

MedROV: Towards Real-Time Open-Vocabulary Detection Across Diverse Medical Imaging Modalities

Tooba Tehreem Sheikh

Jean Lahoud

Rao Muhammad Anwer

Fahad Shahbaz Khan

Salman Khan

Hisham Cholakkal

Mohamed Bin Zayed University of Artificial Intelligence (MBZUAI)

{tooba.sheikh, jean.lahoud, rao.anwer, fahad.khan, salman.khan, hisham.cholakkal}

@mbzuai.ac.ae

Abstract

Traditional object detection models in medical imaging operate within a closed-set paradigm, limiting their ability to detect objects of novel labels. Open-vocabulary object detection (OVOD) addresses this limitation but remains underexplored in medical imaging due to dataset scarcity and weak text-image alignment. To bridge this gap, we introduce MedROV, the first Real-time Open Vocabulary detection model for medical imaging. To enable open-vocabulary learning, we curate a large-scale dataset, Omnis, with 600K detection samples across nine imaging modalities and introduce a pseudo-labeling strategy to handle missing annotations from multi-source datasets. Additionally, we enhance generalization by incorporating knowledge from a large pre-trained foundation model. By leveraging contrastive learning and cross-modal representations, MedROV effectively detects both known and novel structures. Experimental results demonstrate that MedROV outperforms the previous state-of-the-art foundation model for medical image detection with an average absolute improvement of 40 mAP50, and surpasses closed-set detectors by more than 3 mAP50, while running at 70 FPS, setting a new benchmark in medical detection. Our source code, dataset, and trained model are available at [MedROV](#).

1. Introduction

Object detection in medical imaging plays an important role in identifying abnormalities such as tumors, fractures, and diseased cells across diverse modalities, including CT scans, X-rays, MRIs, and histopathology slides. Unlike natural images, medical imaging poses unique challenges due to its multi-modal nature, where each modality has distinct visual and semantic characteristics.

Traditional object detection models, including single-

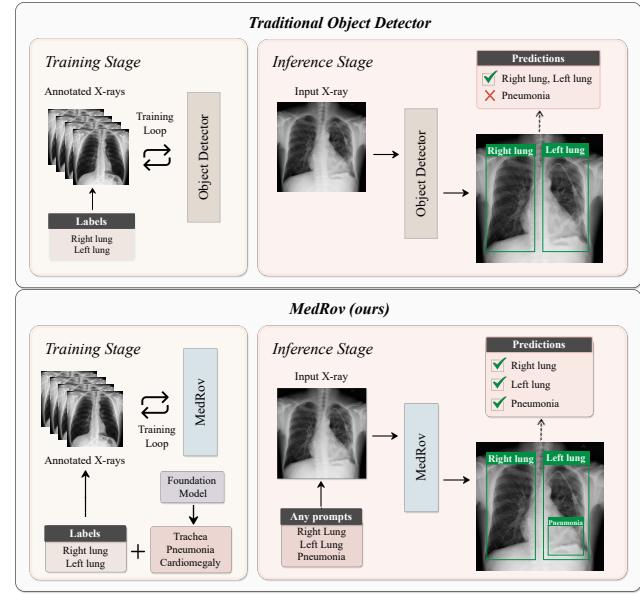


Figure 1. **Comparison of Traditional Object Detection and MedROV:** Traditional detectors are trained on a fixed set of categories and cannot recognize unseen classes. For example, the model shown here detects the left and right lungs but fails to detect pneumonia, which is present in the image. In contrast, MedROV is a Real-time Open Vocabulary detection model for medical imaging that leverages the BioMedCLIP foundation model to enable detection of both seen and unseen classes. At inference, it can detect any class described by a text prompt, if present in the image.

stage [40, 54] and two-stage [17, 41] detectors, have been adapted for medical tasks like tumor detection and organ localization. However, these models are limited to detecting only predefined categories (closed-set detection), making them ineffective in real-world medical scenarios where new, critical abnormalities may emerge and require immediate detection. On the other hand, Open Vocabulary Object Detection (OVOD) addresses the limitations of closed-

set detection models by enabling the detection of novel objects through vision-language alignment [10] or region-level vision-language pre-training [25]. While OVOD has shown significant potential in natural images, its application to medical imaging remains largely unexplored. Adapting OVOD methods, such as YOLO-World, to the medical domain presents challenges due to the scarcity of large, diverse, and well-annotated image-text datasets necessary for learning meaningful visual-language associations. Additionally, medical images exhibit complex variations in object size and shape, overlapping objects, and imbalances in class distribution, making the application of OVOD in the medical field more challenging.

To address these challenges, we introduce MedROV, the first Real-time Open Vocabulary detection model for medical imaging, designed to detect both known and novel structures across nine imaging modalities (Fig. 1). Trained on a large-scale dataset of over 600K samples, our model outperforms existing OVOD methods, previous state-of-the-art medical detection methods, and closed-set detectors by leveraging the BioMedCLIP [52] foundation model to enhance open vocabulary detection. In summary, our main contributions are as follows:

- We introduce the first open-vocabulary object detector for medical images, capable of detecting both known and unknown structures, by curating Omnis, a large-scale dataset of over 600K detection samples across 9 imaging modalities (CT, MRI, X-ray, Ultrasound, Histopathology, Dermoscopy, Fundoscopy, Endoscopy, and Microscopy).
- We adapt YOLO-World, an open vocabulary detector originally designed for natural images, to the medical domain by training it on our dataset and addressing the challenge of missing annotations when integrating multiple datasets across different modalities.
- We improve the model’s detection and generalization performance by incorporating information from the BioMedCLIP foundation model, leveraging its vision-language features for improved performance.
- We conduct extensive experiments comparing our method with existing OVOD approaches, the previous state-of-the-art model [53], and closed-set detectors. Our model achieves significant improvement in zero-shot detection performance on the Omnis test set compared to the baseline YOLO-World, while maintaining comparable speed (YOLO-World: 72 FPS, Ours: 70 FPS). Additionally, it outperforms the BioMedParse foundation model by an average absolute improvement of 40 mAP50 and surpasses closed-set methods by more than 3 mAP50.

2. Related Work

Recent advances in deep learning have significantly improved performance across various medical imaging tasks, including classification, segmentation, and detection. Ob-

ject detection, in particular, plays a crucial role in identifying anatomical structures and pathological abnormalities within medical scans. While traditional object detectors have been successfully adapted from natural images to medical domains, they largely operate under closed-set assumptions, limiting their ability to generalize to unseen or rare categories. Meanwhile, the emergence of large-scale foundation models and open-vocabulary object detection techniques has enabled the development of more flexible and generalizable systems. In this section, we review prior work in three relevant areas: object detection in medical imaging, the application of foundation models to medical tasks, and recent progress in open-vocabulary object detection.

