

Bridging the Dimensionality Gap: A Scoping Review of 3D Vision-Language Models in Medical Imaging

Anonymous ACL submission

Abstract

While large multimodal models have achieved remarkable success in general domains, adapting them to medical imaging faces a fundamental dimensionality gap: transition from 2D snapshots to 3D volumetric data (e.g., CT, MRI). This review systematically examines 37 studies on 3D Vision-Language Models in healthcare, capturing the rapid research surge of 2024-2025 in this emerging field. We further propose a categorical framework that classifies these studies by multimodal fusion mechanisms (projection, attention, and adapters), volumetric encoding strategies (slice-based vs. native 3D), and language processing (encoder vs. foundation model). Our analysis highlights the growing adoption of parameter-efficient fine-tuning due to computational constraints, but significant challenges remain, including hallucinations, a lack of spatial grounding, and misalignment between evaluation metrics and clinical utility. This survey aims to clarify current methodologies and identify future directions in volumetric medical AI.

1 Introduction

The integration of Computer Vision (CV) and Natural Language Processing (NLP) has fundamentally reshaped artificial intelligence research, enabling models to jointly reason over visual and language modalities (Karpathy and Fei-Fei, 2015; Lu et al., 2019; Tan and Bansal, 2019; Radford et al., 2021). General-domain Foundation Models (FMs) and Vision-Language Models (VLMs), such as CLIP (Radford et al., 2021) and GPT-4V (OpenAI et al., 2023), have demonstrated exceptional performance in cross-modal representation learning and generation. Despite this progress, a critical bottleneck limits their adoption in clinical settings: the dimensionality mismatch between general-purpose VLMs and medical imaging data. Most existing VLMs are designed for 2D images (Lu et al., 2019; Tan and Bansal, 2019; Radford et al., 2021),

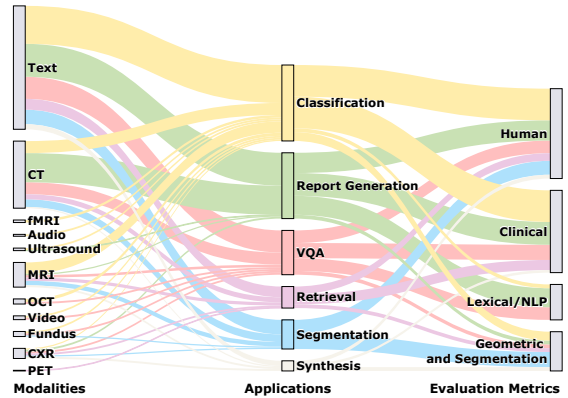


Figure 1: The relationships between data modality, downstream vision-language task, and evaluation metrics examined in this review.

whereas key diagnostic modalities, such as Computed Tomography (CT) and Magnetic Resonance Imaging (MRI), produce three-dimensional volumes (Milletari et al., 2016). Modeling volumetric data exponentially increases computational complexity and requires preserving spatial and depth information, which is often lost in 2D pre-trained encoders (Yang et al., 2021; Zhang et al., 2021). Consequently, simply extending 2D architectures to 3D medical imaging often yields degraded performance or unjustifiable computational costs.

In response to this challenge, recent research activity has rapidly expanded around 3D Medical VLMs (hereafter 3D-MedVLMs) (Lu et al., 2025; Chang et al., 2025; Xin et al., 2025) (Figure 1). Proposed solutions range from decomposing 3D volumes into pseudo-2D sequences to developing native volumetric encoders and utilizing Large Language Models (LLMs) via lightweight adaptation mechanisms (Dack et al., 2023; Kakkar et al., 2024; Liu et al., 2025; Luo et al., 2025). However, these efforts remain fragmented across design choices, such as different feature fusion, pre-training objectives, and evaluation metrics, making systematic comparison difficult.

068 Unlike general or 2D-focused surveys (Caffagni
 069 et al., 2024; Sim et al., 2025), this review fo-
 070 cuses on the unique challenges of 3D-MedVLMs.
 071 Furthermore, we distinguish our work from the
 072 clinically-oriented roadmap of Wu et al. (2025) by
 073 identifying a broader set of studies and by focusing
 074 on the unique methodological and practical chal-
 075 lenges posed by volumetric data in clinical settings.

076 Our contributions are threefold: (i) Systematic
 077 Classification: We provide a framework to categor-
 078 ize 3D-MedVLMs based on vision encoding and
 079 multimodal fusion techniques (Figure 4). (ii) Rea-
 080 soning Analysis: We analyze how recent models
 081 capture intrinsic medical dependencies. (iii) Clini-
 082 cal Alignment: We identify persistent challenges
 083 in clinical evaluation, highlighting the discrepancy
 084 between standard NLP metrics and clinical factual
 085 correctness for 3D-MedVLMs.

086 2 Scope and Literature Selection

087 Our scoping review adheres to the Preferred Re-
 088 porting Items for Systematic reviews and Meta-
 089 Analyses extension for Scoping Reviews (Fig-
 090 ure 2) (Tricco et al., 2018). We searched 4
 091 databases: PubMed¹, IEEE Xplore², ACM Digital
 092 Library³, and ACL Anthology⁴, covering studies
 093 published from 2021 to 2025. Our search strategy
 094 was designed to systematically identify studies at
 095 the intersection of volumetric imaging and mul-
 096 timodal learning. We employed Boolean queries
 097 that combined three primary concept groups: (i) 3D
 098 imaging modalities, (ii) multimodal architectures,
 099 and (iii) clinical text or downstream tasks. The
 100 complete list of database-specific search queries is
 101 detailed in Appendix A.

102 We then screened the articles based on their ti-
 103 tles and abstracts according to predefined criteria.
 104 Inclusion criteria required studies (i) were pub-
 105 lished in English, (ii) underwent peer review, and
 106 (iii) presented and benchmarked a model. Exclu-
 107 sion criteria filtered out (i) studies unrelated to 3D-
 108 MedVLMs, (ii) non-English publications, and (iii)
 109 secondary literature such as systematic reviews,
 110 case studies, or descriptive studies lacking experi-
 111 mental results. After screening, we extracted meta-
 112 data from each included study, such as models,
 113 datasets, applications, results, and reported limi-
 114 tations. To ensure accuracy, two annotators cross-

¹<https://pubmed.ncbi.nlm.nih.gov/>

²<https://ieeexplore.ieee.org/>

³<https://dl.acm.org/>

⁴<https://aclanthology.org/>

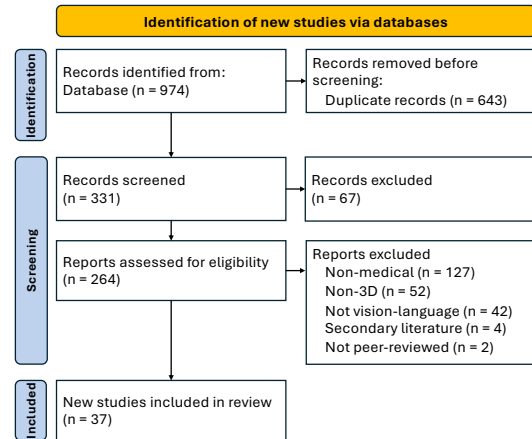


Figure 2: PRISMA-ScR flow diagram.

115 verified the study selection and metadata extraction,
 116 with a third one resolving any disagreements.

117 From an initial set of 974 papers retrieved from
 118 databases, we removed 643 duplicates. During the
 119 initial screening, 6 papers were excluded, and 3
 120 additional papers were excluded for misaligned ob-
 121 jectives or lack of relevance to 3D medical VLMs.
 122 Ultimately, 37 studies met the inclusion criteria
 123 and form the basis of this review. Important meta-
 124 data, including resource links, architectures, train-
 125 ing strategies, and downstream tasks for the studies,
 126 are presented in Table 1.

127 3 Multimodal Fusion

128 The core challenge in developing 3D-MedVLMs
 129 lies in aligning 3D volumetric data with textual rep-
 130 resentations. Existing approaches can be broadly
 131 categorized into three interaction paradigms, dis-
 132 tinguished by the level at which visual and textual
 133 modalities intersect (Figure 3).

134 3.1 Global Alignment (Dual-Encoder)

135 Global alignment strategies prioritize computa-
 136 tional efficiency by encoding visual and textual
 137 modalities in separate streams that interact only
 138 at the final embedding stage. This Dual-Encoder
 139 approach is well-suited to zero-shot classification
 140 and retrieval. Within this paradigm, methods differ
 141 in how cross-modal interaction is realized.

