051

052

053

054

055

056

057

058

059

060

061

LNTransformer: Lung Nodule Transformer for Sparse CT Segmentation

Anonymous CVPR - CVMI submission

Paper ID 32

Abstract

001 Accurate segmentation of lung nodules in computed tomography (CT) scans is challenging due to extreme class im-002 003 balance, where nodules appear sparsely among healthy tissue. We introduce a novel two-stage approach for lung 004 nodule segmentation, framing it as an anomaly detection 005 problem. The method consists of two stages: Stage 1 em-006 007 ploys a custom Detection Transformer architecture with de-008 formable attention and focal loss to generate region proposals, addressing class imbalance and localizing sparse nod-009 ules. In Stage 2, the predicted bounding boxes are refined 010 into segmentation masks using a fine-tuned variant of the 011 Segment Anything Model (SAM). To address sparsity and 012 013 enhance spatial context, a 5mm Maximum Intensity Projec-014 tion is applied to improve differentiation between nodules, bronchioles, and vascular structures. The model achieves a 015 stage-2 DiceC of 91.4%, with stage-1 yielding an F1 score 016 of 94.2%, 95.2% sensitivity, and 93.3% precision on the 017 018 LUNA16 dataset despite extreme sparsity, where only 5% 019 of slices contain a nodule, outperforming existing state-ofthe-art methods. The model was additionally validated on 020 a privately procured test dataset of 30 patients with signifi-021 cantly different characteristics, achieving a Dice coefficient 022 023 of 78.3% despite significant distribution drift, demonstrat-024 ing strong generalization to clinical variability and establishing our approach as the new state-of-the-art for lung 025 nodule segmentation. 026

027 1. Introduction

Lung cancer is a leading cause of cancer-related deaths 028 029 worldwide [3], with early detection and accurate assessment being crucial for improving outcomes. Tumor boards, com-030 031 prising oncologists, radiologists, surgeons, pathologists, 032 and other specialists, collaboratively review complex lung 033 cancer cases to determine the best treatment plan. Accurate segmentation of lung nodules in CT scans is essential to 034 provide critical information about size, location, and spread. 035 Tumor board evaluations typically involve manual segmen-036 037 tation, which slows decision-making and increases resource

demands.

Implementing a clinical decision support (CDS) system 039 for auto-segmentation of lung nodules can enhance work-040 flow efficiency, improve patient outcomes, and reduce costs 041 [26]. Despite the benefits, current models struggle with the 042 extreme class imbalance that appears in CT data, as lung 043 nodules appear infrequently among healthy tissue. Nodule 044 volumes are much smaller than the overall lung volume, 045 vary widely in size and location, and often have similar 046 shape and density to vasculature on an axial CT slice. These 047 challenges underscore the need for a customized architec-048 ture that addresses these limitations to ensure a reliable tool 049 for clinicians. 050

We present a two-stage framework for automated lung nodule segmentation tailored for tumor boards. The first stage performs region proposal to detect sparse nodules in lung CT scans, while the second stage refines these regions for precise pixel-level segmentation. Maximum Intensity Projection (MIP) is applied to enhance nodule visibility, and custom Focal Loss is used to address class imbalance. Our aim is to provide an architecture that can handle the sparsity of lung nodules effectively, making it ready for deployment in clinical applications.

2. Related Work

Thoracic Computed Tomography (CT) involves a series of 062 2D cross-sectional greyscale images that when combined, 063 form a detailed 3D representation of the patient's thorax. 064 The LUNA16 dataset, derived from the LIDC-IDRI dataset, 065 consists of 888 thoracic CT scans containing 1,186 anno-066 tated lung nodules, annotated by four radiologists with nod-067 ules larger than 3mm considered relevant. The challenge 068 lies in the sparse occurrence of nodules, only 0 to 5% of 069 slices contain a nodule [24], and in the variability of voxel 070 sizes and scan resolutions across patients. This dataset 071 serves as a critical benchmark with numerous studies train-072 ing architectures such as CNNs, 3D-CNNs, and U-Net [10]. 073 However, transformer architectures remain underexplored 074 in this domain. 075

