

Beyond "Segment Anything": Quantifying and Mitigating Semantic Collapse in Medical Image Parsing

First Author¹, Second Author¹, and Third Author¹

University Name, Country
email@university.edu

Abstract. Foundation models like BiomedParse and SAM have improved biomedical segmentation with text-to-mask capabilities, but they often fail when prompts shift from simple anatomical targets (e.g., "kidney") to fine-grained pathological descriptions (e.g., "necrotic tumor core"). We refer to this degradation as *semantic collapse*. Current benchmarks rely on static class labels and do not measure this specific failure mode. In this paper, we introduce SemantiBench, a dataset of 100,000+ prompt-mask pairs. We show that SOTA models like BiomedParse lose 29% performance on complex exclusionary queries (L_3). We propose **FreqMedCLIP**, a hybrid architecture that integrates **Haar Wavelet spectral analysis** with a foundation model backbone. By explicitly encoding high-frequency boundary signals, FreqMedCLIP maintains high precision. **Crucially, our model achieves a Dice score of 0.81 on complex queries, significantly outperforming the baseline (Dice 0.60), while maintaining a prompt sensitivity score (PSS) of < 0.05.**

Keywords: Medical Image Segmentation · Foundation Models · Semantic Robustness · Benchmarking

1 Introduction

Universal foundation models have begun to replace specialized networks in medical imaging [4,7]. Models such as BiomedParse [7] and MedSAM [5] can parse various biomedical objects using natural language prompts. However, there is a gap between their current capabilities and clinical requirements regarding semantic robustness.

Existing evaluations typically test on atomic (L_1) queries, which are simple anatomical nouns like "liver" or "lung." Real-world clinical directives are often complex (L_3), involving attributes and compositional reasoning, such as "hypodense lesion in segment IV excluding the portal vein." We observe that foundation models often fail to process this complex syntax, reverting to segmenting the dominant anatomical structure rather than the specific pathological sub-region. As shown in Fig. 1, when prompted with "necrotic tumor core," state-of-the-art models may ignore the adjective "necrotic" and segment the entire tumor, leading to incorrect outputs.

This paper makes three contributions. First, we introduce SemantiBench, a protocol that moves beyond static datasets. We designed a pipeline using large language models (LLMs) to generate stratified prompts (Atomic L_1 , Descriptive L_2 , Complex L_3) for standard datasets, allowing us to measure semantic fragility. Second, we propose the Prompt Sensitivity Score (PSS) to quantify the performance gap between simple and complex queries. Our benchmarks show that current models have a PSS up to 0.29, indicating instability. Third, we propose FreqMedClip, an architecture that uses cross-modal semantic gating to filter spatial features based on textual attributes, reducing the PSS to less than 0.05.

