

LEARNING FROM ADVERSITY: SEMANTIC-AWARE MASK REFINEMENT THROUGH ADVERSARIAL PERTURBATION

Anonymous authors

Paper under double-blind review

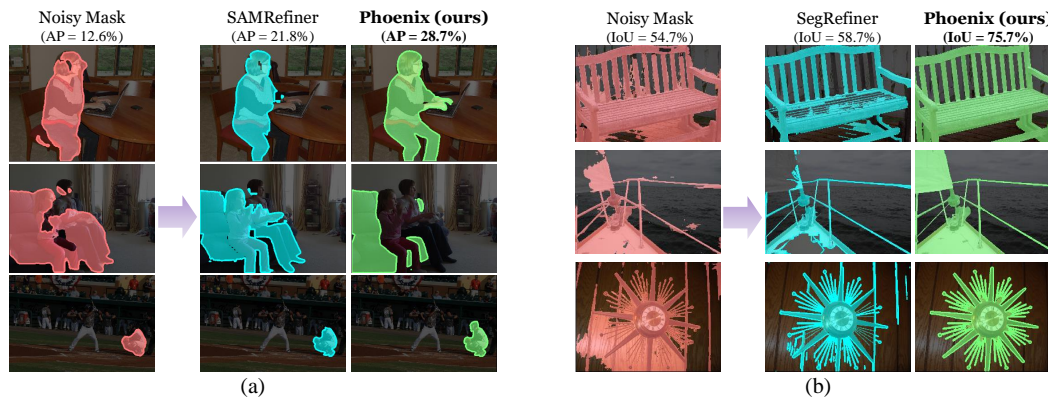


Figure 1: **Qualitative Comparison** of our mask refinement method with existing approaches. The (a) instance segmentation and (b) fine-grained segmentation examples show our superior mask refinement performance on complex structures and boundary details.

ABSTRACT

Despite significant advances in image segmentation, even state-of-the-art models produce masks with imperfect boundaries, semantic inconsistencies, and structural errors. Mask refinement addresses these limitations, yet current approaches rely on simplistic synthetic noise that fails to capture the complex error patterns of real segmentation models. We introduce Phoenix, a novel framework that leverages adversarial learning to generate semantically meaningful noise patterns and contrastive learning to model refinement relationships. Our approach consists of two key innovations: (1) Adversarial Mask Perturbation, which employs embedding attacks to create semantic-aware noise that mimics real segmentation errors, and (2) Contrastive Mask Refinement Learning, which establishes a tri-directional framework that ensures feature consistency within semantic regions while maintaining separation between classes. Experiments demonstrate that Phoenix significantly outperforms existing methods across diverse tasks, while consistently enhancing state-of-the-art segmentation models with substantial improvements.

1 INTRODUCTION

Image segmentation provides pixel-level understanding crucial for numerous applications. Despite remarkable progress in segmentation architectures, even state-of-the-art models [Cheng et al. \(2022\)](#); [Kirillov et al. \(2023\)](#); [Xie et al. \(2021a\)](#) exhibit persistent limitations: imprecise boundaries at complex contours, semantic confusion between similar objects, and structural inconsistencies that violate object integrity. These errors stem from fundamental challenges that remain despite extensive model scaling and data collection efforts.

This circumstance establishes mask refinement as a distinct and versatile complementary approach to conventional segmentation methods, directly addressing their systematic limitations. It enhances model performance without architectural changes, making it valuable in cases where model updates are infeasible, such as due to proprietary restrictions, computational limitations, or data

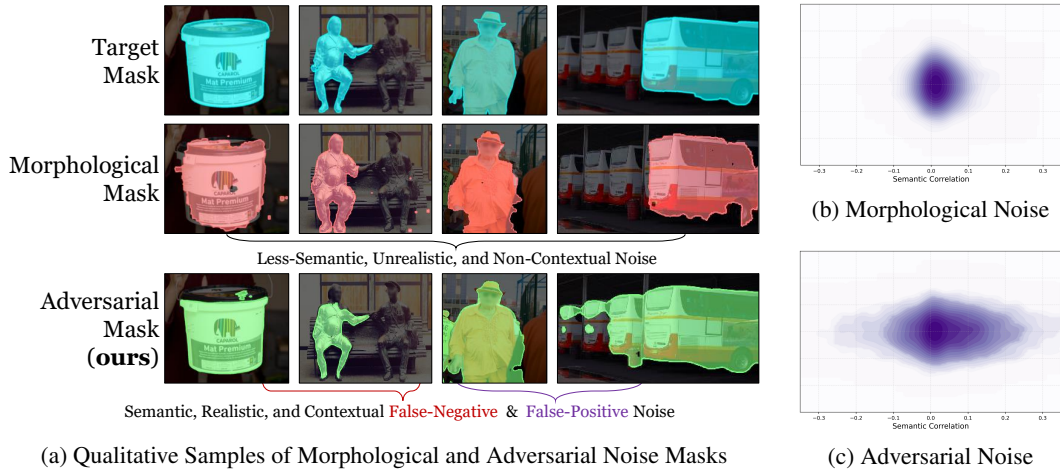


Figure 2: **Qualitative Comparison** of noise patterns (a) between morphological and our adversarial noise masks. (b) Semantic correlation distribution of morphological noise showing a narrow distribution, and (c) our adversarial noise showing a broader distribution and more diverse error types.

collection costs. Furthermore, mask refinement is pivotal in label-efficient learning, including semi-supervised Kim et al. (2023); Wang et al. (2022) and weakly-supervised settings Kim et al. (2022); Tian et al. (2021), transforming low-quality pseudo-labels into reliable supervision, which maximizes learning outcomes even with limited annotations.

The construction of realistic noisy and clean mask pairs is central to effective refinement learning. Recent efforts Lin et al. (2025); Tang et al. (2021) have advanced mask refinement through various approaches, yet significant limitations persist. Methods like SegFix Yuan et al. (2020) and SegRefiner Wang et al. (2023) rely on synthetic noise generated through morphological operations, producing simplistic, spatially random perturbations that often fail to capture the structured, context-dependent errors of real neural networks. This fundamental limitation restricts their ability to address the complex challenges encountered in real-world segmentation tasks, as shown in Figure 1.

"What doesn't kill you makes you stronger." This wisdom reflects how adversity can become a catalyst for growth, a principle we translate from philosophy to algorithm design. In natural systems, adaptation occurs in response to meaningful challenges. Similarly, the effective learning system emerges from facing realistic, meaningful obstacles rather than artificial ones. Drawing inspiration from this, we introduce a novel framework for semantic-aware mask refinement through adversarial perturbation. We dub our framework *Phoenix* since it transforms challenging noise patterns into superior refinement capabilities, mirroring how the mythical bird rises stronger from the ashes.

The Phoenix framework builds upon two key innovations. First, Adversarial Mask Perturbation (AMP) employs adversarial embedding attacks to generate semantically meaningful, contextually aware noise patterns. By optimizing against the model's learned representations, AMP creates perturbations that concentrate precisely where real segmentation models struggle most, such as at semantically challenging boundaries and ambiguous regions, as shown in Figure 2. Notably, this approach offers controllability over noise generation patterns while maintaining computational efficiency. Second, our Contrastive Mask Refinement Learning (CMRL) explicitly models the relationships between ground truth, noisy input, and current prediction. Unlike conventional pixel-wise objective functions, this tri-directional contrastive mechanism promotes clear separation between foreground and background features while ensuring consistency within semantic regions, which is tailored to the mask refinement task.

Our experiments demonstrate that Phoenix significantly outperforms existing refinement methods Lin et al. (2025); Tang et al. (2021); Wang et al. (2023); Yuan et al. (2020) across diverse tasks. When refining masks from semi-supervised and weakly-supervised models, Phoenix achieves absolute gains of up to +16.1% in AP^{mask} . Applied to state-of-the-art segmentation models, it consistently improves performance across diverse architectures. Furthermore, we demonstrate Phoenix's effectiveness in fine-grained mask refinement tasks and its potential for self-supervised refinement without ground-truth annotations.

2 RELATED WORK

Image Segmentation Image segmentation has evolved significantly with the advent of deep learning. Early approaches like FCN Long et al. (2015) and U-Net Ronneberger et al. (2015) established fully-convolutional architectures that became foundational for subsequent research. Performance improvements followed through innovations in feature extraction (*e.g.*, DeepLab Chen et al. (2018), PSPNet Zhao et al. (2017)) and, more recently, with transformer-based architectures (*e.g.*, SegFormer Xie et al. (2021a), Mask2Former Cheng et al. (2022)) that leverage global context modeling. The recent Segment Anything Model (SAM) Kirillov et al. (2023) represents a significant advance in general-purpose segmentation through its prompt-based inference and zero-shot capabilities.

Despite these advances, generating high-quality segmentation masks remains challenging, particularly at object boundaries and in complex scenes. This challenge is magnified in the semi-supervised Wang et al. (2022) and weakly semi-supervised Kim et al. (2023) settings, where initial segmentation predictions are often coarse and contain substantial noise.

Mask Refinement Mask refinement addresses the limitations of primary segmentation models by enhancing mask quality through post-processing. Early approaches employed traditional techniques like conditional random fields (CRF) Krähenbühl & Koltun (2011) and graph cuts Boykov & Jolly (2001) to improve boundary adherence. These methods often rely on handcrafted features and struggle with semantically complex scenes. Learning-based refinement has gained traction with various approaches. SegFix Yuan et al. (2020) employs residual correction for boundaries, BPR Tang et al. (2021) introduces a boundary patch refinement that focuses on local regions around object boundaries. SegRefiner Wang et al. (2023) adopts a diffusion process for mask refinement modeling with a two-stage approach combining coarse prediction and fine-grained refinement. Recently, SAMRefiner Lin et al. (2025) leverages SAM’s zero-shot capabilities through visual prompt engineering for mask refinement. Despite its effectiveness with a training-free setup, it relies solely on fixed pre-trained representations that weren’t optimized for the refinement task. In contrast, Phoenix utilizes SAM’s architecture with efficient fine-tuning techniques specifically designed to address key refinement challenges, such as correction pattern learning and boundary precision.

The primary challenge in effective mask refinement modeling lies in generating representative training data that captures realistic error patterns. Existing approaches Wang et al. (2023); Yuan et al. (2020) predominantly rely on morphological perturbations of ground-truth masks to simulate noise, such as random boundary perturbations with dilation and erosion operations and region modifications. These methods produce synthetic noise patterns that fail to capture the semantic nature of errors in real segmentation models, which are highly structured and context-dependent.

Adversarial Learning Adversarial learning has primarily focused on attack and defense mechanisms for model robustness Goodfellow et al. (2015); Madry et al. (2018). Adversarial attacks on segmentation models Arnab et al. (2018) and black-box perturbation methods Huang & Zhang (2020) demonstrate how adversarial techniques can expose model vulnerabilities through carefully crafted perturbations. However, existing work treats adversarial perturbations as destructive tools for testing rather than constructive mechanisms for training data generation. We introduce a paradigm shift by repurposing adversarial attack techniques as data augmentation for mask refinement, operating in the embedding space to generate semantically meaningful noise patterns that mimic realistic segmentation errors instead of merely maximizing prediction failures.

Contrastive Learning Contrastive learning has achieved remarkable success in representation learning Chen et al. (2020); He et al. (2020) with recent extensions to dense prediction tasks. Wang & Isola (2020) demonstrate that contrastive objectives in feature space can improve dense prediction quality through alignment and uniformity principles. Pixel-wise contrastive methods Wang et al. (2021); Xie et al. (2021b) have been developed for semantic segmentation pre-training and unsupervised representation learning. However, existing approaches primarily focus on learning general visual representations or establishing class boundaries. In contrast, our Contrastive Mask Refinement Learning introduces a novel tri-directional framework specifically designed for modeling the refinement relationship between noisy inputs, predictions, and ground truth masks, explicitly capturing the transformation from incorrect to correct predictions through self-improvement regularization.

