / (2 \times \text{True positive} + \text{False positive} + \text{False negative})$$

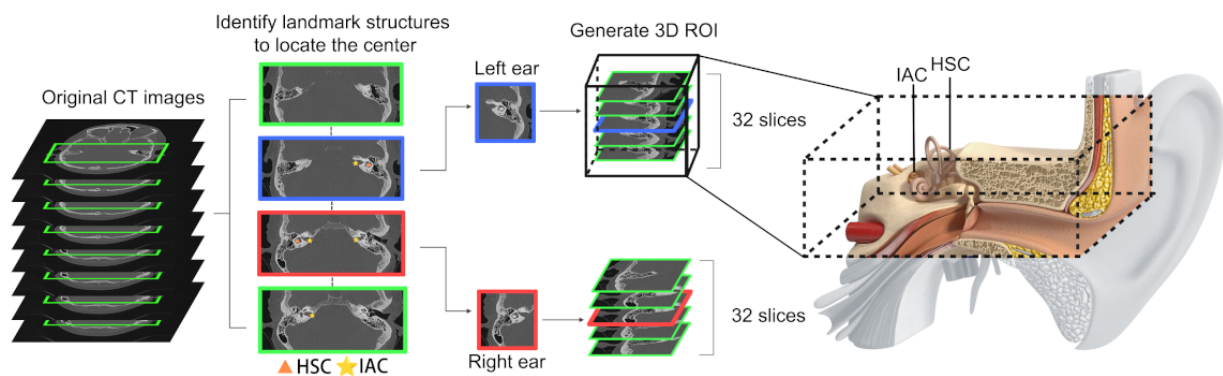
Results

ROI Extraction

The region proposal network successfully extracted the 3D ROI containing the critical anatomies on each side, including the

tympanic cavity and sinus tympani (Figure 3). This has been confirmed by manual inspection of the generated images in all cases from both data sets.

Figure 3. Generation of the 3D ROI. The region proposal network identifies landmark structures in each of the full-sized sequential CT slices and determines the center of the middle ear on each side. A 3D image comprising 32 stacks of axial slices in 150×150 pixels is subsequently segmented. This ROI encompasses an extensive range of critical anatomies within the temporal bone for the evaluation of COM. CT: computed tomographic; HSC: horizontal semicircular canal; IAC: internal auditory canal; ROI: region of interest.



Task 1

Our model exhibited decent performance in identifying pathological changes in the middle ear, achieving a mean accuracy of 87.8%, recall of 85.3%, specificity of 91.3%, and

precision of 93.3% on the internal data set (Table 3). It also demonstrated a near-perfect AUROC score of 0.96. These performance metrics remained generally consistent on the external data set, with a comparable AUROC score of 0.93, indicating reasonable generalizability (Figure 4).

Table 3. Performance of the baseline 2D and the proposed 3D models.

Task and model	Size (MB)	Data set	Accuracy, mean (SD)	Recall, mean (SD)	Specificity, mean (SD)	Precision, mean (SD)	F ₁ -score, mean (SD)	AUROC ^a , mean (SD)	P value
1									
3D	14.2	EENT ^b	0.878 (0.017)	0.853 (0.032)	0.913 (0.067)	0.933 (0.045)	0.89 (0.012)	0.00959 (0.00011)	.003
2D	274	EENT	0.861 (0.019)	0.845 (0.028)	0.883 (0.052)	0.909 (0.036)	0.875 (0.016)	0.00939 (0.00013)	N/A ^c
3D	14.2	WU ^d	0.843 (0.015)	0.756 (0.047)	0.934 (0.021)	0.924 (0.018)	0.83 (0.022)	0.00933 (0.0001)	.05
2D	274	WU	0.821 (0.023)	0.744 (0.078)	0.901 (0.046)	0.891 (0.039)	0.808 (0.036)	0.00918 (0.00012)	N/A
2									
3D	14.2	EENT	0.783 (0.04)	0.808 (0.025)	0.77 (0.054)	0.652 (0.06)	0.721 (0.042)	0.00853 (0.0003)	<.001
2D	274	EENT	0.67 (0.037)	0.716 (0.144)	0.646 (0.119)	0.523 (0.044)	0.596 (0.036)	0.00744 (0.00025)	N/A
3D	14.2	WU	0.812 (0.033)	0.614 (0.085)	0.878 (0.031)	0.626 (0.078)	0.618 (0.069)	0.00826 (0.00055)	<.001
2D	274	WU	0.676 (0.103)	0.479 (0.224)	0.741 (0.185)	0.41 (0.086)	0.411 (0.096)	0.00714 (0.00049)	N/A

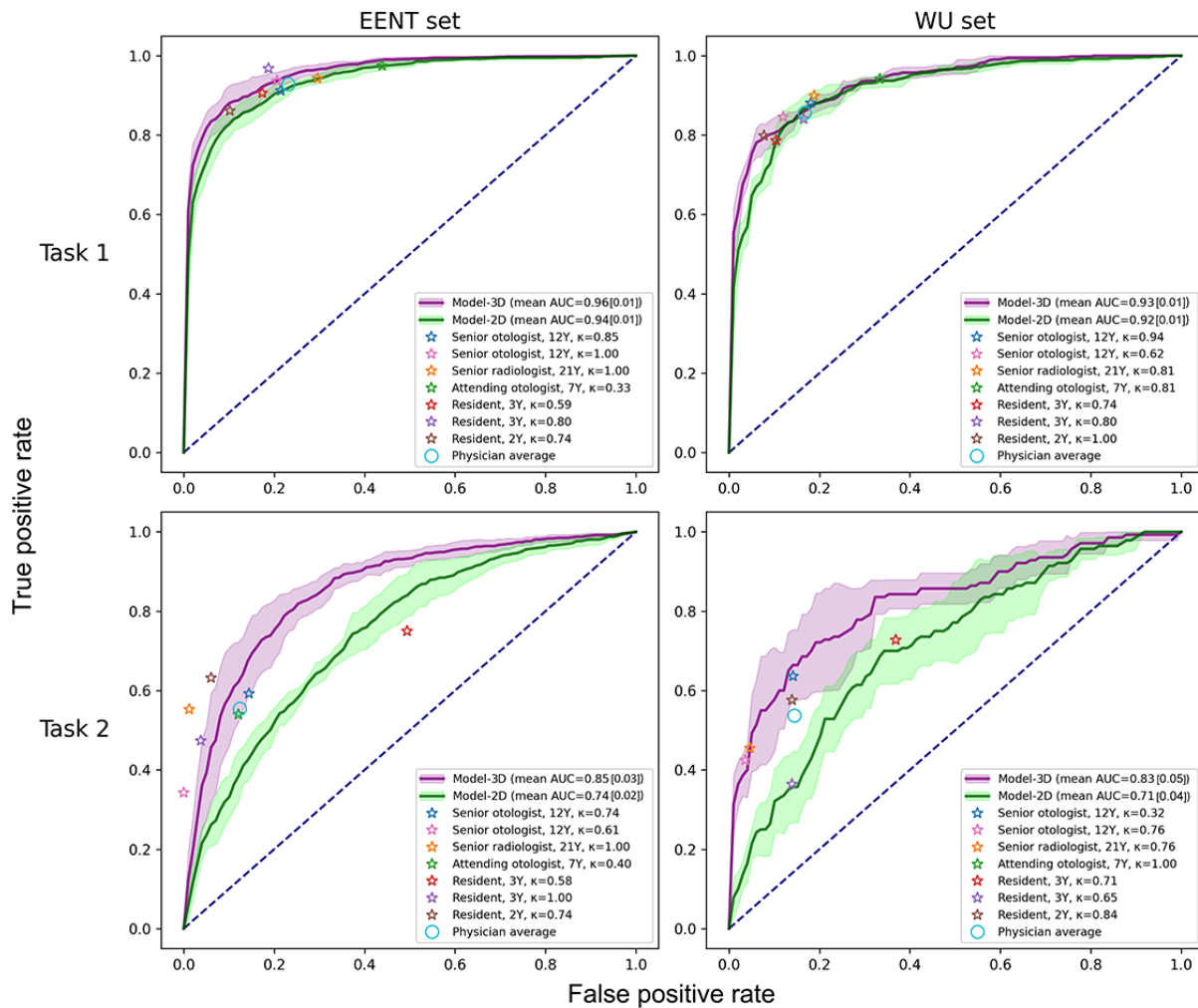
^aAUROC: area under the receiver operating characteristic curve.