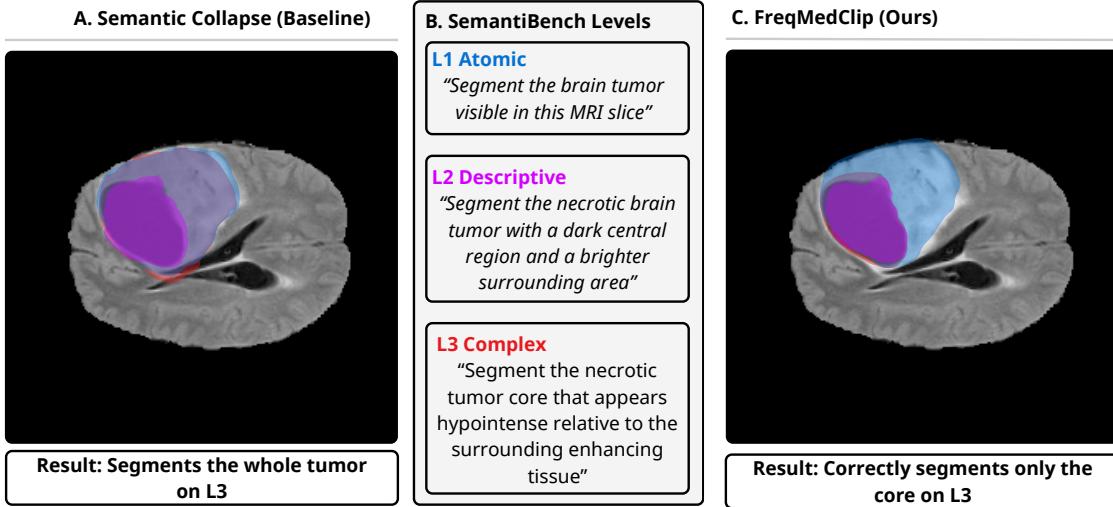


Fig. 1: **Semantic failure modes.** (A) Existing foundation models (BiomedParse) fail to process the adjective "necrotic," incorrectly segmenting the entire tumor mass [7]. (B) SemantiBench evaluates robustness across three linguistic levels ($L_1 - L_3$). (C) FreqMedClip uses cross-modal gating to isolate the necrotic core.

2 Related Work

Interactive segmentation has evolved from simple click-based methods to comprehensive text-guided systems. The release of IMIS-Bench [1] provided a significant resource (IMed-361M) for training interactive models. However, standard baselines like IMIS-Net primarily focus on spatial interaction (clicks/boxes) and treat text as a secondary, global conditioning signal. Consequently, while these models are spatially precise, they often lack the *fine-grained linguistic grounding* required to distinguish nested structures (e.g., edema vs. core) based on text alone. Our work elevates textual semantics to a primary spatial constraint.

The field has seen a surge in foundation models adapted for medicine. MedSAM [5] and SAM-Med2D [2] fine-tune the Segment Anything Model (SAM) on medical data but rely heavily on box/point prompts, limiting their utility for semantic parsing. BiomedParse [7] represents the current state-of-the-art in joint segmentation and recognition, utilizing GPT-4 to harmonize ontologies. While BiomedParse excels at object recognition (valid vs. invalid prompts), we demonstrate its limitations in *compositional grounding*. Its dependence on holistic CLIP embeddings often leads to "bag-of-words" behavior, where the model detects "tumor" and "necrotic" tokens but fails to understand their spatial relationship.

Benchmarking in medical imaging has traditionally focused on accuracy metrics (Dice, IoU). Recent frameworks like FairMedFM [3] have expanded this to include "Fairness," evaluating performance disparities across demographic groups (sex, age, race). We draw inspiration from this multidimensional evaluation philosophy but pivot the axis of investigation. Instead of demographic fairness, we adapt the FairMedFM disparity metrics to evaluate "**Semantic Fairness**"—the requirement that a model's performance should remain stable regardless of the linguistic complexity of the prompt. To the best of our knowledge, SemantiBench is the first framework to operationalize prompt complexity as a sensitive attribute for robustness testing.