Class imbalance often biases models toward the majority class, often leading to high accuracy but poor detection ca-077

pabilities. To mitigate this several strategies exist: oversam-078 079 pling increases the minority class but misrepresents real-080 world prevalence [20], while class weighting forces models to focus on underrepresented cases but can increase false 081 082 positives [5, 8]. Focal loss mitigates these issues by dynamically adjusting the loss based on prediction confidence, 083 down-weighting well-classified examples and emphasizing 084 hard-to-classify, often maintaining precision while increas-085 086 ing accuracy [15].

Maximum Intensity Projection (MIP) is a widely used 087 radiology technique [9] that enhances nodule visibility by 088 089 combining adjacent CT slices into a single 2D image, pro-090 jecting the highest attenuation voxel from a volume onto a 2D plane to preserve 3D spatial information [9]. MIP helps 091 distinguish nodules, which appear as blobs, from vessels, 092 elongated tube-like structures, and improves detection of 093 3–10mm nodules [11]. 094

Previous methods predominantly relied on Convolu-095 tional Neural Networks (CNNs) and variations such as 3D 096 U-Nets, MRUNet-3D, and V-Nets, which utilize hierarchi-097 cal convolutional layers to segment nodules. While CNNs 098 have proven effective at learning local features, their per-099 formance deteriorates significantly when capturing long-100 range spatial relationships, which are critical in differentiat-101 ing nodules from similarly dense structures such as bronchi-102 oles and vessels. Recent hybrid approaches like SW-UNet 103 and DB-Net attempted to combine CNN architectures with 104 105 attention mechanisms but remained limited by their fundamentally convolutional base. Pure transformer approaches 106 107 such as Detection Transformer (DETR) and Deformable-DETR have shown promise in general object detection tasks 108 but require significant architectural adaptations and custom 109 loss functions to be effective for medical imaging chal-110 lenges. 111

Transformer architectures have emerged as a powerful 112 113 alternative to CNNs in medical imaging. While CNNs excel at capturing local features, they struggle with long-114 115 range dependencies, relationships between distant regions in an image [21]. Transformer self-attention effectively 116 models these dependencies, making it particularly valu-117 able for distinguishing nodules from vessels [21]. DETR, 118 a vision transformer, directly predicts object locations via 119 120 self-attention, replacing traditional region proposal methods but struggles with slow convergence and small object 121 detection [7, 27]. Deformable-DETR improves this by in-122 troducing a deformable attention mechanism which is spa-123 124 tially adaptive and computationally efficient. Unlike stan-125 dard self-attention, which attends to all pixels in an image, 126 Deformable Attention selectively focuses on a small set of dynamically learned sampling points around a reference lo-127 cation. This allows the model to adaptively refine its recep-128 tive field and capture fine-grained details of small objects 129 130 while significantly reducing computational overhead [27].

Segment Anything Model (SAM), trained on 1 billion131masks, enables promptable segmentation and is shown to132effectively transfer knowledge to new datasets with fine-133tuning [12, 17]. MedSAM fine-tunes SAM on 1.5 million134medical image-mask pairs to focus on anatomical complex-135ities in clinical settings [17].136

3. Methodology

We present a novel approach to lung nodule segmentation 138 for tumor boards by framing the task as anomaly detec-139 tion. Our method splits the task into two stages: Stage 1 140 serves as a region proposal phase to localize sparse nodules, 141 while Stage 2 refines these bounding boxes into pixel-wise 142 segmentation masks. While the building blocks are well-143 known, our novelty lies in unification of key architectural 144 components such as DETR, SAM along with strategies such 145 as deformable attention, focal loss and MIP into a special-146 ized framework. Our training dataset is based on LUNA16 147 and consists of 9,676 CT slices, preprocessed to enhance 148 nodule visibility through CLAHE, Otsu's thresholding, and 149 a customized training regimen to improve convergence. Our 150 model is validated on an independent test set from the Uni-151 versity Health Network (UHN), featuring 30 patients with 152 diverse imaging protocols to provide a more robust evalua-153 tion. 154