^bEENT: Eye, Ear, Nose, and Throat Hospital of Fudan University.

^cN/A: not applicable.

^dWU: Wuhan Union Hospital.

Figure 4. Receiver operating characteristic plots for the benchmark tests. The curve and the shaded area indicate the mean (1 SD) of a model, respectively. Clinical experts are marked by colored asterisks for individual performance and by an open circle for averaged performance. The dotted diagonal line represents a random classifier. AUC: area under the curve; EENT: Eye, Ear, Nose, and Throat Hospital of Fudan University; WU: Wuhan Union Hospital.



Task 2

This model also demonstrated satisfactory predictive capabilities in differentiating between cholesteatoma and noncholesteatomatous cases. On both data sets, the model managed to correctly identify whether a case involved cholesteatoma in approximately 4 out of 5 instances (with accuracies of 78.3% and 81.3%). Generalizability was further supported by the comparable AUROC scores of 0.85 and 0.83 on the internal and the external data set, respectively (Table 3).

Ablation Study

This model exhibited a reasonable balance between predictability and efficiency (Table 4). Compared with models 1 and 3, it achieved significantly better performance in both tasks ($P < .01$). In addition, despite having approximately 60% fewer parameters, the proposed model demonstrated equivalent performance to model 2 in both tasks ($P = .26$ and $.91$, respectively), indicating its enhanced computational efficiency.

Table 4. Ablation study on the 3D classification network.

Task and model	Size (MB)	Accuracy, mean (SD)	Recall, mean (SD)	Specificity, mean (SD)	Precision, mean (SD)	F ₁ -score, mean (SD)	AUROC ^a , mean (SD)	P value
1								
Proposed	14.2	0.878 (0.017)	0.853 (0.032)	0.913 (0.067)	0.933 (0.045)	0.89 (0.012)	0.00959 (0.00011)	N/A ^b
Model 1	4.0	0.858 (0.03)	0.827 (0.046)	0.901 (0.058)	0.921 (0.043)	0.87 (0.028)	0.00947 (0.00019)	<.001
Model 2	34.5	0.884 (0.014)	0.862 (0.021)	0.914 (0.041)	0.933 (0.03)	0.895 (0.012)	0.00961 (0.00009)	.26
Model 3	64.8	0.864 (0.022)	0.851 (0.062)	0.887 (0.074)	0.914 (0.053)	0.878 (0.019)	0.0095 (0.00019)	.003
2								
Proposed	14.2	0.783 (4.0)	0.808 (0.025)	0.77 (0.054)	0.652 (0.06)	0.721 (0.042)	0.00853 (0.0003)	N/A
Model 1	4.0	0.758 (0.048)	0.712 (0.118)	0.783 (0.065)	0.636 (0.064)	0.668 (0.075)	0.00817 (0.0006)	.006
Model 2	34.5	0.782 (0.036)	0.795 (0.071)	0.775 (0.074)	0.659 (0.071)	0.716 (0.032)	0.00862 (0.00031)	.91
Model 3	64.8	0.756 (0.056)	0.76 (0.059)	0.754 (0.109)	0.634 (0.088)	0.685 (0.037)	0.00826 (0.000047)	.003

^aAUROC: area under the receiver operating characteristic curve.

^bN/A: not applicable.

Benchmarks

Compared with the 2D approach, the 3D network demonstrated significantly superior performance in both tasks across data sets ($P \leq .05$). In particular, the proposed model exhibited a substantial performance gain in differentiating between cholesteatoma and noncholesteatomata, with an increase of more than 10% in all outcome metrics on both data sets (Table 3).

This model also matched or even surpassed the diagnostic capabilities of human experts in both tasks (Figure 4). It exhibited marginally superior performance compared with

human eyes in the first task ($P=.05$) and significantly outperformed them in the visually challenging task 2 ($P<.001$). Post hoc pairwise comparisons revealed that the model excelled over the attending otologist in task 1 and 2 resident fellows in task 2, rivaling all senior clinicians (Table 5). Similar results were shown across the breakdown of data sources, with a notable finding that the model outperformed a senior otologist in task 2 on the EENT subset (Table S4 in Multimedia Appendix 1). Moreover, the proposed model demonstrated perfect consistency, surpassing all human experts who exhibited higher SDs in all outcome metrics and lower scores of intrarater reliability.

Table 5. Benchmark performance against human experts.

Task and rater	Accuracy	Recall	Specificity	Precision	F ₁ -score	Kappa values	P value
1							
The 3D model, mean (SD)	0.878 (0.017)	0.853 (0.032)	0.913 (0.067)	0.933 (0.045)	0.89 (0.012)	0.01 (0.00)	N/A ^a
Expert average, mean (SD)	0.857 (0.022)	0.898 (0.042)	0.804 (0.094)	0.86 (0.05)	0.876 (0.013)	0.0082 (0.0009)	.05
Senior otologist A: 12 Y ^b	87.3%	89.8%	84.1%	87.9%	88.8%	0.75	.79
Senior otologist B: 12 Y	85.7%	89.9%	80.4%	85.6%	87.7%	0.92	.49
Senior radiologist: 21 Y	85.4%	92.4%	76.3%	83.4%	87.7%	0.87	.37
Attending otologist: 7 Y	81.1%	96.0%	61.9%	76.5%	85.2%	0.73	.002
Resident A: 3 Y	85.9%	85.5%	86.4%	89.1%	87.2%	0.71	.56
Resident B: 3 Y	87.4%	91.2%	82.5%	87.1%	89.1%	0.81	.74
Resident C: 2 Y	86.8%	83.5%	91.2%	92.4%	87.7%	0.92	.96
2							
The 3D model, mean (SD)	0.843 (0.015)	0.756 (0.047)	0.934 (0.021)	0.924 (0.018)	0.83 (0.022)	0.01 (0.00)	N/A
Expert average, mean (SD)	0.741 (0.052)	0.549 (0.123)	0.865 (0.139)	0.772 (0.135)	0.622 (0.061)	0.0072 (0.0012)	<.001
Senior otologist A: 12 Y	73.8%	36.7%	98.2%	93.0%	52.6%	0.70	.07
Senior otologist B: 12 Y	75.8%	60.6%	85.7%	73.3%	66.3%	0.47	.25
Senior radiologist: 21 Y	79.5%	52.3%	97.0%	91.9%	66.7%	0.86	.82
Attending otologist: 7 Y	74.5%	55.0%	87.0%	73.2%	62.8%	0.74	.11
Resident A: 3 Y	63.8%	74.3%	56.9%	52.9%	61.8%	0.67	<.001
Resident B: 3 Y	72.3%	44.0%	90.9%	76.2%	55.8%	0.77	.02
Resident C: 2 Y	78.8%	61.5%	89.9%	79.8%	69.4%	0.80	.96

^aN/A: not applicable.