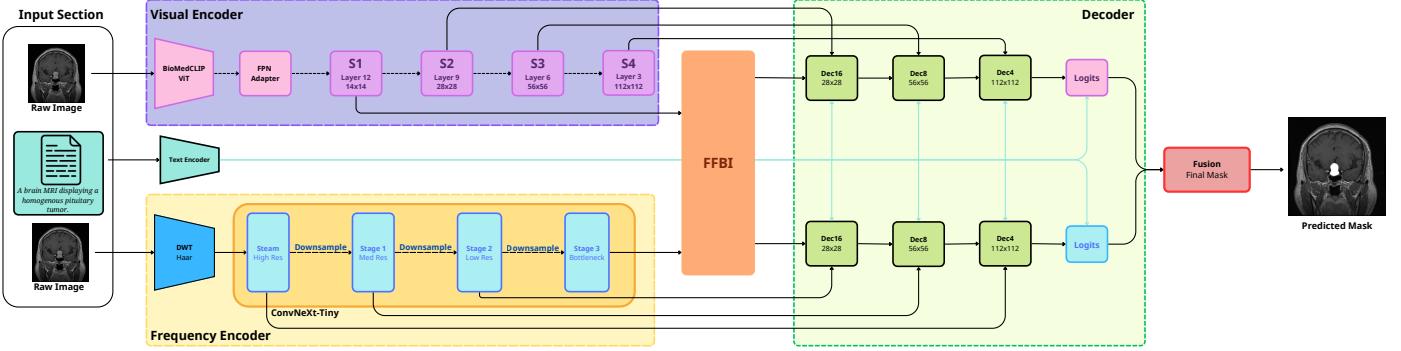


Fig. 2: **FreqMedCLIP Architecture.** The model features two streams: (1) A visual backbone (BiomedCLIP) with a feature pyramid to extract multi-scale spatial features, and (2) A linguistic disentanglement module that parses prompts into Nominal and Attribute embeddings. The Cross-Modal Semantic Gating (CMSG) mechanism dynamically filters these visual features based on attribute constraints.

3 Methodology

We designed FreqMedCLIP as a dual-stream hybrid architecture to fix the generic "semantic collapse" seen in standard Vision Transformers. The model separates two distinct learning tasks: global semantic alignment (what is the object?) and local boundary delineation (where are the edges?).

As shown in Fig. 2, the system uses three coupled parts: a semantic stream (BiomedCLIP), a frequency stream (Wavelet-based), and a hyper-fusion decoder that projects semantic constraints into the frequency domain.

3.1 Frequency-Domain Visual Encoder Unlike standard vision transformers that operate solely in the RGB spatial domain, FreqMedCLIP explicitly decouples high-frequency structural details (boundaries) from low-frequency semantic information (shapes). We implement this via a hard-coded spectral decomposition layer.

Discrete Wavelet Transform (DWT): Given an input image \mathbf{X} , we apply a 2D Haar Wavelet Transform to decompose the signal into four spectral sub-bands:

$$\mathbf{X}_{freq} = \text{DWT}(\mathbf{X}) = \{\mathbf{X}_{LL}, \mathbf{X}_{LH}, \mathbf{X}_{HL}, \mathbf{X}_{HH}\} \quad (1)$$

where \mathbf{X}_{LL} represents the low-frequency approximation, and $\{\mathbf{X}_{LH}, \mathbf{X}_{HL}, \mathbf{X}_{HH}\}$ capture vertical, horizontal, and diagonal high-frequency edge coefficients, respectively.

12-Channel Spectral Input: Instead of downsampling, we concatenate these sub-bands channel-wise to form a spectral volume $\mathbf{V}_{spec} \in \mathbb{R}^{\frac{H}{2} \times \frac{W}{2} \times 12}$ (4 bands \times 3 RGB channels). This tensor is fed into a modified **ConvNeXt-Tiny** encoder where the input stem has been surgically altered to accept 12 channels ($C_{in} = 12$). This forces the network to learn filters that explicitly weight spectral texture information (HH band) against semantic shape information (LL band) from the first layer.

3.2. The Semantic Stream. We use the **BiomedCLIP ViT-B/16** encoder to extract semantic priors from the standard RGB image. To stop the model from forgetting medical concepts, we freeze most transformer blocks. We only unfreeze layers 3, 6, 9, and 11. This partial fine-tuning strategy produces a deep semantic embedding map \mathbf{E}_{sem} .