3.1. Data Preprocessing & Datasets

Our preprocessing pipeline prepares CT scan data for in-156 put into our Stage 1 network shown in Figure 1. We 157 first standardize anatomical structures by resampling CT 158 slices to a consistent voxel spacing of $1 \times 1 \times 1$ mm, en-159 suring uniformity. In order to isolate lung tissue from 160 the surrounding background, we employ Otsu's method for 161 thresholding[19]. This is followed by morphological oper-162 ations, including connected component analysis and region 163 erosion to obtain cleanly cropped lung regions. Slices at 164 the superior and inferior cranio-caudal extremes, which pro-165 vide minimal diagnostic value, are removed based on non-166 zero area size. This reduces the model's search space from 167 15M to 5.25M pixels per patient improving focus on rel-168 evant areas. Post segmentation Contrast Limited Adaptive 169 Histogram Equalization (CLAHE) is applied to improve the 170 visibility of subtle features like small nodules [23]. Images 171 are cropped to dimensions of 256x256, and a 5mm MIP is 172 finally applied to improve visibility of lung nodules by pro-173 ducing a 2D image that highlights the densest features. This 174 follows the expression 175

$$I_{\text{MIP}}(x, y) = \max_{z} \{ I(x, y, z) \}$$
(1) 176

. $I_{\text{MIP}}(x, y)$ is the 2D image intensity at each (x, y) location, and I(x, y, z) is the intensity of the original 3D image at voxel location (x, y, z). A slab thickness of 5mm 179

136 137

155



Figure 1. Data processing pipeline with nodule visible at the top left of lung. a) Original CT slice b) Post Otsu segmentation and CLAHE c) Post 5mm MIP

180 was chosen as a compromise between differentiating the shapes of vessels and nodules and limiting overlap be-181 tween structures. The final training dataset consists of 9,676 182 183 MIP CT slices, with 1,226 containing nodules, split 70%-20%-10% before augmentation. Non-nodule slices were 184 slightly undersampled during training resulting 12.7% pos-185 itive class. This adjustment was necessary because lower 186 nodule rates made training less effective due to excessive 187 sparsity. The test set maintained a 5% nodule rate to re-188 flect real-world conditions, with the adjustment leading to 189 higher test accuracy. Additionally, the training set contained 190 augmentations to enhance variability including flips, rota-191 tions ($\pm 15^{\circ}$), brightness shifts ($\pm 15\%$), and Gaussian noise 192 (0.001-0.18% SD). 193

194 A supplementary test dataset was obtained in collabo-195 ration with University Health Network (UHN), comprising 400 treated patients imaged with a TOSHIBA Aquil-196 ion scanner. The CT images have a 3mm slice thickness, 197 0.781mm pixel spacing, and patient ages ranging from 29 198 to 90 years (median: 68). All images were segmented 199 by a radiation oncologist. A subset of 30 patients (5,610 200 201 CT slices) was randomly selected as a second validation set. Figure 2 highlights nodule diameter variations between 202 UHN and LUNA16. LUNA16 primarily contains small 203 nodules (3-10mm), whereas UHN includes a broader dis-204 205 tribution, with many tumors between 25–55mm. UHN's 206 3mm slice thickness exceeds LUNA16's 2.5mm cutoff, potentially reducing nodule visibility. Pixel spacing also dif-207 fers, with UHN at 0.781mm and LUNA16 varying from 208 0.46-0.98mm. These differences introduce significant dis-209 tribution drift, making UHN a strong test for model gener-210 211 alizability.