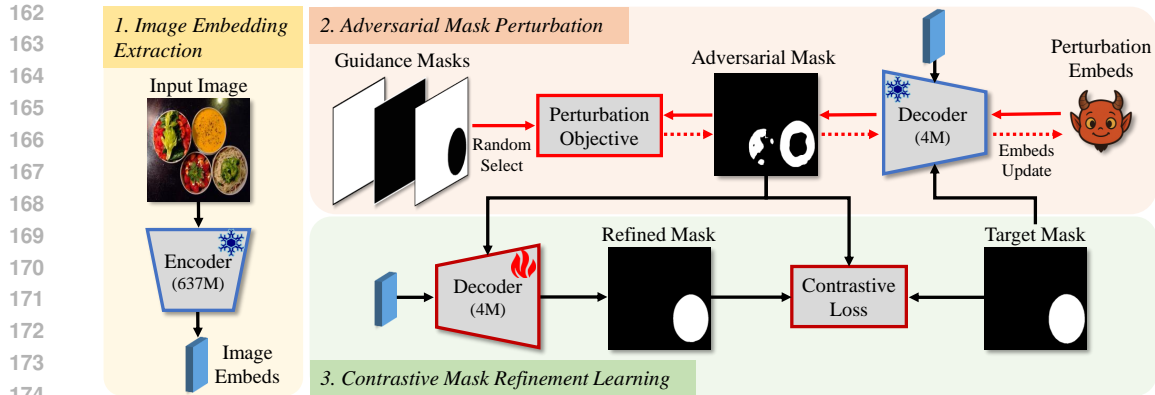


Figure 3: **Overview** of our Phoenix framework. The pipeline consists of three main components: (1) Image Embedding Extraction using SAM’s encoder, (2) Adversarial Mask Perturbation that generates realistic noise patterns through adversarial embedding attacks, and (3) Contrastive Mask Refinement Learning that uses the tri-directional relationships between masks to improve refinement quality.

3 METHODOLOGY

3.1 PROBLEM FORMULATION

Mask refinement aims to transform a noisy or coarse segmentation mask into a high-quality mask that accurately delineates object boundaries and semantic regions. Formally, given an input image $\mathcal{I} \in \mathbb{R}^{H \times W \times 3}$ and a corresponding noisy mask $\mathcal{M}_n \in \{0, 1\}^{H \times W}$, the mask refiner f generates a refined mask $\mathcal{M}_r \in \{0, 1\}^{H \times W}$ such that: $\mathcal{M}_r = f(\mathcal{I}, \mathcal{M}_n; \theta)$, where θ represents the parameters of the refiner model, and H and W denote the height and width of the image.

3.2 FRAMEWORK OVERVIEW

Figure 3 illustrates our Phoenix framework consisting of two key innovations: (1) Adversarial Mask Perturbation (AMP) and (2) Contrastive Mask Refinement Learning (CMRL). Phoenix builds upon pre-trained SAM Kirillov et al. (2023) with fine-tuning only the lightweight decoder. The encoder f_{enc} takes an input image \mathcal{I} and extracts image embeddings \mathbf{E}_{img} . The decoder f_{dec} then processes these embeddings along with the noisy mask \mathcal{M}_n to produce the refined mask \mathcal{M}_r . In addition, we extract point and box prompts from the noisy mask and incorporate them into the decoder’s input as visual prompt embeddings \mathbf{E}_v , following the prompt sampling strategy used in Lin et al. (2025).

3.3 ADVERSARIAL MASK PERTURBATION (AMP)

Limitations of Morphological Noise Approaches. Existing mask refinement methods Wang et al. (2023); Yuan et al. (2020) predominantly rely on morphological perturbations (erosion, dilation, boundary modifications) to create synthetic noise from ground-truth masks. These approaches have several limitations, as evident in Figure 2a: (1) They produce unrealistic noise patterns that fail to capture the semantic nature of errors in real segmentation models, as morphological operations create structurally simplistic and spatially random noise, (2) They generate noise with limited diversity, unable to represent the wide range of failure modes in modern segmentation models, and (3) They are contextually blind, with perturbations operating independently of image content, making them unable to simulate errors from semantic confusion between similar objects.

Proposed Adversarial Approach. We introduce Adversarial Mask Perturbation (AMP), which leverages adversarial embedding attacks to generate semantically meaningful, contextually aware noise patterns. Given a target (ground-truth) mask \mathcal{M}_t , we inject learnable perturbation embeddings $\mathbf{E}_p \in \mathbb{R}^{P \times C}$ into the decoder f_{dec} alongside visual prompt embeddings \mathbf{E}_v derived from \mathcal{M}_t , where P is the number of embeddings and C is the embedding dimension. Importantly, the perturbation embeddings \mathbf{E}_p are used exclusively for generating noisy masks during data augmentation without affecting the pretrained decoder parameters. The decoder f_{dec} remains frozen during perturbation, and the \mathbf{E}_p are not involved in actual training or inference of the decoder.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

Algorithm 1 Adversarial Mask Perturbation

Require: Image embeds \mathbf{E}_{img} , perturbation embeds \mathbf{E}_p , visual prompt embeds \mathbf{E}_v , target mask \mathcal{M}_t , guidance mask \mathcal{M}_g , adversarial criterion \mathcal{D} , initial step size α_0 , IoU threshold τ , margin ϵ , maximum inner iterations N .

Ensure: Noisy mask \mathcal{M}_n with IoU $[\tau, \tau + \epsilon]$

```

1:  $\alpha = \alpha_0$ 
2: repeat
3:   for  $i = 1$  to  $N$  do
4:      $\mathcal{M}_n = f_{dec}(\mathbf{E}_{img}, [\mathbf{E}_p; \mathbf{E}_v])$ 
5:      $iou = \text{compute\_IoU}(\mathcal{M}_n, \mathcal{M}_t)$ 
6:     if  $iou < \tau + \epsilon$  then
7:       break
8:     end if
9:     Save current state:  $\mathbf{E}_p^{prev} = \mathbf{E}_p$ 
10:     $\mathcal{L}_{adv} = -\mathcal{D}(\mathcal{M}_n, \mathcal{M}_g)$ 
11:    Compute gradient  $\nabla_{\mathbf{E}_p} \mathcal{L}_{adv}$ 
12:     $\mathbf{E}_p = \mathbf{E}_p + \alpha \cdot \text{sign}(\nabla_{\mathbf{E}_p} \mathcal{L}_{adv})$ 
13:  end for
14:  if  $iou \geq \tau$  then
15:    return  $\mathcal{M}_n$  {Target IoU achieved}
16:  end if
17:   $\alpha = \alpha/10$  {Decay step size}
18:  Restore previous state:  $\mathbf{E}_p = \mathbf{E}_p^{prev}$ 
19: until maximum decay steps reached
20: return  $\mathcal{M}_n$ 

```

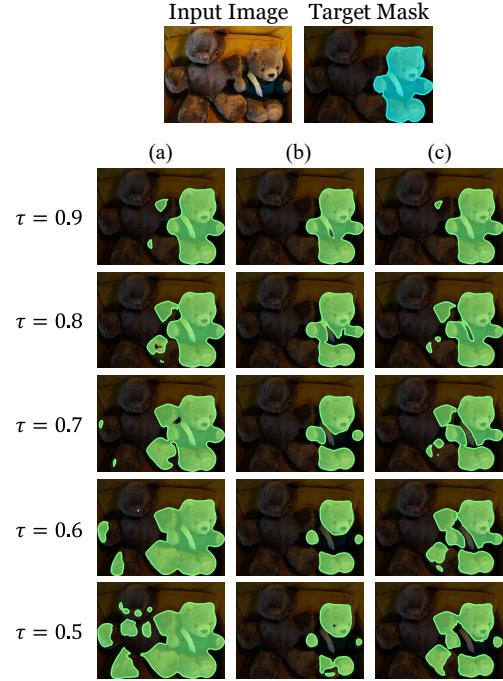


Figure 4: **Qualitative Samples** of generated noisy masks according to the IoU threshold τ and guidance mask (a) expansion guide, (b) contraction guide, and (c) inversion guide.

Unlike conventional adversarial attacks [Arnab et al. \(2018\)](#); [Goodfellow et al. \(2015\)](#); [Madry et al. \(2018\)](#) that aim to maximize classification error by perturbing input images, our approach repurposes adversarial techniques as constructive tools for noise generation. Specifically, we adapt the FGSM [Goodfellow et al. \(2015\)](#) to operate in the embedding space rather than image pixel space: $\mathbf{E}_p \leftarrow \mathbf{E}_p + \alpha \cdot \text{sign}(\nabla_{\mathbf{E}_p} \mathcal{L}_{adv})$, where α controls the perturbation magnitude and \mathcal{L}_{adv} is an adversarial objective. By operating in the embedding space rather than input image space, we achieve both higher computational efficiency and high-level semantic perturbation [Huang & Zhang \(2020\)](#).

Controllable Noise Generation. The *Guidance Mask* (\mathcal{M}_g) determines the semantic direction of perturbation by modifying the adversarial objective: $\mathcal{L}_{adv} = -\mathcal{D}(f_{dec}(\mathbf{E}_{img}, [\mathbf{E}_p; \mathbf{E}_v]), \mathcal{M}_g)$, where \mathcal{D} represents an adversarial criterion (e.g., Dice or MSE loss). Three primary configurations yield distinct noise patterns, as shown in Figure 4: (1) *Expansion Guide* ($\mathcal{M}_g = \mathbf{1}$, all ones) generates false-positive errors by pushing boundaries outward, (2) *Contraction Guide* ($\mathcal{M}_g = \mathbf{0}$, all zeros) creates false-negative errors that mimic under-segmentation behaviors, (3) *Inversion Guide* ($\mathcal{M}_g = \mathbf{1} - \mathcal{M}_t$) produces a balanced distribution of both error types.

Furthermore, to automatically calibrate the noise magnitude, Algorithm 1 implements an adaptive threshold-guided noise generation approach that adjusts the perturbation strength based on the IoU threshold parameter τ . Given initial step size α_0 , IoU threshold τ , margin ϵ , and maximum inner iterations N , it iteratively updates \mathbf{E}_p according to the FGSM rule, the process continues until the IoU between the generated noisy mask and the target mask falls within the range $[\tau, \tau + \epsilon]$. If this condition is not met until N iterations, it reduces the step size $\alpha \leftarrow \alpha/10$ and restarts the process. As shown in Figure 4, the IoU threshold τ enables control over noise intensity: lower thresholds generate more aggressive noise patterns, while higher thresholds produce subtle noise.

Theoretical Analysis of Semantic Distribution. For a pretrained decoder f_{dec} with fixed parameters θ_{dec} , the embedding gradient magnitude $\|\nabla_{\mathbf{E}_p} f_{dec}(\mathbf{E}_{img}, [\mathbf{E}_p; \mathbf{E}_v]; \theta_{dec})\|$ varies across spatial locations based on the model’s uncertainty. Formally, this gradient magnitude relates to the local classification uncertainty: $\|\nabla_{\mathbf{E}_p} f_{dec}(\mathbf{E}_{img}, [\mathbf{E}_p; \mathbf{E}_v]; \theta_{dec})\|_2 \propto -\log p(y|\mathbf{E}_{img}, [\mathbf{E}_p; \mathbf{E}_v]; \theta_{dec})$, where $p(y|\mathbf{E}_{img}, [\mathbf{E}_p; \mathbf{E}_v]; \theta_{dec})$ represents the decoder’s confidence in its prediction y at each spatial location [Kendall & Gal \(2017\)](#). Consequently, our FGSM-based update intrinsically amplifies

270 perturbations in regions of high uncertainty, precisely where segmentation models typically struggle,
 271 producing semantically meaningful noise patterns.

272 We quantify the semantic nature of our generated noise through semantic correlation analysis between
 273 the perturbation magnitude and image feature gradients using the LVIS dataset Gupta et al. (2019).
 274 This analysis computes Pearson correlation Pearson (1895) between noise spatial distribution and
 275 semantic features extracted from edge detection and texture maps to measure how perturbations relate
 276 to image semantics (detailed in Appendix B.5). As shown in Figure 2c, adversarial noise exhibits
 277 a broad distribution of semantic correlations spanning $[-0.6, 0.8]$, indicating its ability to capture
 278 both semantically important regions (positive values) and homogeneous areas (negative values). In
 279 contrast, morphological noise (Figure 2b) shows a narrow distribution concentrated around zero,
 280 confirming its semantic-agnostic nature. This quantitative analysis validates that AMP produces noise
 281 patterns that align with the complex error distributions.