2.1. Object Detection in Medical Images

Object detection is essential in medical imaging, allowing for the accurate detection of tumors, lesions, and other abnormalities. Deep learning-based detectors in natural images have been adapted for medical imaging. RT-DETR with multi-scale feature extraction has been applied to diabetic retinopathy detection [18], while BGF-YOLO [21] and SOCR-YOLO [27] have introduced enhancements to the YOLO architecture for brain tumor detection and lesion detection, respectively. To further enhance detection accuracy across varying object sizes and complex backgrounds, recent work by [49] proposed a cross-scale attention and multi-layer feature fusion method based on YOLOv8 for skin disease detection. However, these approaches remain limited to closed-set object detection, restricting their ability to identify novel or previously unseen structures.

2.2. Foundation Models in Medical Images

Recent foundation models have advanced medical imaging through multi-modal learning and self-supervised training. MedSAM [31] adapts SAM for medical image segmentation, while BioMedCLIP [52], trained on the PMC15M dataset, enhances CLIP for biomedical image understanding. MEDITRON [9], a large-scale medical language model, outperformed previous state-of-the-art models on USMLE-style questions. MedPaLM-2 [43], an adaptation of PaLM-2 for healthcare, demonstrates strong performance in clinical question answering. Additionally, BioMedParse [53], a biomedical foundation model for image parsing, introduces a unified approach for joint segmentation, detection, and recognition. However, most of these models remain limited to classification and question answering tasks, restricting their ability to tackle more complex and open-ended challenges in medical imaging.

2.3. Open-Vocabulary Object Detection

Open-Vocabulary Object Detection (OVOD) extends traditional detection by identifying both known and novel objects through semantic understanding. Significant progress

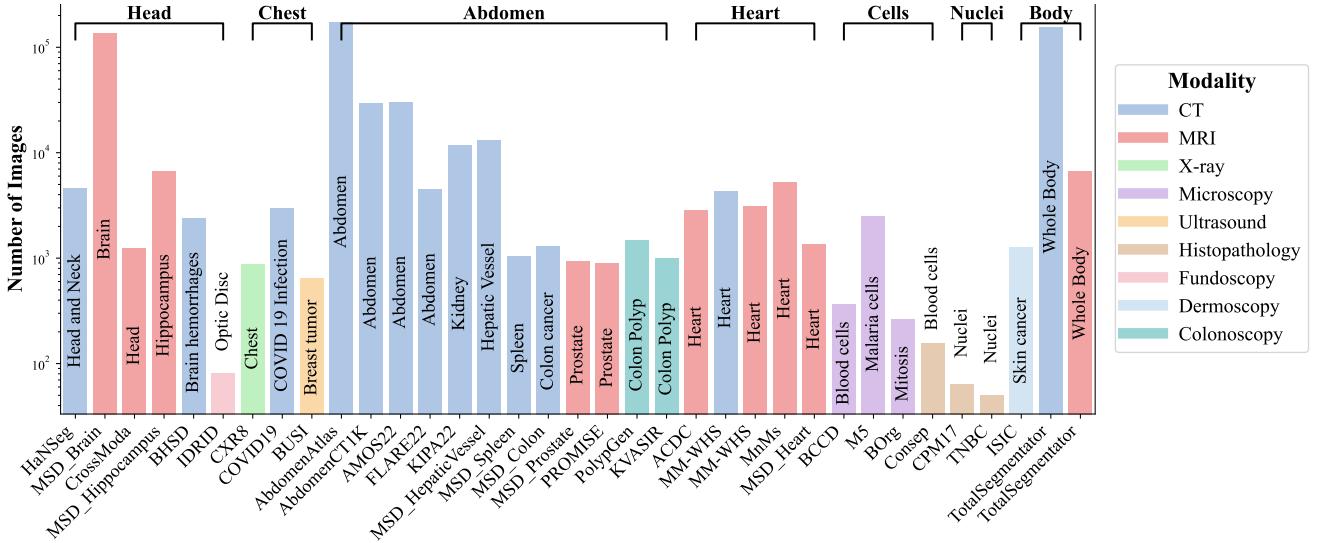


Figure 2. **An overview of the Omnis 600K dataset.** We curate a large-scale object detection dataset for medical imaging by incorporating 35 datasets with diverse modalities (represented in different colors), anatomical regions (displayed at the top), and target areas (indicated within the bars).

has been made on natural images with models such as DINO-X [42], which enhances open-world detection using a Transformer-based encoder-decoder architecture, and DetCLIPv3 [50], which integrates open-set detection with captioning for detailed descriptions. Detic [55] leverages image-level labels and large vocabularies to perform open-vocabulary detection and classification. GLIP [24] enables zero-shot generalization by formulating detection as a matching task between image regions and text. Additionally, YOLO-World [10] integrates the lightweight YOLO framework with CLIP for real-time detection. Despite these advancements, OVOD in medical imaging remains largely unexplored. Our approach builds upon YOLO-World by incorporating the BioMedCLIP foundation model to enhance vision-language alignment and improve OVOD performance in medical imaging.

3. Methodology

We introduce MedROV, a Real-time Open Vocabulary detection model for medical imaging. In our work, we build upon the baseline YOLO-World [10] (Section 3.1) and incorporate domain-specific adaptations. To this end, we curate a large-scale dataset of 600K detection samples spanning nine imaging modalities, detailed in Section 3.2. We address the missing annotation challenge through pseudo-labeling, and make use of a large-scale medical foundation model, BioMedCLIP [52], to improve generalizability (Section 3.3).

3.1. Baseline

YOLO-World is an open-vocabulary object detection framework that extends the YOLOv8 architecture [40] by integrating vision-language modeling to enable detection beyond a fixed set of categories. It incorporates a CLIP-based text encoder to extract text embeddings from user-defined prompts and a Re-parameterizable Vision-Language Path Aggregation Network (RepVL-PAN) for cross-modal fusion of text and image features.

During training, YOLO-World constructs region-text pairs by assigning both positive and negative text labels to image regions. Positive labels correspond to objects present in the image, while negatives are randomly sampled from the remaining dataset classes that are not in the image, forming a dynamic vocabulary of size M labels per image. A region-text contrastive loss \mathcal{L}_{Con} is employed to optimize alignment between predicted object embeddings and text embeddings, where the object-text similarity is computed via an L2-normalized dot product followed by an affine transformation: $s_{k,j} = \alpha \cdot \text{L2-Norm}(e_k) \cdot \text{L2-Norm}(w_j)^T + \beta$, where e_k is the predicted object embedding, w_j is the text embedding, and α, β are learnable scaling and shifting parameters. The model also incorporates an IoU loss (\mathcal{L}_{IoU}) and a Distributed Focal Loss (\mathcal{L}_{DFL}) for bounding box regression. The total loss is given by: $\mathcal{L}(I) = \mathcal{L}_{\text{Con}} + \lambda_I \cdot (\mathcal{L}_{\text{IoU}} + \mathcal{L}_{\text{DFL}})$, where λ_I is an indicator set to 1 for samples with reliable box annotations (e.g., detection or grounding data) and 0 for weakly labeled image-text data. While YOLO-World excels in natural image detection, its performance on medical images is limited

