

MedLayBench-V: A Large-Scale Benchmark for Expert-Lay Semantic Alignment in Medical Vision Language Models

Anonymous ACL submission

Abstract

Medical Vision-Language Models (Med-VLMs) have achieved expert-level proficiency in interpreting diagnostic imaging. However, current models are predominantly trained on professional literature, limiting their ability to communicate findings in the lay register required for patient-centered care. While text-centric research has actively developed resources for simplifying medical jargon, there is a critical absence of large-scale multimodal benchmarks designed to facilitate lay-accessible medical image understanding. To bridge this resource gap, we introduce **MedLayBench-V**, the first large-scale multimodal benchmark dedicated to expert-lay semantic alignment. Unlike naive simplification approaches that risk hallucination, our dataset is constructed via a Structured Concept-Grounded Refinement (SCGR) pipeline. This method enforces strict semantic equivalence by integrating Unified Medical Language System (UMLS) Concept Unique Identifiers (CUIs) with micro-level entity constraints. MedLayBench-V provides a verified foundation for training and evaluating next-generation Med-VLMs capable of bridging the communication divide between clinical experts and patients.

1 Introduction

Enhancing the linguistic accessibility of clinical documentation has emerged as a paramount objective in biomedical Natural Language Processing (NLP). Driven by the imperative to facilitate patient-centered care, recent research has coalesced around tasks such as Biomedical Lay Summarization (BioLaySumm) and Neural Text Simplification (NTS) (Shardlow and Nawaz, 2019; Yao et al., 2024). Collectively framed as Medical Lay Language Generation (MLLG), these efforts aim to translate highly specialized medical jargon into the accessible lay register. This paradigm shift is epitomized by initiatives like the BioLaySumm shared

Input (Image Caption):

"Thoracic CT scan showing perihilar lymphadenomegaly."

Problem: Naive Text Simplification (LLM only)

→ "The scan shows signs of lung cancer..."

(Hallucination)

→ "There is swelling in the chest..."

(Vague, Loss of Modality)

Ours: Structured Concept-Grounded Refinement

Concept Mappings (C):

Thoracic CT [CUI:C0040405] → "Chest CT scan"

perihilar [Entity] → "near lung center"

Lymphadenomegaly [CUI:C0024265] → "enlarged lymph nodes"

↓ + LLM Refinement (Grammar and Fluency)

→ "The Chest CT scan shows enlarged lymph nodes

near the center of the lungs."

Figure 1: **Motivation.** Our method prevents hallucinations by enforcing Structured Constraints: It explicitly maps extracted Concepts and Entities (e.g., *lymphadenomegaly*) to lay terms, ensuring diagnostic accuracy while preserving specific details.

tasks (Xiao et al., 2025; Goldsack et al., 2024) and recent benchmarks like MedAgentBoard (Zhu et al., 2025a), where MLLG is established as a core competency for medical artificial intelligence (AI). Recent studies attribute success in this domain to the advanced semantic reasoning of Large Language Models (LLMs), which allows them to modify lexical complexity while maintaining semantic invariance, thereby ensuring that core medical facts are preserved despite the stylistic shift (Liao et al., 2025).

While the text-to-text simplification landscape has advanced significantly, the integration of this lay perspective into multimodal systems remains an open challenge. Medical Vision-Language Models (Med-VLMs), such as those trained on RO-COV2 (Rückert et al., 2024) or PMC-OA, have achieved expert-level proficiency in interpreting diagnostic imaging (Lozano et al., 2025). However,

063 a critical limitation persists in their current training
064 paradigm. Unlike text-centric LLMs that are be-
065 coming increasingly adaptable to the lay register,
066 current Med-VLMs are predominantly optimized
067 for the rigid clinical jargon found in professional
068 literature. As illustrated in Figure 1, this domain-
069 specific optimization creates a significant barrier to
070 usability; while models successfully encode visual
071 features into technical tokens like ‘Pneumothorax’,
072 their ability to ground the same visual evidence in
073 natural language equivalents like ‘Collapsed lung’
074 remains unsupported due to the lack of parallel lay
075 data. This suggests that without a dedicated bench-
076 mark to facilitate expert-to-lay alignment, Med-
077 VLMs will remain confined to a specialized lexi-
078 con, severely limiting their applicability in patient-
079 centered care.

080 Overcoming this resource scarcity, however,
081 presents significant methodological challenges. Ex-
082 isting multimodal benchmarks are exclusively pop-
083 ulated with expert-level reports and offer no ground
084 truth for lay-accessible descriptions. Furthermore,
085 relying on standard lexical metrics like BLEU (Pa-
086 pineni et al., 2002) is insufficient for validation as
087 they inherently penalize the vocabulary shifts re-
088 quired for simplification (Zhao et al., 2024a). More-
089 over, constructing a benchmark via naive LLM
090 generation carries the risk of hallucination or the
091 omission of vital quantitative details, which com-
092 promises the factual integrity required for medical
093 AI (Liao et al., 2025).

094 To bridge this divide, we introduce
095 **MedLayBench-V**, the first multimodal benchmark
096 designed to facilitate patient-centric medical
097 image understanding. Drawing inspiration
098 from recent text-centric approaches that lever-
099 age structured medical knowledge to enhance
100 summary relevance (Ming et al., 2025), we
101 extend this philosophy to the multimodal domain
102 via a novel **Structured Concept-Grounded**
103 **Refinement** (SCGR) pipeline. Our approach
104 synergizes macro-level conceptual mapping from
105 the Unified Medical Language System (UMLS)
106 with micro-level entity constraints extracted via
107 Named Entity Recognition (NER) (Bodenreider,
108 2004). This hybrid strategy ensures that the
109 generated lay captions maintain strict semantic
110 equivalence with the original expert reports
111 while effectively transitioning to the lay register.
112 Using this verified dataset, we establish the first
113 comprehensive baselines for expert-lay alignment,
114 providing a standardized foundation for future

research in accessible medical AI. 115

Our contributions are summarized as follows: 116

- To the best of our knowledge, we intro- 117
duce **MedLayBench-V**, the first foundational 118
benchmark encompassing diverse medical 119
imaging modalities specifically curated to 120
bridge the linguistic divide between clinical 121
experts and laypersons. 122
- We propose the SCGR pipeline, a veri- 123
fiable framework that extends knowledge- 124
guided text simplification principles to vision- 125
language tasks, ensuring high clinical correct- 126
ness and hallucination control. 127
- We establish a comprehensive evaluation pro- 128
tocol for Expert-Lay semantic alignment and 129
provide standardized baselines, offering a ro- 130
bust foundation for future research in patient- 131
centered medical AI. 132

2 Related Works 133

2.1 Patient-Centered Clinical Reporting 134

135 The complexity of medical documentation creates
136 significant barriers to patient understanding, driv-
137 ing the need for automated systems that can trans-
138 late clinical narratives into accessible language.
139 To address this, the field has evolved from early
140 Neural Text Simplification (NTS) efforts into the
141 broader paradigm of Medical Lay Language Gen-
142 eration (MLLG) (Shardlow and Nawaz, 2019; Yao
143 et al., 2024). This transition is marked by large-
144 scale community initiatives such as the BioLay-
145 Summ shared tasks and the MedAgentBoard bench-
146 mark, which provide standardized tasks to bridge
147 the communication gap between experts and layper-
148 sons (Xiao et al., 2025; Zhu et al., 2025a).

149 Within this text-centric landscape, LLMs have
150 achieved remarkable proficiency, effectively bal-
151 ancing lexical simplification with semantic invari-
152 ance as demonstrated by frameworks (Liao et al.,
153 2025). However, this progress has yet to perme-
154 ate the multimodal domain. Unlike the thriving
155 domain for text-only models, there is a critical ab-
156 sence of benchmarks designed to evaluate Med-
157 VLMs leaving it unclear whether current SOTA
158 models can successfully ground visual findings in
159 lay-accessible language without compromising fac-
160 tual accuracy.