282 **Computational Efficiency.** Our implementation maintains high computational efficiency by reusing
 283 image embeddings from the encoder’s single forward pass (637M parameters for ViT-H). Since AMP
 284 operates only on the lightweight decoder (4M parameters, 1.01 GFlops), each perturbation update
 285 requires only 6 ms on a V100 GPU, enabling efficient noise generation.

287 3.4 CONTRASTIVE MASK REFINEMENT LEARNING (CMRL)

288 **Motivation.** Traditional refinement approaches Tang et al. (2021); Wang et al. (2023); Yuan et al.
 289 (2020) primarily rely on pixel-wise classification, which treats each pixel independently. The funda-
 290 mental challenge in mask refinement lies in modeling complex error patterns and their relationship to
 291 correct image contexts. Building on advances in contrastive learning, we propose a Contrastive Mask
 292 Refinement Learning (CMRL) approach that explicitly models the relationships between ground
 293 truth, noisy input, and model output masks. While conventional contrastive methods Chen et al.
 294 (2020); He et al. (2020) focus on representation learning by contrasting positive and negative pairs,
 295 our tri-directional approach establishes a more complex and unique relationship structure across three
 296 different masks and is designed to explicitly guide the mask refinement learning, which is tailored
 297 specifically for the mask refinement task.

298 **Formulation.** CMRL is designed with a tri-directional contrastive framework to address three
 299 objectives: (1) maintaining clear separation between foreground and background features, (2) ensuring
 300 consistency among features within the same semantic region, and (3) enabling self-improvement
 301 through bootstrapping from successful refinements. First, CMRL categorizes pixels into six distinct
 302 regions based on their classification in target (\mathcal{M}_t), noisy (\mathcal{M}_n), and refined (\mathcal{M}_r) masks:

$$\begin{aligned} \mathcal{T}_{fg} &= (\mathcal{M}_t = 1) \wedge (\mathcal{M}_n = 1) \wedge (\mathcal{M}_r = 1), & \mathcal{T}_{bg} &= (\mathcal{M}_t = 0) \wedge (\mathcal{M}_n = 0) \wedge (\mathcal{M}_r = 0) \\ \mathcal{S}_{fg} &= (\mathcal{M}_t = 1) \wedge (\mathcal{M}_n = 0) \wedge (\mathcal{M}_r = 1), & \mathcal{S}_{bg} &= (\mathcal{M}_t = 0) \wedge (\mathcal{M}_n = 1) \wedge (\mathcal{M}_r = 0) \\ \mathcal{F}_{fg} &= (\mathcal{M}_t = 1) \wedge (\mathcal{M}_n = 0) \wedge (\mathcal{M}_r = 0), & \mathcal{F}_{bg} &= (\mathcal{M}_t = 0) \wedge (\mathcal{M}_n = 1) \wedge (\mathcal{M}_r = 1) \end{aligned}$$

303 where \mathcal{T} , \mathcal{S} , and \mathcal{F} denote true, success, and failure regions, with fg and bg indicating foreground and
 304 background. This categorization creates a natural curriculum where the model progressively refines
 305 its predictions by learning from both initially correct regions and its own successful refinements.

306 **Notation.** Let \mathbf{F} denote upsampled image embeddings. We apply projector g consisting of a 3-
 307 layer MLP to obtain projection feature maps $\mathbf{p} = g(\mathbf{F}) \in \mathbb{R}^{c \times h \times w}$. For each position $i \in \Omega =$
 308 $\{1, \dots, h \times w\}$, $\mathbf{p}_i \in \mathbb{R}^c$ denotes the feature vector at position i . The expectation $\mathbb{E}_{i \in \mathcal{R}}[\cdot]$ denotes
 309 uniform sampling over pixels in region \mathcal{R} . The similarity function $\text{sim}(\mathbf{p}_i, \mathbf{p}_j) = \mathbf{p}_i^\top \mathbf{p}_j$ computes
 310 cosine similarity between L2-normalized features. Unlike conventional segmentation losses that
 311 operate in the pixel space, our contrastive framework operates in the feature space, enabling the
 312 model to learn rich feature representations that capture contextual information Wang & Isola (2020).

313 **Intra-Class Feature Consistency** aims to create coherent feature representations within each seman-
 314 tic class, ensuring that all parts of the same object share similar feature characteristics regardless of
 315 their visual appearance. For foreground features, the intra-class loss is defined as:

$$\mathcal{L}_{intra}^{fg} = -\mathbb{E}_{i \in \mathcal{F}_{fg}} \left[\log \frac{\sum_{j \in \mathcal{S}_{fg} \cup \mathcal{T}_{fg}} \exp(\text{sim}(\mathbf{p}_i, \mathbf{p}_j)/\tau)}{\sum_{k \in \Omega} \exp(\text{sim}(\mathbf{p}_i, \mathbf{p}_k)/\tau)} \right] \quad (1)$$

316 where τ is a temperature parameter. Following the InfoNCE contrastive loss Oord et al. (2018), this
 317 loss maximizes the similarity between foreground failure features and correct foreground features

Table 1: Performance of refined masks on COCO train5K using LVIS annotations. We denote full mask annotations as \mathcal{F} , unlabeled data as \mathcal{U} , and (object center) point annotations as \mathcal{P} .

(a) Semi-Supervised, NB Wang et al. (2022) (b) Weakly Semi-Supervised, PointWSSIS Kim et al. (2023)

Methods	Annotations	AP ^{mask}	AP ^{boundary}	Methods	Annotations	AP ^{mask}	AP ^{boundary}
NB		5.1	1.9	PointWSSIS		12.6	6.3
+SegRefiner	\mathcal{F} 1% + \mathcal{U} 99%	6.2	3.8	+SegRefiner	\mathcal{F} 1% + \mathcal{P} 99%	14.7	9.5
+SAMRefiner		8.1	6.0	+SAMRefiner		21.8	16.4
+Phoenix (ours)		9.8 (+4.7)	8.1 (+6.2)	+Phoenix (ours)		28.7 (+16.1)	23.6 (+17.3)
NB		18.5	9.7	PointWSSIS		25.7	16.4
+SegRefiner	\mathcal{F} 5% + \mathcal{U} 95%	20.4	14.0	+SegRefiner	\mathcal{F} 5% + \mathcal{P} 95%	27.6	20.2
+SAMRefiner		23.8	18.5	+SAMRefiner		32.8	26.2
+Phoenix (ours)		26.7 (+8.2)	22.2 (+12.5)	+Phoenix (ours)		36.3 (+10.6)	30.2 (+13.8)
NB		22.6	12.7	PointWSSIS		30.2	20.4
+SegRefine	\mathcal{F} 10% + \mathcal{U} 90%	25.0	17.9	+SegRefiner	\mathcal{F} 10% + \mathcal{P} 90%	31.7	23.8
+SAMRefiner		28.2	22.1	+SAMRefiner		36.6	29.7
+Phoenix (ours)		30.8 (+8.2)	25.7 (+13.0)	+Phoenix (ours)		38.9 (+8.7)	32.6 (+12.2)

Table 2: Performance of refined masks on COCO validation set using LVIS annotations.

(a) Results on Mask R-CNN

(b) Results on State-of-the-art Segmentation Models

Method	AP ^{mask}	AP ^{boundary}	Method	AP ^{mask}	AP ^{boundary}	Method	AP ^{mask}	AP ^{boundary}
MRCNN(RN50)	39.8	27.3	SOLO	37.4	24.7	CondInst	39.8	29.2
+SegFix	40.6	29.1	+SegRefiner	40.5	31.3	+SegRefiner	41.1	32.2
+BPR	41.0	30.4	+SAMRefiner	44.1	34.2	+SAMRefiner	45.2	35.8
+SegRefiner	41.9	32.6	+Phoenix (ours)	46.2 (+8.8)	37.7 (+13.0)	+Phoenix (ours)	46.7 (+6.9)	38.6 (+9.4)
+SAMRefiner	45.3	35.9	RefineMask	41.2	30.5	Mask2Former	46.8	37.0
+Phoenix (ours)	46.9 (+7.1)	38.8 (+11.5)	+SegRefiner	41.9	33.0	+SegRefiner	47.4	38.8
MRCNN(RN101)	41.6	29.0	+SAMRefiner	44.7	35.3	+SAMRefiner	49.0	39.0
+SegFix	42.2	30.6	+Phoenix (ours)	46.0 (+4.8)	37.9 (+7.4)	+Phoenix (ours)	50.6 (+3.8)	42.1 (+5.1)
+BPR	42.8	32.0	ViTDet	54.6	42.5	MaskDINO	56.8	46.5
+SegRefiner	43.6	34.1	+SegRefiner	55.5	46.0	+SegRefiner	57.0	47.7
+SAMRefiner	46.6	36.9	+SAMRefiner	55.8	46.2	+SAMRefiner	57.0	47.4
+Phoenix (ours)	48.1 (+6.5)	39.8 (+10.8)	+Phoenix (ours)	56.3 (+1.7)	46.7 (+4.2)	+Phoenix (ours)	58.0 (+1.2)	48.6 (+2.1)

while implicitly pushing away features from other regions. A similar loss \mathcal{L}_{intra}^{bg} is computed for background features, and the total intra-class loss is $\mathcal{L}_{intra} = \mathcal{L}_{intra}^{fg} + \mathcal{L}_{intra}^{bg}$.

Inter-Class Feature Contrast enforces clear separation between foreground and background features, particularly in error-prone regions, by maximizing the distance between the two features:

$$\mathcal{L}_{inter}^{fg \rightarrow bg} = \mathbb{E}_{i \in \mathcal{F}_{fg}} \left[\log \sum_{j \in \mathcal{F}_{bg} \cup \mathcal{S}_{bg} \cup \mathcal{T}_{bg}} \exp(\text{sim}(\mathbf{p}_i, \mathbf{p}_j) / \tau) \right] \quad (2)$$

It pushes foreground failure features away from all background features. Likewise, $\mathcal{L}_{inter}^{bg \rightarrow fg}$ is computed. By maximizing the feature distance between different classes, this bidirectional repulsion reshapes the feature space to create clearer decision boundaries between foreground and background. The total inter-class loss is $\mathcal{L}_{inter} = \mathcal{L}_{inter}^{fg \rightarrow bg} + \mathcal{L}_{inter}^{bg \rightarrow fg}$.

Self-Improvement Regularization leverages successfully refined regions to guide the improvement of currently unrefined errors, creating a learning pathway from failure to success within the same image. This component enables the model to learn from its own successful corrections, addressing the unique challenge of transforming incorrect predictions into correct ones.

$$\mathcal{L}_{self} = -\mathbb{E}_{i \in \mathcal{F}_{fg} \cup \mathcal{F}_{bg}} \left[\log \frac{\sum_{j \in \mathcal{S}_{fg} \cup \mathcal{S}_{bg}} \exp(\text{sim}(\mathbf{p}_i, \mathbf{p}_j) / \tau)}{\sum_{k \in \Omega} \exp(\text{sim}(\mathbf{p}_i, \mathbf{p}_k) / \tau)} \right] \quad (3)$$

This loss encourages features from current failure regions ($\mathcal{F}_{fg} \cup \mathcal{F}_{bg}$) to resemble those from successfully refined regions ($\mathcal{S}_{fg} \cup \mathcal{S}_{bg}$) within the same image. Unlike the other components that focus on static correctness, this loss explicitly models the transformation from incorrect to correct predictions, creating a bootstrapping mechanism where the model learns from its own successful refinements to improve regions that are still incorrectly classified.

Loss Function. Our final CMRL loss combines all three objectives with weighting parameters:

$$\mathcal{L}_{CMRL} = \lambda_{intra} \cdot \mathcal{L}_{intra} + \lambda_{inter} \cdot \mathcal{L}_{inter} + \lambda_{self} \cdot \mathcal{L}_{self} \quad (4)$$

where λ_{intra} , λ_{inter} , and λ_{self} balance the respective objectives. This contrastive loss is combined with a traditional segmentation loss to produce the final training objective.