142 Projection-based alignment, popularized by
 143 CLIP (Radford et al., 2021), employs lightweight
 144 linear layers to map modalities into a shared la-
 145 tent space. Extensions such as RadCLIP (Lu et al.,
 146 2025) introduce slice-pooling adapters that aggre-
 147 gate 2D slice embeddings into a unified volumetric
 148 representation, enabling global image-text match-

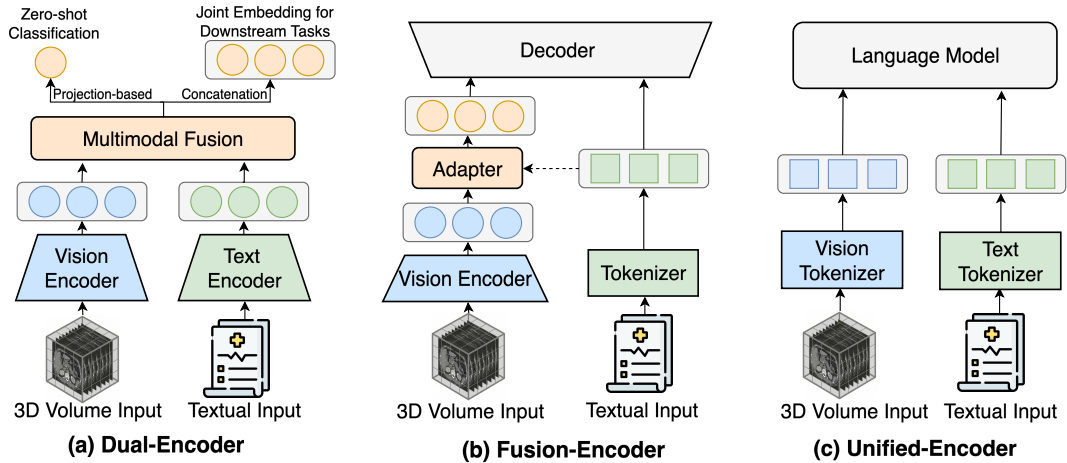


Figure 3: Overview of vision-language model architectures. This figure illustrates three architectural paradigms for studies integrating 3D medical imaging and clinical text. a, **Dual-Encoder** employs separate encoders to process visual and textual data in parallel, followed by multimodal fusion to achieve global alignment (e.g., CLIP (Radford et al., 2021) and ALIGN (Jia et al., 2021)). b, **Fusion-Encoder** utilizes a vision encoder and adapter module to connect vision embeddings with tokenized text for unified decoding (e.g., LLaVA (Liu et al., 2023), BLIP-2 (Li et al., 2023), and Flamingo (Alayrac et al., 2022)). c, **Unified-Encoder** bypasses complex fusion mechanisms by projecting both visual and textual inputs into a unified language model token space (e.g., Med3DVLM (Xin et al., 2025) and RadGenome-ChestCT (Zhang et al., 2025)).

149 ing. Similarly, HCL-AL (Gao et al., 2025) aligns
 150 3D features with text through learned projections
 151 for anatomical retrieval. Concatenation-based
 152 fusion, on the other hand, directly concatenates
 153 features from independent encoders before a task-
 154 specific head. Zhao et al. (2025) follows this ap-
 155 proach by joining ultrasound video features with
 156 clinical text embeddings to predict thyroid nodule
 157 invasiveness, while MultiModalGAN (Ghadekar
 158 et al., 2025) concatenates CNN image features with
 159 BERT embeddings to condition generative adver-
 160 sarial networks for image synthesis.

161 While Dual-Encoder offers computational effi-
 162 ciency, it often sacrifices fine-grained contextual
 163 information (Wang et al., 2022).

164 3.2 Dense Interaction (Fusion-Encoder)

165 Dense interaction paradigms address the above lim-
 166 itation by integrating visual and textual modal-
 167 ities at intermediate layers, enabling discriminative
 168 tasks such as segmentation and grounding. These
 169 Fusion-Encoder architectures leverage attention-
 170 based deep fusion to capture token-level interac-
 171 tions between visual and textual sequences.

172 Such a strategy is essential for tasks that re-
 173 quire precise semantic grounding. For example,
 174 Ye et al. (2024) employs bi-attention blocks that
 175 combine self-attention over instructions with cross-
 176 attention over MRI features to drive text-prompted

177 segmentation. Some architectures further adopt
 178 the Encoder-Decoder framework, pairing a fusion
 179 encoder with a shared language decoder (e.g., VL-
 180 BERT (Liu et al., 2021b)).

181 Despite their effectiveness in localization and
 182 grounding, Fusion-Encoder incurs higher computa-
 183 tional overhead than Dual-Encoder (Li et al., 2022).

184 3.3 Generative Alignment (Unified-Encoder)

185 Unlike previous paradigms that coordinate sepa-
 186 rate encoders, Generative Alignment harnesses
 187 the unified processing capabilities of foundation
 188 models to jointly handle visual and textual modal-
 189 ities. These Unified-Encoder models project 3D
 190 visual features directly into the token space of an
 191 LLM, thereby eliminating explicit fusion mecha-
 192 nisms and capitalizing on pre-trained foundation
 193 models.

194 In this approach, the “architecture” is defined by
 195 the adapter-based mechanism that bridges modal-
 196 ities. For example, Med3DVLM (Xin et al.,
 197 2025) maps 3D visual features into a Qwen-based
 198 LLM (Qwen et al., 2024), enabling joint reason-
 199 ing over 3D visual tokens and textual instructions.
 200 Given the high cost of full fine-tuning, frameworks
 201 like RadGenome-ChestCT (Zhang et al., 2025) em-
 202 ploy Parameter-Efficient Fine-Tuning (PEFT) by
 203 inserting trainable bottleneck adapters into frozen
 204 Llama (Grattafiori et al., 2024) backbones, thereby

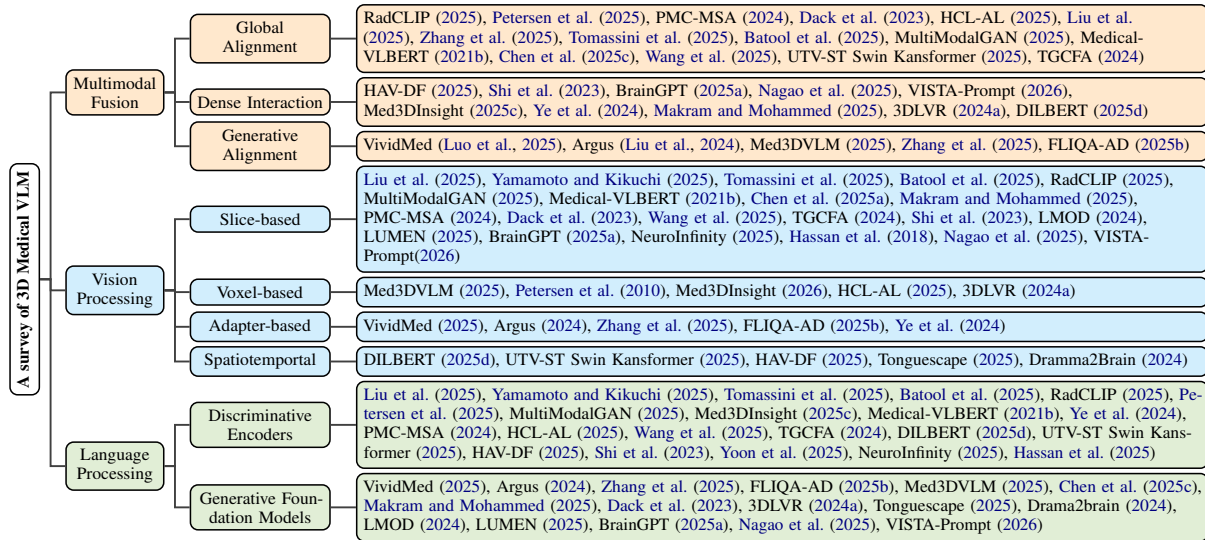


Figure 4: Categorical framework of 3D medical VLM research and representative studies included in this work.

preserving general-world knowledge in LLM while adapting to 3D medical data. Hybrid approaches such as FLIQA-AD (Chen et al., 2025b) integrate 3D adapters directly within Vision Transformers before feeding features into text generators such as FLAN-T5 (Chung et al., 2022).

While Unified-Encoder offers strong flexibility and performance, it remains susceptible to hallucination when the LLM relies on parametric knowledge rather than visual input (Li et al., 2023).

4 Single Modal Representation

While multimodal fusion enables cross-modal reasoning, the efficacy of 3D-MedVLMs relies on robust feature processing. The primary challenge lies in the vision processing, where models must efficiently process high-dimensional 3D data. In parallel, language encoders require strategies tailored to specific downstream tasks. This section examines the encoding mechanisms for them.

4.1 Vision Processing

The dimensionality mismatch between 2D models and 3D volumetric medical data (CT, MRI, PET) has led to four main encoding strategies.

Slice-based Feature Extraction decomposes volumes into 2D slices processed by 2D convolutional or Transformer pretrained on large datasets (Deng et al., 2009; Mei et al., 2022). It is computationally efficient at the cost of full-volumetric context, or when aggregating 2D slices provides sufficient contextual information (Lu et al., 2025; Kunhimon et al., 2024).

Native 3D architectures operate directly on voxel patches to capture the intricate spatial dependencies inherent in volumetric organs. This strategy is essential for tasks requiring precise anatomical localization along the Z-axis but entails high memory and computational demands. To overcome these challenges, Chen et al. (2025c); Xin et al. (2025) incorporate specialized transformer-based models to efficiently capture volumetric features.

Adapter-Based Strategies extends 2D encoders for volumetric data by combining 2D pretraining with 3D spatial awareness (Chen et al., 2025b). These approaches “inflate” pretrained 2D backbones by inserting lightweight adapters, allowing 3D input processing while retaining large-scale 2D pretraining, as seen in FLIQA-AD (Chen et al., 2025b) and VividMed (Luo et al., 2025).

Spatiotemporal Modeling treats depth or time as a sequence, and applies specialized attention or convolutional mechanisms to capture temporal continuity (Chen et al., 2025d). For example, Zhao et al. (2025) applies a Video Swin (Liu et al., 2021c) backbone to ultrasound videos to learn temporal consistencies for predicting thyroid nodule invasiveness. However, spatiotemporal modeling remains relatively unexplored in medicine.

