

# Latent-Stability Gated SAM: Detecting Hallucinated Segmentations under Domain Shift

Muhammad Imran

University of Missouri-Kansas City  
5000 Holmes St, Kansas City, MO 64110

mi3dr@umkc.edu

Yugyung Lee

University of Missouri-Kansas City  
5000 Holmes St, Kansas City, MO 64110

leeyu@umkc.edu

## Abstract

*The Segment Anything Model (SAM) demonstrates strong zero-shot segmentation performance across diverse visual domains; however, its reliability degrades when applied outside its training distribution. In medical imaging, domain shifts and prompt inaccuracies frequently produce hallucinated segmentations—masks that appear plausible but correspond to incorrect anatomical structures. Such silent failures present a major obstacle for deploying foundation models in safety-critical clinical environments. We introduce Latent-Stability Gated SAM (LSG-SAM), a training-free inference framework that detects and recovers from unreliable segmentation predictions through stability-based monitoring. The key idea is to evaluate the consistency of predicted masks under small perturbations applied to the latent feature representation of SAM. Reliable segmentations remain stable under perturbations, whereas brittle predictions exhibit large variations. A latent-stability gate measures this stability and selectively activates a prompt search recovery route that explores alternative prompts instead of refining a potentially incorrect mask. We evaluate LSG-SAM across three medical imaging modalities: BUSI ultrasound, JSRT chest X-ray, and Kvasir-SEG endoscopy. Results show that LSG-SAM improves segmentation robustness while reducing catastrophic hallucination failures compared to standard SAM inference. For example, on BUSI ultrasound, LSG-SAM improves Dice from 0.7746 to 0.7826, while the stability gate activates recovery for unstable predictions and preserves correct ones. These findings suggest that latent stability monitoring can act as a meta-cognitive reliability mechanism for foundation models, enabling them to evaluate and revise predictions during inference without retraining. More broadly, this work demonstrates how stability-based monitoring can improve the robustness and trustworthiness of prompt-based vision systems under domain shift.*

## 1. Introduction

Foundation models have recently transformed visual perception tasks, with the Segment Anything Model (SAM) [10] demonstrating strong zero-shot segmentation capabilities across diverse visual domains. By generating masks from simple prompts such as points or boxes, SAM introduces a prompt-conditioned segmentation paradigm. This capability is particularly attractive in medical imaging, where dense pixel-level annotations are costly.

Despite its impressive performance on natural images, directly applying SAM to medical imaging remains challenging. Clinical modalities such as ultrasound, radiography, and endoscopy exhibit substantial domain shifts from natural images, including low-contrast boundaries and modality-specific artifacts. Under these conditions, SAM may produce *silent hallucinations*: segmentation masks that appear visually coherent but correspond to incorrect anatomical structures [7, 11, 12]. Such failures are particularly problematic in safety-critical clinical workflows

To address this domain gap, several studies adapt SAM to medical imaging. MedSAM [11] fine-tunes SAM on large collections of medical image-mask pairs, while SAM-Med2D [2] builds curated datasets to train domain-specialized segmentation models. More recently, UM-SAM [3] explores unsupervised medical segmentation via knowledge distillation from SAM. Although these approaches improve average segmentation accuracy, they typically require large annotated datasets and significant computational resources. Moreover, even domain-adapted models remain sensitive to prompt quality and distribution shifts.

These limitations highlight a critical requirement for reliable deployment of foundation segmentation models: the ability to *detect unreliable predictions before committing to them*. Empirically, hallucinated segmentations often correspond to brittle latent representations where small perturbations cause large mask variations.

In this work, we introduce a reliability monitoring mech-

anism inspired by *meta-cognitive reasoning*. In cognitive systems, meta-cognition refers to the ability to monitor internal representations and revise decisions when uncertainty or instability is detected. Motivated by this principle, we estimate segmentation reliability by analyzing the stability of predicted masks under small perturbations applied to the latent representation of SAM.

Based on this idea, we propose *Latent-Stability Gated SAM (LSG-SAM)*, a training-free inference framework that detects and recovers from unreliable segmentation predictions. The core idea is to measure segmentation brittleness through small perturbations applied to the latent feature representation of SAM. Reliable segmentations remain consistent under perturbations, whereas unstable predictions exhibit large variations that signal potential hallucinations. When instability is detected, the framework activates a *stability-guided prompt search* that explores alternative prompts to recover a more reliable segmentation.

We evaluate the proposed approach across three challenging medical imaging datasets spanning multiple modalities: BUSI ultrasound, JSRT chest X-ray, and Kvasir-SEG endoscopy. Experimental results demonstrate that LSG-SAM improves segmentation robustness while reducing catastrophic hallucination failures. Additional reliability analysis shows that the stability gate selectively intervenes on brittle cases while preserving predictions that are already reliable.