Table 3: Performance of refined masks on the DIS task using coarse masks from 4 different models. SAMRefiner[†] means using the HQ-SAM Ke et al. (2023) backbone for finer mask prediction.

Methods	DIS-VD		DIS-TE1		DIS-TE2		DIS-TE3		DIS-TE4		Average	
	IoU	Boundary \mathcal{F}	IoU	Boundary \mathcal{F}	IoU	Boundary \mathcal{F}	IoU	Boundary \mathcal{F}	IoU	Boundary \mathcal{F}	IoU	Boundary \mathcal{F}
U-Net	54.8	67.0	44.1	59.2	54.4	66.1	59.3	72.0	61.1	77.8	54.7	68.4
+SAMRefiner [†]	63.7	72.4	56.5	74.3	67.1	76.3	69.3	74.4	64.8	64.8	64.3	72.4
+SegRefiner	58.7	71.0	47.0	63.4	58.1	70.7	63.4	75.9	66.3	81.8	58.7	72.6
+Phoenix (ours)	74.9 (+20.1)	84.2 (+17.2)	69.6 (+25.5)	83.3 (+24.1)	79.3 (+24.9)	85.8 (+19.7)	79.4 (+20.1)	85.8 (+13.8)	75.2 (+14.1)	82.9 (+5.1)	75.7 (+21.0)	84.4 (+16.0)
PSPNet	56.4	69.3	47.6	66.8	57.7	70.5	59.9	71.8	57.6	70.4	55.8	69.8
+SAMRefiner [†]	64.2	71.7	58.8	76.4	67.2	76.6	68.0	72.5	63.2	62.9	64.3	72.0
+SegRefiner	61.9	73.7	50.4	69.0	62.0	74.3	65.3	77.5	66.7	80.3	61.3	75.0
+Phoenix (ours)	75.7 (+19.3)	84.8 (+15.5)	70.3 (+22.7)	84.0 (+17.2)	79.5 (+21.8)	86.0 (+15.5)	79.8 (+19.9)	86.2 (+14.4)	74.4 (+16.8)	82.8 (+12.4)	75.9 (+20.1)	84.8 (+15.0)
HRNet	61.0	74.2	51.0	67.1	61.2	73.6	64.4	77.9	64.7	82.0	60.5	74.9
+SAMRefiner [†]	67.5	74.2	59.7	76.6	68.1	77.9	72.6	76.0	65.8	65.0	66.8	73.9
+SegRefiner	64.4	76.1	52.9	68.7	64.7	75.7	67.6	79.8	70.5	84.6	64.0	77.0
+Phoenix (ours)	76.5 (+15.5)	85.8 (+11.6)	69.6 (+18.6)	83.4 (+16.3)	78.1 (+16.9)	86.4 (+12.8)	80.7 (+16.3)	86.6 (+8.7)	74.9 (+10.2)	83.4 (+1.4)	76.0 (+15.5)	85.1 (+10.2)
ISNet	67.1	79.8	56.2	74.9	67.1	78.4	69.9	81.2	70.1	83.8	66.1	79.6
+SAMRefiner [†]	68.1	75.2	60.6	78.4	69.4	78.2	71.7	74.8	66.0	64.0	67.1	74.1
+SegRefiner	68.3	80.2	56.9	75.0	67.8	79.0	70.9	81.9	72.2	85.2	67.2	80.3
+Phoenix (ours)	76.7 (+9.6)	86.1 (+6.3)	71.4 (+17.7)	85.3 (+10.4)	79.3 (+12.2)	86.7 (+8.3)	80.2 (+10.3)	87.0 (+5.8)	75.3 (+5.2)	84.2 (+0.4)	76.6 (+11.0)	85.9 (+6.3)

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Datasets. For general object mask refinement, we train Phoenix on the LVIS Gupta et al. (2019) dataset following the same setup as SegRefiner Wang et al. (2023). We evaluate on masks from both low-quality and high-quality segmentation models, refining coarse masks generated by semi-supervised (NB Wang et al. (2022)) and weakly semi-supervised (PointWSSIS Kim et al. (2023)) methods by measuring the quality of pseudo labels on the COCO Train5K dataset Kim et al. (2023). We also evaluate on masks produced by various instance segmentation models on the COCO Lin et al. (2014) validation set. Separately, for the fine-grained segmentation task, we train Phoenix on DIS5K Qin et al. (2022) and ThinObject-5K Liew et al. (2021) datasets and evaluate it on the DIS task Qin et al. (2022), which demands precise boundary delineation for thin structures, complex topologies, and fine-grained details.

Evaluation Metrics. We use Average Precision (AP) and boundary AP (AP^{boundary}) Cheng et al. (2021) to evaluate instance segmentation quality, and Intersection over Union (IoU) and Boundary \mathcal{F} measure Perazzi et al. (2016) for fine-grained segmentation tasks.

Implementation Details. Our mask refiner builds upon the pre-trained SAM Kirillov et al. (2023) with the ViT-H Dosovitskiy et al. (2021) backbone. During training, we freeze the encoder (637M parameters) and fine-tune only the decoder (4M parameters). We optimize the model using AdamW Loshchilov & Hutter (2017) with an initial learning rate of 1×10^{-4} and cosine decay scheduling. Training proceeds for 10K iterations with a total batch size of 16 on 8 V100 GPUs, requiring less than 10 hours of total training time. For the adversarial mask perturbation process, we use Dice loss as the perturbation objective with the maximum inner iterations N of 10, the initial step size α_0 of 0.01, and the margin ϵ of 5%. The IoU threshold τ is randomly sampled from a uniform distribution between 0.3 and 0.9, and the guidance mask \mathcal{M}_g is randomly chosen among expansion, contraction, and inversion guides. For contrastive loss weightings, we set $\lambda_{\text{intra}} = 0.4$, $\lambda_{\text{inter}} = 0.4$, and $\lambda_{\text{self}} = 0.2$ by default. The contrastive loss \mathcal{L}_{CMRL} is added to the loss function used in SAM, which is the combination of Dice and Focal Lin et al. (2017) losses.

4.2 RESULTS ON INSTANCE SEGMENTATION

Table 1 presents the performance of various refinement methods on coarse masks generated by semi-supervised Wang et al. (2022) and weakly semi-supervised Kim et al. (2023) approaches. Our method consistently outperforms previous state-of-the-art refiners across all settings, largely due to our semantic-aware noise modeling and tri-directional contrastive learning. Phoenix achieves significant improvements when refining extremely noisy masks produced with minimal supervision (1% fully-annotated and 99% point labels), with absolute gains of up to +16.1% in AP^{mask} and +17.3% in AP^{boundary}. Moreover, Table 2 shows our effectiveness in refining masks from various modern instance segmentation models He et al. (2017); Kirillov et al. (2020); Tian et al. (2020); Cheng et al. (2022); Li et al. (2022; 2023a). Phoenix achieves superior performance to existing refiners in all settings, demonstrating the versatility of our refinement approach across model architectures.

Table 4: **Ablation Study** using instance segmentation results on 1% PointWSSIS (AP¹) and MRCNN with a ResNet-50 backbone (AP²) and fine-grained segmentation results on DIS-UNet (IoU¹) and DIS-ISNet (IoU²), averaged across all DIS tasks. We denote morphological perturbation as *morp* and adversarial perturbation as *adv*.

(a) τ (IoU Threshold)				(b) \mathcal{D} (Adversarial Criterion)			(c) Contrastive Loss Component				
τ	AP ¹	AP ²		Adv Func	AP ¹	AP ²	\mathcal{L}_{intra}	\mathcal{L}_{inter}	\mathcal{L}_{self}	AP ¹	AP ²
0.3	28.2	45.5		L1	27.7	46.5	✗	✗	✗	26.9	46.1
0.5	27.3	46.4		MSE	28.7	46.8	✓	✗	✗	28.1	46.6
0.7	25.3	46.6		BCE	28.4	46.8	✗	✓	✗	28.2	46.6
0.9	24.2	46.2		Focal	28.2	46.6	✓	✓	✗	28.5	46.8
$\mathcal{U}(0.3, 0.9)$	28.7	46.9		Dice	28.7	46.9	✓	✓	✓	28.7	46.9

(d) Mask Perturbation				(e) Noise from Model			(f) Component Analysis					
Method	Perturb	IoU ¹	IoU ²	Noisy Mask	IoU ¹	IoU ²	Perturb	CMRL	AP ¹	AP ²	IoU ¹	IoU ²
SegRefiner	Morp	58.7	74.2	UNet	67.6	71.6	Morp	✗	22.3	45.2	65.3	69.9
	Adv	71.8	76.6		ISNet	65.6	70.2	Adv	✗	26.9	46.1	71.5
Phoenix	Morp	67.8	71.0	Adv (ours)	75.7	79.3	Morp	✓	23.8	45.5	67.8	71.0
	Adv	75.7	77.1		Adv	✓	28.7	46.9	75.7	77.1		

4.3 RESULTS ON FINE-GRAINED SEGMENTATION

Table 3 presents the performance of Phoenix when applied to the challenging DIS Qin et al. (2022) task across multiple test datasets and backbone models Ronneberger et al. (2015); Shen et al. (2022); Wang et al. (2020); Zhao et al. (2017). The results demonstrate that our method consistently outperforms both SAMRefiner and SegRefiner by substantial margins, with average improvements ranging from 11% to 21% in IoU, highlighting the effectiveness of our realistic and contextual noise perturbation modeling in the mask refinement task. Figure 1 provides visual comparisons between our method and existing refiners. Phoenix produces masks with significantly improved boundary adherence and structural integrity compared to both noisy inputs and competing refiners, successfully recovering challenging scenarios, such as thin structures and intricate boundaries.

4.4 ABLATION STUDIES

We conduct extensive ablation studies to analyze the contribution of individual components and hyperparameters. For these experiments, we use low-quality masks from PointWSSIS with 1% supervision (denoted as AP¹) and high-quality masks from Mask R-CNN with ResNet-50 He et al. (2017) backbone (denoted as AP²). Additionally, we evaluate fine-grained segmentation performance on DIS-UNet (denoted as IoU¹) and DIS-ISNet (denoted as IoU²), averaged across all the DIS test splits. **Additional ablation studies and in-depth analyses can be found in the appendix.**

Effect of Adversarial Mask Perturbation Parameters. Table 4a shows the impact of the IoU threshold τ on refinement performance. Higher τ values lead to better performance on high-quality masks (higher AP²) but lower performance on low-quality masks (lower AP¹). For robustness against various noise patterns, we randomly sample τ from a uniform distribution $\mathcal{U}(0.3, 0.9)$ for each noise mask generation step during training, achieving the best balanced performance. Moreover, Table 4b presents the performance with different adversarial criteria, demonstrating our method is relatively robust to the choice of objective, with Dice loss providing the best overall balance.

Effect of Contrastive Loss Components. Table 4c illustrates the performance impact of different contrastive loss components. We observe that both intra-class consistency (\mathcal{L}_{intra}) and inter-class contrast (\mathcal{L}_{inter}) contribute significantly to performance, with their combination yielding substantial improvements. Including all three components achieves the optimal performance.

Effect of Mask Perturbation Method. Table 4d compares different mask perturbation approaches across models. When applying adversarial perturbation to SegRefiner, its performance improves substantially. Conversely, when our Phoenix is trained with simple morphological perturbations, performance drops significantly compared to our full approach. These results highlight the critical importance of sophisticated noise generation in mask refinement learning.

Comparison with Real Model Errors. Table 4e compares our adversarially generated noise patterns with noise obtained directly from real segmentation models. While using output masks from the pre-trained model (UNet [Ronneberger et al. \(2015\)](#) or ISNet [Shen et al. \(2022\)](#)) as noisy masks for training provides adequate performance, our adversarial approach significantly outperforms them. This superiority stems from two key advantages: (1) a wider diversity of error patterns, and (2) controllable perturbation magnitude.