4.2 Language Processing

Clinical text processing in 3D-MedVLMs has progressed from simple semantic-embedding extraction to complex autoregressive reasoning, reflecting a shift from static image-text alignment to instruction-following and generative capabilities. Existing studies identify two primary paradigms.

Discriminative Encoders are suitable for segmentation, and classification tasks. These encoder-only models map clinical text to a fixed latent space (Liu et al., 2025), primarily using transformer-based encoders such as BERT and CLIP (Liu et al., 2021b, 2025). Beyond free-text, Petersen et al. (2025) encodes tabular health data into natural-language descriptions, thereby enabling alignment with 3D MRI features.

Generative Foundation Models have recently been adopted to address complex language understanding and reasoning tasks, such as report generation and visual question answering (VQA). Med3DVLM (Xin et al., 2025) exemplifies this line of work by integrating Qwen2.5-7B (Qwen et al., 2024), allowing the system to process 3D visual tokens alongside textual instructions. Similarly, Argus (Liu et al., 2024), RadGenome-ChestCT (Zhang et al., 2025), and BrainGPT (Li et al., 2025a) adopt Llama models (Grattafiori et al., 2024) as their linguistic backbone, leveraging their broad real-world knowledge to generate detailed, contextually accurate radiology reports.

5 Training Strategies

Pre-training. 3D-MedVLMs’s capacity to interpret clinical text is established through pretraining objectives. Two primary pretraining strategies are commonly used. (1) Contrastive Learning aligns the global image features with corresponding text reports. Standard CLIP uses softmax loss, which requires large batch sizes to provide sufficient negative samples. However, this is often impractical for 3D-MedVLMs due to GPU memory constraints (Li et al., 2019). To address this, approaches such as SigLIP employ a pairwise Sigmoid loss that decouples performance from batch size (Xin et al., 2025), enabling effective alignment even with small batch sizes. However, such global alignment methods may associate a medical term (e.g., “pneumonia”) with broad visual patterns (e.g., general hazy texture), without accurately localizing the pathology. (2) Masked Signal Modeling (MSM) complements contrastive learning by learning dense, local dependencies through reconstructing masked portions of input. This encourages the encoder to model the underlying structure of 3D volumes or clinical text rather than solely on high-level semantics (Chen et al., 2025d; Ye et al., 2024). However, MSM can be computationally expensive for volumetric data and may overfit to low-level

noise if not combined with a discriminative loss. Still, it provides the fine-grained spatial precision that is lost in global contrastive methods.

Fine-Tuning. As 3D-MedVLMs scales to billions of parameters, full model fine-tuning becomes computationally burdensome and prone to forgetting. Parameter-Efficient Fine-Tuning (PEFT) addresses this by only updating a fraction of the model’s weights. Additive methods introduce new trainable parameters to fixed backbones. Common approaches include Adapters, which inject bottleneck layers (Chen et al., 2025b), and Head Tuning, which fine-tunes only the lightweight decision heads (Luo et al., 2025). However, these methods rely heavily on the relevance of the frozen features (Liu et al., 2021b). Selective methods fine-tune a subset of existing parameters. Approaches such as Medical-VLBERT (Liu et al., 2021b) employ alternate learning to selectively optimize the language decoder separate from joint multimodal updates. However, the performance of this strategy is highly sensitive to the training schedule, requiring a balanced ratio of pretraining to transferring iterations to prevent performance degradation. Finally, Reparameterization methods (e.g., LoRA) constrain optimization to a low-rank subspace to reduce memory consumption while maintaining fine-tuning efficiency (Liu et al., 2024). However, shared low-rank constraints often fail to capture unique intra-group heterogeneity, leading to demographic fairness disparities (Li et al., 2025b).

6 Reasoning

3D-MedVLMs require advanced reasoning capabilities. Here, we summarize how included studies capture the intrinsic dependencies of medical data.

Spatial and Topological Reasoning. Standard global embeddings often discard the anatomical hierarchy and connectivity required for advanced reasoning. Recent methods address this by explicitly encoding 3D geometry and enforcing fine-grained visual-textual alignment. Du et al. (2024a) introduces a Topological Feature Extractor based on persistent homology to capture structural invariance beyond pixel-level features. Gao et al. (2025) proposes Hierarchical Contrastive Learning to enforce semantic alignment across coarse-to-fine anatomical hierarchies and establishes a broad anatomical context (e.g., body region) before refining its focus to specific substructures. VividMed (Luo et al.,

2025) further bridges semantic understanding and localization by integrating a VLM with a localization module that utilizes the hidden states of special enclosure tokens to generate segmentation masks and bounding boxes for both 2D and 3D inputs.

Temporal and Longitudinal Reasoning. Diagnostic reasoning also requires capturing both the longitudinal progression of pathology and the spatiotemporal continuity of real-time imaging. 3DVLR (Du et al., 2024a) models “evolve anatomical landscapes” by using optimal transport theory to align feature distributions across different time stamps, thus distinguishing stable anatomical structures from changing pathological conditions (e.g., tumor growth) and enabling the generation of narratives that reason about disease progression. For dynamic imaging, Chen et al. (2025d) propose DILBERT, which aggregates visual geometric features with named entity linguistic features across video frames and employs multimodal distillation to align dynamic visual evidence with diagnostic text, ensuring consistent interpretation of real-time scans such as ultrasound.

Clinical Chain-of-Thought (CoT) and Causal Reasoning. Medical diagnosis also involves a deductive process of hypothesis generation and verification, rather than single-step classification, motivating a shift towards CoT reasoning to mirror clinical workflows. Gai et al. (2025) trains models to output intermediate reasoning steps using datasets (e.g., R-RAD, R-SLAKE) annotated with reasoning traces. Holistic Narrative Synthesis extends reasoning to 3D data by structuring visual findings, such as lesion location and size, into logically flowing narratives. For example, BrainGPT (Li et al., 2025a) produced narratives rated indistinguishable from human-written reports in 74% of cases. Making up for the lack of reasoning-focused metrics, RaTEScore (Zhao et al., 2024) assesses factual correctness of key medical entities (e.g., diagnostic outcomes, anatomical details), and FORTE (Li et al., 2025a) captures report completeness.

7 Benchmark Dataset

In the 3D medical domain, datasets have historically lagged behind 2D counterparts. Our review highlights that 3D datasets are typically smaller, less accessible, and often lack associated reports in prior studies. To address this gap, increasing efforts have focused on curating large-scale public 3D image datasets with text annotations (Zhang

et al., 2025) or employing synthetic data augmentation (Ghadekar et al., 2025; Chen et al., 2025a) to increase both dataset size and diversity (Table 2).

Among these, CT-RATE (Hamamci et al., 2025a) is particularly noteworthy. With more than 50,000 volumes from over 20,000 patients, it offers a substantial scale advantage over other datasets. Additionally, text reports were authored by multiple radiologists with varied writing styles, further enhancing the dataset’s robustness.

Other datasets, such as Ichinose et al. (2023) and RadGenome-ChestCT (Zhang et al., 2025), address a different need: image grounding. While many datasets provide image-text pairs, such as a slice of the lungs with the caption “nodule in both lungs”, they often lack precise localization. Grounding-focused datasets address this gap, enabling models like VividMed (Luo et al., 2025) to produce both textual reports and corresponding bounding boxes or segmentation masks, thereby bridging the semantic gap between “what” and “where”.

Another group of datasets moves beyond simple answers to include rationales and explanations. While traditional datasets like VQA-RAD (Lau et al., 2018) and SLAKE (Liu et al., 2021a) included short-form answers, newer initiatives such as R-RAD and R-SLAKE have annotated the intermediate reasoning steps required for medical decision-making (Gai et al., 2025). Similarly, instruction-response datasets like M3D (Bai et al., 2024) have emerged to support complex reasoning.

In contrast to these text-rich initiatives, the visual backbones often rely on anatomical understanding through segmentation masks. Leading this category is TotalSegmentator (Wasserthal et al., 2023), which provides 5,000 CT volumes with comprehensive masks for 117 anatomical classes. Similarly, AbdomenAtlas (Li et al., 2025c), a large-scale dataset comprising over 8,000 volumes, and AbdomenCT-1K (Ma et al., 2022) exemplify ongoing efforts to scale segmentation benchmarks.

8 Evaluation Metrics

Evaluating 3D-MedVLMs requires a comprehensive assessment of technical performance, clinical utility, and robustness. Accordingly, metrics should jointly capture vision-task accuracy and text quality, including fidelity to visual evidence, linguistic coherence, and clinical relevance (Liu et al., 2024; Xin et al., 2025; Zhang et al., 2025).

Geometric and Segmentation Metrics. For

anatomical segmentation and visual grounding, geometric metrics quantify spatial accuracy. Volumetric overlap between predicted and ground-truth masks is commonly measured using the Dice Similarity Coefficient (DSC) and the Intersection over Union (IoU) (Ye et al., 2024; Wang et al., 2025). Boundary accuracy is assessed using the 95th-percentile Hausdorff Distance (HD95) and Mean Surface Distance (MSD), which capture worst-case and average deviations between predicted and reference contours, respectively (Ye et al., 2024). In addition, pixel-level metrics such as Root Mean Square Error (RMSE) assess decomposition precision (Wang et al., 2025).