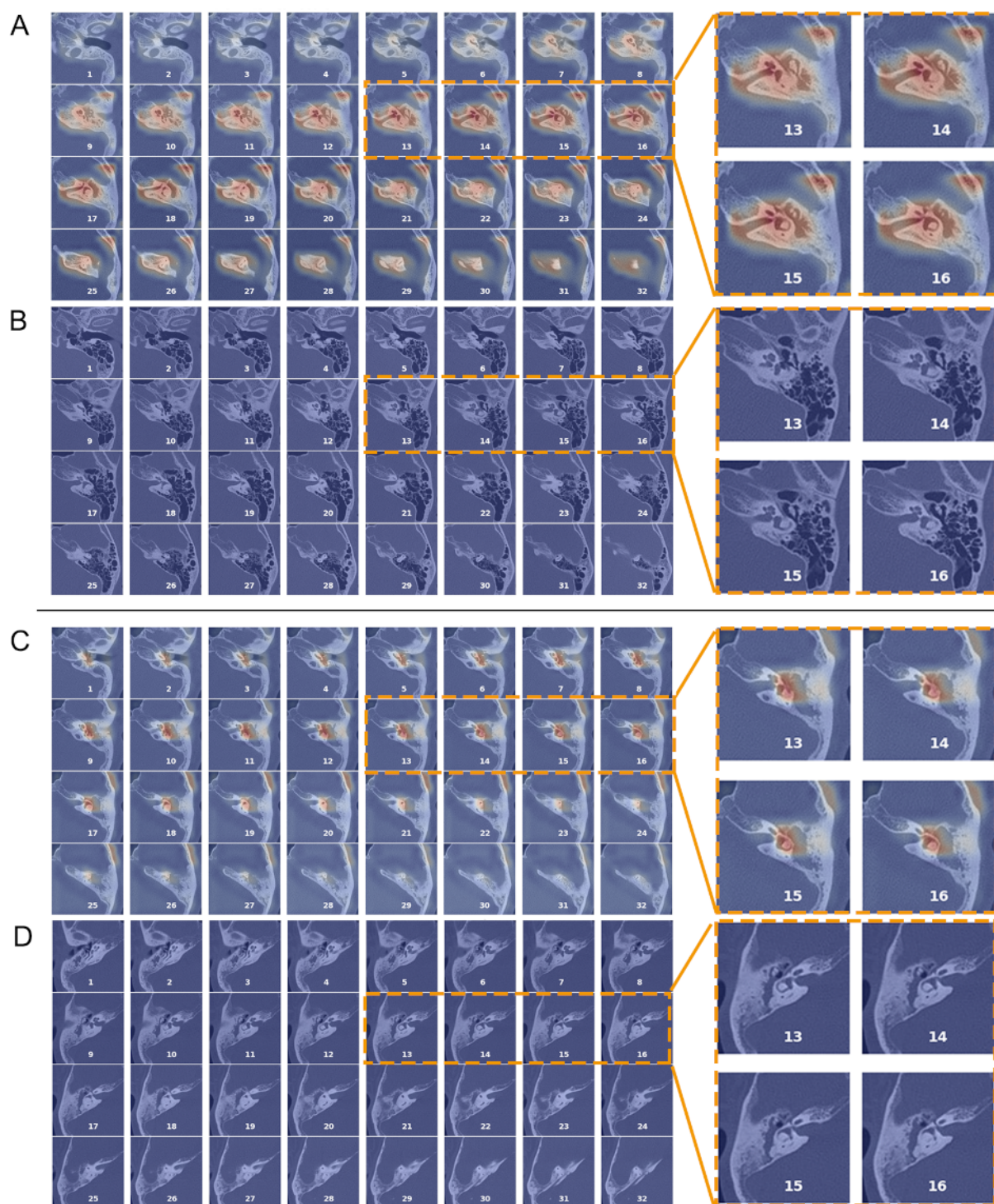
^bY: years of experience in clinical practice.

Visual Assessment of Heat Maps

Heat Maps from both models consistently highlighted the tympanic cavity and mastoid that manifested pathological findings characteristic of the target condition (Figure 5). Specifically, the first model generated a hot signal indicative of soft tissue density in an affected middle ear (Figure 5A), while the signal remained subdued in a normal ear (Figure 5B). Similarly, the second model revealed a distinct hot spot in a

cholesteatomatous ear exhibiting the classic patterns of tympanum widening and ossicular destruction [17,32,33] (Figure 5C). In contrast, a case of CSOM showing intact ossicles surrounded by soft tissue shadows in a normal-sized tympanic cavity did not exhibit a corresponding hot spot (Figure 5D). These observations reflect that the AI’s decision-making strategy aligns reasonably well with established human knowledge for both tasks.

Figure 5. Examples of heat maps. The heat maps, generated in 3D fashion, are superimposed on the original computed tomographic scans and flattened to a series of 2D images for demonstration purpose. (A-B) A pathological and a normal ear, respectively. (C-D) A cholesteatoma and a noncholesteatoma case, respectively. Area marked by hot signals indicate the presence of graphic patterns contributing to a “positive” prediction (ie, a pathological ear in task 1 and a cholesteatoma in task 2).



Clinical Use

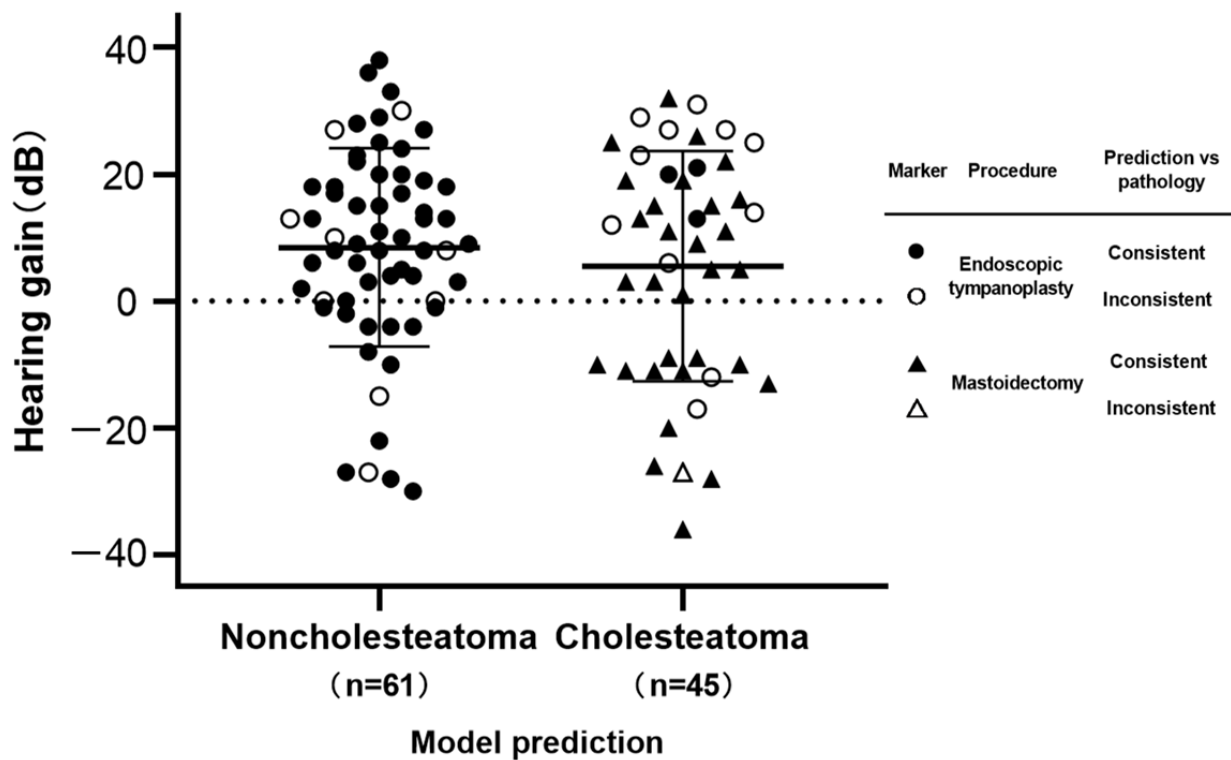
The automatic evaluation system, incorporating the validated 3D model and the heatmap visualization technique, was evaluated for its viability in aiding preoperative assessment in 121 patients with COM (mean age 46.8, SD 16.1 years, 40.5% male). This system achieved an overall accuracy of 81.8% in

distinguishing between cholesteatoma and noncholesteatoma cases. Sixty-nine ears were identified as free of cholesteatoma by the model, all of which received minimally invasive tympanoplasty under endoscopy. During the procedure, 9 ears (13.0%) revealed signs of cholesteatoma, and 5 of them required additional bone-grinding technique for complete removal of the mass. Cholesteatoma was initially predicted in 52 ears, with 37

(71.2%) of them undergoing canal-wall-down mastoidectomy. In the remaining 15 ears, the treating surgeons opted for endoscopic tympanoplasty, overriding the conventional technique for the model’s predicted diagnosis. Clinicians reported that the model predictions aligned with their initial judgment or helped with their decision-making in 90.1%

(109/121) cases. Postoperative hearing results were obtained in 87.6% (106/121) patients who maintained follow-up. Both groups of ears showed normal recovery, with a mean hearing gain of 8.5 (SD 15.6) and 5.5 (SD 18.1) dB, respectively (Figure 6).