161	2.2 Medical Vision-Language Models and Dataset Scarcity		211
162			212
163	In the multimodal domain, Med-VLMs have	modal benchmarks remain limited in scope, pre-	213
164	achieved expert-level proficiency in interpreting	dominantly focusing on specific modalities like	214
165	diagnostic imaging (Zhang et al., 2023; Li et al.,	Chest X-rays (CXR) with restricted dataset sizes.	215
166	2023; Sellergren et al., 2025). These capabilities	Furthermore, these datasets typically rely on end-	216
167	are predominantly driven by large-scale datasets	to-end LLM generation for creating lay captions,	217
168	such as ROCOv2 (Rückert et al., 2024) and	which can perpetuate the very hallucinations they	218
169	BIOMEDICA (Lozano et al., 2025). However,	aim to resolve without rigorous concept-level ver-	219
170	these datasets are exclusively curated from pro-	fication. To facilitate the training of robust, general-	220
171	fessional biomedical literature, thereby optimizing	purpose Med-VLMs, there is a critical need for a	221
172	models strictly for the rigid clinical jargon.	large-scale, diverse benchmark that extends beyond	
173	A critical limitation in existing multimodal	single modalities.	
174	datasets is the scarcity of parallel multimodal data		
175	that pairs medical images with patient-friendly de-	2.4 Evaluation Metrics for Medical Text	222
176	scriptions. While models can successfully align	Generation	223
177	visual features with technical concepts (e.g., “Pneu-	Evaluating the quality of MLLG systems remains a	224
178	mothorax”), the lack of ground truth for natural	persistent challenge due to the inadequacy of exist-	225
179	language equivalents (e.g., “Collapsed lung”) pre-	ing metrics. Traditional n-gram based metrics such	226
180	vents them from learning the lay register. Unlike	as BLEU (Papineni et al., 2002), ROUGE (Lin,	227
181	the text domain where lay benchmarks exist, the	2004), and METEOR (Banerjee and Lavie, 2005)	228
182	vision-language field suffers from this fundamental	measure surface-level overlap. However, they in-	229
183	resource gap, which hinders the development of	herently penalize the vocabulary shifts required	230
184	expert-lay alignment capabilities in VLMs.	for simplification, making them unsuitable for	231
185		expert-to-lay translation tasks (Zhao et al., 2024a;	232
186	2.3 Limitations of Current Benchmarks	Zhang et al., 2019). Conversely, medically-oriented	233
187	To bridge the expert-lay divide, prior research has	metrics like Green (Ostmeier et al., 2024) and	234
188	predominantly focused on text-to-text simplifica-	RaTEScore (Zhao et al., 2024b) focus on clinical	235
189	tion strategies. Early approaches relied on rule-	factuality and entity extraction.	236
190	based methods or phrase tables to substitute medi-	While effective for expert reports, they do not	237
191	cal jargon with simpler synonyms (Shardlow and	assess whether the generated text is understand-	238
192	Nawaz, 2019). With the advent of LLMs, recent	able to a lay audience. Finally, standard readability	239
193	studies have shifted towards generative rewriting,	metrics rely on heuristic formulas (e.g., sentence	240
194	employing models such as GPT-4o to translate clin-	length) rather than actual comprehensibility, often	241
195	ical notes into patient-friendly language (Yao et al.,	failing to capture the semantic nuances required	242
196	2024). However, LLMs frequently generate plausi-	for patient education (Yao et al., 2024). Therefore,	243
197	ble yet factually incorrect descriptions or omit vital	effective MLLG evaluation requires a comprehen-	244
198	quantitative details to satisfy readability constraints,	sive framework that simultaneously assesses visual	245
199	thereby compromising patient safety in clinical set-	grounding, factual correctness, and lay accessibil-	246
200	tings (Moor et al., 2023; Zhu et al., 2025b). For	ity. However, performing such multi-dimensional	247
201	instance, a recent prospective trial demonstrated	evaluation is unfeasible with current VLM datasets	248
202	that while LLM-based simplification significantly	due to the critical absence of lay-aligned references.	249
203	reduces cognitive workload, it introduced factual	To bridge this gap, we introduce MedLayBench-	250
204	errors and omissions in approximately 6–7% of	V, a unified benchmark designed to facilitate this	251
205	reports, necessitating rigorous verification mecha-	holistic evaluation.	252
206	nisms (Prucker et al., 2025).		
207	Recent initiatives, such as the BioLaySumm	3 Methodology	253
208	2025 Shared Task (Xiao et al., 2025) and Layman’s	We introduce MedLayBench-V , a large-scale mul-	254
209	RRG (Zhao et al., 2024a), have begun to incor-	timodal benchmark designed to bridge the gap be-	255
210	porate visual modalities to address these ground-	tween expert clinical jargon and patient-accessible	256
	ing issues. Despite these advances, current multi-	language. To ensure the high semantic fidelity	257
		of this benchmark, we propose the Structured	258
		Concept-Grounded Refinement (SCGR) pipeline.	259

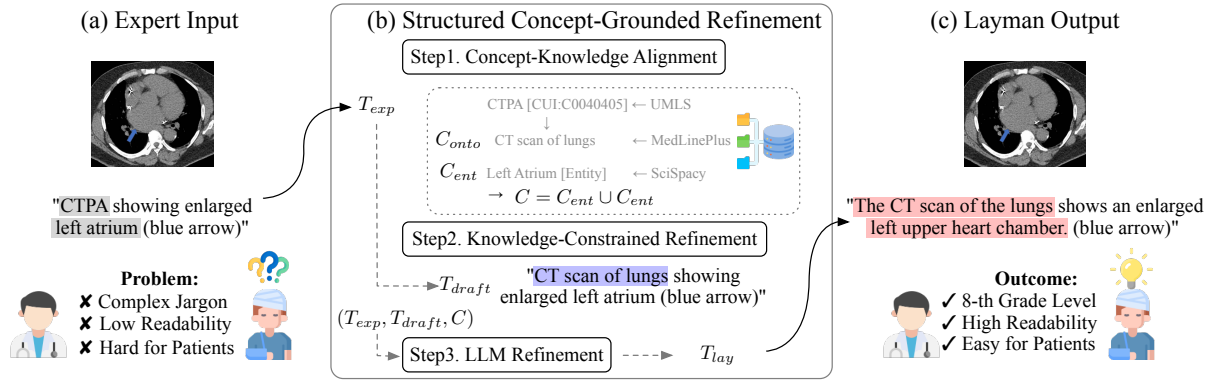


Figure 2: **Overview of the SCGR Framework.** (a) Expert Input extracts technical concepts from the initial jargon-heavy reports. (b) Structured Concept-Grounded Refinement maps terms to lay definitions and employs Llama-3.1-8B to synthesize the final caption, optimizing for syntax and fluency while strictly adhering to factual constraints (Detailed prompt in Appendix). (c) Layman Output provides a clinically accurate and accessible description.

Crucially, our framework explicitly decouples semantic extraction from stylistic refinement. This separation ensures strict *Semantic Equivalence* between the expert and lay registers, mitigating the hallucinations common in end-to-end generation. The pipeline consists of three distinct stages, **corresponding to Steps 1–3 in Figure 2(b)**: (i) Concept-Knowledge Alignment, (ii) Knowledge-Constrained Refinement, and (iii) LLM Refinement.

3.1 Data Source and Task Definition

We utilize the ROCov2 dataset (Rückert et al., 2024)¹ as our seed corpus. Derived from the PubMed Central Open Access (PMC-OA) subset (Lin et al., 2023)², ROCov2 is uniquely advantageous for our task as it provides not only diagnostic captions (T_{exp}) but also pre-computed UMLS Concept Unique Identifiers (CUIs) extracted via the MedCAT toolkit (Kraljevic et al., 2021)³. These existing annotations serve as a critical foundation for our semantic extraction phase.

Despite the richness of this clinical metadata, the expert descriptions in ROCov2 remain inherently unintelligible to non-specialists. Our objective is to augment these pairs with layman-accessible descriptions (T_{lay}), creating the first dual-register medical benchmark optimized for patient-centric VLM training and testing.

¹<https://huggingface.co/datasets/eltorio/ROCOv2-radiology>

²<https://pmc.ncbi.nlm.nih.gov/tools/openftlist/>

³<https://github.com/CogStack/MedCAT2>

3.2 Concept-Knowledge Alignment

To guarantee that the simplified captions retain the diagnostic precision of T_{exp} , we first extract a set of semantic constraints C . This process integrates high-level ontology mapping with fine-grained entity recognition.

Ontology-Based CUI Mapping. We utilize the UMLS Metathesaurus API (Bodenreider, 2004)⁴ to ground clinical terms to CUIs. In contrast to heuristic string matching, direct API querying guarantees precise alignment with standard medical ontologies. This step captures core medical concepts (e.g., C0040405 → “CTPA”). We denote the set of identified CUIs as C_{onto} , ensuring that the pathology is rigorously anchored to standardized terminology.

Fine-Grained Entity Extraction. We supplement CUIs with a biomedical Named Entity Recognition (NER) model, SciSpacy (Neumann et al., 2019)⁵. This module explicitly extracts quantitative attributes (e.g., lesion sizes) and spatial descriptors (C_{ent}) often missed by high-level mapping. We integrate these two sources to establish the final semantic constraint set C . Formally, this is defined as:

$$C = C_{onto} \cup C_{ent} \quad (1)$$

where C_{onto} represents the high-level ontological constraints anchored to UMLS, and C_{ent} denotes the fine-grained entity constraints extracted via NER.

⁴<https://www.nlm.nih.gov/research/umls/>

⁵<https://allenai.github.io/scispacy/>

Table 1: **Linguistic Complexity and Readability Analysis.** Our refinement consistently reduces reading difficulty, improves accessibility, and standardizes vocabulary across the entire dataset.