Clinical (Retrieval, Classification, and Synthesis) Metrics. In image–text retrieval, Recall@K measures whether the correct reference appears among the top- K results (Xin et al., 2025). Classification performance is commonly assessed using Accuracy, F1-score, and AUC (Lu et al., 2025; Zhao et al., 2025). For text-guided image synthesis (e.g., generating tumor-bearing CT volumes from clinical descriptions), Root Mean Square Error (RMSE) and Pearson correlation are used to verify anatomical consistency between generated and reference images (Chen et al., 2025a). Additionally, in speech-oriented models, Medical Term Recognition Rate (MTRR) is used to ensure critical medical terminology is preserved (Lv et al., 2025).

Lexical/NLP Metrics. For report generation and VQA, standard automatic metrics such as BLEU, ROUGE, METEOR, and CIDEr are commonly used to quantify lexical similarity between generated text and reference reports. However, these surface-level metrics often fail to reflect the clinical correctness and diagnostic validity (Liu et al., 2024; Luo et al., 2025). To address these limitations, clinically-oriented evaluation measures have been increasingly adopted. Early work computes precision, recall, and F1 scores over extracted clinical findings using clinical labelers such as RadBERT (Yan et al., 2022). More recent metrics like GREEN (Ostmeier et al., 2024) and RaTEScore (Zhao et al., 2024) provide a finer-grained evaluation of radiological accuracy and report completeness (Liu et al., 2024).

Human Evaluation. Beyond automated metrics, human evaluation remains the gold standard for assessing generated clinical text, with experts rating generated reports for correctness, completeness, and conciseness (Batool et al., 2025). However, evaluation protocols vary: some employ Likert-

scale scoring to quantify clinical accuracy (Liu et al., 2024; Batool et al., 2025), whereas others utilize pairwise ranking in which annotators select the best report among competing systems (Liu et al., 2021b). More recently, LLMs have been proposed as automated judges, but evidence suggests that they may systematically favor AI-generated text over human reference, highlighting the continued need for human evaluation (Batool et al., 2025).

9 Applications

3D-MedVLMs have been adapted for a diverse array of clinical downstream tasks, which we categorize into five primary domains (Table 1).

Diagnostic Classification and Anatomical Identification is a primary application of 3D-MedVLMs in medical imaging. These models extend beyond binary detection to recognize fine-grained recognition of disease subtypes and spatial localization (Gao et al., 2025; Lu et al., 2025). For example, RadCLIP (Lu et al., 2025), a contrastive-based model, demonstrates strong zero-shot generalization. To address data scarcity, Petersen et al. (2025) shows that scaling the number of negative pairs enables training with 62 scans. In contrast, task-specific architectures prioritize fine-grained optimization and yield higher performance on targeted tasks (Zhao et al., 2025; Wang et al., 2025).

Automated Report Generation converts visual findings to clinical narratives. Models such as 3DLVR (Du et al., 2024a) and Medical-VLBERT (Liu et al., 2021b) align visual and linguistic features across modalities. Recent work has extended report generation to dynamic or time-sensitive applications, including spine CT (Batool et al., 2025) and emergency head CT (Tomassini et al., 2025). A central challenge is hallucination, where the model generates clinically incorrect statements. Consequently, evaluation practices have shifted from surface-level lexical metrics to clinically grounded measures (Liu et al., 2024).

Semantic Segmentation and Visual Grounding seeks to align textual findings with anatomical coordinates in volumetric images, which can improve spatial localization and downstream task performance when incorporated in models (Luo et al., 2025). When traditional visual prompts (e.g., points or bounding boxes) are ineffective, particularly for anatomies with sparse or fragmented tissue, text-prompted pretraining has been proposed to unify heterogeneous datasets and guide segmen-

571 tation via semantic descriptions (Ye et al., 2024).

572 **Image Synthesis** generates anatomically accu- 621
573 rate visual pathologies from textual clinical descrip- 622
574 tions. By conditioning generative models on struc- 623
575 tured metadata, these approaches help mitigate data 624
576 scarcity while maintaining adherence to medical 625
577 guidelines. For example, diffusion-based methods 626
578 have been used to synthesize renal tumors from 627
579 nephrometry scores, allowing granular control over 628
580 tumor morphology (Chen et al., 2025a). Similarly, 629
581 adversarial frameworks have been applied to cross- 630
582 modality tasks; MultiModalGAN synthesizes chest 631
583 X-rays directly from text reports, while CycleGAN 632
584 has been utilized to convert MRI scans into CT-like 633
585 images, providing diagnostic alternatives when CT 634
586 is unavailable (Ghadekar et al., 2025). 635

587 **Retrieval Tasks** have recently shifted from 636
588 purely technical optimization to enabling clini- 637
589 cally robust applications and overcoming data 638
590 scarcity. Med3DVLM (Xin et al., 2025) demon- 639
591 strated that using pairwise sigmoid loss during con- 640
592 trastive learning removed the need for large neg- 641
593 ative batches and outperformed the SOTA model 642
594 by 19%. Complementing these model-centric ad- 643
595 vances, Yamamoto and Kikuchi (2025) tackled the 644
596 bottleneck of dataset construction by incorporating 645
597 lesion- and organ-aware supervision. Their method 646
598 achieved a Top-1 accuracy of 51.7%, demonstrat- 647
599 ing the feasibility of slice-based retrieval for au- 648
600 tomating the curation of large-scale medical VQA 649
601 datasets from routine clinical archives. 650

602 10 Challenges and Future Direction 651

603 While 3D medical VLMs have demonstrated im- 652
604 pressive capabilities, they face challenges before 653
605 they can be widely adopted as clinical assistants. 654
606 We identify five key areas where future efforts need 655
607 be directed: reliability, interpretability, data ecosys- 656
608 tems, clinical alignment, and scalability. 657

609 The major barrier to deploying 3D-MedVLMs 658
610 is ensuring validity and reliability. Current 659
611 models remain prone to hallucinations, where 660
612 models generate plausible but incorrect medical as- 661
613 sertions (Qin et al., 2024). This issue is exacerbated 662
614 in 3D-MedVLMs because radiologists often rely 663
615 on external information (e.g., prior reports), and 664
616 models trained on such data are more prone to hallu- 665
617 cination (Luo et al., 2025). Moreover, maintaining 666
618 precise feature alignment in complex anatomical 667
619 scenarios is also challenging, as simple approaches 668
620 often struggle with large 3D batches or rich cross- 669

621 modal interactions (Xin et al., 2025). Furthermore, 622
623 models often degrade under domain shift, with 624
625 variations in scanner protocols or imaging artifacts 626
627 significantly altering model behavior. Future work 628
629 should focus on robustness to these shifts, ensuring 630
631 models can generalize from curated datasets to the 632
633 real clinical environments (Hamamci et al., 2025b). 634

635 Clinical deployment requires transparency 636
637 and explainability, as trust relies on the ability 638
639 to audit the decision-making process. Future archi- 640
641 tectures may emphasize glass-box designs to 642
643 offer transparent, guideline-aligned reasoning (Gai 644
645 et al., 2025). Promising directions include chain-of- 646
647 radiology-thought prompting combined with visual 647
648 grounding that links textual findings to voxel-level 648
649 evidence (Luo et al., 2025). Equally important is in- 649
650 tegrating uncertainty to identify low-confidence 650
651 predictions and reduce overconfident errors. 651

652 The third priority is building healthy data 653
654 ecosystems. Although hospitals generate abun- 654
655 dant raw data, high-quality, instruction-tuned 3D 655
656 pairs remain scarce. Future efforts should focus on 656
657 self-sustaining automatic data pipelines that 657
658 reduce reliance on manual annotation, for example, 658
659 using LLMs to transform noisy clinical reports into 659
660 machine-readable formats at scale (Zhang et al., 660
661 2025). Simultaneously, data access and diver- 661
662 sity constraints must be addressed, and synthetic 662
663 data generation may be leveraged to augment 663
664 rare pathologies, thereby improving model robust- 664
665 ness and fairness for long-tail conditions (Ghadekar 665
666 et al., 2025; Chen et al., 2025a; Lv et al., 2025). 666

667 Another direction is clinical alignment, as 667
668 a gap exists between standard technical evaluation 668
669 and clinical utility (Liu et al., 2024; Luo et al., 669
670 2025). Qin et al. (2024) shows that high scores on 670
671 lexical metrics (e.g., BLEU, ROUGE) do not guar- 671
672 antee clinically accurate reports, obscuring true 672
673 model readiness. To address this, the field should 673
674 establish standardized, clinically rigorous bench- 674
675 marks that reflect real-world workflows (e.g., diag- 675
676 nostic efficiency and clinician acceptance). 676

677 Finally, scalability remains a major bottle- 677
678 neck. The high computational cost of 3D data pro- 678
679 cessing limits the deployment of 3D-MedVLMs 679
680 in resource-constrained clinical environments (Xin 680
681 et al., 2025). Strategies such as downsampling 681
682 reduce computational overhead but inevitably sac- 682
683 rifice spatial details (Xin et al., 2025). Future work 683
684 should explore parameter-efficient mechanisms to 684
685 maintain diagnostic fidelity while lowering the 685
686 computational overhead (Chen et al., 2025b). 686

673 Limitations

674 While this survey provides a broad and systematic
675 overview of recent VLLM research, our analysis
676 is necessarily qualitative rather than exhaustively
677 quantitative. The rapid evolution and wide scope
678 of the field make detailed head-to-head benchmark-
679 ing difficult, particularly given the heterogeneity
680 of evaluation metrics, datasets, and experimental
681 settings across studies. As a result, performance
682 comparisons are presented primarily in terms of
683 high-level trends rather than rigorous quantitative
684 meta-analysis. In addition, our review is limited
685 to English-language publications and prioritizes
686 peer-reviewed work from major conferences such
687 as ACL, CVPR, and MICCAI. Also, influential
688 arXiv preprints and some emerging or non-English
689 contributions may have been overlooked. Finally,
690 the fast pace of industrial and academic progress
691 means that certain cutting-edge methods may not
692 yet be fully captured. In future work, we plan to
693 incorporate more quantitative meta-analysis where
694 feasible, integrate insights from large-scale indus-
695 trial deployments, and continuously update the sur-
696 vey to reflect ongoing advances in this rapidly de-
697 veloping field.