Figure 6. Postoperative hearing gain for the operated ears with available audiometry outcomes (n=106). Data are categorized according to model predictions. Predictions that agree with the pathological results are denoted by close symbols, while open symbols indicate disagreements. Circles and triangles represent the treatment of endoscopic tympanoplasty and mastoidectomy, respectively. The error bars indicate ±1 SD from the mean.



Discussion

Principal Results

This study demonstrates the robustness and generalizability of an AI model based on 3D CNN for the detection and differential diagnosis of COM using temporal bone CT scans. This model leverages multidimensional diagnostic information from the middle ear, resulting in a significant performance improvement compared with the traditional 2D approach. The framework exhibits comparable or even superior performances to human experts in otologic tasks with clinical significance and visual challenges, especially for classifying between cholesteatoma and noncholesteatomatous cases. In addition, the novel heatmap technique allows inspection of the AI’s logic for decision-making, thereby enhancing the transparency of this model. The resulting AI system serves to automate summarization of critical radiologic findings and enables efficient evaluation of COM with minimum manual input. It provides tangible benefit in assisting otologists during preoperative assessment and results in favorable clinical outcomes that are comparable with historical results [34-37]. These findings further support the clinical viability and advantages of AI technology, which is expected to improve

efficiency, reduce errors, and facilitate precision medicine in health care in the new era of big data.

Comparison With Prior Work

A few AI models have recently been developed to classify common middle ear conditions, such as CSOM, otitis media with effusion, and cholesteatoma [38-41]. However, these models were primarily based on traditional otoscopic images, which are potentially limited by a narrow field of view and insufficient diagnostic information. Temporal bone CT scans, which are increasingly used in otologic workup by virtue of its accessibility, rich amount of anatomical information, and adequate sensitivity in revealing pathological changes, have also been explored in a limited number of studies [22,42-45]. Although these AI models demonstrated decent AUROC scores (eg, >0.9) in common classification tasks, they were all trained to generate predictions based on 2D single-layer CT scans. A potential drawback is the increased likelihood of missing small or peripheral pathological changes (eg, an attic cholesteatoma) and the resultant false negatives.

Efforts were made in this study to establish a 3D approach to take full advantage of all available anatomical information and achieve a better coverage of the tympanum and the mastoid.

Inspection of the extracted ROI suggests that all critical anatomies are visible. Results from the benchmark test indicate that the proposed 3D model outperforms the state-of-the-art 2D approach by a modest performance gain in the detection of COM and by a much larger extent in differentiating between cholesteatoma and noncholesteatoma. This finding has several implications. First, both models are generally adequate in identifying common abnormal patterns from the CT, which are graphically characterized by increased opacification or soft tissue shadows in the middle ear cavity and indicative of pathological conditions in general. This is a relatively simple visual task, during which diagnostic information obtained from a single 2D CT slice is likely sufficient for the purpose and extra findings from other layers provide only minimal contribution to the decision-making. Second, the 3D model has huge advantage over the 2D approach in differentiating cholesteatoma from other types of COM. This task is known to be more visually challenging for humans, often requiring detection of subtle osseous erosions from multiple CT slices, as quite a few pathological changes caused by cholesteatoma are peripheral or noncharacteristic [32,33]. A substantial increase in each outcome measure justifies the advantage of the current 3D model for this task. Moreover, this 3D model has only a simple network structure with a small size (14.5 MB) as opposed to a complex and large-sized 2D network (274 MB), suggesting both higher computational efficiency and performance of the 3D approach. Finally, the AUROC of 0.92-0.94 and accuracy scores of 82.1%-86.1% achieved by the 2D network in this study in detecting COM were equivalent to historical results (0.92%-86%, respectively) in our previous study [22], further indicating the reliability of these findings and potentially the intrinsic limit of using single-layer CT scan for this task. To the best of our knowledge, this is the first study showing quantitative evidence to support the advantage of a 3D CNN model in 2 common otologic tasks based on temporal bone CT scans. It also advances beyond prior retrospective research by showcasing the practicality and benefits of the model in a clinical environment.

Clinical Implications

Cholesteatoma exhibits distinct histology marked by local invasiveness and a propensity for recurrence. The imperative for successful outcomes necessitates complete removal of the mass, particularly because recurrent cholesteatoma complicates revision surgery [46]. Suspected cases often require a canal-wall-down mastoidectomy to expose the tympanum, resulting in an open cavity and a permanently altered sound conduction pathway [46]. Accumulating evidence suggests that noncholesteatoma may spare from mastoidectomy and benefit from minimally invasive procedures such as endoscopic tympanoplasty [47,48]. Therefore, the current AI system holds potential value for otologists in surgical planning. Ears with a low risk of cholesteatoma, as identified by the model, could potentially be treated by less invasive procedures that retain the integrity of canal wall, leading to reduced procedural time and enhanced recovery [6,7,9,49,50]. This clinical merit is supported by the superior benchmark performance in identifying cholesteatoma and the favorable outcomes observed in the prospective study.

While detecting COM in task 1 involves spotting any pathological patterns on CT, which may not fully capture the differences between models in diagnostic capabilities, the increased visual challenges in identifying cholesteatoma substantiate the advantages of the proposed 3D approach for this task. In this study, the 3D model outperformed junior clinicians and demonstrated equivalent or superior performance to senior experts in identifying cholesteatoma based on CT. Notably, the 3D model achieved outcomes that were on par with or better than those based on human interpretation of MRI, which, despite its higher sensitivity, is a more expensive diagnostic method [22,43,45,51-53]. These findings underscore the 3D model's potential as a reliable and cost-effective alternative, offering sufficient COM evaluation with CT alone, thereby reducing the need for the pricier MRI.

The findings from the prospective study indicate that the model is efficacious in clinical environments, especially in distinguishing cholesteatoma from noncholesteatoma. Feedback from our clinical team highlights that the system serves as a reliable and streamlined source for a second opinion. Before surgery, the treating physician can rapidly identify essential details such as the lesion's location and properties, using the model's diagnostic output, and heatmaps. Concordance between the model's predictions and the physician's initial assessment bolsters confidence in surgical planning, thereby streamlining the diagnostic and therapeutic process. In contrast, discrepancies between the model's results and the physician's judgment prompt a detailed case reassessment or team consultation, aiding in the validation of a suitable treatment plan or preparing for intraoperative modifications. This process provides timely advisory support for complex cases, encouraging meticulous evaluation by the physician, minimizing errors, and keeping the clinician's cognitive load in check without compromising their autonomy in decision-making.

It should be noted that even for seasoned otologists and radiologists, who are adept at quickly and accurately reading temporal bone CT scans, a second opinion can add an extra layer of confidence to their assessments. For novice clinicians, who may find the diagnostic process more challenging and time-intensive (Table 5), the model may offer substantial improvements in both the accuracy and the speed of diagnosing and managing COM. This is particularly beneficial for physicians in smaller medical facilities or those early in their careers. Looking ahead, the integration of this model into electronic medical systems or cloud-based servers stands to streamline the provision of immediate second opinions or enable physicians from diverse locations to upload imaging data for dependable diagnostic insights. Such technological progress is poised to advance individualized COM treatments in the big data era, boosting efficiency, reducing costs, and enhancing the quality of patient care.