Linguistic Metric	Train Set ($N = 59,962$)		Validation Set ($N = 9,904$)		Test Set ($N = 9,927$)		Overall (Total) ($N = 79,793$)	
	Expert	Layman	Expert	Layman	Expert	Layman	Expert	Layman
Readability Metrics								
FKGL (Kincaid et al., 1975) ↓	13.05	10.29	13.29	10.50	13.21	10.53	13.10	10.35
CLI (Coleman and Liau, 1975) ↓	15.73	9.82	16.12	10.06	16.02	10.04	15.82	9.88
DCRS (Dale and Chall, 1948) ↓	14.02	11.73	14.09	11.80	14.02	11.77	14.03	11.74
SMOG (Mc Laughlin, 1969) ↓	13.71	12.21	13.85	12.35	13.88	12.41	13.75	12.25
FRE (Flesch, 1948) ↑	26.44	56.15	24.85	55.00	25.64	55.09	26.14	55.88
Lexical Statistics								
Average Sentence Length	23.73	27.81	25.45	28.86	25.61	29.34	24.17	28.13
Vocab Size ↓	36,875	20,589	14,877	9,191	14,865	9,238	44,673	24,085

3.3 Knowledge-Constrained Refinement

Leveraging the semantic constraint set C , we synthesize the lay caption T_{lay} . This phase shifts the linguistic register while strictly adhering to the extracted medical facts.

Lexical Alignment and Draft Synthesis. For each concept in C_{onto} , we retrieve patient-friendly definitions by querying the MedlinePlus vocabulary within the UMLS Metathesaurus. Curated by the National Library of Medicine (NLM), MedlinePlus serves as the authoritative bridge between rigorous clinical ontologies and public health literacy (Miller et al., 2000)⁶. By aligning UMLS CUIs directly with MedlinePlus definitions, we ensure that the terminology is not merely simplified but standardized to a trusted lay register. We then construct an intermediate noisy lay draft (T_{draft}) via deterministic dictionary-based substitution. While grammatically noisy, T_{draft} serves as a reliable lexical basis for the subsequent refinement.

Constraint-Guided Linguistic Refinement. To generate the final accessible caption, we employ Llama-3.1-8B-Instruct (Dubey et al., 2024) within a constrained generation framework. Our structured prompt incorporates: (1) the source text T_{exp} ensuring factual grounding, (2) a strict constraint set C for hallucination mitigation, and (3) the initial draft T_{draft} to steer vocabulary selection. The objective is to downscale linguistic complexity from a college-level register to a high school level, ensuring the output remains semantically faithful to the clinical findings through explicit constraints. Figure 3 demonstrates qualitative examples of our refinement across different modalities.

⁶<https://medlineplus.gov/>

Table 2: **Dataset Statistics and Quality Consistency.**

We report consistency across Train ($N=59,962$), Validation ($N=9,904$), and Test ($N=9,927$). The **Overall** column represents the weighted average ($N=79,793$). High clinical correctness (RaTEScore, GREEN) and consistent simplification scores (LENS) across all splits confirm the robust quality of our refinement pipeline.

Metric	Train	Val	Test	Overall
Relevance				
BLEU-4 (Papineni et al., 2002) ↑	20.99	22.32	22.45	21.34
ROUGE-L (Lin, 2004) ↑	49.33	50.13	50.40	49.56
METEOR (Banerjee and Lavie, 2005) ↑	53.00	53.40	53.56	53.12
Readability				
LENS (Maddela et al., 2023) ↑	63.28	62.91	62.94	63.19
Radiological Factualty				
RaTEScore (Zhao et al., 2024b) ↑	64.66	64.57	65.09	64.70
GREEN (Ostmeier et al., 2024) ↑	69.03	70.14	70.03	69.29

4 Experiments

We demonstrate the value of **MedLayBench-V** through a comprehensive analysis of its linguistic properties and quality consistency, followed by a zero-shot downstream benchmark to evaluate current VLMs' capability in handling both expert and layman medical concepts.

4.1 Evaluation Metrics

To ensure a comprehensive assessment, we employ metrics across four dimensions: textual similarity, linguistic readability, clinical factuality, and downstream utility.

- **Relevance:** We use standard n-gram metrics to measure the structural similarity and lexical overlap between expert and layman captions. Specifically, we report BLEU-4 (Papineni et al., 2002), ROUGE-L (Lin, 2004), and METEOR (Banerjee and Lavie, 2005).

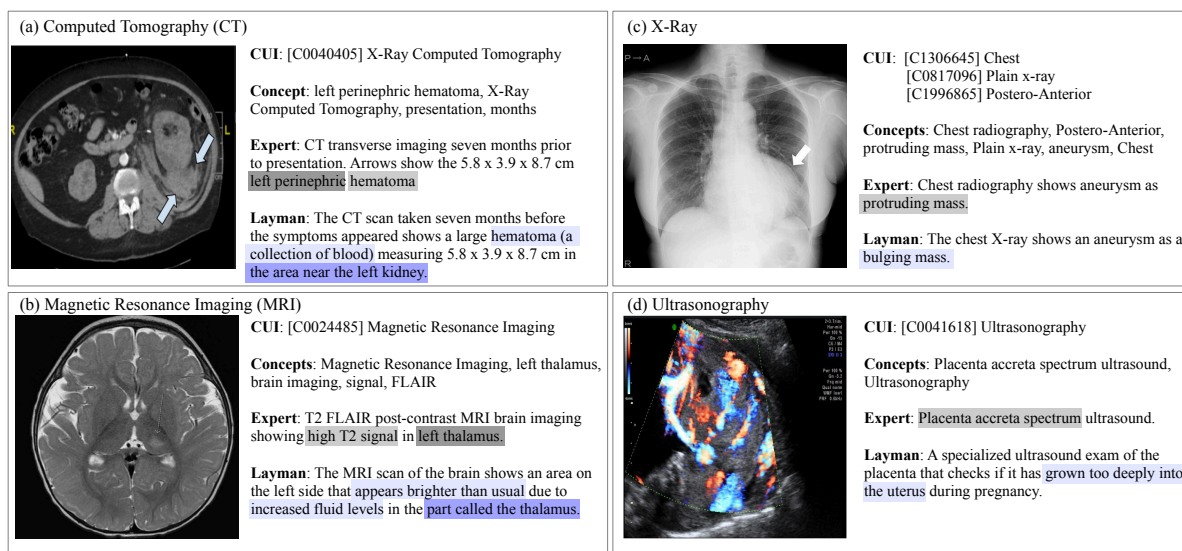


Figure 3: **Qualitative Comparison of Jargon Refinement across Modalities.** The figure illustrates example cases from CT, MRI, X-Ray, and Ultrasound. Highlights indicate the transformation from medical jargon (Original expert-level caption) to patient-friendly language (Layman-level caption). Our method successfully simplifies anatomical terms, structural definitions, and visual descriptions while preserving core medical information.

370 • **Readability:** To quantify the accessibil- 371
372 ity of the text, we utilize Flesch-Kincaid 373
374 Grade Level (FKGL) (Kincaid et al., 375
376 1975), Coleman-Liau Index (CLI) (Cole- 377
378 man and Liau, 1975), Dale-Chall Readability 379
380 Score (DCRS) (Dale and Chall, 1948), Sim- 381
382 ple Measure of Gobbledygook (SMOG) In-
383 dex (Mc Laughlin, 1969), and Flesch Read-
384 ing Ease (FRE) (Flesch, 1948). Additionally,
385 we incorporate LENS (Maddela et al., 2023),
386 a learnable metric specifically optimized for
387 text simplification.

382 • **Radiological Factuality:** Evaluating the clin- 383
384 ical integrity of simplified text is critical. 385
386 We employ Radiological Report Text Eval- 387
388 uation (RaTEScore) (Zhao et al., 2024b) and 389
390 Generative Radiology Report Evaluation and
391 Error Notation (GREEN) (Ostmeier et al.,
392 2024). These model-based metrics are de-
393 signed to detect hallucinations and ensure clin-
394 ical correctness in radiology reports.

391 • **Downstream Performance:** To assess 392
393 whether the simplified text preserves essential 394
395 semantic information for automated analysis, 396
397 we evaluate zero-shot text-to-image retrieval 398
399 performance. We report Recall@K (R@1,
400 R@5, R@10) to measure retrieval accuracy
401 using the generated captions.

4.2 MedLayBench-V Statistics and Quality Analysis

402 We analyze the linguistic characteristics and seman- 403
404 tic consistency of MedLayBench-V, which com- 405
406 prises 79,789 image-text pairs across 7 modal- 407
408 ities, maintaining the original ROCov2 configu-
409 ration (Rückert et al., 2024).

405 **Linguistic Complexity and Accessibility.** As 406
407 presented in Table 1, our refinement pipeline 408
409 successfully standardizes the linguistic complex-
410 ity of medical captions.

409 • **Vocabulary Reduction:** The unique vocabu- 410
411 lary size is reduced by 46.1% in the layman 412
413 version compared to the expert version. This 414
415 indicates a significant removal of long-tail 416
417 medical jargon and noisy tokens, streamlining 418
419 the dataset for generalizable learning.

415 • **Improved Readability:** We observe a con- 416
417 sistent drop in grade-level metrics across the 418
419 entire dataset. Notably, the FKGL drops from 420
421 13.10 to 10.35, and the Coleman-Liau Index 422
423 decreases from a graduate level of 15.82 to 424
425 9.88, aligning with the recommended read-
426 ing level for patient education materials.

422 • **Enhanced Accessibility:** The FRE score 423
424 more than doubles from 26.14 to 55.88. This 425
426 shift in text difficulty from very confusing to

Table 3: Overall top-K retrieval performance on MedLayBench-V across four modalities (X-Ray, CT, MRI, Ultrasound). **Bold** indicates best performance, underline indicates second best performance. Values are presented as Jargon / Layman. All values are in percentage (%).