698 References

699 Michael J Ackerman. 2022. [The visible human project](#).
700 *Inf. Serv. Use*, 42(1):129–136.

701 Parnian Afshar, Shahin Heidarian, Nastaran Enshaei,
702 Farnoosh Naderkhani, Moezedin Javad Rafiee, Anas-
703 tasia Oikonomou, Faranak Babaki Fard, Kaveh
704 Samimi, Konstantinos N Plataniotis, and Arash Mo-
705 hammadi. 2021. [COVID-CT-MD, COVID-19 com-
706 puted tomography scan dataset applicable in machine
707 learning and deep learning](#). *Sci. Data*, 8(1):121.

708 Walid Al-Dhabyani, Mohammed Gomaa, Hussien
709 Khaled, and Aly Fahmy. 2020. [Dataset of breast
710 ultrasound images](#). *Data Brief*, 28(104863):104863.

711 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, An-
712 toine Miech, Iain Barr, Yana Hasson, Karel Lenc,
713 Arthur Mensch, Katie Millican, Malcolm Reynolds,
714 Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda
715 Han, Zhitao Gong, Sina Samangooei, Marianne Mon-
716 teiro, Jacob Menick, Sebastian Borgeaud, and 8 oth-
717 ers. 2022. [Flamingo: A visual language model for
718 few-shot learning](#). In *Proceedings of the 36th Interna-
719 tional Conference on Neural Information Processing
720 System*, pages 23716–23736.

721 Michela Antonelli, Annika Reinke, Spyridon Bakas,
722 Keyvan Farahani, Annette Kopp-Schneider, Ben-
723 nett A Landman, Geert Litjens, Bjoern Menze, Olaf

Ronneberger, Ronald M Summers, Bram van Gin- 724
neken, Michel Bilello, Patrick Bilic, Patrick F Christ, 725
Richard K G Do, Marc J Gollub, Stephan H Heckers, 726
Henkjan Huisman, William R Jarnagin, and 39 oth- 727
ers. 2022. [The medical segmentation decathlon](#). *Nat.*
Commun., 13(1):4128. 728

Fan Bai, Yuxin Du, Tiejun Huang, Max Q-H Meng, and 730
Bo Zhao. 2024. [M3D: Advancing 3D medical image
analysis with multi-modal large language models](#).
arXiv [cs.CV]. 731
732
733

Humaira Batool, Asmat Mukhtar, Sajid Gul Khawaja, 734
Norah Saleh Alghamdi, Asad Mansoor Khan, Adil 735
Qayyum, Ruqqayia Adil, Zawar Khan, Muhammad 736
Usman Akram, Muhammad Usman Akbar, and An- 737
ders Eklund. 2025. [Knowledge distillation and
transformer-based framework for automatic spine CT
report generation](#). *IEEE Access*, 13:42949–42964. 738
739
740

Asma Ben Abacha, Sadid A Hasan, Vivek V Datla, Joey 741
Liu, Dina Demner-Fushman, and Henning Muller. 742
2024. [VQA-med: Overview of the medical visual
question answering task at ImageCLEF 2019](#). 743
744

Olivier Bernard, Alain Lalande, Clement Zotti, Freder- 745
ick Cervenansky, Xin Yang, Pheng-Ann Heng, Irem 746
Cetin, Karim Lekadir, Oscar Camara, Miguel Angel 747
Gonzalez Ballester, Gerard Sanroma, Sandy Napel, 748
Steffen Petersen, Georgios Tziritas, Elias Grinias, 749
Mahendra Khened, Varghese Alex Kollerathu, Gana- 750
pathy Krishnamurthi, Marc-Michel Rohe, and 18 751
others. 2018. [Deep learning techniques for automatic
MRI cardiac multi-structures segmentation and di-
agnosis: Is the problem solved?](#) *IEEE Trans. Med.*
Imaging, 37(11):2514–2525. 752
753
754
755

Patrick Bilic, Patrick Christ, Hongwei Bran Li, Eugene 756
Vorontsov, Avi Ben-Cohen, Georgios Kaissis, Adi 757
Szeskin, Colin Jacobs, Gabriel Efrain Humpire Ma- 758
mani, Gabriel Chartrand, Fabian Lohöfer, Julian Wal- 759
ter Holch, Wieland Sommer, Felix Hofmann, Alexan- 760
dre Hostettler, Naama Lev-Cohain, Michal Drozdal, 761
Michal Marianne Amitai, Refael Vivanti, and 90 oth- 762
ers. 2023. [The liver tumor segmentation benchmark
\(LiTS\)](#). *Med. Image Anal.*, 84(102680):102680. 763
764

Davide Caffagni, Federico Cocchi, Luca Barsellotti, 765
Nicholas Moratelli, Sara Sarto, Lorenzo Baraldi, 766
Lorenzo Baraldi, Marcella Cornia, and Rita Cuc- 767
chiara. 2024. [The revolution of multimodal large
language models: A survey](#). In *Findings of the As-
sociation for Computational Linguistics ACL 2024*,
pages 13590–13618, Stroudsburg, PA, USA. Associ- 768
ation for Computational Linguistics. 769
770
771
772

Victor M Campello, Polyxeni Gkontra, Cristian 773
Izquierdo, Carlos Martin-Isla, Alireza Sojoudi, Pe- 774
ter M Full, Klaus Maier-Hein, Yao Zhang, Zhiqiang 775
He, Jun Ma, Mario Parreno, Alberto Albiol, Fan- 776
wei Kong, Shawn C Shadden, Jorge Corral Acero, 777
Vaanathi Sundaresan, Mina Saber, Mustafa Elattar, 778
Hongwei Li, and 28 others. 2021. [Multi-centre,
multi-vendor and multi-disease cardiac segmentation:
The M&Ms challenge](#). *IEEE Trans. Med. Imaging*,
40(12):3543–3554. 779
780
781
782