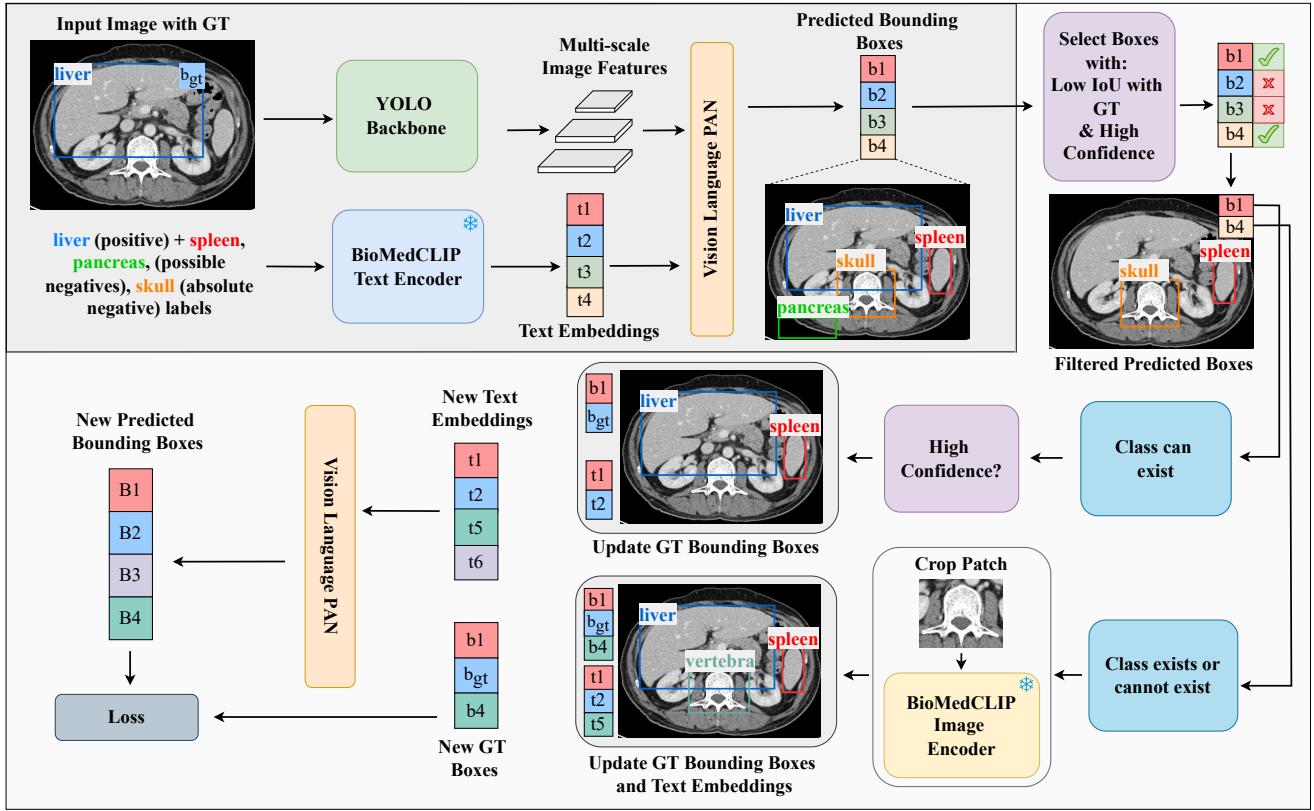


Figure 3. Overall architecture of MedROV. The model takes image and text labels as input. During training, positive and negative labels are used, whereas during testing, free-form text labels can be passed. The YOLO backbone extracts image features, while the BioMedCLIP text encoder generates text embeddings. These features are fused using the Vision-Language PAN (VL-PAN) to obtain bounding box predictions. During training (outer box), predictions are first filtered based on an IoU threshold. If a predicted class is missing in the dataset but can exist, the high-confidence bounding box is added to the ground truth as pseudo-label. Otherwise, the cropped region is passed through the BioMedCLIP image encoder for feature extraction. The extracted features replace one of the negative text label embeddings, updating the ground truth. The updated text embeddings and bounding boxes are passed through VL-PAN again to generate refined predictions. Finally, the loss is computed between the new ground truth and the updated predictions.

due to the domain gap, multi-modality nature of medical imaging, and a lack of large annotated medical datasets.

3.2. Omnis Dataset

Medical imaging spans multiple modalities with significant domain gaps. To enable open-vocabulary detection in MedROV, we curate Omnis 600K, a large-scale medical detection dataset covering nine imaging modalities and diverse anatomical and pathological targets (Fig. 2). It integrates public datasets from CT [3, 20, 28–30, 36, 39, 47, 48, 56], MRI [6, 8, 14], Ultrasound [2], Histopathology [15, 34], X-ray [46], Colonoscopy [26, 37], Microscopy [1, 4, 45], Dermoscopy [11], and Fundoscopy [38].

For 3D modalities like CT and MRI, each volume is processed into in-plane 2D slices for consistency. Segmentation masks are converted to detection bounding boxes by identifying the minimum and maximum coordinates of non-

zero regions per slice. To handle varying intensity distributions, we apply modality-specific normalization. Following MedSAM [31], we clip CT intensities to [-500, 1000] and MRI intensities to the 0.5th–99.5th percentile range. All images are normalized to a [0, 255] scale. Grayscale images are converted to 3-channel by replicating the single channel.

Omnis is specifically designed for open-vocabulary object detection. It comprises 157 training classes and incorporates the BioMedCLIP foundation model to enable generalization beyond fixed categories. Omnis consists of 577k training and 28k validation images. Following MedSAM [31] and BioMedParse [53], we split the data at the volume level to prevent data leakage, ensuring slices from the same 3D scan remain within a single split. We allocate 95% of the data for training and 5% for validation. We also hold out some classes within Omnis, ensuring that entire volumes containing these classes are excluded from training for zero-

shot evaluation. This held-out set is included in our test set (see Section 4 for details). Unlike datasets such as BioMedParse [53], which is limited to 82 predefined categories and follows an image–mask–label triplet format, constraining its ability to generalize to broader tasks, or MedSAM, which lacks detection capabilities and a scalable vocabulary, Omnis is specifically designed for medical OVOD.

3.3. MedROV Architecture

Overview: MedROV, illustrated in Fig. 3, leverages BioMedCLIP [52] as its text encoder, as it is specifically designed for medical applications and trained on a large-scale dataset of medical image–text pairs, making it well-suited for medical image understanding. Our method introduces open-vocabulary capability and enhances performance by training the model on the large-scale Omnis 600K detection dataset and addressing the missing annotations problem in medical datasets through pseudo-labeling. Additionally, we integrate BioMedCLIP features into MedROV for feature alignment (detailed next).