Research Insights

Efforts were undertaken in this study to demystify the criticized nontransparency of DL models, characterized by intricate decision-making strategies within multilayer architectures [30,54]. The nonlinear interactions among these components can yield incomprehensible logic and untraceable predictions

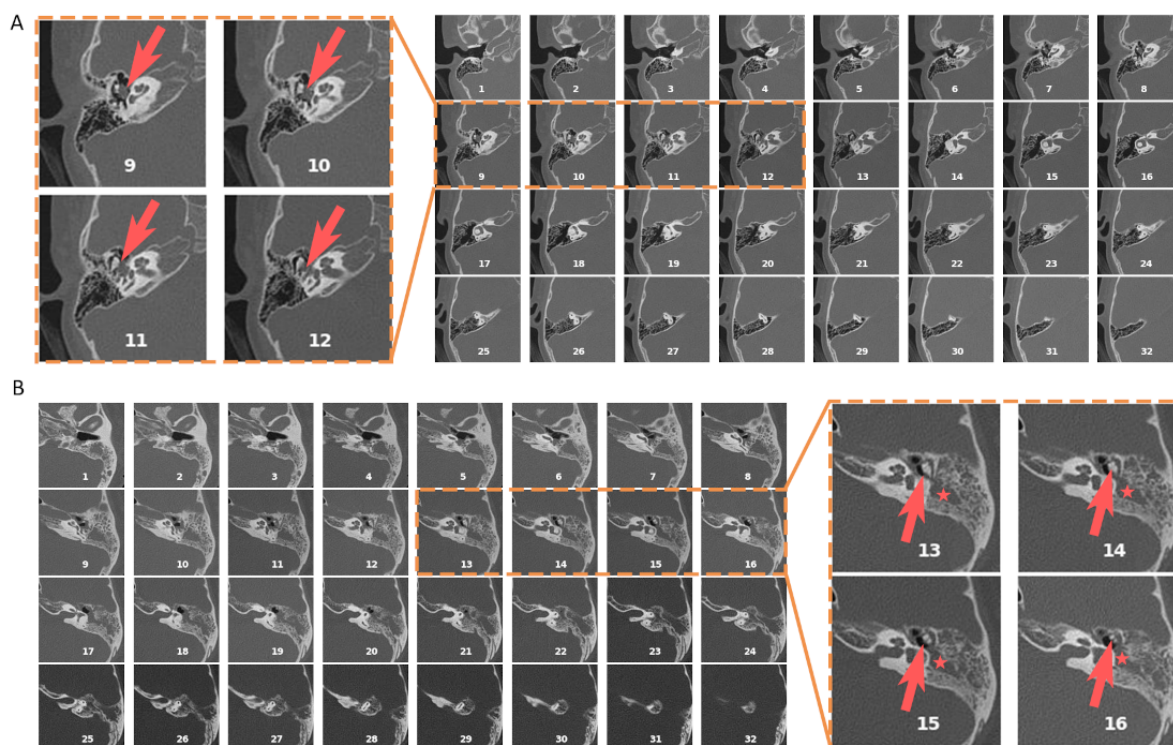
vulnerable to bias or errors, posing a significant challenge to the widespread application of AI in health care. To address this issue, heatmaps, and specifically, the Gradient-Weighted Class Activation Mapping technique, have been used as a method to inspect AI's strategy and enhance human interpretation in a parsimonious manner [55-57]. In this study, the strategy learned by our models to focus on the middle ear and mastoid regions appeared reasonable and aligned with human knowledge in interpreting CT for COM, reinforcing the reliability of this framework. These informative heatmaps can aid clinicians in understanding and validating AI predictions for specific cases, or serve as educational tools for training medical students or junior residents in reading temporal bone CT scans. Ultimately, this approach presents a viable solution for developing explainable AI models for clinical tasks.

Overfitting is a common concern with DL models, especially when data are limited or sourced from a single institute. It can lead to poor performance on new data despite promising results on the original data set. Previous DL models were trained on monocentric CT scans with participant counts ranging from 61 to 562. Lack of external validation and small sample sizes may raise concern about potential overfitting of these models [22,42,43]. Several approaches were used in this study to enhance the generalizability of our framework. First, our models underwent cross-validation on a major data set comprising more than 3000 ears, the largest sample size reported to date. Second, these models were evaluated on external data with different patient origins and image properties. Third, several machine learning methods were applied to minimize the risk of models being tuned to the random features, including early termination of training and the use of a dropout function to decrease the interdependency among network nodes [27]. Consistent performance metrics across data sets in both tasks substantiated the generalizability of this framework. Moreover, the region proposal method proved applicable to CT scans from both sources, demonstrating adaptability despite differences in CT scanner, scan settings, and image quality.

Limitations

This study has several limitations. First, although an external data set was obtained from a hospital in a different city, patients in both data sets shared a common racial background. Further validation on data collected from patients with diverse origins may be necessary to ensure the generalizability of these models. Second, the research was constrained to 2 binary classification tasks relevant to COM. Incorporating additional diagnostic tasks, such as assessing the ossicular chain's integrity and forecasting auditory outcomes, may enrich the diagnostic toolkit. Third, the models were exclusively trained to analyze CT scans, potentially not leveraging AI's full potential in COM evaluation. Comprehensive diagnostics often involve synthesizing information from patient history, clinical symptoms, ear examinations, audiological testing, otoscopy, and various imaging techniques. Overreliance on CT scans alone may introduce limitations in performance and may not always lead to conclusive diagnoses (Figure 7). Fourth, the ablation study examined a limited array of model alternatives. Despite achieving notable performance through initial model structure refinement, future endeavors should include ongoing optimization of the model architecture and detailed analysis of network component functions to optimize the trade-off between model efficacy and computational demands. In addition, this study did not place extensive emphasis on exploring common ethical issues, such as patient privacy, data security, and human autonomy, which are critical considerations in the clinical application of AI and warrant ongoing attention. Finally, this study reported initial findings from the clinical application of the AI system in a small, prospective cohort without a control group. Although the main objective was to show that the current model is ready for clinical implementation, a thorough assessment of the model's clinical benefits will be conducted in an upcoming clinical trial with a more rigorous research design.

Figure 7. Examples of misclassified cases. (A) A pathological ear showing a small-sized soft tissue density near the ossicles (arrows) with no evident sign of osseous erosion or mastoid opacification. (B) A case of cholesteatoma showing soft tissue density (asterisks) but with a visually intact ossicular chain (arrows) and a normal-sized tympanic cavity.



Future Research

Future studies will focus on leveraging novel techniques to enhance model performance and evaluate the effectiveness in larger-scale controlled trials. For example, new models will be trained to perform additional tasks, including evaluation of ossicular chain and forecasting postoperative hearing, which may enhance features of the current AI framework. A broader data set will be compiled from hospitals worldwide to assess and refine the generalizability of these models. Moreover, future models will potentially incorporate multiple sources of clinical information with a fusion layer for generating predictions, mimicking human decision-making strategies, and potentially enhancing model robustness. Ongoing efforts will also be made to refine model architectures and to address ethical issues associated with the use of AI in health care. An active learning framework may be established to integrate feedback loops, allowing clinicians to provide input to the model. This approach is expected to support ongoing model enhancement and reinforcement learning based on human feedback. In the next stage, multicenter, prospective human trials will be conducted

to assess the practical benefits of implementing these AI models in clinical contexts. The ultimate goal of this research line is to establish a robust AI system that can assist clinicians with reliability, efficiency, and transparency in the evaluation and management of ear diseases.

Conclusions

This study presents a 3D CNN model trained to detect pathological changes and identify cholesteatoma based on temporal bone CT scans. The model's performance significantly surpasses the baseline 2D approach, reaching a level comparable with or even exceeding that of human experts in both tasks. The model also exhibits decent generalizability and enhanced comprehensibility through the gradient heatmaps. The resulting AI system allows automatic assessment of COM and shows promising viability in real-world clinical settings. These findings imply the potential of AI as a valuable tool for aiding clinicians in the evaluation of COM. Future research will involve enhancing models with additional source of diagnostic information to perform various clinical tasks and evaluating the benefits of AI models in large-scale controlled trials.

Acknowledgments

This study is supported by the National Natural Science Foundation of China (grants 81970889 to Fanglu Chi and 82271166, 81970880, and 81771017 to Dongdong Ren, U22A20249, 52188102, 52027806 to Jiajie Guo); Natural Science Foundation of Shanghai (grant 22ZR1410100 to Dongdong Ren); the "Zhuo-Xue Plan" of Fudan University (Dongdong Ren); Heng-Jie special technical support plan (Dongdong Ren); and the Shanghai Outstanding Young Medical Talent Program (Dongdong Ren). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Data Availability

The complete source code and trained models are available in [Multimedia Appendices 1](#) and [2](#) as well as in a public repository, which can be accessed in GitHub [58]. The automatic evaluation system is available for individual use with a detailed instruction manual and a walk-through tutorial. The data sets generated during this study may be obtained from the corresponding authors upon reasonable request.