1123	for COVID-19 CT report generation with alternate learning. <i>IEEE Trans. Neural Netw. Learn. Syst.</i> , 32(9):3786–3797.	
1124		
1125		
1126	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning . <i>arXiv [cs.CV]</i> , pages 34892–34916.	
1127		
1128		
1129	Jingxin Liu, Xinran Zhu, Zhangzhen Shi, Donghong An, Lihui Zu, Kailiang Cheng, and Xiaopeng Guo. 2025. Text-guided multimodal deep learning in magnetic resonance imaging for spinal structures segmentation and lumbar abnormalities identification . <i>Quant. Imaging Med. Surg.</i> , 15(10):9710–9728.	
1130		
1131		
1132		
1133		
1134		
1135	Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. 2021c. Video swin transformer . <i>arXiv [cs.CV]</i> , pages 3202–3211.	
1136		
1137		
1138	Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. ViLBERT: Pretraining task-agnostic vision-linguistic representations for vision-and-language tasks . <i>arXiv [cs.CV]</i> , pages 13–23.	
1139		
1140		
1141		
1142	Zhixiu Lu, Hailong Li, Nehal A Parikh, Jonathan R Dillman, and Lili He. 2025. RadCLIP: Enhancing radiologic image analysis through contrastive language-image pretraining . <i>IEEE Trans. Neural Netw. Learn. Syst.</i> , 36(10):17613–17622.	
1143		
1144		
1145		
1146		
1147	Lingxiao Luo, Bingda Tang, Xuanzhong Chen, Rong Han, and Ting Chen. 2025. VividMed: Vision language model with versatile visual grounding for medicine . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 1800–1821, Stroudsburg, PA, USA. Association for Computational Linguistics.	
1148		
1149		
1150		
1151		
1152		
1153		
1154		
1155		
1156	Xiangde Luo, Wenjun Liao, Jianghong Xiao, Jieneng Chen, Tao Song, Xiaofan Zhang, Kang Li, Dimitris N Metaxas, Guotai Wang, and Shaoting Zhang. 2022. WORD: A large scale dataset, benchmark and clinical applicable study for abdominal organ segmentation from CT image . <i>Med. Image Anal.</i> , 82(102642):102642.	
1157		
1158		
1159		
1160		
1161		
1162		
1163	Jianhui Lv, Wadii Boulila, Shalli Rani, and Huamao Jiang. 2025. Enhanced multimodal speech processing for healthcare applications: A deep fusion approach . <i>IEEE J. Sel. Top. Signal Process.</i> , 19(4):600–612.	
1164		
1165		
1166		
1167		
1168	Jun Ma, Yao Zhang, Song Gu, Cheng Zhu, Cheng Ge, Yichi Zhang, Xingle An, Congcong Wang, Qiyuan Wang, Xin Liu, Shucheng Cao, Qi Zhang, Shangqing Liu, Yunpeng Wang, Yuhui Li, Jian He, and Xiaoping Yang. 2022. AbdomenCT-1K: Is abdominal organ segmentation a solved problem? <i>IEEE Trans. Pattern Anal. Mach. Intell.</i> , 44(10):6695–6714.	
1169		
1170		
1171		
1172		
1173		
1174		
1175	Manal Makram and Ammar Mohammed. 2025. Enhancing large vision language models for liver CT scans in medical reports . In <i>2025 International Mobile, Intelligent, and Ubiquitous Computing Conference (MIUCC)</i> , pages 88–95. IEEE.	
1176		
1177		
1178		
1179		
	Ian B Malone, David Cash, Gerard R Ridgway, David G MacManus, Sebastien Ourselin, Nick C Fox, and Jonathan M Schott. 2013. MIRIAD—public release of a multiple time point alzheimer’s MR imaging dataset . <i>Neuroimage</i> , 70:33–36.	1180 1181 1182 1183 1184
	Daniel S Marcus, Anthony F Fotenos, John G Csernansky, John C Morris, and Randy L Buckner. 2010. Open access series of imaging studies: longitudinal MRI data in nondemented and demented older adults . <i>J. Cogn. Neurosci.</i> , 22(12):2677–2684.	1185 1186 1187 1188 1189
	Daniel S Marcus, Tracy H Wang, Jamie Parker, John G Csernansky, John C Morris, and Randy L Buckner. 2007. Open access series of imaging studies (OASIS): cross-sectional MRI data in young, middle aged, nondemented, and demented older adults . <i>J. Cogn. Neurosci.</i> , 19(9):1498–1507.	1190 1191 1192 1193 1194 1195
	Kenneth Marek, Danna Jennings, Shirley Lasch, Andrew Siderowf, Caroline Tanner, Tanya Simuni, Chris Coffey, Karl Kieburtz, Emily Flagg, Sohini Chowdhury, and Others. 2011. The parkinson progression marker initiative (PPMI) . <i>Prog. Neurobiol.</i> , 95(4):629–635.	1196 1197 1198 1199 1200 1201
	Xueyan Mei, Zelong Liu, Philip M Robson, Brett Marinelli, Mingqian Huang, Amish Doshi, Adam Jacobi, Chendi Cao, Katherine E Link, Thomas Yang, Ying Wang, Hayit Greenspan, Timothy Deyer, Zahi A Fayad, and Yang Yang. 2022. RadImageNet: An open radiologic deep learning research dataset for effective transfer learning . <i>Radiol. Artif. Intell.</i> , 4(5):e210315.	1202 1203 1204 1205 1206 1207 1208 1209
	Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. 2016. V-net: Fully convolutional neural networks for volumetric medical image segmentation . In <i>2016 Fourth International Conference on 3D Vision (3DV)</i> , pages 79–87. IEEE.	1210 1211 1212 1213 1214
	Asbjørn Munk, Jakob Ambsdorf, Sebastian Llambias, and Mads Nielsen. 2024. AMAES: Augmented masked autoencoder pretraining on public brain MRI data for 3D-native segmentation . <i>arXiv [eess.IV]</i> .	1215 1216 1217 1218
	Maiko Nagao, Atsushi Teramoto, Kaito Urata, Kazuyoshi Imaizumi, Masashi Kondo, and Hiroshi Fujita. 2025. Preliminary study on image-finding generation and classification of lung nodules in chest CT images using vision–language models . <i>Computers</i> , 14(11):489.	1219 1220 1221 1222 1223 1224
	Yuko Nakagi, Takuya Matsuyama, Naoko Koide-Majima, Hiroto Q Yamaguchi, Rieko Kubo, Shinji Nishimoto, and Yu Takagi. 2024. Unveiling multi-level and multi-modal semantic representations in the human brain using large language models . <i>bioRxiv</i> , pages 20313–20338.	1225 1226 1227 1228 1229 1230
	Friska Natalia, Hira Meidia, Nunik Afriliana, Ala S Al-Kafri, Sud Sudirman, Andrew Simpson, Ali Sophian, Mohammed Al-Jumaily, Wasfi Al-Rashdan, and Mohammad Bashtawi. 2018. Development of ground truth data for automatic lumbar spine MRI image	1231 1232 1233 1234 1235

1350	24452–24470, Stroudsburg, PA, USA. Association	Jakob Wasserthal, Hanns-Christian Breit, Manfred T	1408
1351	for Computational Linguistics.	Meyer, Maurice Pradella, Daniel Hinck, Alexan-	1409
1352	Pratul P Srinivasan, Leo A Kim, Priyatham S Mettu,	der W Sauter, Tobias Heye, Daniel T Boll, Joshy Cyr-	1410
1353	Scott W Cousins, Grant M Comer, Joseph A Izatt,	iac, Shan Yang, Michael Bach, and Martin Segeroth.	1411
1354	and Sina Farsi. 2014. Fully automated detection of	2023. TotalSegmentator: Robust segmentation of	1412
1355	diabetic macular edema and dry age-related macu-	104 anatomic structures in CT images. <i>Radiol. Artif.</i>	1413
1356	lar degeneration from optical coherence tomography	<i>Intell.</i> , 5(5):e230024.	1414
1357	images. <i>Biomed. Opt. Express</i> , 5(10):3568–3577.		
1358	Hao Tan and Mohit Bansal. 2019. LXMERT: Learning	Navodini Wijethilake, Marina Ivory, Oscar MacCor-	1415
1359	cross-modality encoder representations from trans-	mac, Siddhant Kumar, Aaron Kujawa, Lorena Garcia-	1416
1360	formers . In <i>Proceedings of the 2019 Conference</i>	Foncillas Macias, Rebecca Burger, Amanda Hitch-	1417
1361	<i>on Empirical Methods in Natural Language Pro-</i>	ings, Suki Thomson, Sinan Barazi, Eleni Maratos,	1418
1362	<i>cessing and the 9th International Joint Conference</i>	Rupert Obholzer, Dan Jian, Fiona McClenaghan,	1419
1363	<i>on Natural Language Processing (EMNLP-IJCNLP)</i> ,	Kazumi Chia, Omar Al-Salihi, Nick Thomas, Steve	1420
1364	Stroudsburg, PA, USA. Association for Computa-	Connor, Tom Vercauteren, and Jonathan Shapey.	1421
1365	tional Linguistics.	2025. Deep learning consensus-based annotation	1422
1366	Tianchi. 2020. Spinal disease dataset. https://tianchi.aliyun.com/dataset/dataDetail?dataId=79463 .	of vestibular schwannoma from magnetic resonance	1423
1367		imaging: An annotated multi-center routine clinical	1424
1368		dataset (vestibular-schwannoma-MC-RC 2).	1425
1369	Selene Tomassini, Damiano Duranti, Abdallah Zeggada,	Fuping Wu and Xiahai Zhuang. 2020. CF distance: A	1426
1370	Carlo Cosimo Quattrocchi, Farid Melgani, and Paolo	new domain discrepancy metric and application to	1427
1371	Giorgini. 2025. Multi-branch CNN-LSTM fusion	explicit domain adaptation for cross-modality cardiac	1428
1372	network-driven system with BERT semantic evalu-	image segmentation. <i>IEEE Trans. Med. Imaging</i> ,	1429
1373	ator for radiology reporting in emergency head CTs .	39(12):4274–4285.	1430
1374	<i>IEEE J. Transl. Eng. Health Med.</i> , 13:61–74.		
1375	Andrea C Tricco, Erin Lillie, Wasifa Zarin, Kelly K	Jing Wu, Yuli Wang, Zhusi Zhong, Weihua Liao, Natalia	1431
1376	O’Brien, Heather Colquhoun, Danielle Levac, David	Trayanova, Zhicheng Jiao, and Harrison X Bai. 2025.	1432
1377	Moher, Micah D J Peters, Tanya Horsley, Laura	Vision-language foundation model for 3D medical	1433
1378	Weeks, Susanne Hempel, Elie A Akl, Christine	imaging. <i>NPJ Artif. Intell.</i> , 1(1):15.	1434
1379	Chang, Jessie McGowan, Lesley Stewart, Lisa		
1380	Hartling, Adrian Aldcroft, Michael G Wilson,	Yu Xin, Gorkem Can Ates, Kuang Gong, and Wei Shao.	1435
1381	Chantelle Garritty, and 9 others. 2018. PRISMA	2025. Med3DVLM: An efficient vision-language	1436
1382	extension for scoping reviews (PRISMA-ScR): Check-	model for 3D medical image analysis . <i>arXiv [cs.CV]</i> ,	1437
1383	list and explanation . <i>Ann. Intern. Med.</i> , 169(7):467–	pages 1–14.	1438
1384	473.		
1385	Santiago Vitale, José Ignacio Orlando, Emmanuel	Kohei Yamamoto and Tomohiro Kikuchi. 2025. Fea-	1439
1386	Iarussi, and Ignacio Larrabide. 2020. Improving re-	sibility study of CLIP-based key slice selection in	1440
1387	alism in patient-specific abdominal ultrasound simu-	CT images and performance enhancement via lesion-	1441
1388	lation using CycleGANs . <i>Int. J. Comput. Assist.</i>	and organ-aware fine-tuning . <i>Bioengineering (Basel)</i> ,	1442
1389	<i>Radiol. Surg.</i> , 15(2):183–192.	12(10):1093.	1443
1390	Xianghong Wang, Jiajun Xiang, Aihua Mao, Jiayi Xie,	An Yan, Julian McAuley, Xing Lu, Jiang Du, Eric Y	1444
1391	Peng Jin, Mingchao Ding, Yixuan Yuan, Yanye Lu,	Chang, Amilcare Gentili, and Chun-Nan Hsu.	1445
1392	Lequan Yu, Hongmin Cai, Baiying Lei, and Tianye	2022. RadBERT: Adapting transformer-based lan-	1446
1393	Niu. 2025. Clip-driven universal model for multi-	guage models to radiology . <i>Radiol. Artif. Intell.</i> ,	1447
1394	material decomposition in dual-energy CT . <i>IEEE</i>	4(4):e210258.	1448
1395	<i>Trans. Comput. Imaging</i> , 11:349–361.		
1396	Yibin Wang, William Duggar, David Caballero, Toms	Jiancheng Yang, Xiaoyang Huang, Yi He, Jingwei Xu,	1449
1397	Vengaloor Thomas, Neha Adari, Eswarakumar	Canqian Yang, Guozheng Xu, and Bingbing Ni. 2021.	1450
1398	Mundra, and Haifeng Wang. 2023. Brain tumor	Reinventing 2D convolutions for 3D images . <i>IEEE J.</i>	1451
1399	recurrence prediction after gamma knife radiother-	<i>Biomed. Health Inform.</i> , 25(8):3009–3018.	1452
1400	apy from MRI and related DICOM-RT: An open		
1401	annotated dataset and baseline algorithm (brain-TR-	Jiancheng Yang, Rui Shi, Donglai Wei, Zequan Liu,	1453
1402	GammaKnife) .	Lin Zhao, Bilian Ke, Hanspeter Pfister, and Bingbing	1454
1403	Zifeng Wang, Zhenbang Wu, Dinesh Agarwal, and Ji-	Ni. 2023. MedMNIST v2 - a large-scale lightweight	1455
1404	meng Sun. 2022. MedCLIP: Contrastive learning	benchmark for 2D and 3D biomedical image classi-	1456
1405	from unpaired medical images and text . <i>Proc. Conf.</i>	fication . <i>Sci. Data</i> , 10(1):41.	1457
1406	<i>Empir. Methods Nat. Lang. Process.</i> , 2022:3876–		
1407	3887.	Huping Ye, Yushan Deng, and Yi Hong. 2024. Epi-	1458
		cardial adipose tissue segmentation in MRIs using	1459
		text-prompted pretraining model . In <i>2024 IEEE</i>	1460
		<i>International Conference on Bioinformatics and</i>	1461
		<i>Biomedicine (BIBM)</i> , pages 7021–7028. IEEE.	1462