Model	Image → Text			Text → Image		
	Recall@1	Recall@5	Recall@10	Recall@1	Recall@5	Recall@10
OpenAI-Base	1.23 / 1.08	3.96 / 3.74	6.56 / 6.32	1.57 / 1.54	4.41 / 4.51	7.03 / 7.12
CoCa-Large	2.10 / 2.15	5.70 / 5.71	8.24 / 8.07	3.56 / 3.64	8.78 / 8.84	11.97 / 12.11
LAION-2B	2.28 / 2.33	6.67 / 6.58	9.88 / 9.78	4.31 / 4.29	9.94 / 9.88	13.76 / 13.74
OpenCLIP-Huge	3.33 / 3.44	8.71 / 8.43	12.58 / 12.28	5.17 / 5.15	11.88 / 12.10	16.59 / 16.70
PubMedCLIP	4.61 / 4.26	13.46 / 13.12	20.93 / 20.66	4.85 / 4.71	14.49 / 14.43	21.94 / 21.73
BMC-CLIP	22.69 / 22.42	40.83 / 40.36	50.33 / 49.65	23.04 / 23.21	42.09 / 42.03	52.09 / 51.71
PMC-CLIP	<u>28.98</u> / <u>28.38</u>	<u>53.12</u> / <u>52.47</u>	<u>64.14</u> / <u>63.60</u>	<u>30.90</u> / <u>30.24</u>	<u>55.66</u> / <u>55.16</u>	<u>66.11</u> / <u>65.55</u>
BioMedCLIP	31.06 / 30.70	58.52 / 58.11	70.31 / 69.41	32.50 / 32.07	59.94 / 59.09	71.07 / 70.44

fairly difficult ensures the content is accessible to a general audience with a standard high school education.

Quality Consistency across Splits. Table 2 reports the semantic quality and consistency of our dataset.

- **Uniformity:** The relevance metrics including BLEU-4, ROUGE-L, and METEOR show minimal variance between the training, validation, and test sets. For instance, the overall METEOR score is maintained at 53.56, confirming that our pipeline generates stylistically consistent samples regardless of the data split.
- **Simplification Quality:** The LENS score, which is specifically designed for text simplification, remains robust across all splits with an overall average of 62.94. This indicates stable rewriting performance throughout the dataset.
- **Clinical Safety:** Most importantly, the clinical correctness scores represented by RaTEScore and GREEN demonstrate that our simplification process preserves the factual integrity of the original medical reports. Specifically, the test set achieves a RaTEScore of 65.09 and a GREEN score of 70.03, indicating high clinical safety despite the reduced linguistic complexity.

4.3 Downstream Task: Zero-Shot Retrieval

To evaluate the utility of MedLayBench-V, we conducted a zero-shot Image-Text Retrieval (ITR) experiment. This task measures how well models can

align visual features with both *Expert* (original) and *Layman* (refined) textual descriptions. We report the Recall@K metrics for both Image-to-Text and Text-to-Image retrieval in Table 3 and visualize the results in Figure 4.

Baseline Models. We benchmarked a diverse array of dual-encoder architectures, categorized into general-domain and medical-domain models. For the general domain, we employed OpenAI-CLIP (Radford et al., 2021) and OpenCLIP (Cherti et al., 2023) (trained on LAION-2B (Schuhmann et al., 2022)), along with CoCa (Yu et al., 2022), which integrates contrastive and generative objectives. For the medical domain, we selected models pre-trained on large-scale biomedical image-text pairs to assess the impact of domain adaptation. These include PubMedCLIP (Eslami et al., 2023), BMC-CLIP, PMC-CLIP (Lin et al., 2023), and BioMedCLIP (Zhang et al., 2023), which utilize domain-specific encoders aligned with biomedical imagery.

Performance of Medical vs. General VLMs. We observe a clear performance hierarchy based on domain adaptation. While general domain models (e.g., OpenAI-CLIP) struggle with medical contexts (Recall@1 < 5%), medical-specific models show significantly improved alignment. Notably, BioMedCLIP achieves state-of-the-art performance, benefiting from its large-scale pre-training on biomedical literature.

Semantic Preservation in Layman Captions. Crucially, our results demonstrate that simplifying the language does not compromise semantic

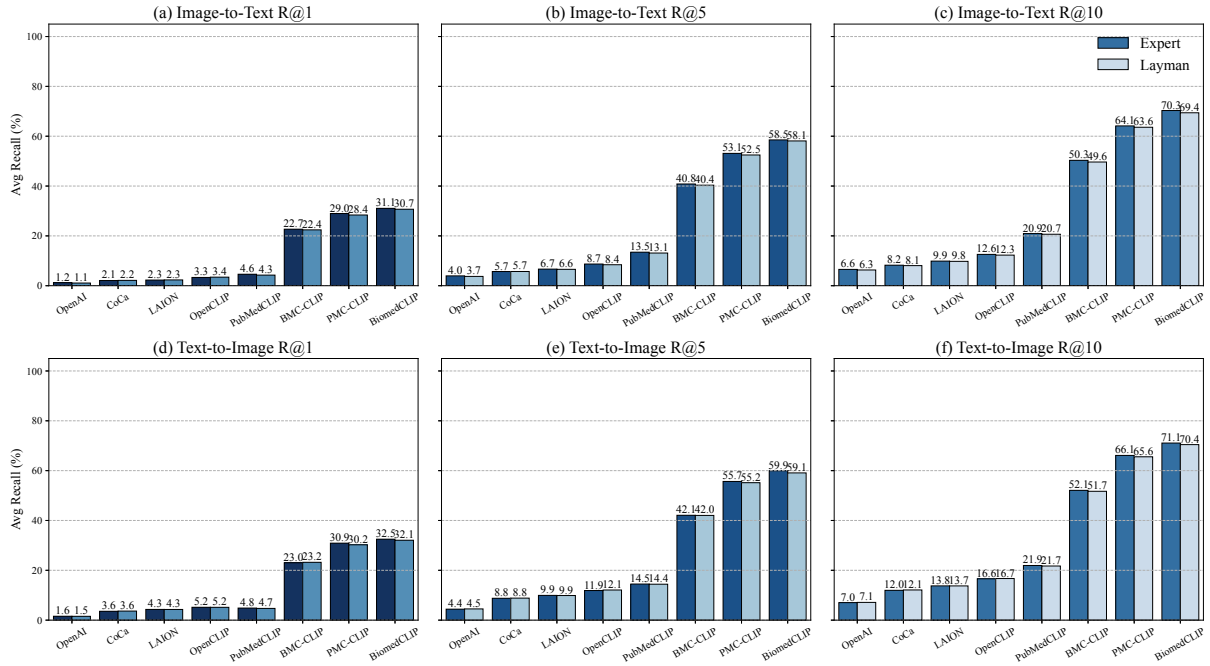


Figure 4: **Zero-Shot Retrieval Performance.** Recall@ K results for Image-to-Text (a-c) and Text-to-Image (d-f) tasks. Dark and light blue bars denote Expert and Layman queries, respectively. Medical VLMs outperform general models with negligible gaps between expert and layman inputs, confirming that our simplification preserves semantic fidelity.

489 fidelity. As evidenced in Table 3, retrieval perfor-
 490 mance remains robust across all medical mod-
 491 els, exhibiting negligible degradation when tran-
 492 sitioning from *Expert* to *Layman* queries. For in-
 493 stance, BioMedCLIP exhibits only a marginal drop
 494 in Image-to-Text Recall@1 (31.06% \rightarrow 30.70%).
 495 This explicitly verifies that MedLayBench-V suc-
 496 cessfully retains the core diagnostic semantics re-
 497 quired for visual alignment, proving that high read-
 498 ability can be achieved without sacrificing medical
 499 accuracy.

500 5 Conclusion

501 In this work, we introduced **MedLayBench-V**, the
 502 first multimodal benchmark designed to quantify
 503 the semantic alignment between clinical jargon
 504 and lay language. By evaluating state-of-the-art
 505 VLMs, we formalized the existence of a represen-
 506 tation alignment gap, revealing that current medical
 507 models are overfitted to the professional register at
 508 the expense of patient accessibility. Our proposed
 509 structured concept-grounded refinement pipeline
 510 provides a foundational framework for developing
 511 next-generation Medical AI that is both clinically
 512 accurate and universally understandable.

513 Limitations

514 While MedLayBench-V establishes a foundation
 515 for patient-centric AI, we acknowledge limitations
 516 regarding the reliance on synthetic data, restriction
 517 to English, and modality imbalances inherited from
 518 the source. Although our pipeline ensures clinical
 519 correctness via structured constraints, synthetic
 520 captions may lack the subtle nuances of human-
 521 authored text, and validation with diverse patient
 522 groups is necessary to fully assess real-world util-
 523 ity.

524 More importantly, we hypothesize that the rep-
 525 resentation alignment gap between clinical jargon
 526 and lay language may have been obscured by the
 527 limited complexity of the current retrieval task. We
 528 posit that a distinct gap exists but requires more
 529 challenging scenarios to be fully exposed. Con-
 530 sequently, our future work will focus on scaling
 531 this benchmark to a wider array of complex down-
 532 stream tasks. By increasing both the scale and
 533 difficulty, we aim to rigorously identify this latent
 534 alignment gap and develop robust methodologies
 535 to effectively bridge the expert-lay divide.