Addressing Missing Annotations: Medical datasets are typically labeled for specific targets, resulting in many object labels being present in images but not annotated. When combining multiple datasets, the number of labels increases, leading to more missing annotations. The absence of these annotations would lead to penalizing the model during training for predicting visible but unannotated objects, causing confusion and reducing learning efficiency. To address this issue, we introduce a dataset-class presence matrix (denoted by M), which categorizes the availability of annotations for each class across datasets. The matrix M assigns a value to each dataset-class pair (d, c) as follows: $M_{d,c} = 1$ if class c is annotated in dataset d , $M_{d,c} = 0$ if class c may be present in d but is not annotated, and $M_{d,c} = -1$ if class c cannot exist in d .

During training, we extract predictions from the detection head and apply non-maximum suppression (NMS) to eliminate redundant boxes. To identify unannotated objects, we compute the Intersection over Union (IoU) between each predicted bounding box and all ground truth boxes, regardless of the class labels (class-agnostic). A prediction is considered a potential missed annotation if its maximum IoU with any ground truth box falls below a predefined threshold T i.e.: $\max_{g \in G} \text{IoU}(p, g) < T$, where p is the predicted bounding box, and G is the set of ground truth boxes. This indicates that the model has detected an object that does not overlap with any annotated ground truth.

For each potentially missed annotation, we check $M_{d,c}$. If $M_{d,c} = 0$, indicating that class c may be present but is not annotated in dataset d , and the model’s prediction confidence exceeds a threshold C , the prediction is added to the ground truth as a pseudo-label to prevent the model from being incorrectly penalized for detecting valid but unanno-

tated objects, improving training stability and performance.

Enhancing Generalization Using Foundation Medical Image Models: To enhance the generalizability and open-vocabulary capability of MedROV, we incorporate knowledge from BioMedCLIP [52], a foundation model trained on a large-scale medical dataset of image–text pairs. Given a prediction, we refer to the dataset-class matrix $M_{d,c}$. If $M_{d,c} = 1$ (class annotated) or $M_{d,c} = -1$ (class cannot exist), we refine predictions by discarding overly small boxes, full-image boxes, or those that predominantly cover background regions (e.g., boxes containing mostly black pixels). The remaining boxes are expanded by a factor $F = 1.3$, and the cropped region is processed through the BioMedCLIP image encoder to align image features with textual representations. For example, if the model predicts a liver in a dataset where liver annotations already exist ($M_{d,c} = 1$), but the actual object is, say, a spleen, BioMedCLIP helps ensure the extracted visual features align more closely with the spleen’s textual features. Similarly, if the model predicts a skull in an abdomen dataset, where such an object cannot exist ($M_{d,c} = -1$), BioMedCLIP encourages alignment with the correct object’s text features, correcting the semantic mismatch. To avoid direct reliance on noisy or incorrect labels, we replace the text features of a negative ground truth sample with the extracted BioMedCLIP image features and add the bounding box to the ground truth, repeating this for up to five boxes per image, sorted by confidence. By training the model on semantically grounded features instead of potentially incorrect labels, this approach mitigates error propagation, enhances robustness, and improves the open-vocab capability and generalizability of MedROV.

The final loss, identical to that of the baseline YOLO-World [10], is computed between the updated ground truth and the model predictions. During inference, MedROV adopts the prompt-then-detect paradigm of YOLO-World [10], enabling real-time detection with comparable speed across diverse medical imaging tasks.

4. Experiments and Results

To validate the effectiveness of MedROV, we conduct extensive experiments across a wide range of medical imaging datasets. We evaluate both detection accuracy and generalization ability under open-vocabulary and cross-modality settings, comparing MedROV against existing state-of-the-art methods and baseline models.

Implementation Details: Starting from the baseline YOLO-World, we fine-tune MedROV on our proposed Omnis 600K dataset for 20 epochs, using 577K training and 28K validation samples. Training is performed on four NVIDIA A100 40GB GPUs with a batch size of 128, a learning rate of 0.0002, and a weight decay of 0.05. To optimize detection performance, we evaluate various combinations of IoU threshold T and confidence threshold C

Models → Dataset ↓	Modality	Classes	YOLO-World		YOLO-World + Our Omnis		MedROV (Ours)	
			mAP50	mAP50:95	mAP50	mAP50:95	mAP50	mAP50:95
BTCV	CT	Base	0.00	0.00	80.7	62.2	83.4	63.1
		Base + Novel	0.03	0.01	74.5	57.5	79.1	59.3
Cervix	CT	Base	0.00	0.00	64.3	39.3	66.9	42.3
		Base + Novel	0.00	0.00	33.3	20.7	33.8	21.3
MSD Liver	CT	Base	0.34	0.09	99.4	94.7	99.4	94.3
		Base + Novel	0.19	0.05	51.4	47.9	58.7	51.7
MSD Pancreas	CT	Base	0.07	0.01	92.5	62.5	92.1	62.5
		Base + Novel	0.04	0.01	46.3	30.8	47.0	31.5
LiTS	CT	Base	2.42	0.61	98.5	91.9	98.2	92.8
		Base + Novel	2.68	0.48	50.7	46.6	57.0	50.9
TotalSegmentator	CT and MRI	Base	0.33	0.12	63.2	50.6	65.7	53.7
		Base + Novel	0.28	0.16	43.5	33.9	46.0	35.5
Omnis Validation Set	All Modalities	Base	0.02	0.01	61.2	45.7	64.1	48.9
Multi-Modality	All Modalities	Base	3.14	1.25	86.5	62.2	89.8	66.7
		Base + Novel	0.59	0.35	38.6	26.7	43.5	31.1

Table 1. Comparison of zero-shot detection performance between MedROV and YOLO-World on BTCV, Cervix, MSD Liver, MSD Pancreas, LiTS, TotalSegmentator, the Omnis validation set, and multi-modality datasets, for both base and novel classes. YOLO-World + Our Omnis denotes the baseline YOLO-World model trained on our Omnis 600K dataset. Top scores are highlighted in **bold**.

values to optimize detection performance. The best results are achieved with an IoU threshold $T = 0.3$, and the confidence threshold of $C = 0.9$, which effectively filters predictions to retain high-quality pseudo-labels and reduce noise.

Evaluation Datasets and Metrics: We evaluate the open-vocabulary detection (OVOD) capability of MedROV through three strategies. First, we test zero-shot transfer on entirely unseen datasets, including the Medical Segmentation Decathlon (Liver and Pancreas) [3] and LiTS [7]. Second, we hold out specific classes and their images from BTCV [22], Cervix [22], TotalSegmentator CT [47] and MRI [13] to assess recognition of novel categories. Third, we evaluate cross-modality generalization using a curated Multi-Modality dataset spanning nine imaging types, created by combining Chest X-ray [19], Breast Lesion [35], DigestPath [12], PH2 [32], CVC ClinicDB [5], LiverHcc-Seg [16], NeurIPS CellSeg [23], and Drishti-GS1 [44]. Performance is measured using mean average precision (mAP) at IoU thresholds of 50 and 50:95.