Effect of Individual Components. Table 4f isolates the contributions of AMP and CMRL. Starting from the baseline (morphological noise with pixel-wise loss), AMP alone provides substantial improvements (+4.6% AP¹, +6.2% IoU¹), demonstrating the importance of semantic-aware noise generation. CMRL alone yields meaningful gains (+1.5% AP¹, +2.5% IoU¹), validating the impact of our tri-directional contrastive framework. Combining both components achieves the best performance (+6.4% AP¹, +10.4% IoU¹), with the combined improvement exceeding the sum of individual contributions (+8.7% IoU¹), demonstrating true synergy. This synergistic effect occurs because CMRL is more effective when encountering realistic noise patterns from AMP; the semantically meaningful error distributions enable more informative tri-directional contrastive relationships, leading to superior refinement learning. Notably, AMP provides larger marginal gains on challenging scenarios, while both components contribute substantially across all settings, confirming their complementary nature.

5 CONCLUSION

We presented Phoenix, a novel framework for mask refinement that leverages adversarial perturbation and contrastive learning to address the limitations of existing approaches. Our method generates semantically meaningful noise patterns that better reflect real-world segmentation errors and employs a tri-directional contrastive learning approach that explicitly models the relationships between ground truth, noisy input, and current prediction. Extensive experiments demonstrate that Phoenix significantly outperforms existing methods across diverse tasks and noise conditions. We believe our work establishes a new foundation for mask refinement and opens up promising directions for future research, including self-supervised refinement and multimodal guidance integration.

REFERENCES

- Anurag Arnab, Ondrej Miksik, and Philip HS Torr. On the robustness of semantic segmentation models to adversarial attacks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 888–897, 2018. 3, 5
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 17
- Yuri Y Boykov and M-P Jolly. Interactive graph cuts for optimal boundary & region segmentation of objects in nd images. In *Proceedings eighth IEEE international conference on computer vision. ICCV 2001*, volume 1, pp. 105–112. IEEE, 2001. 3
- John Canny. A computational approach to edge detection. *IEEE Transactions on pattern analysis and machine intelligence*, (6):679–698, 1986. 19
- Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 801–818, 2018. 3
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PmLR, 2020. 3, 6
- Bowen Cheng, Ross Girshick, Piotr Dollár, Alexander C Berg, and Alexander Kirillov. Boundary iou: Improving object-centric image segmentation evaluation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15334–15342, 2021. 8
- Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 1290–1299, 2022. 1, 3, 8

- 540 Noel Codella, Veronica Rotemberg, Philipp Tschandl, M Emre Celebi, Stephen Dusza, David
541 Gutman, Brian Helba, Aadi Kalloo, Konstantinos Liopyris, Michael Marchetti, et al. Skin lesion
542 analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging
543 collaboration (isic). *arXiv preprint arXiv:1902.03368*, 2019. 16
- 544 Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo
545 Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban
546 scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern
547 recognition*, pp. 3213–3223, 2016. 16
- 548 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
549 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,
550 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale.
551 In *International Conference on Learning Representations*, 2021. 8, 17, 24
- 552 Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The
553 pascal visual object classes (voc) challenge. *International journal of computer vision*, 88:303–338,
554 2010. 15
- 555 Golnaz Ghiasi, Xiuye Gu, Yin Cui, and Tsung-Yi Lin. Scaling open-vocabulary image segmentation
556 with image-level labels. In *European conference on computer vision*, pp. 540–557. Springer, 2022.
557 15, 26
- 558 Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial
559 examples. In *International Conference on Learning Representations (ICLR)*, 2015. 3, 5
- 560 Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance
561 segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern
562 recognition*, pp. 5356–5364, 2019. 6, 8, 15, 16
- 563 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
564 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
565 pp. 770–778, 2016. 20
- 566 Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the
567 IEEE international conference on computer vision*, pp. 2961–2969, 2017. 8, 9, 16, 20
- 568 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for
569 unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on
570 computer vision and pattern recognition*, pp. 9729–9738, 2020. 3, 6
- 571 Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint
572 arXiv:1606.08415*, 2016. 17
- 573 Zhichao Huang and Tong Zhang. Black-box adversarial attack with transferable model-based
574 embedding. In *International Conference on Learning Representations*, 2020. 3, 5
- 575 Lei Ke, Mingqiao Ye, Martin Danelljan, Yu-Wing Tai, Chi-Keung Tang, Fisher Yu, et al. Segment
576 anything in high quality. *Advances in Neural Information Processing Systems*, 36:29914–29934,
577 2023. 8
- 578 Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer
579 vision? *Advances in neural information processing systems*, 30, 2017. 5
- 580 Beomyoung Kim, Youngjoon Yoo, Chae Eun Rhee, and Junmo Kim. Beyond semantic to instance
581 segmentation: Weakly-supervised instance segmentation via semantic knowledge transfer and
582 self-refinement. In *Proceedings of the IEEE/CVF conference on computer vision and pattern
583 recognition*, pp. 4278–4287, 2022. 2
- 584 Beomyoung Kim, Joonhyun Jeong, Dongyoon Han, and Sung Ju Hwang. The devil is in the
585 points: Weakly semi-supervised instance segmentation via point-guided mask representation.
586 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
587 11360–11370, 2023. 2, 3, 7, 8, 20

- 594 Alexander Kirillov, Yuxin Wu, Kaiming He, and Ross Girshick. Pointrend: Image segmentation as
595 rendering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
596 pp. 9799–9808, 2020. [8](#)
- 597 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
598 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings*
599 *of the IEEE/CVF international conference on computer vision*, pp. 4015–4026, 2023. [1](#), [3](#), [4](#), [8](#), [17](#)
- 600 Philipp Krähenbühl and Vladlen Koltun. Efficient inference in fully connected crfs with gaussian
601 edge potentials. *Advances in neural information processing systems*, 24, 2011. [3](#)
- 602 Feng Li, Hao Zhang, Huaizhe Xu, Shilong Liu, Lei Zhang, Lionel M Ni, and Heung-Yeung Shum.
603 Mask dino: Towards a unified transformer-based framework for object detection and segmentation.
604 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp.
605 3041–3050, 2023a. [8](#)
- 606 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image
607 pre-training with frozen image encoders and large language models. In *International conference*
608 *on machine learning*, pp. 19730–19742. PMLR, 2023b. [15](#), [26](#)
- 609 Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer
610 backbones for object detection. In *European conference on computer vision*, pp. 280–296. Springer,
611 2022. [8](#)
- 612 Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yinan Zhao, Hang Zhang, Peizhao Zhang,
613 Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted
614 clip. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp.
615 7061–7070, 2023. [15](#), [26](#)
- 616 Jun Hao Liew, Scott Cohen, Brian Price, Long Mai, and Jiashi Feng. Deep interactive thin object
617 selection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*,
618 pp. 305–314, 2021. [8](#)
- 619 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
620 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer vision–*
621 *ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings,*
622 *part v 13*, pp. 740–755. Springer, 2014. [8](#), [16](#)
- 623 Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object
624 detection. In *Proceedings of the IEEE international conference on computer vision*, pp. 2980–2988,
625 2017. [8](#), [23](#)
- 626 Yuqi Lin, Minghao Chen, Wenxiao Wang, Boxi Wu, Ke Li, Binbin Lin, Haifeng Liu, and Xiaofei
627 He. Clip is also an efficient segmenter: A text-driven approach for weakly supervised semantic
628 segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
629 *Recognition*, pp. 15305–15314, 2023. [15](#), [16](#)
- 630 Yuqi Lin, Hengjia Li, Wenqi Shao, Zheng Yang, Jun Zhao, Xiaofei He, Ping Luo, and Kaipeng Zhang.
631 Samrefiner: Taming segment anything model for universal mask refinement. In *The Thirteenth*
632 *International Conference on Learning Representations*, 2025. [2](#), [3](#), [4](#), [15](#), [18](#), [30](#)
- 633 Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic
634 segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
635 pp. 3431–3440, 2015. [3](#)
- 636 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*
637 *arXiv:1711.05101*, 2017. [8](#), [18](#)
- 638 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
639 Towards deep learning models resistant to adversarial attacks. In *International Conference on*
640 *Learning Representations*, 2018. [3](#), [5](#)
- 641 David Marr and Ellen Hildreth. Theory of edge detection. *Proceedings of the Royal Society of*
642 *London. Series B. Biological Sciences*, 207(1167):187–217, 1980. [19](#)

- 648 Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural
649 networks for volumetric medical image segmentation. In *2016 fourth international conference on*
650 *3D vision (3DV)*, pp. 565–571. Ieee, 2016. 23
- 651 Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive
652 coding. *arXiv preprint arXiv:1807.03748*, 2018. 6, 19
- 653 A Paszke. Pytorch: An imperative style, high-performance deep learning library. *arXiv preprint*
654 *arXiv:1912.01703*, 2019. 18
- 655 Karl Pearson. Vii. note on regression and inheritance in the case of two parents. *proceedings of the*
656 *royal society of London*, 58(347-352):240–242, 1895. 6, 18
- 657 Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc Van Gool, Markus Gross, and Alexander
658 Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation.
659 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 724–732,
660 2016. 8
- 661 Xuebin Qin, Hang Dai, Xiaobin Hu, Deng-Ping Fan, Ling Shao, and Luc Van Gool. Highly accurate
662 dichotomous image segmentation. In *European Conference on Computer Vision*, pp. 38–56.
663 Springer, 2022. 8, 9, 20
- 664 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
665 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
666 models from natural language supervision. In *International conference on machine learning*, pp.
667 8748–8763. PmLR, 2021. 15, 26
- 668 Shenghai Rong, Bohai Tu, Zilei Wang, and Junjie Li. Boundary-enhanced co-training for weakly
669 supervised semantic segmentation. In *Proceedings of the IEEE/CVF conference on computer*
670 *vision and pattern recognition*, pp. 19574–19584, 2023. 15, 16
- 671 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical
672 image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI*
673 *2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III*
674 *18*, pp. 234–241. Springer, 2015. 3, 9, 10, 17, 20
- 675 Tiancheng Shen, Yuechen Zhang, Lu Qi, Jason Kuen, Xingyu Xie, Jianlong Wu, Zhe Lin, and Jiaya
676 Jia. High quality segmentation for ultra high-resolution images. In *Proceedings of the IEEE/CVF*
677 *conference on computer vision and pattern recognition*, pp. 1310–1319, 2022. 9, 10, 20
- 678 Chufeng Tang, Hang Chen, Xiao Li, Jianmin Li, Zhaoxiang Zhang, and Xiaolin Hu. Look closer
679 to segment better: Boundary patch refinement for instance segmentation. In *Proceedings of the*
680 *IEEE/CVF conference on computer vision and pattern recognition*, pp. 13926–13935, 2021. 2, 3, 6
- 681 Zhi Tian, Chunhua Shen, and Hao Chen. Conditional convolutions for instance segmentation. In
682 *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020,*
683 *Proceedings, Part I 16*, pp. 282–298. Springer, 2020. 8
- 684 Zhi Tian, Chunhua Shen, Xinlong Wang, and Hao Chen. Boxinst: High-performance instance
685 segmentation with box annotations. In *Proceedings of the IEEE/CVF Conference on Computer*
686 *Vision and Pattern Recognition*, pp. 5443–5452, 2021. 2
- 687 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
688 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
689 *systems*, 30, 2017. 17
- 690 Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong
691 Mu, Mingkui Tan, Xinggang Wang, et al. Deep high-resolution representation learning for visual
692 recognition. *IEEE transactions on pattern analysis and machine intelligence*, 43(10):3349–3364,
693 2020. 9
- 694 Mengyu Wang, Henghui Ding, Jun Hao Liew, Jiajun Liu, Yao Zhao, and Yunchao Wei. Segrefiner:
695 Towards model-agnostic segmentation refinement with discrete diffusion process. *Advances in*
696 *Neural Information Processing Systems*, 36:79761–79780, 2023. 2, 3, 4, 6, 8, 30

- 702 Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through align-
703 ment and uniformity on the hypersphere. In *International conference on machine learning*, pp.
704 9929–9939. PMLR, 2020. 3, 6
- 705
706 Xinlong Wang, Rufeng Zhang, Chunhua Shen, Tao Kong, and Lei Li. Dense contrastive learning
707 for self-supervised visual pre-training. In *Proceedings of the IEEE/CVF conference on computer
708 vision and pattern recognition*, pp. 3024–3033, 2021. 3
- 709
710 Zhenyu Wang, Yali Li, and Shengjin Wang. Noisy boundaries: Lemon or lemonade for semi-
711 supervised instance segmentation? In *Proceedings of the IEEE/CVF conference on computer
712 vision and pattern recognition*, pp. 16826–16835, 2022. 2, 3, 7, 8
- 713
714 Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer:
715 Simple and efficient design for semantic segmentation with transformers. *Advances in neural
716 information processing systems*, 34:12077–12090, 2021a. 1, 3
- 717
718 Zhenda Xie, Yutong Lin, Zheng Zhang, Yue Cao, Stephen Lin, and Han Hu. Propagate yourself:
719 Exploring pixel-level consistency for unsupervised visual representation learning. In *Proceedings
720 of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16684–16693, 2021b.
721 3
- 722
723 Yuhui Yuan, Jingyi Xie, Xilin Chen, and Jingdong Wang. Segfix: Model-agnostic boundary refine-
724 ment for segmentation. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow,
725 UK, August 23–28, 2020, Proceedings, Part XII 16*, pp. 489–506. Springer, 2020. 2, 3, 4, 6
- 726
727 Zhuoyang Zhang, Han Cai, and Song Han. Efficientvit-sam: Accelerated segment anything model
728 without performance loss. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
729 Pattern Recognition*, pp. 7859–7863, 2024. 24
- 730
731 Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing
732 network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
733 2881–2890, 2017. 3, 9
- 734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

756 APPENDIX

757
758 A ADDITIONAL APPLICATIONS OF PHOENIX

759
760 A.1 SELF-SUPERVISED MASK REFINEMENT LEARNING.