References

- 537 Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An
538 automatic metric for mt evaluation with improved cor-
539 relation with human judgments. In *Proceedings of*
540 *the acl workshop on intrinsic and extrinsic evaluation*
541 *measures for machine translation and/or summariza-*
542 *tion*, pages 65–72.
- 543 Olivier Bodenreider. 2004. The unified medical lan-
544 guage system (umls): integrating biomedical termi-
545 nology. *Nucleic acids research*, 32(suppl_1):D267–
546 D270.
- 547 Mehdi Cherti, Romain Beaumont, Ross Wightman,
548 Mitchell Wortsman, Gabriel Ilharco, Cade Gordon,
549 Christoph Schuhmann, Ludwig Schmidt, and Jenia
550 Jitsev. 2023. Reproducible scaling laws for con-
551 trastive language-image learning. In *Proceedings*
552 *of the IEEE/CVF conference on computer vision and*
553 *pattern recognition*, pages 2818–2829.
- 554 Meri Coleman and Ta Lin Liao. 1975. A computer
555 readability formula designed for machine scoring.
556 *Journal of Applied Psychology*, 60(2):283.
- 557 Edgar Dale and Jeanne S Chall. 1948. A formula for
558 predicting readability: Instructions. *Educational re-*
559 *search bulletin*, pages 37–54.
- 560 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,
561 Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,
562 Akhil Mathur, Alan Schelten, Amy Yang, Angela
563 Fan, and 1 others. 2024. The llama 3 herd of models.
564 *arXiv preprint arXiv:2407.21783*.
- 565 Sedigheh Eslami, Christoph Meinel, and Gerard
566 De Melo. 2023. Pubmedclip: How much does clip
567 benefit visual question answering in the medical do-
568 main? In *Findings of the Association for Computa-*
569 *tional Linguistics: EACL 2023*, pages 1181–1193.
- 570 Rudolph Flesch. 1948. A new readability yardstick.
571 *Journal of applied psychology*, 32(3):221.
- 572 Tomas Goldsack, Carolina Scarton, Matthew Shard-
573 low, and Chenghua Lin. 2024. Overview of the bi-
574 olaysumm 2024 shared task on the lay summariza-
575 tion of biomedical research articles. *arXiv preprint*
576 *arXiv:2408.08566*.
- 577 J Peter Kincaid, Robert P Fishburne Jr, Richard L
578 Rogers, and Brad S Chissom. 1975. Derivation of
579 new readability formulas (automated readability in-
580 dex, fog count and flesch reading ease formula) for
581 navy enlisted personnel. Technical report.
- 582 Zeljko Kraljevic, Thomas Searle, Anthony Shek, Lukasz
583 Roguski, Kawsar Noor, Daniel Bean, Aurelie Mas-
584 cio, Leilei Zhu, Amos A Folarin, Angus Roberts, and
585 1 others. 2021. Multi-domain clinical natural lan-
586 guage processing with medcat: the medical concept
587 annotation toolkit. *Artificial intelligence in medicine*,
588 117:102083.
- Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto
Usuyama, Haotian Liu, Jianwei Yang, Tristan Nau-
mann, Hoifung Poon, and Jianfeng Gao. 2023. Llava-
med: Training a large language-and-vision assistant
for biomedicine in one day. *Advances in Neural In-*
formation Processing Systems, 36:28541–28564.
- Weibin Liao, Tianlong Wang, Yinghao Zhu, Yasha
Wang, Junyi Gao, and Liantao Ma. 2025. Magi-
cal: Medical lay language generation via semantic
invariance and layperson-tailored adaptation. *arXiv*
preprint arXiv:2508.08730.
- Chin-Yew Lin. 2004. Rouge: A package for automatic
evaluation of summaries. In *Text summarization*
branches out, pages 74–81.
- Weixiong Lin, Ziheng Zhao, Xiaoman Zhang, Chaoyi
Wu, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023.
Pmc-clip: Contrastive language-image pre-training
using biomedical documents. In *International Con-*
ference on Medical Image Computing and Computer-
Assisted Intervention, pages 525–536. Springer.
- Alejandro Lozano, Min Woo Sun, James Burgess,
Liangyu Chen, Jeffrey J Nirschl, Jeffrey Gu, Ivan
Lopez, Josiah Aklilu, Anita Rau, Austin Wolf-
gang Katzer, and 1 others. 2025. Biomedica: An
open biomedical image-caption archive, dataset, and
vision-language models derived from scientific litera-
ture. In *Proceedings of the Computer Vision and Pat-*
tern Recognition Conference, pages 19724–19735.
- Mounica Maddela, Yao Dou, David Heineman, and
Wei Xu. 2023. Lens: A learnable evaluation metric
for text simplification. In *Proceedings of the 61st*
Annual Meeting of the Association for Computational
Linguistics (Volume 1: Long Papers), pages 16383–
16408.
- G Harry Mc Laughlin. 1969. Smog grading-a new read-
ability formula. *Journal of reading*, 12(8):639–646.
- Naomi Miller, Eve-Marie Lacroix, and Joyce EB
Backus. 2000. Medlineplus: building and main-
taining the national library of medicine’s consumer
health web service. *Bulletin of the Medical Library*
Association, 88(1):11.
- Shufan Ming, Yue Guo, and Halil Kilicoglu. 2025. To-
wards knowledge-guided biomedical lay summariza-
tion using large language models. In *Proceedings of*
the Second Workshop on Patient-Oriented Language
Processing (CL4Health), pages 285–297.
- Michael Moor, Oishi Banerjee, Zahra Shakeri Hossein
Abad, Harlan M Krumholz, Jure Leskovec, Eric J
Topol, and Pranav Rajpurkar. 2023. Foundation mod-
els for generalist medical artificial intelligence. *Nat-*
ure, 616(7956):259–265.
- Mark Neumann, Daniel King, Iz Beltagy, and Waleed
Ammar. 2019. Scispacy: fast and robust models
for biomedical natural language processing. *arXiv*
preprint arXiv:1902.07669.

644	Sophie Ostmeier, Justin Xu, Zhihong Chen, Maya Varma, Louis Blankemeier, Christian Bluethgen, Arne Edward Michalson Md, Michael Moseley, Curtis Langlotz, Akshay S Chaudhari, and 1 others. 2024. Green: Generative radiology report evaluation and error notation. In <i>Findings of the association for computational linguistics: EMNLP 2024</i> , pages 374–390.	701
645		702
646		703
647		704
648		705
649		706
650		707
651		
652	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pages 311–318.	708
653		709
654		710
655		711
656		
657	Philipp Prucker, Keno K Bressemer, Jan Peeken, Mateo Jukic, Alexander W Marka, Maximilian Strenzke, Su Hwan Kim, Christian J Mertens, Dominik Weller, Tristan Lemke, and 1 others. 2025. A prospective controlled trial of large language model-based simplification of oncologic ct reports for patients with cancer. <i>Radiology</i> , 317(2):e251844.	712
658		713
659		714
660		715
661		716
662		717
663		
664	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, and 1 others. 2021. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pages 8748–8763. PmLR.	718
665		719
666		720
667		721
668		
669		
670		
671	Johannes Rückert, Louise Bloch, Raphael Brüngel, Ahmad Idrissi-Yaghir, Henning Schäfer, Cynthia S Schmidt, Sven Koitka, Obioma Pelka, Asma Ben Abacha, Alba G. Seco de Herrera, and 1 others. 2024. Rocov2: Radiology objects in context version 2, an updated multimodal image dataset. <i>Scientific Data</i> , 11(1):688.	722
672		723
673		724
674		725
675		726
676		
677		
678	Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, and 1 others. 2022. Laion-5b: An open large-scale dataset for training next generation image-text models. <i>Advances in neural information processing systems</i> , 35:25278–25294.	727
679		728
680		729
681		730
682		
683		
684		
685	Andrew Sellergren, Sahar Kazemzadeh, Tiam Jaroensri, Atilla Kiraly, Madeleine Traverse, Timo Kohlberger, Shawn Xu, Fayaz Jamil, Cían Hughes, Charles Lau, and 1 others. 2025. Medgemma technical report. <i>arXiv preprint arXiv:2507.05201</i> .	731
686		732
687		733
688		734
689		735
690		736
691		
692		
693	Chenghao Xiao, Kun Zhao, Xiao Wang, Siwei Wu, Sixing Yan, Tomas Goldsack, Sophia Ananiadou, Noura Al Moubayed, Liang Zhan, William K Cheung, and 1 others. 2025. Overview of the biolaysumm 2025 shared task on lay summarization of biomedical research articles and radiology reports. In <i>Proceedings of the 24th Workshop on Biomedical Language Processing</i> , pages 365–377.	737
694		738
695		739
696		740
697		741
698		742
699		743
700		
	Zonghai Yao, Nandyala Siddharth Kantu, Guanghao Wei, Hieu Tran, Zhangqi Duan, Sunjae Kwon, Zhichao Yang, and Hong Yu. 2024. Readme: Bridging medical jargon and lay understanding for patient education through data-centric nlp. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 12609–12629.	744
		745
		746
		747
		748
		749
		750
		751
		752
		753
		754
	Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. 2022. Coca: Contrastive captioners are image-text foundation models. <i>arXiv preprint arXiv:2205.01917</i> .	755
		756
		757
		758
		759
		760
		761
		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800

A Implementation Details and Prompts

In this section, we provide a comprehensive breakdown of the SCGR pipeline’s implementation. The core of our approach lies in the rigorous separation of semantic extraction and stylistic refinement, as detailed in Algorithm 1. To ensure that the LLM adheres strictly to clinical facts while simplifying the syntax, we engineered a specific prompt template shown in Figure A2. By explicitly defining the system role as a "Medical Text Simplifier" and enforcing a JSON output format, we enable seamless

755 integration into downstream training pipelines. The
756 "Critical Instructions" block serves as a safeguard
757 against common pitfalls such as hallucinations or
758 the use of subjective pronouns (e.g., "your body"),
759 ensuring the output remains objective and profes-
760 sional.