4.1. Results

MedROV is the first open-vocabulary object detection (OVOD) framework tailored for medical imaging. Initial experiments showed that state-of-the-art OVOD methods for natural images—such as OV-DETR [51], OWL-ViT [33], and GLIP [24]—perform poorly on medical data, with near-zero accuracy across benchmarks. This highlights the significant domain gap and the need for medical-specific solutions. Given these limitations, we selected YOLO-World [10] as our starting point due to its strong performance.

While its zero-shot inference was limited on medical data, fine-tuning it on our Omnis 600K dataset led to more meaningful results. We benchmarked MedROV against both traditional closed-set detectors and the previous state-of-the-art medical imaging detection model, BioMedParse [53], under both open- and closed-vocabulary settings. MedSAM [31] was excluded from comparison, as it is a segmentation-only model that outputs binary masks without class labels, making it unsuitable for object detection.

Comparison of MedROV with YOLO-World: MedROV outperforms the baseline YOLO-World [10], which fails at zero-shot detection on medical datasets due to a significant domain gap between natural and medical images. To ensure a fair comparison, we evaluate against YOLO-World trained on our Omnis 600K dataset (YOLO-World + Our Omnis). As shown in Table 1, MedROV consistently surpasses this stronger baseline across both base and novel classes. For example, it improves base + novel mAP50 from 74.5 to 79.1 on the BTCV dataset, from 43.5 to 46.0 on the TotalSegmentator CT and MRI dataset, and from 38.6 to 43.5 on the multi-modality dataset. Additionally, MedROV retains real-time performance (MedROV: 72 FPS vs. YOLO-World: 70 FPS), where real-time refers to inference speeds exceeding a threshold of 30 FPS, making it a practical solution for accurate and efficient medical object detection.

Comparison of MedROV with Closed-set Detectors: Table 2 presents a comparison between MedROV, YOLO-World + Our Omnis (baseline trained on our Omnis 600K dataset), and traditional closed-set detectors (YOLOv8,

Models → Metrics ↓	YOLOv8	Fine-tuned YOLOv9	YOLOv10	YOLO-World + Our Omnis Zero-Shot	Fine-tuned	MedROV (Ours) Zero-Shot	Fine-tuned
mAP50	77.3	74.7	76.2	74.5	72.2	79.1	81.9
mAP50:95	59.1	57.1	58.2	57.5	55.1	59.3	63.1

Table 2. Performance comparison of MedROV with YOLO-World + Our Omnis (baseline trained on our proposed Omnis 600K dataset) and closed-set detectors. All models are fine-tuned on the BTCV training set and evaluated on the BTCV test set. Top scores are highlighted in **bold**.

Models → Dataset ↓	Modality	BioMedParse [53] mAP50	MedROV (Ours) mAP50	Improvement ↑
BTCV	CT	9.09	79.1	70.01
Cervix	CT	6.49	33.8	27.31
MSD Pancreas	CT	24.61	47.0	22.39
MSD Liver	CT	18.68	58.7	40.02
LiTS	CT	12.63	57.0	44.37
TotalSegmentator	CT and MRI	0.52	46.0	45.48
Multimodality	All Modalities	12.84	43.5	30.66

Table 3. Zero-shot detection performance comparison between MedROV (ours) and the foundation model BioMedParse [53] across multiple medical datasets spanning CT, MRI, and multi-modality settings. MedROV consistently outperforms BioMedParse, with improvements reported in the last column.

YOLOv9, YOLOv10), fine-tuned on the BTCV training set for 20 epochs. In the zero-shot setting, MedROV achieves a strong mAP50 of 79.1, outperforming all fine-tuned closed-set models. When fine-tuned, MedROV further improves to 81.9 mAP50 and 63.1 mAP50:95, outperforming the best closed-set model (YOLOv8) by +4.6 mAP50 and +4.0 mAP50:95, demonstrating its superior generalization across both zero-shot and fine-tuned scenarios.

Comparison of MedROV with BioMedParse: To the best of our knowledge, MedROV is the first to explore real-time OVOD in medical imaging. Given the absence of OVOD methods tailored for medical data, we compared MedROV with BioMedParse [53], a recent multi-task foundation model for detection, segmentation, and recognition.

BioMedParse is trained on a large corpus of image–mask–label triplets to recognize 82 predefined medical object types, but lacks generalization to unseen categories. Additionally, it does not produce confidence scores necessary for standard detection evaluation. To address this, we approximated prediction confidence by averaging the positive pixel values within each predicted mask. Its reliance on anatomy-specific preprocessing (e.g., intensity tuning per organ) further limits scalability across diverse object types. Moreover, BioMedParse struggles with multi-object images, often detecting only a subset of entities or assigning the same label to distinct structures—challenges stemming from its fixed vocabulary and training on single-object image–mask pairs. In contrast, MedROV, trained on 157 categories and incorporating BioMedCLIP, supports open-vocabulary detection, applies uniform preprocessing

across all inputs, and handles multi-object images effectively. For a fair comparison, both models were evaluated on our test datasets using the same preprocessed inputs.

As shown in Table 3, MedROV consistently outperforms BioMedParse across all benchmarks. For example, on the BTCV dataset, it achieves an mAP50 of 79.1, compared to 9.09 by BioMedParse, an improvement of 70.01. Similar trends are observed on MSD Pancreas, LiTS, and TotalSegmentator, highlighting MedROV’s strong detection and generalization performance. Moreover, MedROV operates in real time at 70 FPS with CPU support, whereas BioMedParse runs at only 4 FPS and lacks CPU deployment capability. MedROV’s flexibility makes it better suited for real-world clinical scenarios, where broad category coverage and adaptability are essential.

Qualitative Results: Fig. 4 illustrates MedROV’s zero-shot detection performance across four datasets. The model accurately identifies novel classes such as liver and breast lesions, and liver cancer, while maintaining strong performance on base classes, as seen in the BTCV dataset. Low-confidence detections (e.g., 3% confidence for breast lesion) were retained based on observations from YOLO-World [10], where novel objects can appear even with confidence scores as low as 1%. For qualitative visualization, we use a data-driven, image-specific threshold rather than a fixed one. Specifically, we sort the detection scores in descending order and find the point where the score curve sharply bends (“elbow point”), which reflects a separation between strong and weak detections. This method is only used for visualization and does not affect quantitative evaluation, where