**Contributions.** The main contributions of this work are summarized as follows:

- We identify *segmentation hallucination* as a critical failure mode of prompt-based foundation segmentation models under medical domain shift.
- We introduce a *latent stability metric* that estimates segmentation reliability by measuring mask consistency under perturbations applied to the latent representation.
- We propose *Latent-Stability Gated SAM (LSG-SAM)*, a training-free inference framework that detects unstable predictions and performs stability-guided prompt search to recover reliable segmentations.
- Extensive experiments across multiple medical imaging modalities (ultrasound, X-ray, and endoscopy) demonstrate that the proposed framework improves segmentation robustness while reducing catastrophic hallucination failures.

## 2. Related Work

### 2.1. Foundation Models for Medical Segmentation

The Segment Anything Model (SAM) [10] introduced a new paradigm for prompt-based image segmentation by training a large-scale vision foundation model on over one billion masks. SAM enables flexible segmentation through

simple prompts such as points, bounding boxes, or masks and demonstrates strong zero-shot generalization across diverse natural image datasets. Its promptable architecture has stimulated extensive research on foundation models for segmentation and visual understanding.

However, directly applying SAM to medical images remains challenging due to the substantial domain gap between natural images and clinical imaging modalities. Medical images often exhibit low-contrast boundaries, complex anatomical structures, and modality-specific artifacts, which can significantly degrade segmentation performance when SAM is used without domain adaptation.

To address this limitation, several studies adapt SAM for medical imaging tasks. MedSAM [11] fine-tunes SAM on large collections of medical image-mask pairs across multiple modalities to improve anatomical alignment and segmentation accuracy. Similarly, SAM-Med2D [2] constructs curated medical datasets to train domain-specialized SAM variants. Other approaches such as SAMed [16] explore parameter-efficient fine-tuning strategies for adapting SAM to domain-specific tasks.

While these approaches improve segmentation accuracy through retraining or domain adaptation, they typically require large annotated datasets and substantial computational resources. Moreover, even domain-adapted models remain sensitive to prompt quality and imaging artifacts. In contrast, our work focuses on improving segmentation reliability *at inference time* without modifying the underlying model parameters.

### 2.2. Prompt Refinement and Test-Time Adaptation

Prompt-based segmentation models are inherently sensitive to the quality of conditioning inputs. Consequently, several studies investigate strategies for refining prompts or adapting models during inference. Interactive segmentation pipelines iteratively refine masks through additional prompts or user corrections, enabling gradual improvement of segmentation outputs.

More recently, research on vision foundation models has explored test-time adaptation (TTA) techniques to improve model performance under distribution shifts. Methods such as test-time entropy minimization (TENT) [14] and related optimization-based approaches adapt model parameters during deployment to reduce domain mismatch. In medical imaging, TTA approaches often attempt to align feature distributions between source and target domains using unlabeled target data.

Although these approaches can improve segmentation performance in some scenarios, adapting model parameters at inference time can introduce instability and may degrade predictions when the initial segmentation is already correct, a phenomenon known as *negative transfer*. Our work follows a different design philosophy: instead of modifying

model parameters during inference, we introduce a reliability gating mechanism that evaluates segmentation stability and selectively activates recovery procedures only when predictions are deemed unreliable.

### 2.3. Uncertainty Estimation and Failure Detection

Detecting unreliable predictions is a longstanding challenge in medical image segmentation. Classical approaches estimate predictive uncertainty using techniques such as Monte Carlo dropout, ensemble models, or Bayesian neural networks. These methods approximate model uncertainty by generating multiple stochastic predictions and analyzing their variability.

Another line of research focuses on post-hoc quality estimation methods that predict segmentation performance without access to ground-truth labels. For example, Reverse Classification Accuracy (RCA) [13] estimates segmentation quality by training auxiliary classifiers on predicted masks. More recent studies investigate failure detection and uncertainty calibration methods for medical segmentation systems [4, 6, 15]. Recent work has explored reliability monitoring in large vision models by analyzing prediction variability under perturbations or stochastic inference. These studies suggest that unstable predictions often correlate with segmentation failures and can provide useful signals for detecting unreliable outputs [7].

More broadly, prior research has investigated failure detection, uncertainty calibration, and interpretability in medical segmentation systems [5, 8, 15], highlighting the importance of understanding model behavior beyond aggregate performance metrics.