761 B Detailed Dataset Statistics

762 MedLayBench-V encompasses a diverse range of
763 medical imaging modalities, mirroring real-world
764 clinical distributions. As summarized in Table 4,
765 Computed Tomography (CT) and X-Ray consti-
766 tute the majority of the dataset, reflecting their
767 prevalence in diagnostic radiology. Table 5 fur-
768 ther breaks down the top co-occurring concepts
769 for each modality, confirming that our extrac-
770 tion pipeline correctly identifies modality-specific
771 anatomical structures (e.g., "left ventricle" in Ultra-
772 sound, "coronary artery" in Angiography). Addi-
773 tionally, Figure A3 illustrates the long-tail distribu-
774 tion of both UMLS concepts and raw terms. This
775 indicates that while common terms (head, tail) are
776 frequent, the dataset also preserves a vast array of
777 rare, specific medical conditions, which is crucial
778 for training robust medical VLMs.

779 C Semantic Preservation Analysis

780 To empirically validate that our simplification
781 process preserves the underlying medical seman-
782 tics, we analyzed the embedding space of various
783 Vision-Language Models. Figure A4 visualizes
784 the t-SNE projections of image-text embeddings
785 for both Expert (original) and Layman (refined)
786 captions. Across different architectures (OpenAI-
787 CLIP, BioMedCLIP, PMC-CLIP), we observe that
788 the distributions of Expert and Layman embed-
789 dings are nearly isomorphic. Furthermore, the high
790 cosine similarity (≈ 0.99) and low Euclidean dis-
791 tance distributions confirm that the transition to
792 lay language does not shift the semantic vector
793 significantly. This serves as strong evidence that
794 MedLayBench-V successfully lowers the linguis-
795 tic barrier without compromising the diagnostic
796 information required for model training.

797 D Extended Qualitative Analysis

798 To further demonstrate the robustness and versatil-
799 ity of the SCGR pipeline, we provide an extended
800 set of qualitative examples across diverse imaging
801 modalities. Figure A5 and Figure A6 illustrate

802 how our pipeline handles specific linguistic chal-
803 lenges, ranging from simplifying complex vascular
804 anatomy in CT/MRI to interpreting acoustic
805 artifacts in ultrasound. Each example highlights
806 the transformation from the original expert report
807 (**Expert**) to the generated patient-friendly caption
808 (**Layman**). Key medical terms are highlighted in
809 gray/red while their simplified explanations are
810 highlighted in blue to visualize the semantic align-
811 ment.

812 E Ablation Study: Impact of Naive LLM 813 Simplification

814 To rigorously validate the necessity of our Struc-
815 tured Concept-Grounded Refinement (SCGR)
816 pipeline, we conducted an ablation study compar-
817 ing our method against a standard "Naive LLM"
818 baseline. In this baseline setting, we prompted the
819 same backbone model (Llama-3.1-8B) to simplify
820 the expert reports using a generic instruction (e.g.,
821 "*Rewrite this medical report for a patient*") without
822 the intermediate concept mapping or entity con-
823 straints.

824 **Catastrophic Performance Drop.** As illustrated
825 in Figure A1, relying solely on the LLM's para-
826 metric knowledge results in a catastrophic degra-
827 dation of downstream retrieval performance. For
828 instance, the Image-to-Text Recall@1 for BioMed-
829 CLIP plummets from **31.06%** (using our SCGR
830 data) to approximately **5.3%** (using Naive LLM
831 data). Similar trends are observed across all evalu-
832 ated VLMs and modalities.

833 **Analysis of Failure Modes.** Qualitative inspec-
834 tion reveals that while the Naive LLM generates flu-
835 ent and readable text, it suffers from severe **seman-**
836 **tic drift** and **hallucination**. Without structured
837 constraints, the model tends to: (1) over-simplify
838 specific pathologies into vague terms (e.g., convert-
839 ing "pneumothorax" into generic "lung problem"),
840 losing the discriminative features required for re-
841 trieval; or (2) hallucinate plausible but incorrect
842 details to fill narrative gaps. This ablation confirms
843 that high-quality medical lay language generation
844 cannot be achieved by LLMs alone and necessitates
845 the explicit knowledge grounding provided by our
846 SCGR framework.

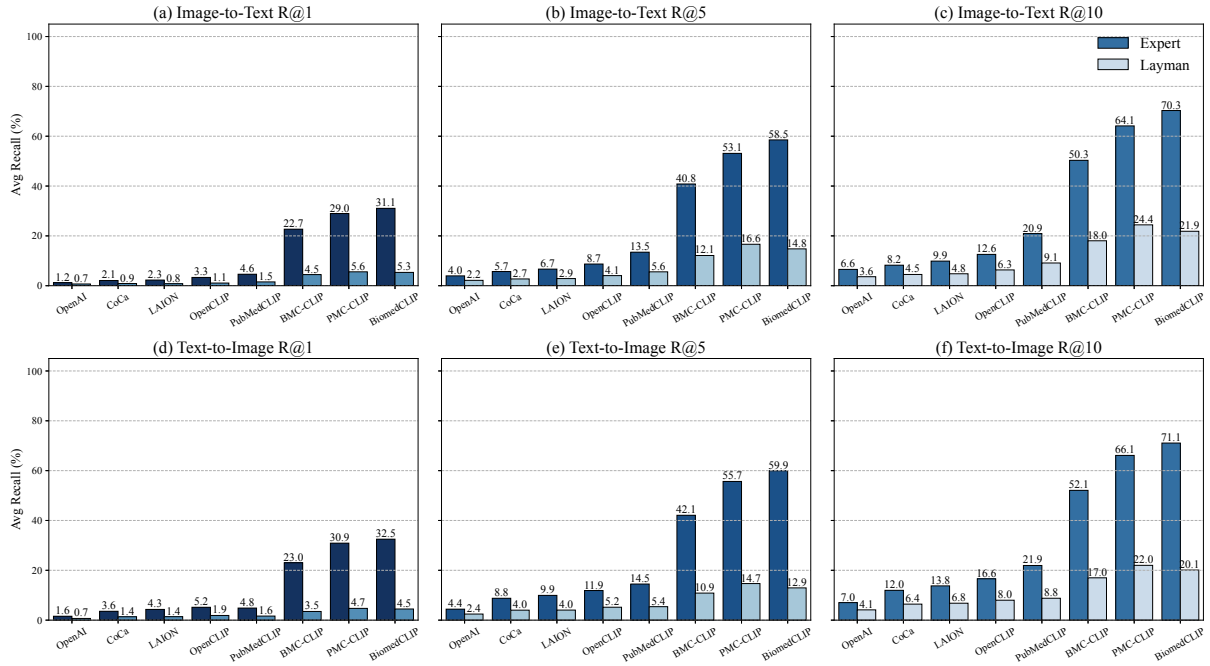


Figure A1: **Impact of Naive LLM-only Simplification.** Performance analysis using layman captions generated by a standard LLM without structured grounding. Unlike our SCGR method (Main text Figure 4), the naive approach causes severe semantic drift, leading to a drastic drop in retrieval accuracy (e.g., BioMedCLIP R@1 drops from 31.1% to 5.3%). This highlights the necessity of our concept-grounded pipeline.