3.2. Stage One and Two Model Architecture

Figure 3 overviews our two-stage approach, where Stage 1 213 generates region proposals to localize potential lung nod-214 ules in CT scans. Input images pass through a ResNet-50 215 CNN backbone for multi-scale feature extraction, then aug-216 mented with 2D sine-cosine positional encodings and pro-217 cessed by the encoder's DSA layers, which refine features 218 by attending to a sparse set of learnable sampling points 219 around each nodule. DSA aggregates features as 220

$$\mathbf{y}_q = \sum_{m=1}^M W_m \left(\sum_{k=1}^K A_{mqk} \cdot \mathbf{x} \Big(\mathbf{p}_q + \Delta \mathbf{p}_{mqk} \Big) \right) \quad (2) \quad \mathbf{221}$$

where M is the attention head count, K the sampled points 222 per head, W_m the projection matrices, and A_{mqk} the at-223 tention weight. The learnable offsets allow the model 224 to dynamically adjusts its receptive field for small and 225 irregular nodules. Regular self-attention has complexity 226 $O(H^2W^2C)$, but DSA reduces this to O(HWKC) where 227 height H, width W, channels C and K representing the num-228 ber of sampled points per attention head. Final decoder 229 heads refine object queries into bounding boxes and con-230 fidence scores, which serve as inputs for Stage 2 segmen-231 tation. Stage 1 was trained for 15 epochs using AdamW 232 with a learning rate of 10^{-4} , scheduled to reduce every 10 233 epochs. Training used an L4 GPU with mixed precision 234 (16-bit), batch size of 4, gradient clipping (0.1), and accu-235 mulation over 6 batches. A grid search optimized hyperpa-236 rameters for stability and performance. 237

We fine-tuned MedSAM's pretrained weights on a 238 dataset of 1,400 CT slices, where ground truth bounding boxes were used as prompts for SAM to simulate Stage 240

CVPR - CVMI 2025 Submission #32. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 2. Comparison of Size Distributions in LUNA16 (left) and UHN Test Dataset (right). LUNA16 contains primarily nodules, while UHN includes nodules and tumors.



Figure 3. Overview of the proposed LN-Transformer architecture, illustrating a sample CT MIP with a nodule in the top left corner. MIP deconstructed into its 5 associated slices for Stage 2.

241 1 predictions. These acted as attention cues enabling 242 targeted refinement while maintaining full-image context. The dataset included 1000 slices with ground truth nod-243 ule bounding boxes and 400 non-nodule slices to enhance 244 false positive discrimination from Stage 1. Ground truth 245 246 segmentation masks served as labels, and training was con-247 ducted using Dice-CrossEntropy loss with the Adam opti-248 mizer. Once trained, the Stage 2 auto-segmentation model processed MIP slices and associated bounding boxes from 249 250 Stage 1, reconstructing individual CT slices to generate precise pixel-wise segmentation masks as illustrated in Fig-251 252 ure 3.

3.3. Implementation Details

We trained our Stage 1 Deformable-DETR and Stage 2 finetuned SAM on a single L4 GPU with mixed precision (16-

bit). For Stage 1, we used the AdamW optimizer with an 256 initial learning rate of 1e-4 for the main parameters and 1e-257 5 for the ResNet-50 backbone, reducing by a factor of 10 258 every 10 epochs. A batch size of 4 was used, and gradi-259 ents were accumulated over 6 steps for stable learning. In 260 practice, processing each 256×256 MIP slice required ap-261 proximately 6 GB of GPU memory, with an average infer-262 ence speed of about 50 ms per slice. These hardware con-263 straints motivated our choice of deformable attention, which 264 is more memory-efficient than naive attention while retain-265 ing fine-grained focus on small objects. 266

3.4. Focal Loss

To address class imbalance, we incorporated focal loss into268the DETR loss function to enhance nodule detection by
down-weighting easy samples and emphasizing hard-to-269

336

337

338

339

340

341

342

343

344

345

346

347

348

349

361

classify cases [14]. The focal loss is defined as

$$FL(p_t) = -\alpha_t (1 - p_t)^{\gamma} \log(p_t) \tag{3}$$

273, where p_t is the predicted probability of the correct class,274 α_t balances positive and negative examples, and γ adjusts275focus towards challenging samples. Hyperparameter tuning276found $\gamma = 2$ and $\alpha_t = 0.25$ to provide an optimal balance277between precision and recall.