Authors' Contributions

YL conceptualized and designed the study, reviewed and analyzed the data, performed computer programming, developed and evaluated the AI models, wrote and edited the manuscript, prepared the figures, and supervised the entire project; BC retrieved, validated and analyzed the data, drafted the manuscript, and prepared the figures; YS provided data resources and validated the data. HS, YW, and JL retrieved data and evaluated the models. XN and LY retrieved and validated the data. JG acquired funding support, validated the model, and edited the manuscript. SB performed model validation and deployment. YC, JX, and JY evaluated the models. YH and BL provided data resources; FC provided data resources and funding support; and DR conceptualized the study, provided funding support and data resources, and edited the manuscript. All authors have reviewed, discussed, and approved the manuscript. No generative artificial intelligence tool was used during the preparation of this manuscript. BC, YL, and YS have contributed equally to this study and are credited as co-first authors. YL, BL, and DR are designated as the corresponding authors. Contact details for correspondence are as follows: Yike Li, Department of Otolaryngology-Head and Neck Surgery, Vanderbilt University Medical Center, 1215 21st Avenue South, Rm. 10410, Medical Center East, Nashville, TN 37232, USA. E-mail: yike.li.1@vumc.org. ORCID: 0000-0001-8465-130X; Bo Liang, Department of Radiology, Union Hospital, Tongji Medical College, Huazhong University of Science and Technology, Jiefang Avenue #1277, Wuhan, Hubei, 430022, China. E-mail: xiehelb@163.com. ORCID: 0000-0002-3494-4187; DongDong Ren, Department of Otorhinolaryngology, Eye, Ear, Nose and Throat Hospital, 83 Fenyang Road, Shanghai, 200031, China. E-mail: dongdongren@fudan.edu.cn. ORCID: 0000-0002-2889-9375.

Conflicts of Interest

YL served as an associate editor for the Journal of Medical Internet Research at the time of manuscript submission. YL has abstained from participating in any peer-reviewed or editorial decision-making processes related to this article.

Multimedia Appendix 1

Additional tables.

[\[DOCX File, 24 KB-Multimedia Appendix 1\]](#)

Multimedia Appendix 2

Source code for model development and validation.

[\[ZIP File \(Zip Archive\), 69 KB-Multimedia Appendix 2\]](#)

References

1. Schilder AGM, Chonmaitree T, Cripps AW, Rosenfeld RM, Casselbrant ML, Haggard MP, et al. Otitis media. *Nat Rev Dis Primers*. 2016;2(1):16063. [doi: [10.1038/nrdp.2016.63](https://doi.org/10.1038/nrdp.2016.63)] [Medline: [27604644](https://pubmed.ncbi.nlm.nih.gov/27604644/)]
2. World Health Organization. Chronic Suppurative Otitis Media: Burden of Illness and Management Options. Geneva, Switzerland. World Health Organization; 2004.
3. Lustig L, Limb C, Baden R, LaSalvia M. Chronic Otitis Media, Cholesteatoma, and Mastoiditis in Adults. Waltham, MA (citirano 145 2019). UpToDate; 2018.
4. Takahashi M, Motegi M, Yamamoto K, Yamamoto Y, Kojima H. Endoscopic tympanoplasty type I using interlay technique. *J Otolaryngol Head Neck Surg*. 2022;51(1):45. [doi: [10.1186/s40463-022-00597-3](https://doi.org/10.1186/s40463-022-00597-3)] [Medline: [36397175](https://pubmed.ncbi.nlm.nih.gov/36397175/)]
5. Ohki M, Kikuchi S, Tanaka S. Endoscopic type 1 tympanoplasty in chronic otitis media: comparative study with a postauricular microscopic approach. *Otolaryngol Head Neck Surg*. 2019;161(2):315-323. [doi: [10.1177/0194599819838778](https://doi.org/10.1177/0194599819838778)]
6. Hsu Y, Kuo C, Huang T. A retrospective comparative study of endoscopic and microscopic tympanoplasty. *J Otolaryngol Head Neck Surg*. 2018;47(1):44. [doi: [10.1186/s40463-018-0289-4](https://doi.org/10.1186/s40463-018-0289-4)]
7. Yang Q, Wang B, Zhang J, Liu H, Xu M, Zhang W. Comparison of endoscopic and microscopic tympanoplasty in patients with chronic otitis media. *Eur Arch Otorhinolaryngol*. 2022;279(10):4801-4807. [doi: [10.1007/s00405-022-07273-2](https://doi.org/10.1007/s00405-022-07273-2)]
8. Tsetsos N, Vlachtsis K, Stavarakas M, Fyrmpas G. Endoscopic versus microscopic ossiculoplasty in chronic otitis media: a systematic review of the literature. *Eur Arch Otorhinolaryngol*. 2020;278(4):917-923. [doi: [10.1007/s00405-020-06182-6](https://doi.org/10.1007/s00405-020-06182-6)]
9. Tarabichi M, Ayache S, Nogueira JF, Al Qahtani M, Pothier DD. Endoscopic management of chronic otitis media and tympanoplasty. *Otolaryngol Clin North Am*. 2013;46(2):155-163. [doi: [10.1016/j.otc.2012.12.002](https://doi.org/10.1016/j.otc.2012.12.002)]