Algorithm 1: SCGR framework

```

Input : Set of Expert Captions
           $\mathcal{T}_{exp} = \{T_{exp}^{(1)}, \dots, T_{exp}^{(N)}\}$ 
Output: Set of Layman Captions  $\mathcal{T}_{lay}$ 

1 Initialize  $\mathcal{T}_{lay} \leftarrow \emptyset$ 
2 foreach  $T_{exp} \in \mathcal{T}_{exp}$  do
   // Step 1: Hybrid Concept Extraction
3    $C_{onto} \leftarrow \text{ExtractCUIs}(T_{exp})$  // MedCAT
4    $C_{ent} \leftarrow \text{ExtractEntities}(T_{exp})$  // SciSpacy
5    $C \leftarrow C_{onto} \cup C_{ent}$ 

   // Step 2: Knowledge Retrieval & Drafting
6    $T_{draft} \leftarrow T_{exp}$ 
7   foreach  $c \in C$  do
8      $def \leftarrow \text{RetrieveLayDef}(c)$ 
   // MedlinePlus
9      $T_{draft} \leftarrow \text{Substitute}(T_{draft}, c, def)$ 
10  end

   // Step 3: Constrained Refinement (LLM)
11   $P \leftarrow \text{ConstructPrompt}(T_{exp}, C, T_{draft})$ 
12   $T_{lay} \leftarrow \text{Generate}(P)$  // Llama-3

   // Step 4: Quality Verification
13  if  $\text{CheckFactuality}(T_{lay}, T_{exp})$  then
14     $\mathcal{T}_{lay}.add(T_{lay})$ 
15  end
16 end
17 return  $\mathcal{T}_{lay}$ 

```

SCGR Instruction Prompt Template

[System Role]

You are a precise **Medical Text Simplifier**. Rewrite the report for a high school student using the provided Concepts.

[Critical Instructions]

- Source of Truth:** Trust the Original Caption completely. Ignore hallucinations in the Draft.
- Objective Tone:** No 'you'/'your'. Use 'the patient' or 'the body'.
- Strict Format:** Return **ONLY** the refined sentence. No "Note:" or explanations.
- No Hallucinations:** Do not invent words. Keep unclear terms in parentheses.

[User Input Template]

Original (Fact): "{Expert Caption (T_{exp})}"
Concepts: [{"Verified UMLS Concepts (C)}]"
Draft (Ref): "{Noisy Layman Draft (T_{draft})}"
[Structured Output]

```

{
  "layman_caption": "The CT scan shows an enlarged heart..."
}

```

Figure A2: **Prompt Construction for SCGR.** The prompt enforces strict adherence to the *Original Caption* as the source of truth while utilizing the *Draft* only for stylistic reference. The output is constrained to an objective, third-person tone.

Table 4: Distribution of Imaging Modalities. The number of image-caption pairs for each modality as reported in the original ROCov2 dataset (Rückert et al., 2024).

Code	Modality Name	Count
DRCT	Computed Tomography (CT)	27,747
DRXR	X-Ray (Plain Radiography)	21,997
DRMR	Magnetic Resonance Imaging (MRI)	12,657
DRUS	Ultrasonography	11,429
DRAN	Angiography	4,799
DRCO	Combined Modality	728
DRPE	Positron Emission Tomography (PET)	432
Total		79,789

Table 5: Detailed Top 5 Concepts Distribution per Modality. The frequency of the top 5 co-occurring concepts extracted from the text context for each major imaging modality.

(a) Computed Tomography (CT)			(b) Magnetic Resonance Imaging (MRI)		
Rank	Concept	Freq	Rank	Concept	Freq
1	X-Ray Computed Tomography	27,747	1	Magnetic Resonance Imaging	12,659
2	CT scan	3,474	2	arrow	1,246
3	abdomen	2,869	3	image	957
4	arrow	2,802	4	patient	712
5	image	1,669	5	brain	661

(c) Ultrasonography			(d) Plain X-Ray		
Rank	Concept	Freq	Rank	Concept	Freq
1	Ultrasonography	11,422	1	Plain x-ray	21,936
2	arrow	892	2	Anterior-Posterior (AP)	9,606
3	image	633	3	Chest	7,196
4	left ventricle	610	4	Postero-Anterior (PA)	4,302
5	left atrium	564	5	Bone structure of cranium	3,973

(e) Angiography			(f) Positron-Emission Tomography (PET)		
Rank	Concept	Freq	Rank	Concept	Freq
1	angiogram	4,766	1	Positron-Emission Tomography	432
2	arrow	435	2	uptake	98
3	Coronary angiography	228	3	increased	64
4	right coronary artery	219	4	image	40
5	stenosis	195	5	patient	34

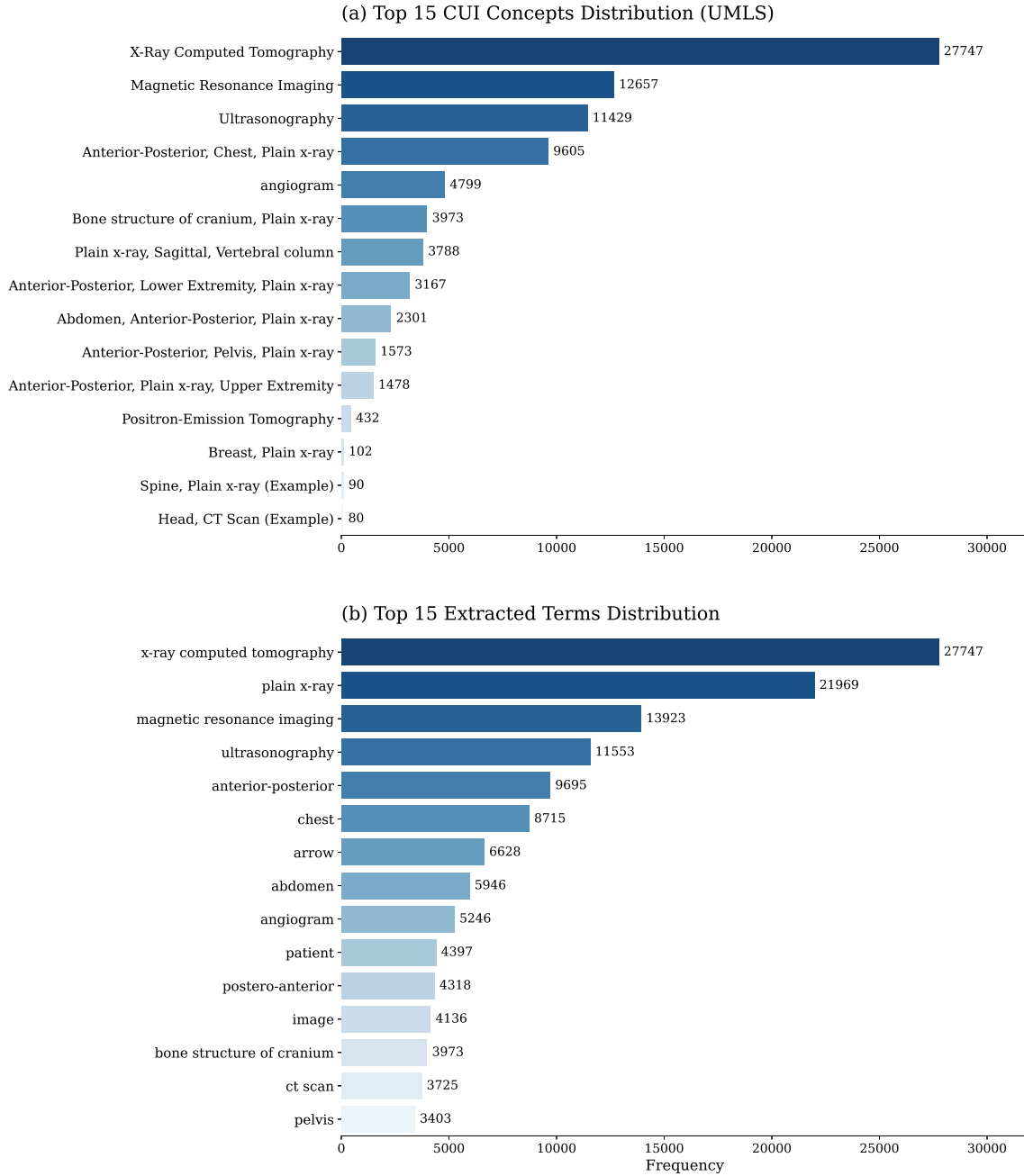


Figure A3: **Distribution of Top 15 Concepts and Terms.** (a) The frequency of Unique Medical Language System (UMLS) Concept Unique Identifiers (CUIs) mapped from the dataset. (b) The frequency of raw extracted terms directly from the captions. Both distributions illustrate the long-tail nature of medical findings in the dataset.