4. Results

Table 1 summarizes the performance metrics with nodules 279 280 size categories: small (up to 7mm), medium (7-15mm), and 281 large (over 15mm). Precision measures the proportion of 282 correctly identified nodules among predictions, while sensitivity captures the percentage of actual nodules detected. 283 The F1 score balances these metrics, and the Dice Coef-284 285 ficient evaluates the overlap between predicted and actual segmentation masks. Slice accuracy quantifies the propor-286 tion of CT slices correctly classified as containing or not 287 containing a nodule, providing a high-level assessment of 288 the ability to distinguish nodule-present and absent slices. 289 For medium and large nodules, the model achieves high pre-290 cision (96.7% and 97.8%) and recall (97.0% and 99.2%). 291 Stage 2 yields a 91.4% DiceC, with only a 3% drop from 292 the Stage 1 F1-score, indicating strong segmentation accu-293 racy. On the UHN test set, the model attains 86.8% slice-294 wise accuracy, 79.1% F1-score, and 78.3% Dice, reflecting 295 robustness despite greater tumor heterogeneity (25-55mm) 296 and distribution drift. The model maintains high precision 297 (74.9%) and sensitivity (83.5%), confirming generalization 298 across diverse clinical conditions. 299

Table 2 presents a comparison of our proposed DETR-300 SAM approach against comparable models on LUNA16. 301 Our DiceC of 91.4% outperforms next best models such 302 as MRUNet-3D, a multi-resolution U-Net with 3D convo-303 lutions (89.0%), and DB-NET, a dual-branch CNN with at-304 tention mechanisms (88.9%). For sensitivity and specificity, 305 our scores of 95.2% and 93.3% exceed prior state-of-the-306 art models such as ConvLSTM (92.2% sensitivity) and SW-307 UNet, a sliding window U-Net (89.0% specificity). 308

309 Figure 4 illustrates full pipeline results, showing Stage 1 bounding box predictions (LUNA16: red, UHN: blue) 310 and Stage 2 segmentation masks from MIP-reconstructed 311 slices. The top slices highlight complex vascular struc-312 tures and bronchioles that mimic or obscure small nodules, 313 314 while UHN tumors are larger and more numerous. Bound-315 ing boxes exhibit high IoU, leading to accurate Stage 2 segmentation, where most tumors are nearly perfectly delin-316 eated. However, the model occasionally misclassifies con-317 nected pulmonary vessels as part of the nodule, as seen in 318 319 the bottom-left mask.

5. Discussion

Our two-stage approach outperforms CNN and U-Net ar-321 chitectures. Models like MRUNet-3D [1], DB-Net [4] en-322 hance feature extraction but are constrained by fixed recep-323 tive fields. Hybrid models such as 3D-MSViT [18] improve 324 specificity (97.8%) and sensitivity but focus more on de-325 tection than segmentation. Bi-FPN and MV-DCNN empha-326 size sensitivity but lack a balanced trade-off with specificity. 327 Unlike prior hybrid methods, we train end-to-end trans-328 former models that make predictions without intermediary 329 processing [7, 27]. Other models oversample LUNA16 nod-330 ule slices during training, distorting real-world prevalence, 331 and lack external validation [6]. By preserving natural nod-332 ule sparsity and validating on independent clinical data, our 333 method surpasses previous models across all metrics and 334 ensures greater generalization to variability. 335

While Stage 1 struggles slightly with nodules under 7mm, clinical significance is limited as nodules ;6mm rarely warrant follow-up [13]. Stage 1's ability to sift through highly sparse data and still detect nodules and tumors at state-of-the-art rates is the key contribution of this work. However, errors in Stage 1 propagate to Stage 2, highlighting dependency on precise region proposals. Future work could explore adaptive confidence thresholds to reduce error propagation. Unlike balanced classification, our approach follows an anomaly detection paradigm where accuracy holds greater significance due to the rarity of positive instances. Despite lower performance on UHN, the model generalizes despite large nodule variability, imaging protocols, and distribution shifts.