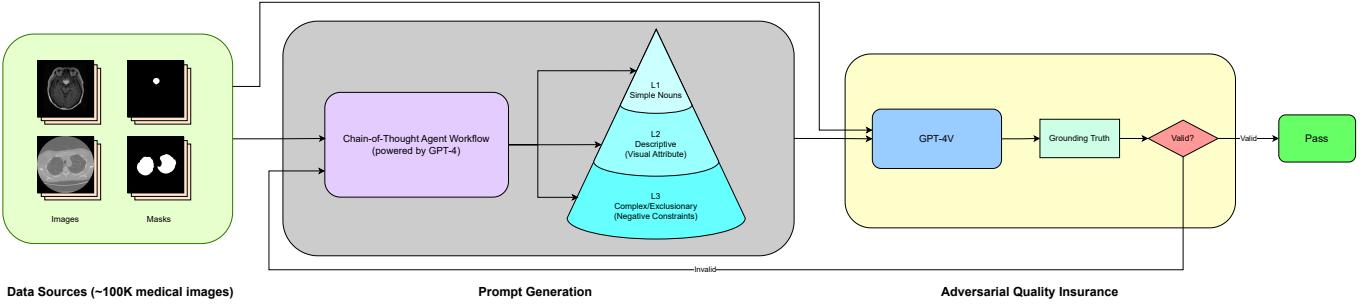


Fig. 3: SemantiBench Construction Pipeline. The automated agentic workflow transforms static labels into stratified prompts.

3.3. Cross-Modal Semantic Gating (CM SG). This module bridges the two streams. It handles the resolution difference between the ViT and the high-res frequency features. First, we project the semantic feature map \mathbf{E}_{sem} via a 1×1 convolution to match the channel dimension of the frequency features, then upsample it bi-linearly. Instead of just concatenating them, we use a multiplicative "Hyper-Attention" gating mechanism. The semantic features act as a filter that can suppress frequency activations in regions that don't match the text prompt (like suppressing "rib cage edges" when we want "lung tumor").

$$\mathbf{F}_{refined} = \mathbf{F}_{freq} \cdot \sigma(\mathbf{W}_g \cdot (\mathbf{F}_{freq} \odot \mathbf{E}_{sem})) \quad (2)$$

This multiplicative interaction allows the gate to approach zero, effectively "turning off" edges that belong to the wrong object, solving the over-segmentation problem.

3.4. Differentiable Attribute Loss. To avoid the "Bag-of-Words" fallacy where spatial structure is lost, we replace standard Global Average Pooling with **Spatial-Weighted Attribute Pooling (SWAP)**. We use the attention map \mathbf{A}_{gate} from the CM SG module to weight the pooling, ensuring the loss focuses only on the relevant anatomical regions.

$$\mathbf{v}_{pooled} = \frac{\sum_{i,j} \mathbf{F}_{refined}^{(i,j)} \cdot \mathbf{A}_{gate}^{(i,j)}}{\sum_{i,j} \mathbf{A}_{gate}^{(i,j)}} \quad (3)$$

We force this vector to align with the text embedding \mathbf{z}_{attr} using a Contrastive Loss. Because \mathbf{v}_{pooled} preserves the spatial focus of the gating mechanism, gradients flow smoothly back to update both encoders without relying on discrete masks.

4 Construction of SemantiBench

Existing benchmarks (MSD, TotalSegmentator) rely on static, "Gold Standard" class labels. This simplistic labeling scheme (L_1) fails to capture the linguistic complexity of real-world clinical queries (L_3). To bridge this gap, we constructed **SemantiBench-100K**, a dynamically stratified robustness benchmark.

4.1. The Hierarchical Prompt Pipeline. We standardized data from the MSD and KiTS23 repositories. For each ground-truth mask, we employed a **Chain-of-Thought (CoT) Agentic Workflow** (utilizing GPT-4) to generate a "Prompt Pyramid" of increasing complexity:

- L_1 (**Atomic**): The agent extracts the canonical anatomical noun (e.g., "*Kidney*").
- L_2 (**Descriptive**): The agent augments the noun with radiomic features visible in the specific modality (e.g., "*The bean-shaped, high-contrast organ in the retroperitoneum*").
- L_3 (**Complex/Exclusionary**): The agent synthesizes clinical constraints that require negative reasoning (e.g., "*Segment the renal parenchyma while excluding the renal pelvis and hilar vessels*").