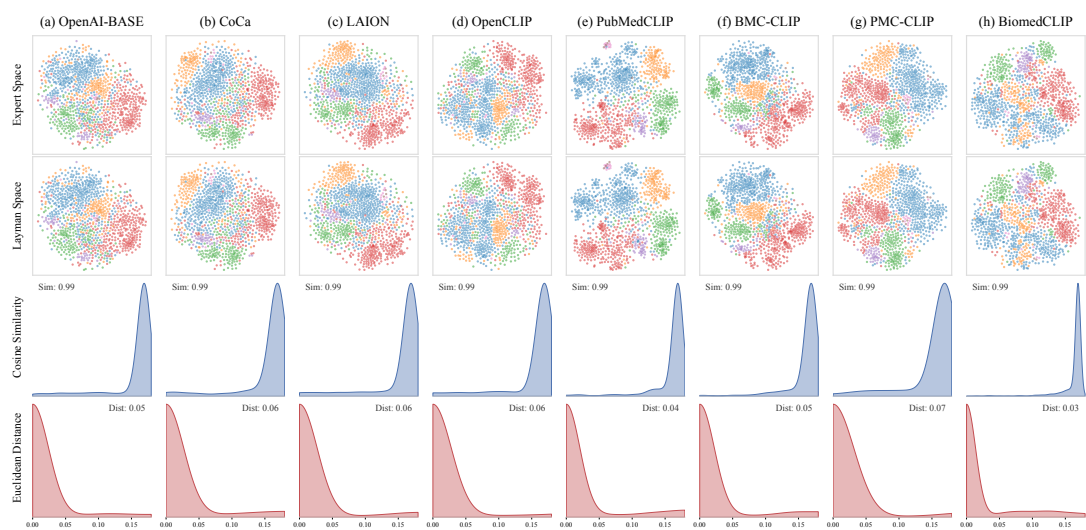
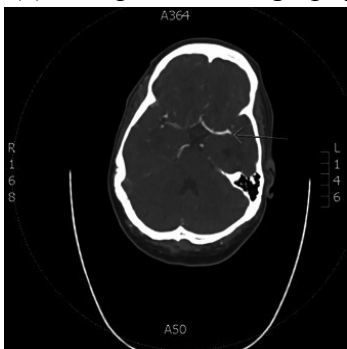


Figure A4: **Embedding space visualization across different CLIP models.** Each column represents a different vision-language model. Row 1-2: t-SNE projections of expert (original) and layman (simplified) text embeddings, colored by imaging modality. Row 3: Distribution of cosine similarity between paired expert-layman embeddings. Row 4: Distribution of Euclidean distance between paired embeddings. High cosine similarity ($\text{Sim} \approx 0.99$) and low Euclidean distance ($\text{Dist} \approx 0.05\text{-}0.07$) indicate that semantic meaning is largely preserved after text simplification across all models.

(a) Computed Tomography (CT)

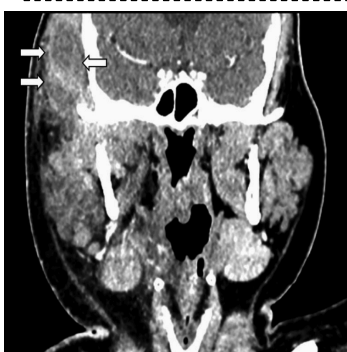


CUI: [C0040405] X-Ray Computed Tomography

Concept: X-Ray Computed Tomography, anterior cerebral artery, internal carotid artery, occlusive thrombi, M2 segment, angiogram

Expert: CT brain perfusion and angiogram on day 16, showing left MCA and ACA territories with occlusive thrombi demonstrated within the inferior division of the M2 segment of the left MCA and within the callosomarginal branch of the left ACA. There was also a complete left ICA occlusion.

Layman: A CT scan showed that there were blood clots blocking part of two main arteries supplying oxygen-rich blood to one side of the brain - specifically, the lower section of the middle cerebral artery's smaller branches and a small branch off the front part of another major artery called the anterior cerebral artery. Additionally, the main pipeline bringing oxygenated blood from the neck to this area had completely blocked.



CUI: [C0040405] X-Ray Computed Tomography

Concept: X-Ray Computed Tomography, right temporalis muscle, rim-enhancing, intramuscular, white arrows, hypodensity

Expert: On coronal cut, there is presence of rim-enhancing intramuscular hypodensity within the bulky (white arrows) and thickened right temporalis muscle with enhancement of the right temporalis muscle.

Layman: There is an area on a CT scan where part of the right temple muscle appears swollen and has some empty space inside it that shows up as darker than usual, indicated by the thicker-than-normal right temple muscle marked by the white arrows.

(b) Magnetic Resonance Imaging (MRI)

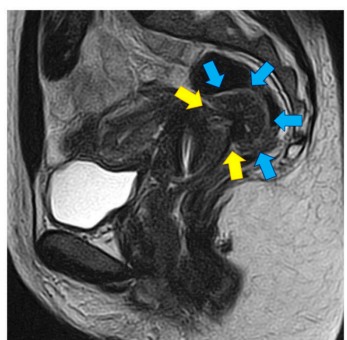


CUI: [C0024485] Magnetic Resonance Imaging

Concept: Inferoseptal late gadolinium enhancement, Magnetic Resonance Imaging, myocardial fibrosis, left ventricle, regional scar

Expert: Inferoseptal late gadolinium enhancement on MRI indicating regional scar and/or myocardial fibrosis. LV, left ventricle.

Layman: The heart's MRI shows signs of scarring or thickening inside one part of the muscle that pumps blood out of the heart, specifically the area between two chambers called the septum. This can be due to damage from past injury or disease affecting the main pumping chamber known as the left ventricle.



CUI: [C0024485] Magnetic Resonance Imaging

Concept: Magnetic Resonance Imaging, T2-weighted sagittal view, anterior rectal wall, posterior part, yellow arrow, infiltrates

Expert: Pelvic MRI with a T2-weighted sagittal view of a DE nodule (blue arrow) invading the anterior rectal wall. The nodule infiltrates the anterior rectal wall at the level of posterior part of the cervix (yellow arrow).

Layman: A pelvic Magnetic Resonance Image shows a growth (DE nodule indicated by the blue arrow) that has invaded the front part of the rectum's outer layer. This growth also affects the area near the back portion of the uterus (cervix), which is located behind it.

Figure A5: **Qualitative Analysis on Cross-Sectional Modalities.** Comparison of expert and layman descriptions for (a) CT and (b) MRI.

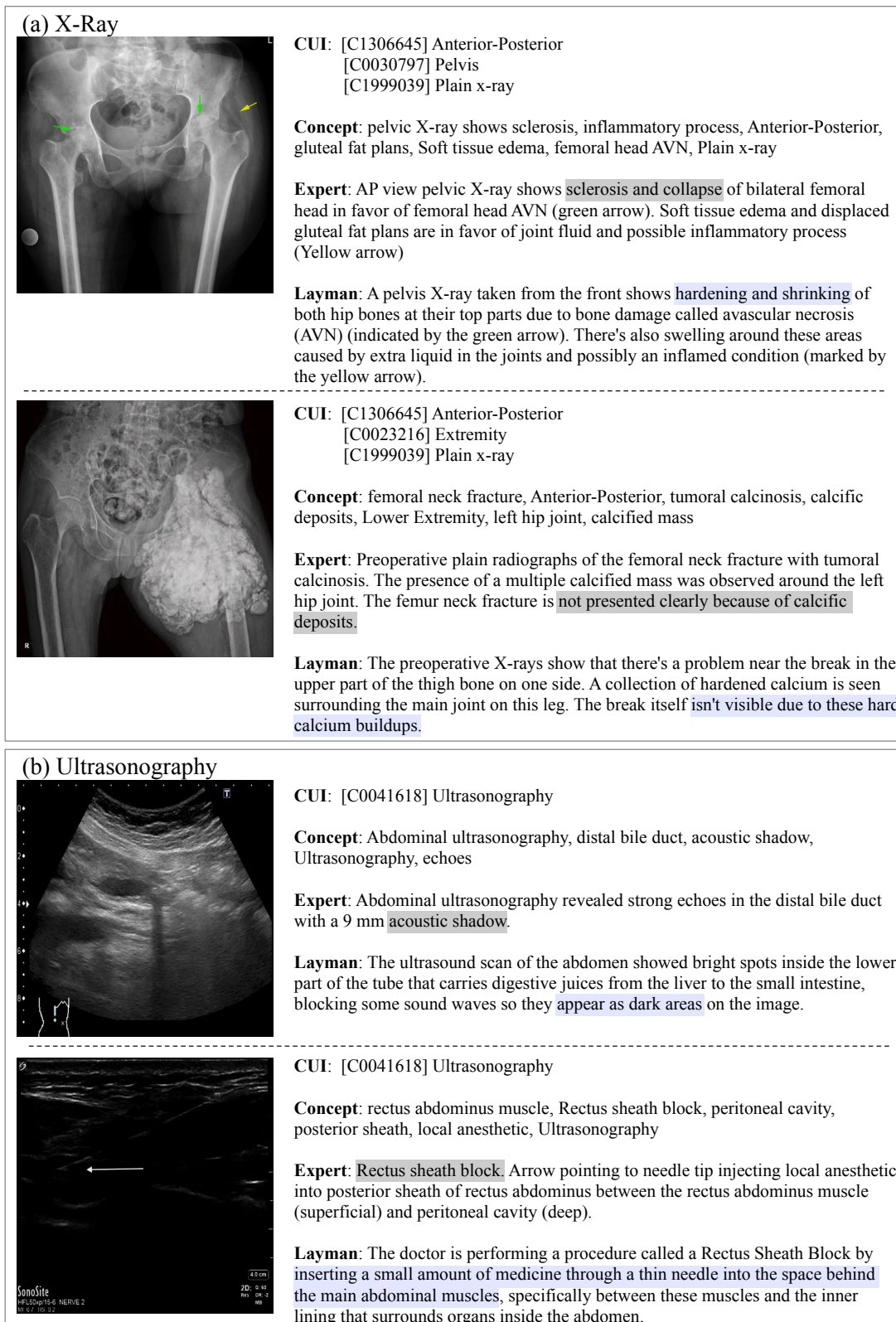


Figure A6: **Qualitative Analysis on Cross-Sectional Modalities.** Comparison of expert and layman descriptions for (a) X-Ray and (b) Ultrasonography.