We acknowledge the limitations of performing all pre-350 processing and modeling steps in 2D. Segment- ing lung 351 nodules in 2D may introduce challenges in differentiating 352 lesions near the thoracic wall and could impact visibility 353 for subsolid and ground-glass opacity (GGO) nodules, par-354 ticularly in thicker MIP slices. Future work will explore 355 extending our framework to 3D volumetric processing to 356 enhance spatial continuity and improve detection of these 357 challenging cases. Additionally, a subclass analysis could 358 provide deeper insights into performance across different 359 nodule types. 360

5.1. Ablation Study

To assess the impact of each component in our pipeline, we 362 conducted an ablation study (Table 3). Applying SAM di-363 rectly to lung CT scans resulted in a Dice coefficient of only 364 3.4%, confirming its inability to segment nodules without 365 guidance. MedSAM, despite being pre-trained on medi-366 cal imaging datasets, exhibited similarly poor performance, 367 achieving 14.1% when applied directly and 26.7% after 368 fine-tuning, indicating that pre-training alone does not en-369 sure generalization to lung nodule segmentation. Introduc-370 ing a detection stage significantly improved results. Using 371



Figure 4. Qualitative results from Stage 1 (top row) region proposals for LUNA16 (red) and UHN (blue), followed by Stage 2 (middle and bottom rows) masks from deconstructed MIP images.

Table 1.	Performance	Metrics for	Region	Proposal	(Stage	1) and	Auto-	Segmentation	(Stage	2) on	LUNA16	and	UHN	Test S	Sets.	IoU
threshold	= 0.5.															

Metric		LUNA16		UHN Test Set			
	F1/Dice	Precision	Sensitivity	F1/Dice	Precision	Sensitivity	
Stage 1: Region Proposal							
F1 Score	94.2%	_	-	79.1%	74.9%	83.5%	
Avg IoU (All)	_	93.3%	95.2%	_	_	_	
IoU (Small)	_	78.4%	83.3%	_	_	_	
IoU (Medium)	_	96.7%	97.0%	_	_	_	
IoU (Large)	_	97.8%	99.2%	_	_	_	
Slice Accuracy	97.1	_	-	86.8%	-	_	
Stage 2: Auto-Segmentation							
Dice Coefficient	91.4%	-	-	78.3%	-	-	

372 DETR in Stage 1 increased the Dice coefficient to 48.5%,373 demonstrating that bounding box proposals help narrow the

search space for SAM. However, DETR struggled significantly with small nodules with detection rates of only about 375

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

434

CVPR - CVMI 2025 Submission #32. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

Author	Architecture	DiceC (%)	Sensitivity (%)	Specificity (%)
Agnes et al. [1]	MRUNet-3D	89.0	94.8	84.2
Bhattacharyya et al. [4]	DB-NET	88.9	90.2	77.9
Song et al. [22]	ConvLSTM	84.0	87.8	81.5
Ma et al. [16]	SW-UNet	84.0	82.0	89.0
Annavarapu et al. [2]	Bi-FPN	82.8	92.2	78.9
Cao et al. [6]	DB-ResNet	82.7	89.4	79.6
Wang et al. [25]	MV-DCNN	77.9	87.0	77.3
Our Method	DETR-SAM	91.4	95.2	93.3

Table 2. Comparison of Nodule Segmentation on DiceC, Sensitivity, and Specificity

Table 3. Ablation Study Results: Performance of Different Configurations in the Proposed Pipeline

Configuration	F1 Score (%)	Dice Coefficient (%)
Direct SAM Application	_	3.4
Finetuned SAM Application	—	19.7
Finetuned MedSAM Application	—	26.7
S1: DETR + S2: Finetuned MedSAM	52.1	48.5
S1: DefDETR + S2: Finetuned MedSAM	94.2	91.4

30%. Replacing DETR with Deformable-DETR yielded
particularly striking results more than doubling segmentation accuracy, underscoring the critical role of adaptive attention in rare object detection.