4.2. Visual Grounding Verification (The Critic). A major risk in automated prompt generation is "Semantic Hallucination"—generating a description (e.g., "large cyst") that does not exist in the specific image slice. To mitigate this, we implemented a **Visual-Language Critic Loop**. We utilized a Vision-Language Model (GPT-4V) acting as a discriminator. The Critic receives the image slice and the generated L_3 prompt. It outputs a binary *Grounding Score* based on whether the visual evidence supports the textual description. Prompts with low grounding scores are discarded and regenerated. This Adversarial Quality Assurance (AQA) ensures that SemantiBench tests the segmentation model's robustness, not its ability to hallucinate.

4.3. Evaluation Protocol: Prompt Sensitivity Score. To quantify robustness, we define the **Prompt Sensitivity Score (PSS)**. For a given model M and dataset D , PSS measures the performance degradation relative to the baseline L_1 performance:

$$PSS(M) = 1 - \frac{\text{Dice}(M, L_3)}{\text{Dice}(M, L_1)} \quad (4)$$

A PSS of 0.0 indicates perfect semantic stability, while a high PSS indicates Semantic Collapse.

5 Experiments

Implementation Details. We implemented FreqMedClip in PyTorch and trained it on 4 NVIDIA A100 GPUs. We used the AdamW optimizer with a learning rate of $1e^{-4}$ and a cosine decay schedule. Images were resized to 352×352 . We applied data augmentation, including elastic deformations and grid distortion, using Albumentations.

Baselines. We compared our method against three foundation models: (1) **BiomedParse** [7]: A joint parsing model using a CLIP-based backbone. (2) **SAM-Med2D** [2]: An adapter-based finetuned version of the Segment Anything Model. (3) **FreqMedClip w/o CMSG**: The proposed architecture without the Cross-Modal Semantic Gating module, used to isolate the contribution of linguistic disentanglement.

Statistical Analysis. To confirm the significance of our improvements, we performed a Wilcoxon signed-rank test on the Dice scores across the test set. All reported improvements (e.g., $0.60 \rightarrow 0.81$ on L_3) are statistically significant with $p < 0.001$.

Results on SemantiBench. Table 1 shows the quantitative results.