1463 Joonhyeok Yoon, Hangeol Park, Minjun Kim, Hwihun
1464 Jeong, Se Young Chun, Sooyeon Ji, and Jongho Lee.
1465 2025. [Clinical dementia rating classification using](#)
1466 [integrated vision and language information](#). *IEEE*
1467 *Access*, 13:184602–184617.

1468 Jianpeng Zhang, Yutong Xie, Yan Wang, and Yong
1469 Xia. 2021. [Inter-slice context residual learning for](#)
1470 [3D medical image segmentation](#). *IEEE Trans. Med.*
1471 *Imaging*, 40(2):661–672.

1472 Xiaoman Zhang, Chaoyi Wu, Ziheng Zhao, Jiayu Lei,
1473 Weiwei Tian, Ya Zhang, Weidi Xie, and Yanfeng
1474 Wang. 2025. [Development of a large-scale grounded](#)
1475 [vision language dataset for chest CT analysis](#). *Sci.*
1476 *Data*, 12(1):1636.

1477 Xiaoman Zhang, Chaoyi Wu, Ziheng Zhao, Weixiong
1478 Lin, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023.
1479 [PMC-VQA: Visual instruction tuning for medical](#)
1480 [visual question answering](#). *arXiv [cs.CV]*.

1481 Weike Zhao, Chaoyi Wu, Xiaoman Zhang, Ya Zhang,
1482 Yanfeng Wang, and Weidi Xie. 2024. [RaTEScore: A](#)
1483 [metric for radiology report generation](#). In *Proceed-*
1484 *ings of the 2024 Conference on Empirical Methods in*
1485 *Natural Language Processing*, pages 15004–15019,
1486 Stroudsburg, PA, USA. Association for Computa-
1487 tional Linguistics.

1488 Yufang Zhao, Yue Li, Yanjing Zhang, Xiaohui Yan,
1489 Guolin Yin, and Liping Liu. 2025. [Enhancing thyroid](#)
1490 [nodule assessment with UTV-ST swin kansformer: A](#)
1491 [multimodal approach to predict invasiveness](#). *IEEE*
1492 *Access*, 13:29081–29090.

1493 Xiahai Zhuang and Juan Shen. 2016. [Multi-scale patch](#)
1494 [and multi-modality atlases for whole heart segmenta-](#)
1495 [tion of MRI](#). *Med. Image Anal.*, 31:77–87.

1496 **A Search Query Variables**

1497 {3d_terms} = "CT"

1498 OR "MRI"

1499 OR "3D"

1500 OR "volumetric"

1501 OR "volume"

1502 OR "computed tomography"

1503 OR "magnetic resonance"

1504 OR "3D image"

1505 OR "biomedical image"

1506 OR "3D reconstruction"

1507 OR "point cloud"

1508
1509 {model_terms} = "multimodal"

1510 OR "vision-language"

1511 OR "foundation model"

1512 OR "self-supervised"

1513 OR "contrastive"

1514 OR "CLIP"

1515 OR "LVM"

1516 OR "image-text"

1517 OR "transformer"

1518 OR "pretrain*"

1519 OR "masked autoencoder"

1520 OR "cross-modal"

1521 OR "report-grounded"

1522 OR "retrieval"

1523
1524 {text_terms} = "radiology report"

1525 OR "clinical text"

1526 OR "report generation"

1527 OR "text"

1528 OR "language"

1529 OR "PACS"

1530 OR "DICOM"

1531 OR "diagnosis"

1532 OR "medical report"

1533 OR "biomedical text"

1534
1535 {video_terms} = "medical video"

1536 OR "surgical video"

1537 OR "ultrasound video"

1538 OR "endoscopy video"

1539 OR "cine MRI"

1540 OR "temporal CT"

Table 1: Overview of Included Studies.

Study	Code	Visionmodel	LLM	3D handling	Inst follow	Pretrained	Fine-tuned	Input	Task
VividMed (2025)	Link	ViT, SAM	Vicuna	Voxel	✓	E2E	E2E	CT, X-ray	Q R S
Argus (2024)	Link	ViT	Llama3	Voxel	✓	Vision	E2E	CT	R
Liu et al. (2025)	✗	ConvNeXt	MedCLIP	Slice	✗	Vision	E2E	MRI	S C
Zhang et al. (2025)	Link	3D ConvNet	Llama3	Voxel	✓	–	Language	CT	R Q
Yamamoto and Kikuchi (2025)	✗	ViT	CLIP	Slice	✗	–	E2E	CT	R
Tomassini et al. (2025)	Link	VGG16	LSTM, BERT	Slice	✗	–	Language	CT	R
Batool et al. (2025)	✗	ViT	BioBERT, MiniLM	Slice	✗	E2E	E2E	CT	R
FLIQA-AD (2025b)	Link	3D ViT	bioClinicalBert, FLAN-T5	Voxel	✓	–	E2E	MRI	C Q
Med3DVLM (2025)	Link	DCFramer	Qwen2.5-7B-Instruct	Voxel	✓	E2E	E2E	CT	Q C R R
RadCLIP (2025)	Link	ViT	CLIP	Slice	✗	E2E	Vision	CT, MRI	C R
Petersen et al. (2025)	✗	Swin-T, MedNeXt, ResNet	BERT	Voxel	✗	Vision	E2E	MRI	C R
MultiModalGAN (2025)	✗	CNN	BERT	Slice	✗	–	E2E	MRI, X-Ray	S
Med3DInsight (2025c)	✗	nnFormer	CLIP	Voxel	✗	E2E	E2E	CT, MRI	S C
Medical-VLBERT (2021b)	✗	DenseNet	BERT	Slice	✗	E2E	E2E	CT	R C
Chen et al. (2025a)	✗	ViT-H SAM	LLaVA-Med	Slice	✓	E2E	E2E	CT	S S
Ye et al. (2024)	✗	ViT	BERT	Slice	✗	Vision	Vision	MRI	S
Makram and Mohammed (2025)	✗	ViT, CLIP	Vicuna, Mistral, MPT, LLaMA, GPT-4o	Slice	✓	–	–	CT, MRI	C R Q
PMC-MSA (2024)	✗	ViT	CLIP	Slice	✗	–	E2E	CT, MRI	C
Dack et al. (2023)	✗	ViT, Swin-T, ResNet	GPT2, BioClinicalBERT, PubMedBERT, RadBERT	Slice	✗	–	E2E	CT	C
HCL-AL (2025)	Link	3D ResNet	CLIP	Voxel	✗	Language	–	PET, CT	R
Wang et al. (2025)	✗	Siamese	CLIP	Slice	✗	Language	–	CT	S
TGCFa (2024)	✗	U-Net	CLIP	Slice	✗	–	Vision	CT, Fundus, MRI	S
3DLVR (2024a)	✗	Transformer	Transformer	Voxel	✗	–	E2E	CT	R
DILBERT (2025d)	✗	3D CNN	BERT	Video	✗	E2E	E2E	Ultrasound	R
Zhao et al. (2025)	✗	Video Transformer	Kansformer	Video	✗	–	E2E	Ultrasound	C
HAV-DF (2025)	✗	3D-Conv, ResNet	LSTM	Video	✗	–	E2E	Audio, Video	C
Tonguescape (2025)	✗	Gemini 1.5 Pro, GPT-4o, LLaVANEXT-Interleave, Phi-3.5-vision-instruct, Qwen2-VL-Instruct, VideoLLaMA2		Video	✓	–	E2E	MRI video	Q
Drama2brain (2024)	Link	DeiT, ResNet, GIT, BridgeTower, LLaVa-v1.5	Word2Vec, BERT, GPT2, OPT, Llama2	Video	Both	–	–	fMRI	C
Shi et al. (2023)	✗	ResNet101	–	Slice	✗	–	E2E	CT	R
LMOD (2024)	✗	LLaVA, InternVL, Yi-6B, Qwen	VILA, GPT-4o	Slice	✓	–	E2E	CFP, LP, OCT, SS, SLO	C Q
LUMEN (2025)	✗	YOLOv9, LLaMA-3-VILA1.5, LLaMA3-VILA-M3	LLaMA-3-VILA1.5, LLaMA3-VILA-M3, LLaMA 3.3	Slice	✓	–	E2E	CT	C S R
Yoon et al. (2025)	✗	ViT	Transformer	Voxel	✗	E2E	E2E	MRI	C
BrainGPT (2025a)	Link	ViT	LLaMA	Slice	✓	–	E2E	CT	R
NeuroInfinity (2025)	✗	ViT	CLIP	Slice	✗	–	E2E	MRI	C