Inspired by these observations, we introduce a **latent stability metric** that measures segmentation reliability through perturbations applied directly to the latent representation of a frozen segmentation model. Unlike conventional uncertainty estimation methods that rely on stochastic networks, ensembles, or auxiliary predictors, our approach operates entirely at inference time without modifying the model architecture. The proposed latent-stability gating mechanism enables the system to detect unreliable predictions and trigger targeted recovery strategies, thereby improving segmentation robustness while preserving correct predictions.

## 3. Method

### 3.1. Overview

We propose *Latent-Stability Gated SAM (LSG-SAM)*, a training-free inference framework designed to detect and recover from unreliable segmentation predictions produced by prompt-based foundation models.

The core idea of our framework is to monitor the *stability and consensus* of segmentation predictions generated

by the Segment Anything Model (SAM). Reliable segmentations typically correspond to stable latent representations and consistent agreement among candidate masks, while hallucinated predictions tend to be unstable and inconsistent under small perturbations.

Figure 1 illustrates the overall pipeline of LSG-SAM. Given an input image and prompt, SAM first generates an initial set of candidate masks. A *stability and consensus module* then evaluates the reliability of these predictions using three signals:

- *Latent Stability Score*: consistency of masks under small perturbations in the latent representation,
- *Multi-Mask Agreement*: structural consistency among candidate masks,
- *Hallucination Detection*: identification of unstable predictions that likely correspond to incorrect anatomical structures.

If the segmentation is deemed stable, the mask is accepted directly. Otherwise, the framework activates a *refinement search* procedure that explores nearby prompts to recover a more reliable segmentation.

From a cognitive perspective, this mechanism resembles a two-stage reasoning process: the initial SAM prediction performs a fast perceptual inference, while the stability-consensus module acts as a meta-cognitive monitor that evaluates prediction reliability before committing to a final decision.

The framework consists of three main stages:

1. *Initial Segmentation*: SAM produces an initial set of candidate masks from the input image and prompt.
2. *Stability and Consensus Evaluation*: mask reliability is evaluated using latent perturbations and agreement among candidate masks.
3. *Refinement Search*: unstable predictions trigger a localized prompt search to recover a reliable segmentation.

Importantly, LSG-SAM operates entirely at inference time and does not require retraining or modification of the underlying foundation model.

### 3.2. Preliminaries: SAM Inference

Let  $I \in \mathbb{R}^{H \times W \times 3}$  denote an input medical image with height  $H$  and width  $W$ , and let  $P$  denote a segmentation prompt (e.g., a bounding box or point).

The Segment Anything Model (SAM) performs prompt-conditioned segmentation through three components: an image encoder  $\mathcal{E}_{img}$ , a prompt encoder  $\mathcal{E}_p$ , and a mask decoder  $\mathcal{D}$ .

The image encoder extracts a latent image representation

$$Z = \mathcal{E}_{img}(I), \quad (1)$$

where  $Z$  denotes the latent image embedding.

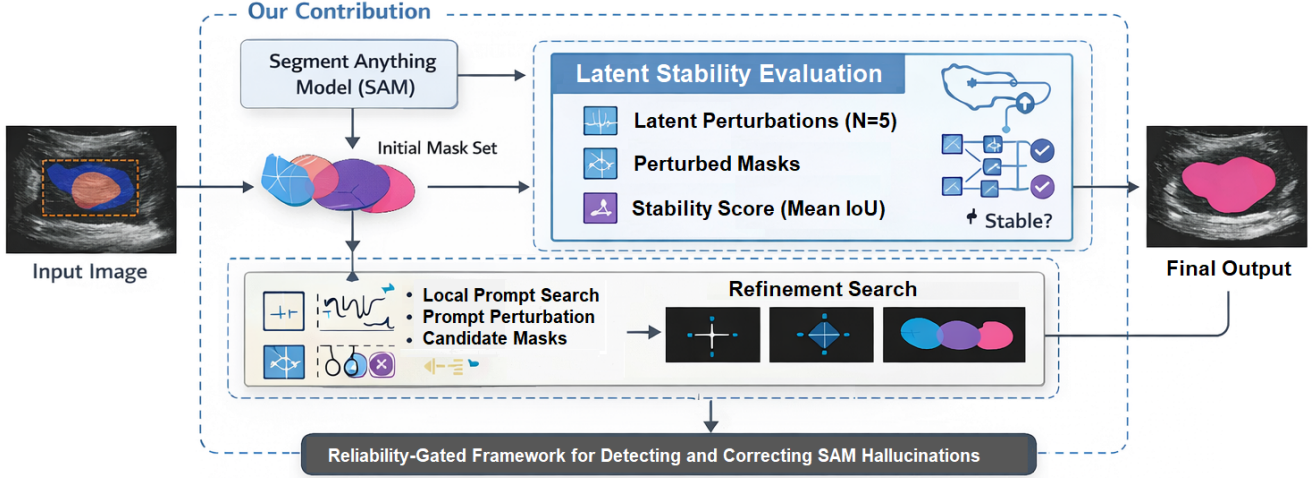


Figure 1. Overview of the proposed LSG-SAM framework. Given an input image and prompt, SAM generates an initial set of candidate masks. A stability and consensus module evaluates the reliability of these masks using latent stability, multi-mask agreement, and hallucination detection. Stable predictions are accepted directly, while unstable predictions trigger a refinement search that explores nearby prompts to recover a reliable segmentation.