380 These results underscore the necessity of a two-stage approach and how a robust detection stage mitigates class 381 imbalance. One possible reason for DETR's superiority is 382 383 its self-attention mechanism enables more flexible, contextaware feature representations of lung anatomy, allowing 384 385 it to distinguish nodules from normal structures even in ambiguous cases, potentially learning hierarchical relation-386 ships between tissue types and structural anomalies in a way 387 that segmentation models alone cannot. 388

389 6. Conclusion

This study introduces LN-Transformer, a two-stage trans-390 former framework tailored for sparse lung nodule segmen-391 tation. Key contributions include integrating Deformable-392 DETR with focal loss, MIP, SAM for mask refinement 393 394 to address class imbalance and segment nodules. Our 395 method achieves state-of-the-art results on LUNA16 (F1: 94.2%, Dice: 91.4%) and demonstrates robust general-396 ization on an independent clinical dataset (F1: 79.1%, 397 Dice: 78.3%), highlighting its potential for clinical appli-398 cation. 399

400 References

[1] S. Angel Agnes and J. Anitha. Appraisal of deep-learning
 techniques on computer-aided lung cancer diagnosis with

computed tomography screening.Journal of Medical403Physics, 45(2):98–106, 2020.5, 7404

- [2] C. S. R. Annavarapu, S. A. B. Parisapogu, N. V. Keetha, P. K.
 Donta, and G. Rajita. A bi-fpn-based encoder-decoder model for lung nodule image segmentation. *Diagnostics (Basel)*, 13
 (8):1406, 2023. 7
- [3] Julie A. Barta, Charles A. Powell, and Juan P. Wisnivesky. Global epidemiology of lung cancer. *Annals of Global Health*, 85(1):8, 2019.
- [4] Debnath Bhattacharyya, N. Thirupathi Rao, Eali Stephen Neal Joshua, and Yu-Chen Hu. A bi-directional deep learning architecture for lung nodule semantic segmentation. *The Visual Computer*, 39(11):5245–5261, 2023. 5, 7
- [5] Mateusz Buda, A. Maki, and M. Mazurowski. A systematic study of the class imbalance problem in convolutional neural networks. *Neural Networks: The Official Journal of the International Neural Network Society*, 106:249–259, 2017.
 2
- [6] Haichao Cao, Hong Liu, Enmin Song, Chih-Cheng Hung, Guangzhi Ma, Xiangyang Xu, Renchao Jin, and Jianguo Lu. Dual-branch residual network for lung nodule segmentation. *Applied Soft Computing*, 86:105934, 2020. 5, 7
- [7] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. *ArXiv*, abs/2005.12872, 2020. 2, 5
- [8] Robin Chan, M. Rottmann, Fabian Hüger, Peter Schlicht, and H. Gottschalk. Metafusion: Controlled false-negative reduction of minority classes in semantic segmentation. *ArXiv*, abs/1912.07420, 2019. 2
 430
- [9] D. Cody. Aapm/rsna physics tutorial for residents: topics in

510

511

ct. image processing in ct. *Radiographics : a review publication of the Radiological Society of North America, Inc, 22*5:1255–68, 2002. 2