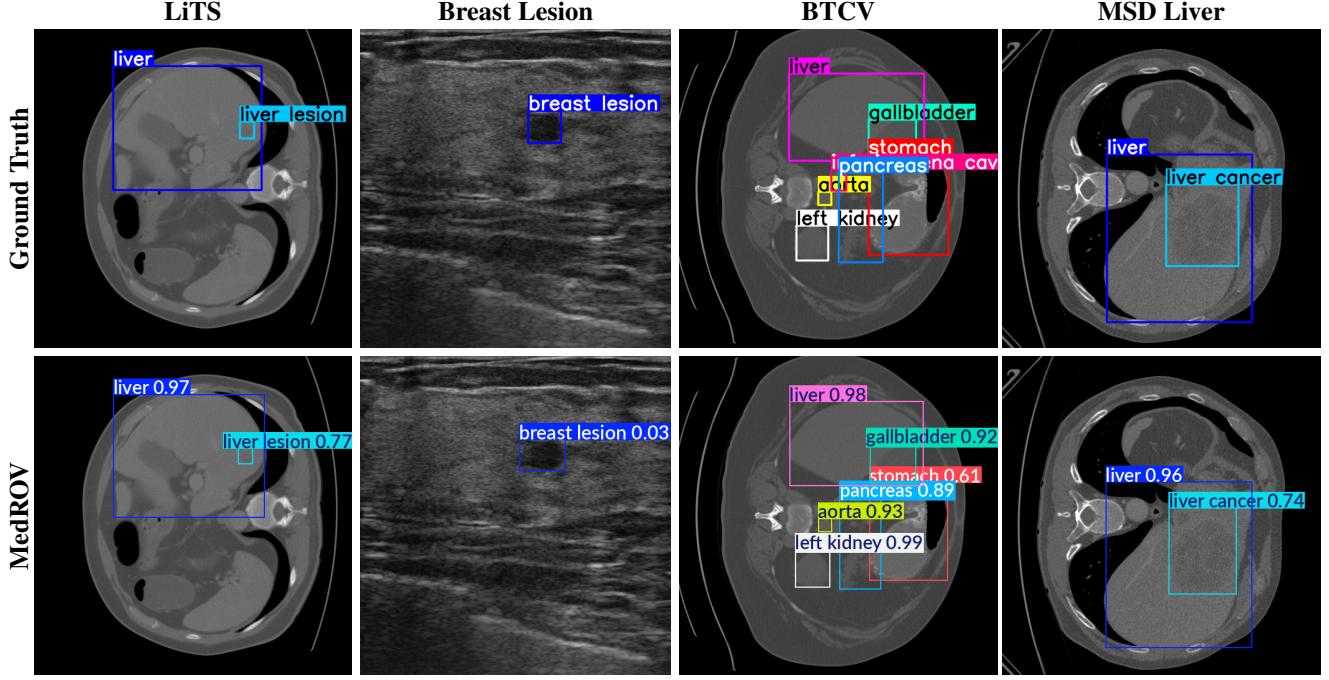


Figure 4. Visual comparison of MedROV’s zero-shot detection performance on four datasets: LiTS, Breast Lesion, BTCV, and MSD Liver. The model successfully detects both known and novel classes, including liver lesion, breast lesion, and liver cancer.

	Omnis 600K Dataset	Missing Annotations	Enhancing Generalization	Base Classes		Base + Novel Classes	
				mAP50	mAP50:95	mAP50	mAP50:95
Baseline	-	-	-	0.02	0.01	0.01	0.00
Baseline	✓	-	-	79.1	61.7	48.0	37.1
Baseline*	✓	-	-	78.8	61.5	48.4	37.5
Baseline*	✓	✓	-	81.1	64.3	50.8	39.5
Baseline*	✓	✓	✓	81.8	66.7	51.3	40.3

Table 4. Ablation study of MedROV with diverse configurations, highlighting the impact of fine-tuning the baseline on our Omnis 600K dataset, addressing missing annotations, and enhancing generalization. Results are reported for base and base + novel classes on the TotalSegmentator CT dataset. Baseline* represents MedROV with BioMedCLIP Text Encoder. Top scores are highlighted in **bold**.

boxes are ranked by raw confidence as in standard AP computation. Moreover, in zero-shot settings, relative confidence is often more informative than absolute values, so even low-confidence detections can be meaningful.

Ablation study on the impact of different configurations: Table 4 illustrates that fine-tuning the baseline YOLO-World on our Omnis-600K dataset significantly improves performance compared to the zero-shot baseline. Replacing the CLIP text encoder with BioMedCLIP (Baseline*) results in a slight decrease in base class performance, with mAP50 dropping from 79.1 to 78.8, but improves novel class performance from 48.0 to 48.4. This reflects better alignment with medical terminology, as BioMedCLIP is trained on a large dataset of medical image–text pairs. Despite the minor tradeoff, the BioMedCLIP text encoder is

used in all experiments due to its domain relevance. Adding pseudo-labeling to address missing annotations further improves the Base + Novel mAP50 from 48.4 to 50.8. Finally, incorporating BioMedCLIP image features to enhance generalization leads to the highest performance, achieving 81.8 mAP50 on base classes and 51.3 mAP50 on base + novel classes.

5. Conclusion

We introduce MedROV, the first real-time open-vocabulary detection method for medical imaging. By adapting YOLO-World to the medical domain, we replace the CLIP text encoder with BioMedCLIP and curate Omnis, a large-scale dataset comprising 600K samples across nine imaging modalities. To address missing annotations when merging

datasets, we employ a pseudo-labeling strategy. We also integrate knowledge from a foundation model to enhance generalization. MedROV outperforms the baseline YOLO-World, recent foundation models, and closed-set detectors in open-vocabulary performance, while maintaining real-time speed. Future work will focus on curating an open-vocabulary test set with a broader range of categories and extending the approach to 2D and 3D segmentation tasks.

Acknowledgments

This work is partially supported by the MBZUAI-WIS Joint Program for AI Research (Grant number WIS-P008) and the NVIDIA Academic Grant 2025.