Table 1: Overview of Included Studies. (Continued)

Study	Code	Visionmodel	LLM	3D handling	Inst follow	Pretrained	Fine-tuned	Input	Task
Hassan et al. (2025)	✗	ResNet-50, PVT, ViT	BERT, Word2Vec, XLNet	Slice	✗	-	E2E	OCT	Ⓢ
Nagao et al. (2025)	✗	BLIP, GiT	BLIP, GiT	Slice	✗	-	E2E	CT	Ⓢ Ⓢ
VISTA-Prompt (2026)	✗	DCFormer	Qwen-2.5	Slice	✓	-	Language	CT	Ⓢ Ⓢ

Note: Ⓢ - Report Generation, Ⓢ - Visual Question Answering, Ⓢ - Diagnostic Classification and Anatomical Identification, Ⓢ - Segmentation and Visual Grounding, Ⓢ - Image Synthesis and Modality Conversion, Ⓢ - Retrieval, Inst Follow - Instruction following, E2E - End-to-end

Table 2: Overview of recent benchmark datasets.

Dataset Name	Data Modality	Annot. Type	Scale (# samples)	Public URL	Institution	Clinical Specialty
TotalSegmentator (2023)	CT	S	1,204	Link	Single	Clinical
VinDr-CXR (2022)	X-Ray	S C	18,000	Link	Multiple	Clinical
VQA-RAD (2018)	CT, MRI, X-Ray	Q	315	Link	Single	Clinical
SLAKE (2021a)	CT, MRI, X-Ray	Q	642	Link	Multiple	Clinical
VQA-Med (2024)	Angiography, CT, Mammograph, MRI, PET, X-Ray, Ultrasound	Q	4,200	Link	—	Clinical
MIMIC-CXR (2019)	X-Ray	C R	377,110	Link	Single	Clinical
CT-RATE (2025a)	CT	C R	25,692	Link	Single	Clinical
BIMCV-R (2024)	CT	R	8,069	Link	Single	Clinical
INSPECT (2023)	CT	C R	23,248	Link	Single	Pulmonary
MRSpineSeg2021 (2021)	MRI	S	215	Link	Single	Orthopedics
RSNA Lumbar Spine Degenerative (2024)	MRI	C	1,975	Link	Single	Orthopedics
Tianchi Spinal Disease (2020)	MRI	S C	201	Link	Single	Orthopedics
Lumbar Spine MRI dataset (2018)	MRI	S R	515	Link	Single	Orthopedics
RadGenome-Chest CT (2025)	CT (Chest)	S R Q	25,692	Link	Single	Clinical
VerSe20 (2021)	CT	—	300	Link	Single	Orthopedics
ADNI (2010)	MRI	C	8,315	Link	Multiple	Neurology
OASIS-2 (2010)	MRI	C	373	Link	Single	Neurology
rtMRIDB	MRI (Video)	C	26K	Link	Single	Otolaryngology
VowelVideo (2025)	MRI (Video)	C	120	Link	Single	Otolaryngology
drama2brain (2024)	fMRI	C	6 subjects	Link	Single	Neurology
LMOD (2024)	CFP, LP, OCT, SLO, SS	S C Q	21,993	Link	Multiple	Ophthalmology
OASIS-3 (2019)	MRI	C	2,842	Link	Single	Neurology
OASIS-4 (2020)	MRI	C	676	Link	Single	Neurology
CQ500 (2018)	CT	R	1,154	Link	Single	Neuroradiology
ImagesOASIS (2007)	MRI	C	80,00	Link	Single	Neurology
MIRIAD (2013)	MRI	C	708	Link	Single	Neurology
Zhang dataset (2018)	OCT	C	109,309	Link	Multiple	Ophthalmology
Duke-1 (2014)	OCT	C	38,300	Link	Multiple	Ophthalmology
Duke-2 (2015)	OCT	C	610	Link	Single	Ophthalmology
Duke-3 (2014)	OCT	C	3,231	Link	Multiple	Ophthalmology
Rabbani (2018)	OCT	C	4,241	Link	Single	Ophthalmology
BIOMISA (2018)	OCT	C	4,163	Link	Single	Ophthalmology
PMC-VQA (2023)	—	Q	149K	Link	Multiple	Clinical
M3D-Cap (2024b)	CT	Q	120k	Link	Multiple	Clinical
M3D-VQA (2024)	CT	Q	660K	Link	Multiple	Clinical
RadCLIP (2025)	CT, MRI, X-Ray	Q	1,210,083	Link	Multiple	Clinical
ChestXpert (2019)	X-Ray	C	224,316	Link	Single	Clinical
Crystal Clean Brain Tumor (2023)	MRI	C	3,264	Link	—	Neurosurgery
IXI Brain	MRI	C	600	Link	Multiple	Neurology
COVID-CT-MD (2021)	CT	C	23,349	Link	Single	Pulmonary

Continued on next page

Table 2: Overview of recent benchmark datasets. (Continued)

Dataset Name	Data Modality	Annot. Type	Scale (# samples)	Public URL	Institution	Clinical Specialty
BRAINS-45K (2024)	MRI	—	44,756	Link	Multiple	Neurology
Brain-TR-GammaKnife (2023)	MRI	Q	17,191	Link	Single	Neurosurgery
3DSeg-8 (2019)	CT, MRI	Q	25K	Link	Multiple	Clinical
MM-WHS (2016)	CT	S	20	Link	Single	Cardiology
MSD-Heart (2022)	CT, MRI	S	30	Link	—	Cardiology
CHAOS (2021)	CT, MRI	S	80	Link	Single	Clinical
AbdomenCT-1K (2022)	CT	S	1,112	Link	Multiple	Clinical
VS (2025)	MRI	S	676	Link	Multiple	Neurosurgery
LiTs (Bilic et al., 2023)	CT	S	201	Link	Single	Hepatology
PPMI (2011)	MRI	C	582	Link	Multiple	Neurology
COVID-19 CT (2021b)	CT	R	1,104	Link	Multiple	Pulmonary
AbdomenAtlas (2025c)	CT	S	3.2m	Link	Multiple	Clinical
WORD (2022)	CT	S	30,495	Link	Single	Clinical
KiTS23 (2023)	CT	S	599	Link	Single	Urology
ACDC (2018)	MRI	S	150	Link	Single	Cardiology
M&Ms (2021)	MRI	S	375	Link	Multiple	Cardiology
Wu and Zhuang (2020)	MRI	S	1,568	Link	Multiple	Cardiology
Emidec (2020)	MRI	S	150	Link	Single	Cardiology
MyoPS2020 (2022)	MRI	S	45	Link	Single	Cardiology
Eurorad	CT, MRI	R	5K	Link	—	Clinical
Visible Human Project (2022)	CT	C	650	Link	—	Research
Auto-PET (2022)	CT, PET	C	1,169	Link	Multiple	Clinical
DL-spectral CT (2024)	CT	S	1,100	Link	Synthetic	Clinical
MedMNIST (2023)	Dermatology, CT, MRI, OCT, Pathology, Retina, X-Ray, Ultrasound	C	718K	Link	Multiple	Clinical
MS-CMR	MRI	S	<1K	Link	—	Clinical
RIGA+ (2022)	Fundus	S	750	Link	Multiple	Ophthalmology
BUSI (2020)	Ultrasound	R	780	Link	Single	Breast Surgery
USSS (2020)	Ultrasound	R	926	Link	Synthetic	Clinical, Simulation
MedDialog (2020)	Audio	D	260K	Link	Multiple	Clinical
MedVidQA	Video	Q	6.4K	Link	Multiple	Clinical, Educational

Note: **S** - Segmentation Mask and Bounding Box, **Q** - Question-Answer Pair, **C** - Classification Label, **R** - Report, **D** - Dialogue