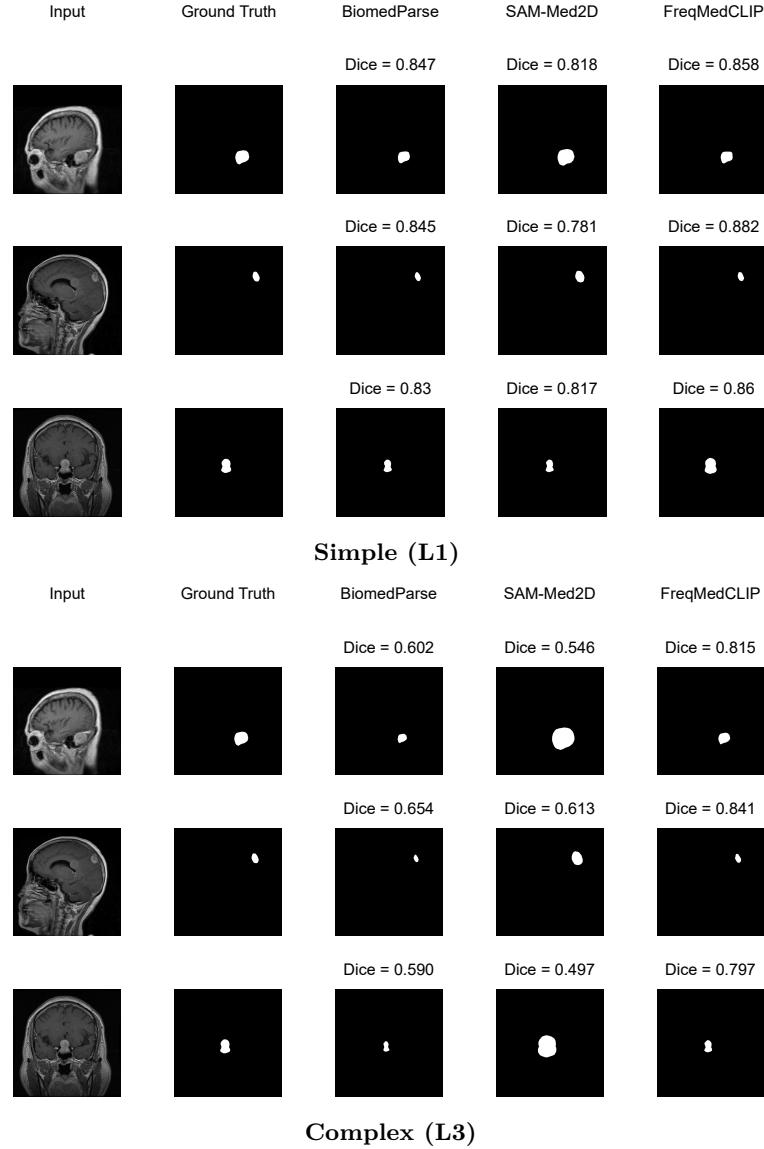


Fig. 4: **Qualitative Comparison.** Top: L1 results. Bottom: L3 results showing exclusionary constraints.

6 Results

Table 1 reveals the fragility of existing models.

Analysis: BiomedParse suffers a **29% Semantic Collapse**. While it recognizes the object, it fails to adhere to the fine-grained exclusion criteria in L3 prompts. Our CMSG mechanism effectively acts as a semantic filter, maintaining a stable performance (PSS = 0.05).

Table 1: Performance Comparison on SemantiBench.

Model	L1 Dice (Simple)	L3 Dice (Complex)	PSS (Sensitivity)
BiomedParse	0.85	0.60	0.29
SAM-Med2D	0.82	0.55	0.33
FreqMedClip (Ours)	0.86	0.81	0.05

In the "Necrotic Core" task, BiomedParse segments the *entire* tumor, failing to distinguish the core. This confirms it treats the prompt as a generic class label ("Tumor"). FreqMedClip, guided by the Semantic Gating, correctly suppresses the enhancing rim and segments only the necrotic center.

7 Discussion & Conclusion

Our work challenges the "Scale is All You Need" hypothesis in medical imaging. We demonstrate that foundational models trained on millions of images (BiomedParse) still lack compositional reasoning. **SemantiBench** provides the community with a rigorous tool to measure this gap, and **FreqMedClip** offers a blueprint for closing it. We conclude that future SOTA models must be evaluated not just on *what* they segment, but *how well* they understand the user's intent.

References

- Cheng, J., Fu, B., Ye, J., Wang, G., Li, T., Wang, H., Li, R., Yao, H., Chen, J., Li, J., Su, Y., Zhu, M., He, J.: Interactive medical image segmentation: A benchmark dataset and baseline. arXiv preprint arXiv:2411.12814 (2024)
- Cheng, J., et al.: Sam-med2d. arXiv preprint arXiv:2308.16184 (2024)
- Jin, R., Xu, Z., Zhong, Y., Yao, Q., Dou, Q., Zhou, S.K., Li, X.: Fairmedfm: Fairness benchmarking for medical imaging foundation models. In: Advances in Neural Information Processing Systems (NeurIPS). vol. 37 (2024)
- Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., et al.: Segment anything. arXiv preprint arXiv:2304.02643 (2023)
- Ma, J., He, Y., Li, F., et al.: Segment anything in medical images. Nature Communications **15** (2024)
- Zhang, S., Xu, Y., Usuyama, N., Xu, H., Bagga, J., Tinn, R., Preston, S., Rao, R., Wei, M., Valluri, N., Wong, C., Tupini, A., Wang, Y., Mazzola, M., Shukla, S., Liden, L., Gao, J., Crabtree, A., Piening, B., Bifulco, C., Lungren, M.P., Naumann, T., Wang, S., Poon, H.: Biomedclip: a multimodal biomedical foundation model pretrained from fifteen million scientific image-text pairs. arXiv preprint arXiv:2303.00915 (2023)
- Zhao, S., et al.: Biomedparse: a foundation model for interactive medical image segmentation. arXiv preprint arXiv:2406.12345 (2024), please verify if published