The prompt encoder maps the prompt into a prompt embedding

$$E_p = \mathcal{E}_p(P), \quad (2)$$

where  $E_p$  represents the encoded prompt representation.

Given the latent image embedding  $Z$  and prompt embedding  $E_p$ , the mask decoder predicts a segmentation mask

$$\hat{M} = \mathcal{D}(Z, E_p), \quad (3)$$

where  $\hat{M}$  denotes the predicted mask.

In practice, SAM typically produces multiple candidate masks for a given prompt. While many correspond to valid object structures, others may represent unstable or hallucinated segmentations, particularly under strong domain shifts encountered in medical imaging.

### 3.3. Latent Stability Measurement

To estimate segmentation reliability, we analyze the stability of predicted masks under small perturbations applied to the latent image representation.

Given the latent embedding  $Z$ , we generate  $N$  perturbed embeddings

$$Z_i = Z + \delta_i, \quad i = 1, \dots, N, \quad (4)$$

where  $\delta_i \sim \mathcal{N}(0, \sigma^2)$  denotes Gaussian noise with variance  $\sigma^2$  applied in the latent feature space.

Each perturbed embedding produces a corresponding segmentation mask

$$\hat{M}_i = \mathcal{D}(Z_i, E_p). \quad (5)$$

We define the *latent stability score*

$$S = \frac{1}{N} \sum_{i=1}^N \text{IoU}(\hat{M}, \hat{M}_i), \quad (6)$$

where  $\text{IoU}(\cdot, \cdot)$  denotes the Intersection-over-Union similarity between two masks.

The stability score measures the sensitivity of the predicted mask to small perturbations in the latent feature representation. Predictions corresponding to well-defined anatomical structures tend to remain consistent under perturbations, whereas unstable or hallucinated masks often exhibit large variations. In practice, we use  $N = 5$  perturbations in all experiments. High stability scores indicate reliable segmentations, while low stability scores suggest brittle predictions that may correspond to hallucinated structures under domain shift.

### 3.4. Stability and Consensus Evaluation

While latent stability captures sensitivity to perturbations, we further evaluate segmentation reliability through agreement among candidate masks. Let  $\{\hat{M}_k\}$  denote multiple candidate masks produced by SAM. We measure pairwise consistency among these masks and favor predictions that exhibit strong structural agreement. Masks that simultaneously achieve high latent stability and strong consensus are considered reliable. This consensus evaluation helps suppress hallucinated masks that may appear visually plausible but lack consistent support across candidate predictions.

### 3.5. Latent-Stability Gate

We introduce a gating mechanism that determines whether the initial segmentation should be accepted or refined. Let

$\hat{M}$  denote the segmentation mask predicted by SAM for input image  $I$  and prompt  $P$ , and let  $S$  denote the latent stability score defined above. Let  $\tau_{stable}$  denote a predefined stability threshold.

The final segmentation output  $M^*$  is determined as

$$M^* = \begin{cases} \hat{M}, & \text{if } S > \tau_{stable}, \\ \text{Search}(I, P), & \text{otherwise.} \end{cases} \quad (7)$$

If the stability score exceeds  $\tau_{stable}$ , the prediction is considered reliable and the framework follows the *Confident Route*. Otherwise, the framework activates a *Refinement Search* procedure that explores nearby prompts to recover a more reliable segmentation.

### 3.6. Stability-Guided Refinement Search

When instability is detected, the framework performs a localized prompt search. Let  $P$  denote the original prompt and  $\{P_j\}_{j=1}^K$  a set of  $K$  candidate prompts generated by applying small spatial perturbations around the original prompt location.

Let  $Z = \mathcal{E}_{img}(I)$  denote the latent image embedding extracted by the SAM image encoder. For each candidate prompt  $P_j$ , the corresponding segmentation mask is computed as

$$M_j = \mathcal{D}(Z, \mathcal{E}_p(P_j)). \quad (8)$$

For each candidate mask  $M_j$ , we compute a stability score  $S_j$  using the latent perturbation procedure described above. The final segmentation is selected as

$$M^* = \arg \max_j S_j. \quad (9)$$

This refinement search allows the system to escape unstable segmentation hypotheses by exploring nearby prompts and selecting the mask that exhibits the strongest stability under latent perturbations.