761
762 Our intriguing finding is the potential for self-supervised mask refinement without ground-truth
763 masks. By leveraging SAM’s zero-shot capability, we can generate pseudo-target masks using a
764 grid of prompting points, enabling training with target and noisy mask pairs without ground-truth
765 annotations. As shown in Table 5a, our self-supervised pipeline achieves competitive performance
766 in AP¹ compared to the fully supervised setting. This success can be attributed to our contrastive
767 learning framework, which focuses on the relationship between mask pairs and learns generalizable
768 noise patterns rather than overfitting to specific annotations. However, this self-supervised approach
769 has limitations. Its effectiveness depends on SAM’s zero-shot capability for the target task. For fine-
770 grained segmentation tasks where SAM struggles to generate high-quality masks, the self-supervised
771 pipeline underperforms. Nevertheless, these results point to promising future directions in reducing
772 annotation requirements for mask refinement.

773
774 A.2 INTEGRATION OF VISUAL CUES.

775 We identify limitations when handling extremely challenging noise patterns, such as complete
776 mislocalization and uncertain target objects. In these scenarios, refinement becomes inherently
777 ambiguous without additional guidance. To address this, we explored integrating (ground-truth)
778 visual prompts (box and points derived from the ground-truth mask) that provide explicit target
779 object information, achieving remarkable performance improvements, as shown in Table 5b. Another
780 important challenge is handling misclassified object masks. Since our approach operates in a class-
781 agnostic manner, it cannot correct semantic class errors inherent from real models. To address this
782 limitation, incorporating open vocabulary models Ghiasi et al. (2022); Liang et al. (2023) or text
783 embeddings Li et al. (2023b); Radford et al. (2021) will be a promising research direction. These
784 insights highlight both the current capabilities and future potential of mask refinement systems,
785 pointing toward more robust, multimodal approaches that can handle increasingly diverse and
786 challenging segmentation scenarios.

787 Table 5: **Ablation Study** using instance segmentation results on 1% PointWSSIS (AP¹) and fine-
788 grained segmentation results on DIS-UNet (IoU¹) and DIS-ISNet (IoU²), averaged across all DIS
789 tasks. We denote the self-supervised mask refinement as *SSL*.

790 (a) Self-Supervised Mask Refinement Learning (b) Integration of Ground-Truth Visual Prompt

SSL	AP ¹	IoU ¹	Prompt	AP ¹	IoU ¹
✗	28.7	75.7	-	28.7	75.7
✓	27.4	58.0	Point	30.3	77.0
			Box	37.3	78.5
			Point+Box	37.4	78.6

791
792
793
794
795
796
797
798 A.3 SEMANTIC SEGMENTATION

799 We evaluate Phoenix on the PASCAL VOC2012 Everingham et al. (2010) semantic segmentation
800 benchmark to demonstrate its effectiveness in the semantic segmentation domain. While our main
801 paper focuses on instance and fine-grained segmentation, this evaluation aims to verify the versatility
802 of our refinement approach across diverse segmentation tasks. For this evaluation, we follow the
803 *split-then-merge* strategy used in SAMRefiner Lin et al. (2025) for mask refinement in the semantic
804 segmentation domain and leverage the Phoenix model trained on the LVIS Gupta et al. (2019)
805 instance segmentation dataset. This approach allows us to refine pseudo semantic labels generated by
806 unsupervised Zhou et al. (2022) or weakly-supervised Lin et al. (2023); Rong et al. (2023) models
807 on the VOC2012 training set, addressing the challenge of noisy pseudo-labels that often plague
808 these learning paradigms. As shown in Table 6a, Phoenix consistently improves the performance
809 of various semantic segmentation methods. Namely, Phoenix improves the unsupervised method,
MaskCLIP Zhou et al. (2022), from 47.8% to 59.1% (+11.3%) in mIoU, effectively refining coarse

Table 6: **Evaluation of Phoenix on additional applications.** (a) Semantic segmentation results on PASCAL VOC2012 training set across different models: \mathcal{U} (unlabeled) and \mathcal{I} (image-level labels) supervisions. (b) Ground truth annotation refinement on COCO2017 validation set. The † denotes that the model is trained using the VOC2012 dataset.

(a) Results on VOC2012.

Method	Annotations	mIoU
MaskCLIP		47.8
+SAMRefiner	\mathcal{U}	57.3
+Phoenix (ours)		59.1 (+11.3)
BECO		66.3
+SAMRefiner	\mathcal{I}	71.8
+Phoenix (ours)		75.0 (+8.7)
CLIP-ES		70.8
+SAMRefiner	\mathcal{I}	79.3
+Phoenix (ours)		81.6 (+10.8)

(b) Refined masks on COCO2017 *val*.

Data	AP^{mask}	$AP^{boundary}$
COCO val	38.3	27.3
+SAMRefiner	41.5	33.0
+Phoenix† (ours)	42.2 (+3.9)	33.6 (+6.3)
+Phoenix (ours)	43.7 (+5.4)	36.1 (+8.8)

pseudo-labels to achieve more accurate semantic boundaries. In addition, Phoenix enhances the weakly-supervised methods, BECO Rong et al. (2023) and CLIP-ES Lin et al. (2023), by +8.7% and +10.8%, respectively. The qualitative results in Figure 6 demonstrate Phoenix’s ability to correct both over-segmentation and under-segmentation errors in pseudo-labels, particularly at object boundaries and in regions with complex semantic transitions. These substantial gains highlight the universal applicability of our approach.

A.4 HUMAN ANNOTATION CORRECTION

We investigate the effectiveness of Phoenix for improving manually annotated masks, which often contain imperfections, particularly at complex boundaries. This application is important because human annotations in datasets like COCO Lin et al. (2014) are frequently coarse due to the sparse polygon points format used during annotation, limiting boundary precision. We use the COCO dataset’s annotations as representative of human-annotated masks and evaluate their quality against the more precise LVIS Gupta et al. (2019) annotations as reference. The LVIS dataset provides a finer version of COCO mask annotations through more detailed annotation protocols, which can be regarded as baseline refined masks. As shown in Table 6b, Phoenix significantly improves the annotation quality of COCO2017 *val* ground-truth masks, enhancing AP^{mask} by +5.4% and $AP^{boundary}$ by +8.8%. The qualitative examples in Figure 7 illustrate how Phoenix refines coarse COCO annotations to better align with precise ground truth boundaries. The improvement is particularly pronounced at boundaries, where human annotators often struggle with precision. Our method effectively identifies and corrects these imperfections, suggesting potential applications in annotation workflow enhancement and dataset quality improvement. This capability could significantly reduce the time and cost associated with creating high-quality segmentation datasets by allowing annotators to focus on object identification rather than precise boundary delineation.

A.5 ZERO-SHOT GENERALIZATION

To assess the generalization capabilities of Phoenix beyond its training domain, we evaluate our method in a zero-shot setting on two challenging datasets: Cityscapes Cordts et al. (2016) for urban scene understanding and the ISIC2018 skin lesion segmentation dataset Codella et al. (2019) for medical imaging. Importantly, Phoenix is trained exclusively on the LVIS dataset and applied directly to these domains without any domain-specific fine-tuning or adaptation.

Cityscapes Evaluation Table 7a shows the performance of Phoenix when applied to refine instance segmentation masks on the Cityscapes validation set. For the evaluation metric, we follow the COCO-style evaluation metrics Lin et al. (2014). Despite the substantial domain gap between LVIS and Cityscapes, Phoenix demonstrates robust generalization. Our method achieves consistent improvements over the baseline Mask R-CNN (R50) He et al. (2017), with gains of +1.1% in AP and +1.0% in AP50. Notably, while SAMRefiner shows degraded performance (-3.6% AP), Phoenix maintains and enhances the original segmentation quality, highlighting the effectiveness of our noise modeling approach across different visual domains.

Medical Imaging Evaluation The results on the ISIC2018 skin lesion segmentation dataset (Table 7b) further demonstrate Phoenix’s domain adaptability. Medical imaging presents unique challenges including varying illumination conditions, texture variations, and irregular lesion boundaries that differ substantially from natural images in LVIS. Despite these challenges, Phoenix achieves substantial improvements over the baseline UNet [Ronneberger et al. \(2015\)](#), with gains of +4.6% in IoU and +6.0% in F1 score. The superior performance compared to SAMRefiner (+3.7% IoU improvement) underscores the robustness of our tri-directional contrastive learning and semantic-aware noise modeling in handling domain-specific segmentation challenges. These zero-shot results validate that Phoenix learns generalizable refinement principles rather than dataset-specific patterns, making it a versatile solution for mask refinement across diverse application domains without requiring additional training or adaptation.

Table 7: **Zero-Shot Generalization of Phoenix** on (a) Cityscapes instance segmentation validation set and (b) Medical imaging, ISIC2018 skin lesion segmentation dataset.

(a) Results on Cityscapes			(b) Results on Medical Imaging Dataset		
Method	AP	AP50	Method	IoU	F1
MRCNN (R50)	36.1	60.9	UNet	54.6	71.2
+SAMRefiner	32.5	58.4	+SAMRefiner	55.5	76.1
+Phoenix (ours)	37.2 (+1.1)	61.9 (+1.0)	+Phoenix (ours)	59.2 (+4.6)	77.2 (+6.0)

B IMPLEMENTATION DETAILS

B.1 NETWORK ARCHITECTURE

Phoenix builds upon the SAM [Kirillov et al. \(2023\)](#) architecture with specific modifications optimized for the mask refinement task. Our architecture leverages key components from SAM while incorporating specialized elements for refinement.

The encoder component consists of the pre-trained ViT-H [Dosovitskiy et al. \(2021\)](#) encoder from SAM, which remains frozen during training to maintain computational efficiency. This encoder generates image embeddings \mathbf{E}_{img} with dimensions $256 \times 64 \times 64$ given an input image of 1024×1024 resolution.

The decoder takes the noisy input mask as a dense prompt through SAM’s prompt encoder. The mask prompt is processed at $4 \times$ lower resolution than the input image, then downscaled an additional $4 \times$ using two convolutional layers with GELU [Hendrycks & Gimpel \(2016\)](#) activations and layer normalization [Ba et al. \(2016\)](#), ultimately mapping to a 256-dimensional embedding. These dense prompt embeddings are added element-wise with the image embeddings \mathbf{E}_{img} to incorporate mask information.

Additionally, we derive point and box coordinates from the noisy input mask to obtain visual prompt embeddings \mathbf{E}_v . For points, we extract positional encodings summed with learnable foreground/background embeddings. For boxes, we encode the top-left and bottom-right corners using positional encodings combined with learned corner-specific embeddings. During the adversarial mask perturbation process, the perturbation embeddings \mathbf{E}_p are concatenated with these visual prompt embeddings, maintaining the standard embedding dimension of 256.