10. Watts S, Flood LM, Clifford K. A systematic approach to interpretation of computed tomography scans prior to surgery of middle ear cholesteatoma. *J Laryngol Otol.* 2000;114(4):248-253. [doi: [10.1258/0022215001905454](https://doi.org/10.1258/0022215001905454)]
11. Selwyn D, Howard J, Cuddihy P. Pre-operative prediction of cholesteatomas from radiology: retrospective cohort study of 106 cases. *J Laryngol Otol.* 2019;133(06):477-481. [doi: [10.1017/s0022215119001154](https://doi.org/10.1017/s0022215119001154)]
12. Songu M, Altay C, Onal K, Arslanoglu S, Balci MK, Ucar M, et al. Correlation of computed tomography, echo-planar diffusion-weighted magnetic resonance imaging and surgical outcomes in middle ear cholesteatoma. *Acta Otolaryngol.* 2015;135(8):776-780. [doi: [10.3109/00016489.2015.1021931](https://doi.org/10.3109/00016489.2015.1021931)]
13. Mahmutoglu AS, Celebi I, Sahinoglu S, Cakmakci E, Sozen E. Reliability of preoperative multidetector computed tomography scan in patients with chronic otitis media. *J Craniofac Surg.* 2013;24(4):1472-1476. [doi: [10.1097/scs.0b013e31829031b1](https://doi.org/10.1097/scs.0b013e31829031b1)]
14. Pandey AK, Bapuraj JR, Gupta AK, Khandelwal N. Is there a role for virtual otoscopy in the preoperative assessment of the ossicular chain in chronic suppurative otitis media? Comparison of HRCT and virtual otoscopy with surgical findings. *Eur Radiol.* 2009;19(6):1408-1416. [doi: [10.1007/s00330-008-1282-5](https://doi.org/10.1007/s00330-008-1282-5)] [Medline: [19153741](https://pubmed.ncbi.nlm.nih.gov/19153741/)]
15. Chee NW, Tan TY. The value of pre-operative high resolution CT scans in cholesteatoma surgery. *Singapore Med J.* 2001;42(4):155-159.
16. Zaman SU, Rangankar V, Muralinath K, Shah V, Pawar R. Temporal bone cholesteatoma: typical findings and evaluation of diagnostic utility on high resolution computed tomography. *Cureus.* 2022;14(3):e22730. [doi: [10.7759/cureus.22730](https://doi.org/10.7759/cureus.22730)] [Medline: [35386487](https://pubmed.ncbi.nlm.nih.gov/35386487/)]
17. Gaurano JL, Joharjy IA. Middle ear cholesteatoma: characteristic CT findings in 64 patients. *Ann Saudi Med.* 2004;24(6):442-447. [doi: [10.5144/0256-4947.2004.442](https://doi.org/10.5144/0256-4947.2004.442)] [Medline: [15646162](https://pubmed.ncbi.nlm.nih.gov/15646162/)]
18. Ardila D, Kiraly AP, Bharadwaj S, Choi B, Reicher JJ, Peng L, et al. End-to-end lung cancer screening with three-dimensional deep learning on low-dose chest computed tomography. *Nat Med.* 2019;25(6):954-961. [doi: [10.1038/s41591-019-0447-x](https://doi.org/10.1038/s41591-019-0447-x)] [Medline: [31110349](https://pubmed.ncbi.nlm.nih.gov/31110349/)]
19. Mikhael PG, Wohlwend J, Yala A, Karstens L, Xiang J, Takigami AK, et al. Sybil: a validated deep learning model to predict future lung cancer risk from a single low-dose chest computed tomography. *J Clin Oncol.* 2023;41(12):2191-2200. [doi: [10.1200/jco.22.01345](https://doi.org/10.1200/jco.22.01345)]
20. Yamashita R, Long J, Longacre T, Peng L, Berry G, Martin B, et al. Deep learning model for the prediction of microsatellite instability in colorectal cancer: a diagnostic study. *Lancet Oncol.* 2021;22(1):132-141. [doi: [10.1016/s1470-2045\(20\)30535-0](https://doi.org/10.1016/s1470-2045(20)30535-0)] [Medline: [33387492](https://pubmed.ncbi.nlm.nih.gov/33387492/)]
21. Li Y, Guo J, Yang P. Developing an image-based deep learning framework for automatic scoring of the pentagon drawing test. *J Alzheimers Dis.* 2022;85(1):129-139. [doi: [10.3233/jad-210714](https://doi.org/10.3233/jad-210714)]
22. Wang Y, Li Y, Cheng Y, He Z, Yang J, Xu J, et al. Deep learning in automated region proposal and diagnosis of chronic otitis media based on computed tomography. *Ear Hear.* 2020;41(3):669-677. [doi: [10.1097/aud.0000000000000794](https://doi.org/10.1097/aud.0000000000000794)]
23. Sundgaard JV, Harte J, Bray P, Laugesen S, Kamide Y, Tanaka C, et al. Deep metric learning for otitis media classification. *Med Image Anal.* 2021;71:102034. [doi: [10.1016/j.media.2021.102034](https://doi.org/10.1016/j.media.2021.102034)] [Medline: [33848961](https://pubmed.ncbi.nlm.nih.gov/33848961/)]
24. Watson DS, Krutzinna J, Bruce IN, Griffiths CE, McInnes IB, Barnes MR, et al. Clinical applications of machine learning algorithms: beyond the black box. *BMJ.* 2019;364:l886. [doi: [10.1136/bmj.l886](https://doi.org/10.1136/bmj.l886)]
25. Castelvechi D. Can we open the black box of AI? *Nature.* 2016;538(7623):20-23. [doi: [10.1038/538020a](https://doi.org/10.1038/538020a)] [Medline: [27708329](https://pubmed.ncbi.nlm.nih.gov/27708329/)]
26. Redmon J, Divvala S, Girshick R, Farhadi A. You only look once: unified, real-time object detection. *IEEE*; 2016. Presented at: Proceedings of the IEEE conference on computer vision and pattern recognition; 2016:779-788; Las Vegas, NV, USA. [doi: [10.1109/CVPR.2016.91](https://doi.org/10.1109/CVPR.2016.91)]
27. Srivastava N, Hinton G, Krizhevsky A, Sutskever I, Salakhutdinov R. Dropout: a simple way to prevent neural networks from overfitting. *J machine Learn Res.* 2014;15(1):1929-1958.
28. Kingma DP, Ba J. Adam: a method for stochastic optimization. *arXiv.* 13:54-29. Preprint published online Dec 2014. [FREE Full text]
29. Tseng C, Lai M, Wu C, Yuan S, Ding Y. Comparison of the efficacy of endoscopic tympanoplasty and microscopic tympanoplasty: a systematic review and meta - analysis. *Laryngoscope.* 2016;127(8):1890-1896. [doi: [10.1002/lary.26379](https://doi.org/10.1002/lary.26379)] [Medline: [27861950](https://pubmed.ncbi.nlm.nih.gov/27861950/)]
30. Selvaraju RR, Cogswell M, Das A, Vedantam R, Parikh D, Batra D. Grad-cam: visual explanations from deep networks via gradient-based localization. 2017. Presented at: Proceedings of the IEEE International Conference on Computer Vision; 2017:618-626; Venice, Italy.
31. Python Core Team. Python: a dynamic, open source programming language. Python Software Foundation. URL: <https://www.python.org/> [accessed 2024-06-01]
32. Baráth K, Huber AM, Stämpfli P, Varga Z, Kollias S. Neuroradiology of cholesteatomas. *AJNR Am J Neuroradiol.* 2011;32(2):221-229. [doi: [10.3174/ajnr.a2052](https://doi.org/10.3174/ajnr.a2052)]
33. Gulati M, Gupta S, Prakash A, Garg A, Dixit R. HRCT imaging of acquired cholesteatoma: a pictorial review. *Insights Imaging.* 2019;10(1):92. [doi: [10.1186/s13244-019-0782-y](https://doi.org/10.1186/s13244-019-0782-y)] [Medline: [31578644](https://pubmed.ncbi.nlm.nih.gov/31578644/)]