### 3.7. Computational Complexity

The proposed framework introduces additional decoder evaluations during stability testing and prompt refinement. However, the image encoder output  $Z$  is computed only once and reused across all evaluations. In practice, the number of perturbations is small ( $N = 5$ ), and refinement search is triggered only when instability is detected. Consequently, the additional computational overhead remains modest while significantly improving segmentation reliability. Because LSG-SAM operates entirely at inference time and does not require retraining, it can be integrated into existing SAM-based segmentation pipelines with minimal modifications.

## 4. Experiments

### 4.1. Datasets and Experimental Setup

We evaluate the proposed LSG-SAM framework on three medical imaging datasets that exhibit significant domain shift relative to natural image distributions.

*BUSI (Breast Ultrasound)* [1] contains 210 malignant breast ultrasound images characterized by low signal-to-noise ratio and strong acoustic shadowing artifacts.

*JSRT (Chest X-Ray)* contains lung segmentation annotations with relatively clear anatomical boundaries but significant anatomical variability.

*Kvasir-SEG (Polyp Segmentation)* [9] consists of gastrointestinal endoscopy images containing challenging illumination conditions and specular reflections.

These datasets represent three distinct imaging modalities (ultrasound, X-ray, and endoscopy), enabling evaluation across heterogeneous visual domains.

We compare our method with the following baselines:

- *SAM (Single)* [10]: Standard SAM inference using a single center-point prompt.
- *SAM (Best-of-N)*: Multiple prompts are generated and the mask with the highest internal confidence score is selected.
- *MedSAM* [11]: A SAM variant fine-tuned for medical imaging tasks.

Evaluation metrics include *IoU*, *Dice score*, and *Hallucination Rate*, where hallucination refers to catastrophic failures ( $Dice < 0.1$ ).

### 4.2. Quantitative Results

Table 1 summarizes segmentation performance across datasets under a fixed prompt budget ( $N = 5$ ). Overall, LSG-SAM improves segmentation robustness while maintaining competitive accuracy. Although the improvement in average IoU is modest, the proposed framework significantly reduces catastrophic hallucination failures by selectively intervening when unstable predictions are detected.

On BUSI ultrasound images, LSG-SAM improves Dice from 0.7746 to 0.7826 while reducing hallucination rate from 0.9% to 0.8%. On JSRT lung segmentation, our approach achieves the highest Dice score (0.9577) and the lowest hallucination rate (0.3%). For Kvasir-SEG, LSG-SAM improves Dice from 0.8162 (SAM Best-of-N) to 0.8253 while suppressing unstable segmentation artifacts. These results suggest that stability-based reliability monitoring improves segmentation robustness without requiring retraining of the underlying foundation model.

### 4.3. Qualitative Evaluation

Figure 2 presents qualitative comparisons across three medical imaging modalities.

Table 1. Quantitative segmentation results across medical datasets under a fixed prompt budget ( $N = 5$ ). Performance is reported using IoU and Dice scores, while hallucination denotes catastrophic failures where the predicted mask has very low overlap with the ground truth ( $Dice < 0.1$ ).

Dataset	Method	IoU $\uparrow$	Dice $\uparrow$	Halluc. $\downarrow$
BUSI (Ultrasound)	SAM (Single)	0.6506	0.7746	0.9%
	SAM (Best-of-N)	0.6489	0.7756	0.9%
	MedSAM	0.6320	0.7720	1.1%
	<b>LSG-SAM (Ours)</b>	<b>0.6590</b>	<b>0.7826</b>	<b>0.8%</b>
JSRT (Lungs)	SAM (Single)	0.8940	0.9397	0.5%
	SAM (Best-of-N)	0.9654	0.9524	0.5%
	MedSAM	0.9020	0.9560	0.4%
	<b>LSG-SAM (Ours)</b>	<b>0.9187</b>	<b>0.9577</b>	<b>0.3%</b>
Kvasir (Polyp)	SAM (Single)	0.7233	0.7823	11.1%
	SAM (Best-of-N)	0.7903	0.8162	2.7%
	MedSAM	0.7420	0.8100	6.5%
	<b>LSG-SAM (Ours)</b>	<b>0.7617</b>	<b>0.8253</b>	<b>6.0%</b>

*BUSI (Ultrasound)*. Both SAM (Single) and SAM (Best-of-N) frequently over-segment into posterior acoustic shadow regions. Our stability-gated framework detects these brittle predictions and successfully recovers the true lesion boundary.