References

- [1] Bccd: Blood cell count and detection, 2018. [4](#)
- [2] Walid Al-Dhabayani, Mohammed Gomaa, Hussien Khaled, and Aly Fahmy. Dataset of breast ultrasound images. *Data in brief*, 28:104863, 2020. [4](#)
- [3] Michela Antonelli, Annika Reinke, Spyridon Bakas, Keyvan Farahani, Annette Kopp-Schneider, Bennett A Landman, Geert Litjens, Bjoern Menze, Olaf Ronneberger, Ronald M Summers, et al. The medical segmentation decathlon. *Nature communications*, 13(1):4128, 2022. [4, 6](#)
- [4] Muhammad Awais, Mehaboobathunnisa Sahul Hameed, Bidisha Bhattacharya, Orly Reiner, and Rao Muhammad Anwer. Borg: A brain organoid-based mitosis dataset for automatic analysis of brain diseases. *arXiv preprint arXiv:2406.19556*, 2024. [4](#)
- [5] Jorge Bernal, F Javier Sánchez, Gloria Fernández-Esparrach, Debora Gil, Cristina Rodríguez, and Fernando Vilariño. Wm-dova maps for accurate polyp highlighting in colonoscopy: Validation vs. saliency maps from physicians. *CMIG*, 43:99–111, 2015. [6](#)
- [6] Olivier Bernard, Alain Lalande, Clement Zotti, Frederick Cervenansky, Xin Yang, Pheng-Ann Heng, Irem Cetin, Karim Lekadir, Oscar Camara, Miguel Angel Gonzalez Ballester, et al. Deep learning techniques for automatic mri cardiac multi-structures segmentation and diagnosis: Is the problem solved? *TMI*, 2018. [4](#)
- [7] Patrick Bilic, Patrick Christ, Hongwei Bran Li, Eugene Vorontsov, Avi Ben-Cohen, Georgios Kaassis, Adi Szeskin, Colin Jacobs, Gabriel Efrain Humpire Mamani, Gabriel Chartrand, et al. The liver tumor segmentation benchmark (lits). *MIA*, 2023. [6](#)
- [8] M Campello and K Lekadir. Multi-centre multi-vendor & multi-disease cardiac image segmentation challenge (m&ms). In *Medical Image Computing and Computer Assisted Intervention*, 2020. [4](#)
- [9] Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, et al. Meditron-70b: Scaling medical pretraining for large language models. *arXiv preprint arXiv:2311.16079*, 2023. [2](#)
- [10] Tianheng Cheng, Lin Song, Yixiao Ge, Wenyu Liu, Xinggang Wang, and Ying Shan. Yolo-world: Real-time open-vocabulary object detection. In *CVPR*, pages 16901–16911, 2024. [2, 3, 5, 6, 7](#)
- [11] Noel Codella, Veronica Rotemberg, Philipp Tschandl, M Emre Celebi, Stephen Dusza, David Gutman, Brian Helba, Aadi Kalloo, Konstantinos Liopyris, Michael Marchetti, et al. Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic). *arXiv preprint arXiv:1902.03368*, 2019. [4](#)
- [12] Qian Da, Xiaodi Huang, Zhongyu Li, Yanfei Zuo, Chenbin Zhang, Jingxin Liu, Wen Chen, Jiahui Li, Dou Xu, Zhiqiang Hu, et al. Digestpath: A benchmark dataset with challenge review for the pathological detection and segmentation of digestive-system. *MIA*, 2022. [6](#)
- [13] Tugba Akinci D’Antonoli, Lucas K Berger, Ashraya K Indrakanti, Nathan Vishwanathan, Jakob Weiß, Matthias Jung, Zeynep Berkarda, Alexander Rau, Marco Reisert, Thomas Küstner, et al. Totalsegmentator mri: Sequence-independent segmentation of 59 anatomical structures in mr images. *preprint arXiv:2405.19492*, 2024. [6](#)
- [14] Reuben Dorent, Aaron Kujawa, Marina Ivory, Spyridon Bakas, Nicola Rieke, Samuel Joutard, Ben Glocker, Jorge Cardoso, Marc Modat, Kayhan Batmanghelich, et al. Cross-modality 2021 challenge: Benchmark of cross-modality domain adaptation techniques for vestibular schwannoma and cochlea segmentation. *Medical Image Analysis*, 83:102628, 2023. [4](#)
- [15] Simon Graham, Quoc Dang Vu, Shan E Ahmed Raza, Ayesha Azam, Yee Wah Tsang, Jin Tae Kwak, and Nasir Rajpoot. Hover-net: Simultaneous segmentation and classification of nuclei in multi-tissue histology images. *Medical image analysis*, 58:101563, 2019. [4](#)
- [16] Moritz Gross, Sandeep Arora, Steffen Huber, Ahmet S Küçükkaya, and John A Onofrey. Liverhcseg: A publicly available multiphasic mri dataset with liver and hcc tumor segmentations and inter-rater agreement analysis. *Data in Brief*, 51:109662, 2023. [6](#)
- [17] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017. [1](#)
- [18] Weijie He, Yuwei Zhang, Ting Xu, Tai An, Yingbin Liang, and Bo Zhang. Object detection for medical image analysis: Insights from the rt-det model. *arXiv preprint arXiv:2501.16469*, 2025. [2](#)
- [19] Stefan Jaeger, Sema Candemir, Sameer Antani, Yí-Xiáng J Wáng, Pu-Xuan Lu, and George Thoma. Two public chest x-ray datasets for computer-aided screening of pulmonary diseases. *Quantitative imaging in medicine and surgery*, 2014. [6](#)
- [20] Yuanfeng Ji, Haotian Bai, Chongjian Ge, Jie Yang, Ye Zhu, Ruimao Zhang, Zhen Li, Lingyan Zhanng, Wanling Ma, Xiang Wan, et al. Amos: A large-scale abdominal multi-organ benchmark for versatile medical image segmentation. *NeurIPS*, 35:36722–36732, 2022. [4](#)
- [21] Ming Kang et al. Bgf-yolo: Enhanced yolov8 with multi-scale attentional feature fusion for brain tumor detection. In *MICCAI*, pages 35–45. Springer, 2024. [2](#)