The mask decoder follows SAM’s transformer [Vaswani et al. \(2017\)](#) decoder architecture with two layers. Each layer performs self-attention on the tokens, cross-attention from tokens to image embeddings, MLP updates to tokens, and cross-attention from image embeddings to tokens. This structure allows bidirectional information flow between prompts and image features. After decoder processing, the updated image embedding is upsampled by $4 \times$ with two transposed convolutional layers. The final mask prediction uses a dot product between the upscaled image embedding with dimensions $32 \times 256 \times 256$ and the output token embeddings.

For the Contrastive Mask Refinement Learning (CMRL), we implement a projector network g that transforms the upsampled image embeddings \mathbf{F} into a space optimized for contrastive learning. This projector consists of 3-layer MLP blocks, generating the projection feature maps $\mathbf{P} = g(\mathbf{F})$ with dimensions $32 \times 256 \times 256$.

B.2 TRAINING PROTOCOL

Our Phoenix is implemented using the Pytorch framework [Paszke \(2019\)](#). Phoenix is trained using the AdamW optimizer [Loshchilov & Hutter \(2017\)](#) with $\beta_1 = 0.9$, $\beta_2 = 0.999$, an initial learning rate of 1×10^{-4} , linear warmup for 500 iterations, and cosine decay scheduling. Training proceeds with a total batch size of 16 (2 samples per GPU) on 8 V100 GPUs, weight decay of 5×10^{-4} , and gradient clipping with a maximum norm value of 0.1.

For each adversarial mask perturbation, we initialize perturbation embeddings \mathbf{E}_p with zeros and add random noise (normal distribution with a mean of 0.0 and a standard deviation of 0.1).

To prevent over-reliance on visual prompts, we randomly omit them with a 30% probability during training by guiding only the noisy mask into the model without the visual prompts. We apply data augmentation, including color jittering and random horizontal flipping, to increase training diversity.

We will release the code upon acceptance.

B.3 INFERENCE STRATEGY

Phoenix inherits SAM’s multimask output capability, generating multiple refined mask candidates and selecting the one with the highest IoU prediction score. This approach leverages the predictive confidence of the model to identify the most accurate refinement.

Following the strategy introduced in [Lin et al. \(2025\)](#), our model implements cascade self-refinement during inference, where each refinement output serves as input for subsequent iterations. This process, repeated for multiple steps, progressively enhances mask quality through iterative improvement. We empirically determine the optimal number of refinement iterations through detailed analysis, as shown in Figure 5a. Our results indicate that 5 iterations provide the most favorable balance between refinement performance and computational efficiency. Each individual refinement step requires only 6 ms for decoder inference on a V100 GPU, resulting in an additional 24 ms total processing time for the complete refinement cascade. This minimal computational overhead maintains highly efficient inference while substantially improving mask quality. In addition, we enhance model training by incorporating this cascade strategy as a curriculum learning mechanism, gradually transitioning from difficult noisy masks to easier self-refined versions.

B.4 DETAILS OF SELF-SUPERVISED MASK REFINEMENT

Our self-supervised approach eliminates the need for ground-truth annotations through a simple yet effective pseudo-supervision strategy, as described in Section A.1. For the pseudo-target generation, we leverage SAM’s automatic mask generation mode by prompting it with a fixed grid of 4×4 points across the image. From the resulting mask candidates, we filter out low-confidence predictions and randomly select one high-confidence mask to serve as our pseudo-target. This approach utilizes SAM’s strong zero-shot capabilities to create high-quality reference masks without human annotation. Once the pseudo-targets are generated, we apply the same adversarial perturbation process and training procedure as in our fully-supervised setting. The model learns to refine synthetically perturbed versions of these pseudo-targets back to their original state, effectively transferring SAM’s segmentation capabilities into our efficient refinement architecture.

B.5 DETAILS OF SEMANTIC CORRELATION ANALYSIS

To quantitatively evaluate the fundamental differences between adversarial and morphological noise patterns, we developed a semantic correlation analysis framework. This analysis aims to measure and characterize how different types of noise patterns relate to semantic image features, providing insights into why our adversarial approach generates more realistic and challenging training examples. This analysis was performed on the LVIS validation set. For each image-mask pair, we compute the Pearson correlation [Pearson \(1895\)](#) between the spatial distribution of noise and semantic features:

$$\text{Corr}(S, N) = \frac{\sum_{x,y} (S(x, y) - \bar{S})(N(x, y) - \bar{N})}{\sqrt{\sum_{x,y} (S(x, y) - \bar{S})^2 \sum_{x,y} (N(x, y) - \bar{N})^2}} \quad (5)$$

where $S(x, y)$ represents semantic feature strength (derived from edge and texture maps extracted using Canny edge detection Canny (1986) and Laplacian of Gaussian filtering Marr & Hildreth (1980)) and $N(x, y)$ represents noise magnitude at position (x, y) .

Figures 2b and 2c in the main paper provided the distribution of semantic correlation values for both noise types. Morphological noise exhibits a compact, near-symmetric distribution centered close to zero (-0.2 to 0.3), indicating little relationship between noise placement and semantic features. This confirms our hypothesis that conventional morphological operations create perturbations based primarily on geometric constraints rather than semantic understanding. In contrast, our adversarial noise shows a substantially broader distribution (-0.6 to 0.8) extending into both positive and negative correlation regions, demonstrating that it captures a diverse spectrum of semantic relationships.

This broader distribution of adversarial noise is highly beneficial for mask refinement learning for several reasons: (1) it provides comprehensive error coverage across both semantically meaningful regions and homogeneous areas, (2) it creates challenging examples at complex boundaries through positive correlations while generating hard negatives in seemingly simple regions through negative correlations, (3) it better mirrors the diverse error patterns produced by real segmentation models, which make mistakes with varying semantic correlations depending on the context.

B.6 DETAILS OF CONTRASTIVE MASK REFINEMENT LEARNING (CMRL)

This section provides practical implementation details for CMRL, addressing how region masks and projection features are computed and batched in practice. Algorithm 2 presents the PyTorch-style pseudo-code.

Region Mask Computation. Given three binary masks, *i.e.*, target \mathcal{M}_t (ground truth), noisy \mathcal{M}_n (input), and refined \mathcal{M}_r (current prediction), we compute six region masks through logical operations. Each mask is binarized using a 0.5 threshold. The six regions are then derived: \mathcal{T}_{fg} identifies pixels classified as foreground in all three masks, \mathcal{F}_{fg} captures false negatives that remain uncorrected (true foreground but predicted as background in both noisy and refined masks), and \mathcal{S}_{fg} represents successful corrections (true foreground, initially misclassified, but corrected in the refined mask). The same logic applies to background regions with subscript *bg*. This categorization creates a pixel-level curriculum where each spatial position belongs to exactly one of the six mutually exclusive regions.

Feature Projection and Sampling. The upsampled image embeddings $\mathbf{F} \in \mathbb{R}^{c \times h \times w}$ are processed through projector g to obtain $\mathbf{p} = g(\mathbf{F}) \in \mathbb{R}^{c \times h \times w}$, where $c = 32$, $h = 256$, and $w = 256$ in our implementation. We then apply L2 normalization along the channel dimension, converting dot products into cosine similarities. For computational efficiency, we sample up to 256 pixels from each region using uniform random sampling. This yields feature vectors $\mathbf{p}_i \in \mathbb{R}^c$ for each sampled position i . The loss components are computed separately for each image in the batch and averaged.

Role of Each Loss Component. The three loss components serve distinct but complementary purposes. The intra-class loss encourages features from failure regions to align with features from correct regions of the same semantic class, achieving within-class consistency by pulling failure foreground features (\mathcal{F}_{fg}) toward correct foreground features ($\mathcal{S}_{fg} \cup \mathcal{T}_{fg}$) through the InfoNCE objective Oord et al. (2018). The inter-class loss enforces separation between foreground and background features by pushing failure foreground features away from all background regions ($\mathcal{F}_{bg} \cup \mathcal{S}_{bg} \cup \mathcal{T}_{bg}$), creating clearer decision boundaries in the feature space. The self-improvement loss enables the model to learn from its own successful corrections by guiding current failures ($\mathcal{F}_{fg} \cup \mathcal{F}_{bg}$) toward successfully refined regions ($\mathcal{S}_{fg} \cup \mathcal{S}_{bg}$), creating a bootstrapping mechanism where the model progressively improves by learning from regions it has already corrected within the same image.

Projector Architecture and Ablation. The projector g consists of three 1×1 convolutional layers with LayerNorm and GELU activations between layers. Table 8g suggests the architecture choice of the projector. Without projection (identity mapping), performance is limited to 27.1% AP¹. Single-layer and two-layer projectors achieve 27.9% and 28.5% respectively, while our three-layer design reaches optimal performance at 28.7%. Adding a fourth layer provides no benefit (28.6%), indicating that three layers provide sufficient capacity.

Algorithm 2 PyTorch-style Pseudo-code for Contrastive Mask Refinement Learning

```

1026
1027
1028 # Input: image features F (Bx32x256x256),
1029 # masks M_t (GT), M_n (Noisy Mask), M_r (Refined Mask) (Bx1xHxW)
1030
1031 # Step 1: Project and normalize features
1032 projector = MLP_layer(in_channels=32, out_channels=32, num_layers=3)
1033 P = F.normalize(projector(F), p=2, dim=1) # (Bx32x256x256)
1034
1035 # Step 2: Define six region masks (binarize and detach)
1036 M_t_bin = (M_t > 0.5).float().detach()
1037 M_n_bin = (M_n > 0.5).float().detach()
1038 M_r_bin = (M_r > 0.5).float().detach()
1039
1040 T_fg = (M_t_bin==1) & (M_n_bin==1) & (M_r_bin==1) # True positive
1041 T_bg = (M_t_bin==0) & (M_n_bin==0) & (M_r_bin==0) # True negative
1042 S_fg = (M_t_bin==1) & (M_n_bin==0) & (M_r_bin==1) # Success FN->TP
1043 S_bg = (M_t_bin==0) & (M_n_bin==1) & (M_r_bin==0) # Success FP->TN
1044 F_fg = (M_t_bin==1) & (M_n_bin==0) & (M_r_bin==0) # Failure: uncorrected FN
1045 F_bg = (M_t_bin==0) & (M_n_bin==1) & (M_r_bin==1) # Failure: uncorrected FP
1046
1047 # Step 3: Compute losses per batch item
1048 for b in range(B):
1049     feat = P[b].view(C, -1).T # (HxW, C) - flattened normalized features
1050
1051     # Sample pixels from each region (max 256 samples per region)
1052     anchor_fg = sample_pixels(feat, F_fg[b], num=256) # Failure foreground
1053     anchor_bg = sample_pixels(feat, F_bg[b], num=256) # Failure background
1054     pos_fg = sample_pixels(feat, T_fg[b] | S_fg[b], num=256) #Correct foreground
1055     pos_bg = sample_pixels(feat, T_bg[b] | S_bg[b], num=256) #Correct background
1056
1057     # Intra-class: Pull failures toward correct same-class features (InfoNCE)
1058     # For foreground: anchor_fg -> pos_fg
1059     sim_pos = (anchor_fg @ pos_fg.T) / tau # (NxM)
1060     sim_all = (anchor_fg @ feat.T) / tau # (NxHxW)
1061     L_intra_fg = -mean(logsumexp(sim_pos, dim=1) - logsumexp(sim_all, dim=1))
1062     # Similarly for background: L_intra_bg
1063
1064     # Inter-class: Push failures away from opposite-class features
1065     # For foreground failures -> background regions
1066     sim_opposite = (anchor_fg @ pos_bg.T) / tau
1067     L_inter_fg = mean(log(1 + exp(sim_opposite)).sum(dim=1))
1068     # Similarly for background: L_inter_bg
1069
1070     # Self-improvement: Guide failures toward success regions
1071     success_feat = sample_pixels(feat, S_fg[b] | S_bg[b], num=512)
1072     failure_feat = sample_pixels(feat, F_fg[b] | F_bg[b], num=512)
1073     sim_success = (failure_feat @ success_feat.T) / tau
1074     sim_all = (failure_feat @ feat.T) / tau
1075     L_self = -mean(logsumexp(sim_success, dim=1) - logsumexp(sim_all, dim=1))
1076
1077 # Step 4: Combine losses
1078 L_CMRL = 0.4*L_intra + 0.4*L_inter + 0.2*L_self
1079

```

C IN-DEPTH ANALYSIS OF PHOENIX COMPONENTS

To provide a comprehensive understanding of Phoenix’s design decisions and component contributions, we conduct detailed analyses that extend beyond the main paper’s ablation studies. Our evaluation protocol maintains consistency with the main paper by employing two distinct noise quality scenarios for instance segmentation: PointWSSIS Kim et al. (2023) with 1% supervision (denoted as AP¹) representing challenging refinement scenarios, and Mask R-CNN He et al. (2017) with ResNet-50 He et al. (2016) backbone (denoted as AP²) representing moderate refinement challenges. For fine-grained segmentation analysis, we evaluate performance on DIS-UNet Ronneberger et al. (2015) outputs (denoted as IoU¹) representing challenging refinement scenarios and DIS-ISNet Shen et al. (2022) outputs (denoted as IoU²) representing moderate refinement challenges, with results averaged across all DIS test splits Qin et al. (2022).