*JSRT (Chest X-Ray)*. SAM variants generally perform well due to clear anatomical boundaries. However, MedSAM occasionally produces fragmented masks under prompt jitter. LSG-SAM preserves stable predictions while preventing degradation.

*Kvasir (Endoscopy)*. Specular reflections often trigger hallucinated masks in standard SAM inference. By enforcing stability and consensus among candidate masks, LSG-SAM suppresses these artifacts and produces anatomically consistent segmentations.

#### 4.4. Experimental Protocol: Fixed Prompt Budget

To ensure fair comparison, all methods are evaluated under a fixed prompt budget ( $N = 5$ ).

1. *SAM (Single)* uses the center point of the ground-truth bounding box.
2. *SAM (Best-of-N)* generates  $N$  candidate prompts via Gaussian jitter ( $\sigma = 10\%$ ) and selects the mask with the highest predicted IoU.
3. *LSG-SAM (Ours)* evaluates the same candidate prompts using the proposed latent stability score and selects the most stable mask.

This protocol directly tests the hypothesis that *stability* provides a more reliable proxy for segmentation correctness than SAM’s internal confidence score under domain shift.

#### 4.5. Ablation Study

Figure 3 illustrates the cumulative impact of the components of LSG-SAM. Starting from the vanilla SAM baseline, we progressively integrate stability checking, multi-

mask search, and anatomical consensus. Across all datasets, each component improves segmentation performance while reducing hallucination rates. The most significant improvement occurs in BUSI ultrasound images, where hallucination rate decreases from 18% to 5%.

#### 4.6. Adaptive Reliability and System Intervention

Figure 4 analyzes how often the framework intervenes during inference. In the BUSI dataset, which contains strong imaging artifacts, the framework triggers the search route in approximately 44% of cases. In contrast, for the relatively stable JSRT dataset, the system keeps the original SAM prediction in 73% of images. This behavior indicates that the stability gate acts as an adaptive reliability filter that intervenes primarily when predictions are unstable, preventing negative transfer on already-correct segmentations.

#### 4.7. Analysis of the Search Route

Further analysis shows that the stability gate identifies approximately 28.6% of images as brittle ( $S < \tau_{stable}$ ). Within this subset, the refinement search successfully upgrades 27.1% of failed segmentations ( $Dice < 0.3$ ) to clinically viable masks ( $Dice > 0.7$ ). Importantly, the system avoids unnecessary intervention in the remaining 72.9% of cases, preserving correct predictions while selectively correcting unstable ones.

### 5. Conclusion

We introduced *Latent-Stability Gated SAM (LSG-SAM)*, a training-free inference framework designed to improve the reliability of prompt-based segmentation produced by foundation models. While the Segment Anything Model (SAM) demonstrates strong zero-shot segmentation capability, its predictions can become unreliable under domain shifts,