- [22] Bennett Landman, Zhoubing Xu, J Igelsias, Martin Styner, Thomas Langerak, and Arno Klein. Miccai multi-atlas labeling beyond the cranial vault–workshop and challenge. In *MICCAI*, 2015. 6
- [23] Kwanyoung Lee, Hyungjo Byun, and Hyunjung Shim. Cell segmentation in multi-modality high-resolution microscopy images with cellpose. In *NeurIPS*, pages 1–11. PMLR, 2023. 6
- [24] Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, et al. Grounded language-image pre-training. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10965–10975, 2022. 3, 6
- [25] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *ECCV*. Springer, 2024. 2
- [26] Shengyuan Liu, Zhen Chen, Qiushi Yang, Weihao Yu, Di Dong, Jiancong Hu, and Yixuan Yuan. Polyp-gen: Realistic and diverse polyp image generation for endoscopic dataset expansion. *arXiv preprint arXiv:2501.16679*, 2025. 4
- [27] Yongjie Liu, Yang Li, et al. Socr-yolo: Small objects detection algorithm in medical images. *IMA*, 34(4):e23130, 2024. 2
- [28] Jun Ma, Yixin Wang, Xingle An, Cheng Ge, Ziqi Yu, Jianan Chen, Qiongjie Zhu, Guoqiang Dong, Jian He, Zhiqiang He, et al. Toward data-efficient learning: A benchmark for covid-19 ct lung and infection segmentation. *Medical physics*, 48(3):1197–1210, 2021. 4
- [29] Jun Ma, Yao Zhang, Song Gu, Cheng Zhu, Cheng Ge, Yichi Zhang, Xingle An, Congcong Wang, Qiyuan Wang, Xin Liu, et al. Abdomenct-1k: Is abdominal organ segmentation a solved problem? *IEEE TPAMI*, 44(10):6695–6714, 2021.
- [30] Jun Ma, Yao Zhang, Song Gu, Xingle An, Zhihe Wang, Cheng Ge, Congcong Wang, Fan Zhang, Yu Wang, Yinan Xu, et al. Fast and low-gpu-memory abdomen ct organ segmentation: the flare challenge. *Medical Image Analysis*, 82:102616, 2022. 4
- [31] Jun Ma, Yuting He, Feifei Li, Lin Han, Chenyu You, and Bo Wang. Segment anything in medical images. *Nature Communications*, 15(1):654, 2024. 2, 4, 6
- [32] Teresa Mendonça, M Celebi, T Mendonca, and J Marques. Ph2: A public database for the analysis of dermoscopic images. *Dermoscopy image analysis*, 2, 2015. 6
- [33] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, et al. Simple open-vocabulary object detection. In *European conference on computer vision*, pages 728–755. Springer, 2022. 6
- [34] Peter Naylor, Marick Laé, Fabien Reyal, and Thomas Walter. Segmentation of nuclei in histopathology images by deep regression of the distance map. *IEEE TMI*, 38(2):448–459, 2018. 4
- [35] Anna Pawłowska, Anna Ćwierz-Pieńkowska, Agnieszka Domalik, Dominika Jaguś, Piotr Kasprzak, Rafał Matkowski, Łukasz Fura, Andrzej Nowicki, and Norbert Żólek. Curated benchmark dataset for ultrasound based breast lesion analysis. *Scientific Data*, 11(1):148, 2024. 6
- [36] Gašper Podobnik, Bulat Ibragimov, Elias Tappeiner, Chanwoong Lee, Jin Sung Kim, Zacharia Mesbah, Romain Modzelewski, Yihao Ma, Fan Yang, Mikolaj Rudecki, et al. Han-seg: The head and neck organ-at-risk ct and mr segmentation challenge. *Radiotherapy and Oncology*, 198:110410, 2024. 4
- [37] Konstantin Pogorelov, Kristin Ranheim Randel, Carsten Grutowicz, Sigrun Losada Eskeland, Thomas de Lange, Dag Johansen, Concetto Spampinato, Duc-Tien Dang-Nguyen, Mathias Lux, Peter Thelin Schmidt, et al. Kvasir: A multi-class image dataset for computer aided gastrointestinal disease detection. In *Proceedings of the 8th ACM on Multimedia Systems Conference*, pages 164–169, 2017. 4
- [38] Prasanna Porwal, Samiksha Pachade, Manesh Kokare, Girish Deshmukh, Jaemin Son, Woong Bae, Lihong Liu, Jianzong Wang, Xinhui Liu, Liangxin Gao, et al. Idrid: Diabetic retinopathy–segmentation and grading challenge. *Medical image analysis*, 59:101561, 2020. 4
- [39] Chongyu Qu, Tiezheng Zhang, Hualin Qiao, Yucheng Tang, Alan L Yuille, Zongwei Zhou, et al. Abdomenatlas-8k: Annotating 8,000 ct volumes for multi-organ segmentation in three weeks. *NeurIPS*, 36, 2024. 4
- [40] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, pages 779–788, 2016. 1, 3
- [41] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE TPAMI*, 39(6):1137–1149, 2016. 1
- [42] Tianhe Ren, Yihao Chen, Qing Jiang, Zhaoyang Zeng, Yuda Xiong, Wenlong Liu, Zhengyu Ma, Junyi Shen, Yuan Gao, Xiaoke Jiang, et al. Dino-x: A unified vision model for open-world object detection and understanding. *arXiv preprint arXiv:2411.14347*, 2024. 3
- [43] Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Mohamed Amin, Le Hou, Kevin Clark, Stephen R Pfahl, Heather Cole-Lewis, et al. Toward expert-level medical question answering with large language models. *Nature Medicine*, pages 1–8, 2025. 2
- [44] Jayanthi Sivaswamy, S Krishnadas, Arunava Chakravarty, G Joshi, A Syed Tabish, et al. A comprehensive retinal image dataset for the assessment of glaucoma from the optic nerve head analysis. *JSM Biomedical Imaging Data Papers*, 2(1):1004, 2015. 6
- [45] Waqas Sultani et al. Towards low-cost and efficient malaria detection. In *CVPR*, pages 20655–20664. IEEE, 2022. 4
- [46] Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammad Bagheri, and Ronald M Summers. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *CVPR*, 2017. 4
- [47] Jakob Wasserthal, Hanns-Christian Breit, Manfred T Meyer, Maurice Pradella, Daniel Hinck, Alexander W Sauter, Tobias Heye, Daniel T Boll, Joshy Cyriac, Shan Yang, et al. Totalsegmentator: robust segmentation of 104 anatomic structures. In *CVPR*, 2024. 4

- tures in ct images. *Radiology: Artificial Intelligence*, 2023.
- [4](#), [6](#)
- [48] Biao Wu, Yutong Xie, Zeyu Zhang, Jinchao Ge, Kaspar Yaxley, Suzan Bahadir, Qi Wu, Yifan Liu, and Minh-Son To. Bhsd: A 3d multi-class brain hemorrhage segmentation dataset. In *MLMI*, pages 147–156. Springer, 2023. [4](#)
- [49] Ting Xu, Yanlin Xiang, Junliang Du, and Hanchao Zhang. Cross-scale attention and multi-layer feature fusion yolov8 for skin disease target detection in medical images. *Journal of Computer Technology and Software*, 4(2), 2025. [2](#)
- [50] Lewei Yao, Renjie Pi, Jianhua Han, Xiaodan Liang, Hang Xu, Wei Zhang, Zhenguo Li, and Dan Xu. Detclipv3: Towards versatile generative open-vocabulary object detection. In *CVPR*, pages 27391–27401, 2024. [3](#)
- [51] Yuhang Zang, Wei Li, Kaiyang Zhou, Chen Huang, and Chen Change Loy. Open-vocabulary detr with conditional matching. In *European conference on computer vision*, pages 106–122. Springer, 2022. [6](#)
- [52] Sheng Zhang, Yanbo Xu, Naoto Usuyama, Hanwen Xu, Jaspreet Bagga, Robert Tinn, Sam Preston, Rajesh Rao, Mu Wei, Naveen Valluri, et al. Biomedclip: a multimodal biomedical foundation model pretrained from fifteen million scientific image-text pairs. *preprint arXiv:2303.00915*, 2023. [2](#), [3](#), [5](#)
- [53] Theodore Zhao, Yu Gu, Jianwei Yang, Naoto Usuyama, Ho Hin Lee, Tristan Naumann, Jianfeng Gao, Angela Crabtree, Jacob Abel, Christine Moung-Wen, et al. Biomedparse: a biomedical foundation model for image parsing of everything everywhere all at once. *arXiv preprint arXiv:2405.12971*, 2024. [2](#), [4](#), [5](#), [6](#), [7](#)
- [54] Yian Zhao, Wenyu Lv, Shangliang Xu, Jinman Wei, Guanzhong Wang, Qingqing Dang, Yi Liu, and Jie Chen. Detrs beat yolos on real-time object detection. In *CVPR*, 2024. [1](#)
- [55] Xingyi Zhou et al. Detecting twenty-thousand classes using image-level supervision. In *ECCV*. Springer, 2022. [3](#)
- [56] Xiahai Zhuang. Multivariate mixture model for myocardial segmentation combining multi-source images. *IEEE TPAMI*, 41(12):2933–2946, 2018. [4](#)