- [10] Shimaa El-bana, A. Al-Kabbany, and M. Sharkas. A two-stage framework for automated malignant pulmonary nodule detection in ct scans. *Diagnostics*, 10, 2020. 1
- [11] J. Gruden, S. Ouanounou, S. Tigges, Shannon D. Norris, and T. Klausner. Incremental benefit of maximum-intensityprojection images on observer detection of small pulmonary nodules revealed by multidetector ct. *AJR. American journal of roentgenology*, 179(1):149–157, 2002. 2
- [12] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao,
 Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and
 Ross Girshick. Segment anything. *arXiv:2304.02643*, 2023.
 2
- [13] Anna Rita Larici, Alessandra Farchione, Paola Franchi,
 Mario Ciliberto, Giuseppe Cicchetti, Lucio Calandriello,
 Annemilia del Ciello, and Lorenzo Bonomo. Lung nodules:
 size still matters. *European Respiratory Review*, 26(146),
 2017. 5
- [14] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and
 Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 2980–2988, 2017. 5
- [15] Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In 2017 IEEE International Conference on Computer Vision (ICCV), pages 2999–3007, 2017. 2
- [16] Jiajun Ma, Gang Yuan, Chenhua Guo, Xiaoming Gang, and
 Minting Zheng. Sw-unet: a u-net fusing sliding window
 transformer block with cnn for segmentation of lung nodules. *Frontiers in Medicine*, 10:1273441, 2023. 7
- [17] J. Ma, Y. He, F. Li, et al. Segment anything in medical images. *Nature Communications*, 15:654, 2024. 2
- [18] Hassan Mkindu, Longwen Wu, and Yaqin Zhao. 3d multiscale vision transformer for lung nodule detection in chest ct
 images. *Signal, Image and Video Processing*, 17:2473–2480,
 2023. 5
- 474 [19] Nobuyuki Otsu. A threshold selection method from gray475 level histograms. *IEEE Transactions on Systems, Man, and*476 *Cybernetics*, SMC-9(1):62–66, 1979. 2
- 477 [20] Wendi Qu, I. Balki, Mauro Mendez, J. Valen, J. Levman, and
 478 P. Tyrrell. Assessing and mitigating the effects of class imbalance in machine learning with application to x-ray imaging. *International Journal of Computer Assisted Radiology* 481 *and Surgery*, 15:2041 – 2048, 2020. 2
- [21] Fahad Shamshad, Salman Khan, Syed Waqas Zamir,
 Muhammad Haris Khan, Munawar Hayat, Fahad Shahbaz
 Khan, and Huazhu Fu. Transformers in medical imaging: A
 survey. *Medical Image Analysis*, 88:102802, 2023. 2
- 486 [22] Guangyu Song, Qian Dai, Yuhang Nie, and Guang Chen.
 487 Differential diagnosis of benign and malignant pulmonary
 488 nodules in ct images based on multitask learning. *Current*489 *Medical Imaging*, 2023. 7
- 490 [23] M. Sundaram, K. Ramar, N. Arumugam, and G. Prabin.
 491 Histogram based contrast enhancement for mammogram images. 2011 International Conference on Signal Processing,

Communication, Computing and Networking Technologies, 493 pages 842–846, 2011. 2 494

- [24] J. Walter, M. Heuvelmans, P. D. de Jong, R. Vliegenthart, 495 P. V. van Ooijen, Robin B. Peters, K. ten Haaf, U. Yousaf-496 Khan, C. van der Aalst, G. D. de Bock, W. Mali, H. Groen, 497 H. D. de Koning, and M. Oudkerk. Occurrence and lung can-498 cer probability of new solid nodules at incidence screening 499 with low-dose ct: analysis of data from the randomised, con-500 trolled nelson trial. The Lancet. Oncology, 17 7:907-916, 501 2016. 1 502
- [25] Shuo Wang, Mu Zhou, Olivier Gevaert, Zhenchao Tang, Di Dong, Zhenyu Liu, and Tian Jie. A multi-view deep convolutional neural networks for lung nodule segmentation. In 2017 39th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), pages 1752–1755, 2017. 7
 508
- [26] David C. Wyatt, Tanya Avery, Alex Quan, and Peter de Waal. Clinical decision support systems and rapid learning in oncology. *BMJ Health & Care Informatics*, 26(1), 2019. 1
- [27] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *ArXiv*, abs/2010.04159, 2020. 2, 5
 513
 514
 515

8