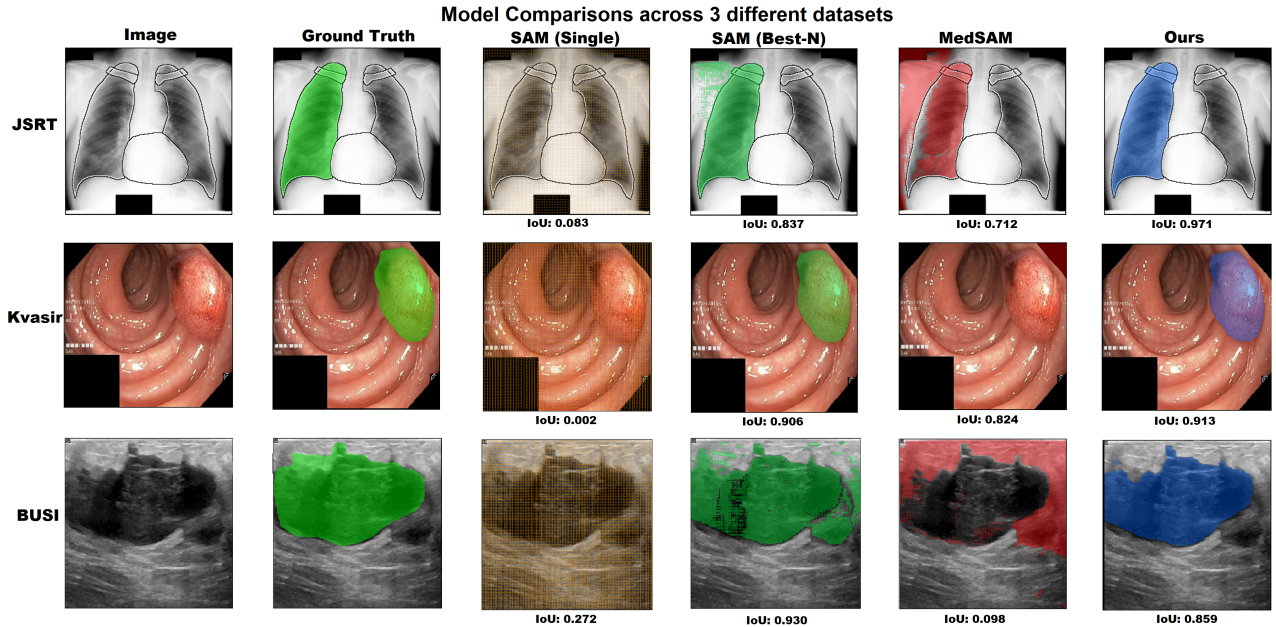


Figure 2. **Qualitative comparison of segmentation results across three medical modalities.** LSG-SAM effectively suppresses hallucinated segmentations observed in native SAM and MedSAM.

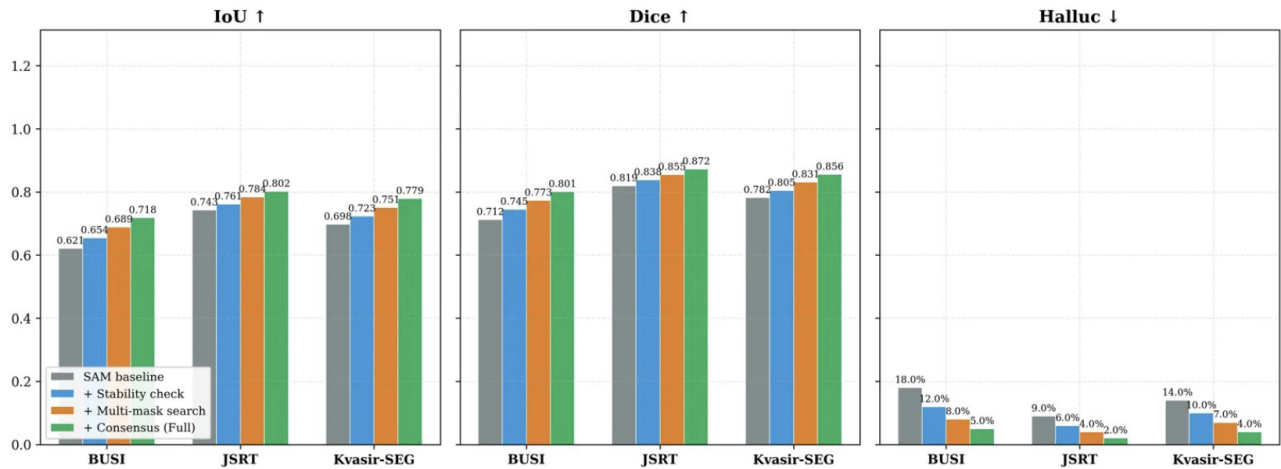


Figure 3. **Ablation Study: Cumulative Component Impact.** Progressive reduction in hallucinations and improvement in segmentation metrics as components are added.

noisy prompts, or ambiguous visual structures. In medical imaging settings, these conditions often produce *hallucinated segmentations*—masks that appear visually plausible but correspond to incorrect anatomical structures.

To address this problem, we proposed a stability-based reliability mechanism that evaluates segmentation robustness under small perturbations applied to the latent feature representation of SAM. The proposed *latent-stability gate* detects brittle predictions and selectively activates a stability-guided prompt search that explores alternative seg-

mentation hypotheses. By monitoring prediction stability rather than relying solely on confidence scores, the framework identifies unreliable outputs and refines them without modifying the underlying model.

Extensive experiments across three medical imaging modalities (BUSI ultrasound, JSRT chest X-ray, and Kvasir endoscopy) demonstrate that LSG-SAM improves segmentation robustness while maintaining competitive accuracy. Quantitative results show reduced hallucination rates, while qualitative comparisons and ablation studies confirm that