34. Daneshi A, Daneshvar A, Asghari A, Farhadi M, Mohebbi S, Mohseni M. Endoscopic versus microscopic cartilage myringoplasty in chronic otitis media. *Iran J Otorhinolaryngol*. 2020;32(112):263-269. [doi: [10.22038/IJORL.2020.44015.2453](https://doi.org/10.22038/IJORL.2020.44015.2453)]
35. Prasad SC, Melia CL, Medina M, Vincenti V, Bacciu A, Bacciu S, et al. Long-term surgical and functional outcomes of the intact canal wall technique for middle ear cholesteatoma in the paediatric population. *Acta Otorhinolaryngol Ital*. 2014;34(5):354-361.
36. Wood CB, O'Connell BP, Lowery AC, Bennett ML, Wanna GB. Hearing outcomes following type 3 tympanoplasty with stapes columella grafting in canal wall down mastoidectomy. *Ann Otol Rhinol Laryngol*. 2019;128(8):736-741. [doi: [10.1177/0003489419841400](https://doi.org/10.1177/0003489419841400)]
37. Chamoli P, Singh CV, Radia S, Shah AK. Functional and anatomical outcome of inside out technique for cholesteatoma surgery. *Am J Otolaryngol*. 2018;39(4):423-430. [doi: [10.1016/j.amjoto.2018.04.008](https://doi.org/10.1016/j.amjoto.2018.04.008)] [Medline: [29748054](https://pubmed.ncbi.nlm.nih.gov/29748054/)]
38. Wu Z, Lin Z, Li L, Pan H, Chen G, Fu Y, et al. Deep learning for classification of pediatric otitis media. *Laryngoscope*. 2020;131(7):E2344-E2351. [doi: [10.1002/lary.29302](https://doi.org/10.1002/lary.29302)]
39. Pichichero ME. Can machine learning and AI replace otoscopy for diagnosis of otitis media? *Pediatrics*. 2021;147(4):e2020049584. [doi: [10.1542/peds.2020-049584](https://doi.org/10.1542/peds.2020-049584)]
40. Tseng CC, Lim V, Jyung RW. Use of artificial intelligence for the diagnosis of cholesteatoma. *Laryngoscope Invest Otolaryngol*. 2023;8(1):201-211. [doi: [10.1002/liv.1008](https://doi.org/10.1002/liv.1008)]
41. Livingstone D, Chau J. Otoscopic diagnosis using computer vision: an automated machine learning approach. *Laryngoscope*. 2020;130(6):1408-1413. [doi: [10.1002/lary.28292](https://doi.org/10.1002/lary.28292)] [Medline: [31532858](https://pubmed.ncbi.nlm.nih.gov/31532858/)]
42. Duan B, Guo Z, Pan L, Xu Z, Chen W. Temporal bone CT-based deep learning models for differential diagnosis of primary ciliary dyskinesia related otitis media and simple otitis media with effusion. *Am J Transl Res*. 2020;14(7):4728-4735.
43. Eroğlu O, Eroğlu Y, Yıldırım M, Karlıdağ T, Çınar A, Akyiğit A, et al. Is it useful to use computerized tomography image-based artificial intelligence modelling in the differential diagnosis of chronic otitis media with and without cholesteatoma? *Am J Otolaryngol*. 2022;43(3):103395. [doi: [10.1016/j.amjoto.2022.103395](https://doi.org/10.1016/j.amjoto.2022.103395)]
44. Khosravi M, Jabbari Moghaddam Y, Esmaeili M, Keshtkar A, Jalili J, Tayefi Nasrabadi H. Classification of mastoid air cells by CT scan images using deep learning method. *J Big Data*. 2022;9(1):1-14. [doi: [10.1186/s40537-022-00596-1](https://doi.org/10.1186/s40537-022-00596-1)]
45. Wang Z, Song J, Su R, Hou M, Qi M, Zhang J, et al. Structure-aware deep learning for chronic middle ear disease. *Expert Syst Appl*. 2022;194:116519. [doi: [10.1016/j.eswa.2022.116519](https://doi.org/10.1016/j.eswa.2022.116519)]
46. Tomlin J, Chang D, McCutcheon B, Harris J. Surgical technique and recurrence in cholesteatoma: a meta-analysis. *Audiol Neurotol*. 2013;18(3):135-142. [doi: [10.1159/000346140](https://doi.org/10.1159/000346140)]
47. Lee S, Lee DY, Seo Y, Kim YH. Can endoscopic tympanoplasty be a good alternative to microscopic tympanoplasty? a systematic review and meta-analysis. *Clin Exp Otorhinolaryngol*. 2019;12(2):145-155. [doi: [10.21053/ceo.2018.01277](https://doi.org/10.21053/ceo.2018.01277)] [Medline: [30674106](https://pubmed.ncbi.nlm.nih.gov/30674106/)]
48. Trinitade A, Page JC, Dornhoffer JL. Therapeutic mastoidectomy in the management of noncholesteatomatous chronic otitis media. *Otolaryngol Head Neck Surg*. 2016;155(6):914-922. [doi: [10.1177/0194599816662438](https://doi.org/10.1177/0194599816662438)] [Medline: [27484233](https://pubmed.ncbi.nlm.nih.gov/27484233/)]
49. Wu L, Liu Q, Gao B, Huang S, Yang N. Comparison of endoscopic and microscopic management of attic cholesteatoma: a randomized controlled trial. *Am J Otolaryngol*. 2022;43(3):103378. [doi: [10.1016/j.amjoto.2022.103378](https://doi.org/10.1016/j.amjoto.2022.103378)] [Medline: [35177254](https://pubmed.ncbi.nlm.nih.gov/35177254/)]
50. Toulouie S, Block - Wheeler NR, Rivero A. Postoperative pain after endoscopic vs microscopic otologic surgery: a systematic review and meta-analysis. *Otolaryngol Head Neck Surg*. 2022;167(1):25-34. [doi: [10.1177/01945998211041946](https://doi.org/10.1177/01945998211041946)] [Medline: [34491858](https://pubmed.ncbi.nlm.nih.gov/34491858/)]
51. Profant M, Sláviková K, Kabátová Z, Slezák P, Waczulíková I. Predictive validity of MRI in detecting and following cholesteatoma. *Eur Arch Otorhinolaryngol*. 2012;269(3):757-765. [doi: [10.1007/s00405-011-1706-8](https://doi.org/10.1007/s00405-011-1706-8)]
52. Lin M, Sha Y, Sheng Y, Chen W. Accuracy of 2D blade turbo gradient- and spin-echo diffusion weighted imaging for the diagnosis of primary middle ear cholesteatoma. *Otol Neurotol*. 2022;43(6):e651-e657. [doi: [10.1097/MAO.0000000000003521](https://doi.org/10.1097/MAO.0000000000003521)] [Medline: [35261384](https://pubmed.ncbi.nlm.nih.gov/35261384/)]
53. Sharifian H, Taheri E, Borgheri P, Shakiba M, Jalali AH, Roshanfekar M, et al. Diagnostic accuracy of non - echo - planar diffusion - weighted MRI versus other MRI sequences in cholesteatoma. *J Med Imaging Radiat Oncol*. 2012;56(4):398-408. [doi: [10.1111/j.1754-9485.2012.02377.x](https://doi.org/10.1111/j.1754-9485.2012.02377.x)] [Medline: [22883647](https://pubmed.ncbi.nlm.nih.gov/22883647/)]
54. Montavon G, Lapuschkin S, Binder A, Samek W, Müller K. Explaining nonlinear classification decisions with deep Taylor decomposition. *Pattern Recognit*. 2017;65:211-222. [doi: [10.1016/j.patcog.2016.11.008](https://doi.org/10.1016/j.patcog.2016.11.008)]
55. Panwar H, Gupta PK, Siddiqui MK, Morales-Menendez R, Bhardwaj P, Singh V. A deep learning and grad-CAM based color visualization approach for fast detection of Covid-19 cases using chest X-ray and CT-scan images. *Chaos Solitons Fractals*. 2020;140:110190. [doi: [10.1016/j.chaos.2020.110190](https://doi.org/10.1016/j.chaos.2020.110190)] [Medline: [32836918](https://pubmed.ncbi.nlm.nih.gov/32836918/)]
56. Cheng C, Ho T, Lee T, Chang C, Chou C, Chen C, et al. Application of a deep learning algorithm for detection and visualization of hip fractures on plain pelvic radiographs. *Eur Radiol*. 2019;29(10):5469-5477. [doi: [10.1007/s00330-019-06167-y](https://doi.org/10.1007/s00330-019-06167-y)] [Medline: [30937588](https://pubmed.ncbi.nlm.nih.gov/30937588/)]

57. He T, Guo J, Chen N, Xu X, Wang Z, Fu K, et al. MediMLP: using grad-CAM to extract crucial variables for lung cancer postoperative complication prediction. *IEEE J Biomed Health Inform.* 2020;24(6):1762-1771. [doi: [10.1109/jbhi.2019.2949601](https://doi.org/10.1109/jbhi.2019.2949601)] [Medline: [31670685](https://pubmed.ncbi.nlm.nih.gov/31670685/)]
58. huntlylee / 3D-Otitis-Media. GitHub. URL: <https://github.com/huntlylee/3D-Otitis-Media> [accessed 2024-08-01]

Abbreviations

AI: artificial intelligence
AUROC: area under the receiver operating characteristic curve
CNN: convolutional neural network
COM: chronic otitis media
CSOM: chronic suppurative otitis media
CT: computed tomography
DL: deep learning
EENT: Eye, Ear, Nose, and Throat Hospital of Fudan University
MRI: magnetic resonance imaging
ROI: region of interest
WU: Wuhan Union Hospital
YOLO: You Only Look Once

Edited by S Ma; submitted 09.08.23; peer-reviewed by Q Chen, SA Javed, CN Hang; comments to author 02.11.23; revised version received 30.11.23; accepted 29.05.24; published 08.08.24

Please cite as:

Chen B, Li Y, Sun Y, Sun H, Wang Y, Lyu J, Guo J, Bao S, Cheng Y, Niu X, Yang L, Xu J, Yang J, Huang Y, Chi F, Liang B, Ren D
A 3D and Explainable Artificial Intelligence Model for Evaluation of Chronic Otitis Media Based on Temporal Bone Computed Tomography: Model Development, Validation, and Clinical Application
J Med Internet Res 2024;26:e51706

URL: <https://www.jmir.org/2024/1/e51706>

doi: [10.2196/51706](https://doi.org/10.2196/51706)

PMID:

©Binjun Chen, Yike Li, Yu Sun, Haojie Sun, Yanmei Wang, Jihan Lyu, Jiajie Guo, Shunxing Bao, Yushu Cheng, Xun Niu, Lian Yang, Jianghong Xu, Juanmei Yang, Yibo Huang, Fanglu Chi, Bo Liang, Dongdong Ren. Originally published in the *Journal of Medical Internet Research* (<https://www.jmir.org>), 08.08.2024. This is an open-access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work, first published in the *Journal of Medical Internet Research* (ISSN 1438-8871), is properly cited. The complete bibliographic information, a link to the original publication on <https://www.jmir.org/>, as well as this copyright and license information must be included.