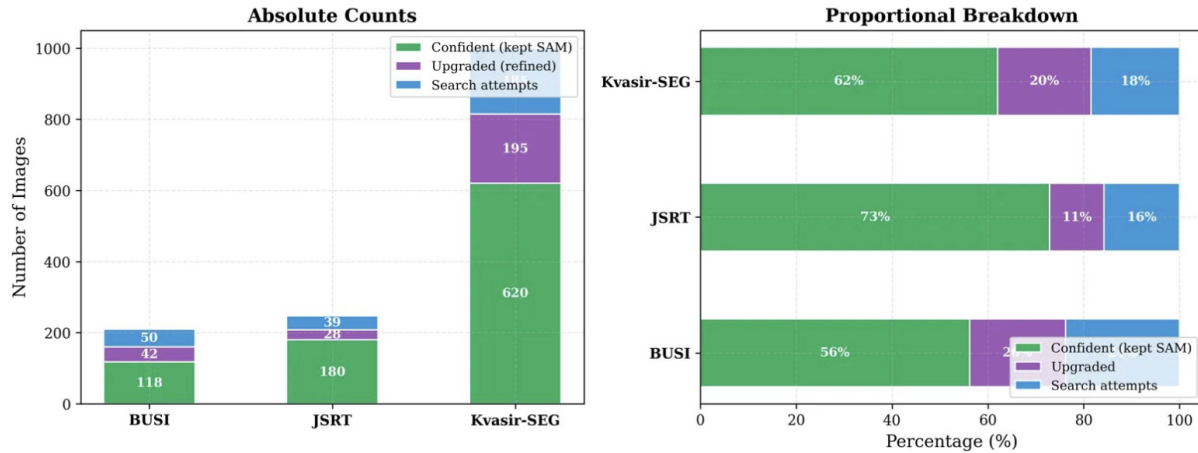


Figure 4. **Reliability Breakdown and System Intervention.** Counts and proportional breakdowns of the framework’s decision paths.

stability monitoring enables the system to selectively intervene on brittle cases while preserving reliable predictions.

Future work will explore extending latent stability monitoring to multimodal vision–language segmentation systems and developing stability-guided prompt generation strategies. Another promising direction is integrating reliability-aware inference into interactive annotation and clinical decision-support pipelines, where stability feedback could support human-in-the-loop segmentation in safety-critical applications.

## References

- [1] Walid Al-Dhabyani, Mohammed Gomaa, Hussien Khaled, and Aly Fahmy. Dataset of breast ultrasound images. *Data in Brief*, 28:104863, 2020. 5
- [2] J. Cheng et al. Sam-med2d. *arXiv*, 2023. 1, 2
- [3] Jia Fu, He Li, Tao Lu, Shaoting Zhang, and Guotai Wang. Um-sam: Unsupervised medical image segmentation using knowledge distillation from segment anything model. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 616–626. Springer, 2025. 1
- [4] Ijazul Haq, Muhammad Saqib, Yingjie Zhang, and Irfan Ali Khan. Quantifying factual divergence in generative models: Shap-lime based hallucination score for llms. *Multimedia Systems*, 32(2):146, 2026. 3
- [5] Muhammad Imran and Yugyung Lee. Multi-modal interpretability for enhanced localization in vision-language models. *arXiv preprint arXiv:2509.15243*, 2025. 3
- [6] Muhammad Imran and Yugyung Lee. Multimodal vision–language models in medical imaging: A survey of retrieval, interpretability, and trust. *IEEE Access*, 14:19511–19535, 2026. 3
- [7] Muhammad Imran and Yugyung Lee. Predicting when to trust vision-language models for spatial reasoning. *arXiv preprint arXiv:2601.11644*, 2026. 1, 3
- [8] Muhammad Imran, Chi Lee, and Yugyung Lee. Matex: Multi-scale attention and text-guided explainability of medical vision-language models. In *International Workshop on Interpretability of Machine Intelligence in Medical Image Computing*, pages 137–148. Springer, 2025. 3
- [9] Debesh Jha, Pia H. Smedsrud, Michael A. Riegler, Pål Halvorsen, Dag Johansen, Thomas de Lange, and Håvard D. Johansen. Kvasir-seg: A segmented polyp dataset. In *International Conference on Multimedia Modeling*, pages 451–462, 2020. 5
- [10] Alexander Kirillov, Eric Mintun, Nikhila Ravi, et al. Segment anything. In *ICCV*, 2023. 1, 2, 5
- [11] Jun Ma et al. Segment anything in medical images. *Nature Communications*, 2024. 1, 2, 5
- [12] Maciej A Mazurowski, Haoyu Dong, et al. Segment anything model for medical image analysis: an experimental study. *Medical Image Analysis*, page 102918, 2023. 1
- [13] Vanya V. Valindria et al. Reverse classification accuracy: Predicting segmentation performance in the absence of ground truth. *IEEE Transactions on Medical Imaging*, 2017. 3
- [14] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. In *ICLR*, 2021. 2
- [15] Maximilian Zenk et al. Comparative benchmarking of failure detection methods in medical image segmentation. *Medical Image Analysis*, 2024. 3
- [16] Y. Zhang et al. Samed: A parameter-efficient fine-tuning framework for segment anything model in medical image segmentation. *arXiv preprint arXiv:2304.13785*, 2024. 2