1080 C.1 COMPUTATIONAL EFFICIENCY ANALYSIS

1081 Table 8 provides a detailed analysis of Phoenix’s computational requirements. Our approach balances
1082 high performance with efficiency across both training and inference stages.

1083 **Training Efficiency:** The training process completes in under 10 hours on 8 V100 GPUs, fine-tuning
1084 only 4M trainable parameters (0.6% of the full model) and consuming approximately 13.5GB of
1085 memory per GPU. This efficiency is achieved by our strategic design choice to fine-tune only the
1086 lightweight decoder while keeping the heavy encoder frozen.

1087 In addition, we investigate the computational cost of our adversarial mask perturbation (AMP) process
1088 during training. As shown in Table 8a, the AMP Time represents the average time required to generate
1089 a single noisy mask through adversarial perturbation. Each perturbation embedding update operation
1090 requires approximately 6ms on a V100 GPU, and we found that an average of 4.6 updates are needed
1091 to satisfy our IoU threshold requirements for each noisy mask generation. Therefore, the complete
1092 adversarial perturbation process takes approximately 27ms ($4.6 \times 6\text{ms}$) per mask, which is remarkably
1093 efficient considering the semantic complexity of the generated noise patterns.

1094 **Inference Efficiency:** The mask refinement time cost shown in Table 8b represents the total refine-
1095 ment (inference) time for processing the COCO train5K dataset, which contains approximately 5K
1096 images and 37K masks. Phoenix demonstrates significant performance advantages while maintaining
1097 computational efficiency comparable to SAM-based methods. Both Phoenix and SAMRefiner require
1098 0.6 hours for processing, as they share the same underlying SAM network architecture, while SegRe-
1099 finer requires 1.4 hours. However, Phoenix achieves substantially higher performance (28.7% AP¹)
1100 compared to both SAMRefiner (21.8% AP¹) and SegRefiner (14.7% AP¹), demonstrating superior
1101 efficiency in terms of performance per computational cost.

1102 This computational efficiency stems from Phoenix’s architectural advantage inherited from SAM,
1103 particularly the ability to reuse image embeddings in the lightweight decoder. Once the heavy
1104 encoder processes an image to generate embeddings, the lightweight decoder can efficiently refine
1105 multiple masks from the same image by reusing these precomputed embeddings. This design is highly
1106 advantageous for scenarios where a single image contains multiple masks, such as conventional
1107 instance and semantic segmentation tasks, where the computational cost scales primarily with the
1108 number of masks rather than images.

1109 C.1.1 EFFECT OF IOU THRESHOLDS

1110 Table 8c examines how different IoU threshold ranges during training affect Phoenix’s performance.
1111 The full range $\mathcal{U}(0.3, 0.9)$ yields the best overall performance, providing exposure to diverse noise lev-
1112 els from severe perturbations (IoU around 0.3) to subtle distortions (IoU around 0.9). Narrower ranges
1113 limit the model’s exposure to the full spectrum of noise patterns, resulting in reduced performance,
1114 particularly on fine-grained segmentation tasks.

1115 C.1.2 EFFECT OF GUIDANCE MASKS

1116 Table 8d analyzes how different guidance mask configurations influence refinement performance. The
1117 results show that each guidance mask type contributes uniquely to overall performance. Individual
1118 expansion or contraction guides provide moderate improvements, while inversion alone performs
1119 better than either. Combining expansion and contraction yields strong results (28.5% AP¹ and 46.8%
1120 AP²), and the full combination achieves optimal performance (28.7% AP¹ and 46.9% AP²). This
1121 analysis confirms the importance of diverse error patterns for robust refinement performance. By
1122 exposing the model to complementary types of errors, we enable it to handle various refinement
1123 scenarios effectively.

1124 C.1.3 IMPACT OF MORPHOLOGICAL NOISE INTEGRATION

1125 We investigate the integration of classical morphological operations with our adversarial perturbation
1126 approach to determine the optimal noise generation strategy for mask refinement training. To this end,
1127 we introduce a probability parameter p_{morph} that controls the random replacement of adversarial
1128 noise masks with morphological noise masks during training, where $p_{morph} = 0.0$ corresponds to
1129 using purely adversarial noise and $p_{morph} = 1.0$ corresponds to using only morphological noise. As

Table 8: **Ablation study** using instance segmentation results on 1% PointWSSIS (AP¹) and MRCNN with a ResNet-50 backbone (AP²) and fine-grained segmentation results on DIS-TE1-UNet (IoU¹) and DIS-TE1-ISNet (IoU²).

(a) Efficiency of Phoenix					(b) Mask refinement time cost				
Specification		Value			Method		AP ¹ Time (h)		
Training Time		≤ 10h			SegRefiner	14.7	1.4		
Training Params		4M (0.6%)			SAMRefiner	21.8	0.6		
GPU Memory Usage		≈ 13.5 GB			Phoenix (ours)	28.7	0.6		
AMP Time		≈ 27 ms							

(c) τ (IoU Threshold)					(d) \mathcal{M}_g (Guidance Mask)				
τ	AP ¹	AP ²	IoU ²	IoU ²	Expansion	Contraction	Inversion	AP ¹	AP ²
$\mathcal{U}(0.3, 0.9)$	28.7	46.9	75.7	77.1	✓	✗	✗	26.9	46.1
$\mathcal{U}(0.5, 0.9)$	28.1	47.0	74.7	77.2	✗	✓	✗	27.1	46.2
$\mathcal{U}(0.3, 0.7)$	28.8	46.4	75.9	75.8	✗	✗	✓	28.0	46.5
$\mathcal{U}(0.5, 0.7)$	28.3	46.6	74.9	76.9	✓	✓	✗	28.5	46.8
					✓	✓	✓	28.7	46.9

(e) p_{morph}		(f) Encoder Study					
p_{morph}	AP ¹	Encoder	AP ¹	AP ²	VRAM (GB)	GFlops (G)	FPS
0.0	28.7	EfficientViT-XL1	28.0	45.9	0.8	323	45.2
0.1	28.6	ViT-B	22.4	42.0	2.8	369	12.4
0.3	27.7	ViT-L	26.6	46.3	4.1	1313	5.6
0.5	26.1	ViT-H	28.7	46.9	4.9	2735	3.5
1.0	23.8						

(g) Projector g		
Projector	AP ¹	AP ²
Identity	27.1	46.2
1-layer MLP	27.9	46.5
2-layer MLP	28.5	46.7
3-layer MLP	28.7	46.9
4-layer MLP	28.6	46.8

shown in the Table 8e, performance consistently degrades with increasing p_{morph} values, declining from 28.7% AP¹ when using purely adversarial noise ($p_{morph} = 0.0$) to 23.8% AP¹ when relying exclusively on morphological noise ($p_{morph} = 1.0$). This demonstrates that traditional morphological operations are insufficient for generating realistic noise patterns. These results underscore the importance of our semantic-aware adversarial noise modeling, which captures more nuanced and contextually relevant perturbations compared to the geometric transformations provided by classical morphological operations, ultimately leading to more effective mask refinement capabilities.

C.1.4 EFFECT OF ADVERSARIAL MASK PERTURBATION PARAMETERS

We analyze how various parameters affect our adversarial perturbation process to provide practical guidance for implementation. Figure 5b examines Phoenix’s sensitivity to the initial step size (α_0) and maximum iteration count (N), revealing remarkable stability across a wide range of values with AP¹ consistently above 28.0% for $N \geq 5$ and $\alpha_0 \in [0.001, 0.1]$. The optimal configuration is achieved at $\alpha_0 = 0.01$ and $N = 10$, while extreme values (very small α_0 or large α_0 with small N) should be avoided as they lead to insufficient or unrealistic perturbations.

Similarly, Figure 5c shows the impact of perturbation embedding count (P), with performance improving substantially as P increases from 1 to 50 (AP¹ rising from 26.7% to 29.0%), then plateauing beyond this point. This suggests that $P = 50$ provides an optimal balance between effectiveness and efficiency, capturing the necessary perturbation patterns without excessive computational overhead.

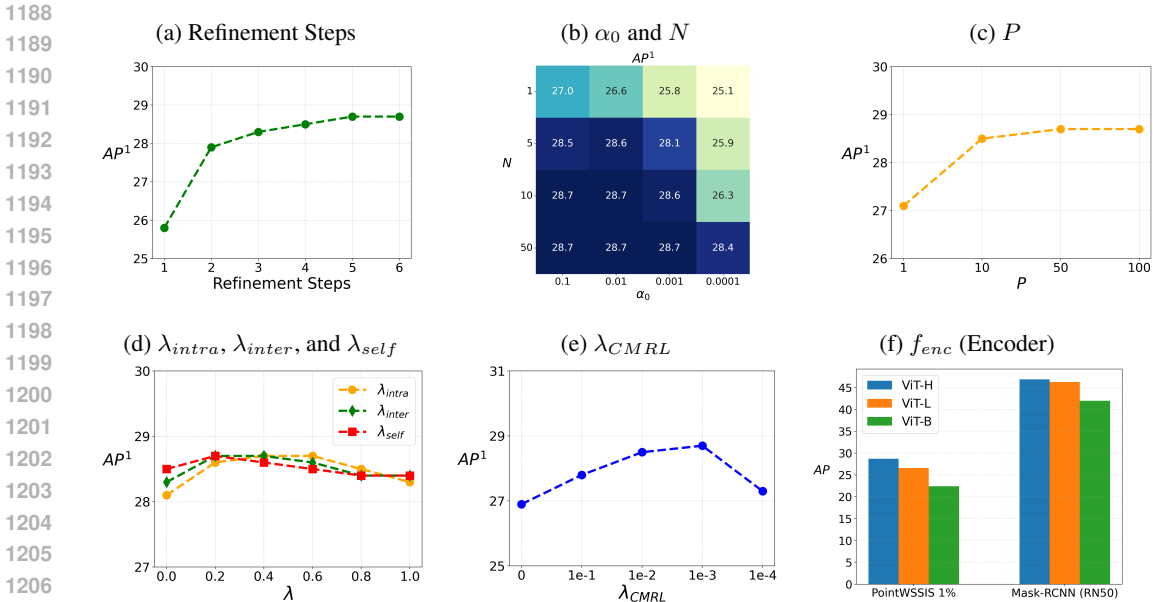


Figure 5: Ablation study (continue) using instance segmentation AP¹ results.

Through this analysis, we establish that once these parameters (*i.e.*, α_0 , N , and P) are set to reasonable values, the IoU threshold parameter τ becomes the primary parameter to control the perturbation intensity without sophisticated parameter tuning.

C.1.5 EFFECT OF CONTRASTIVE MASK REFINEMENT LEARNING (CMRL) LOSS WEIGHTS

We conduct a comprehensive analysis of how CMRL component weights influence Phoenix’s performance to optimize our contrastive learning strategy. The total training loss of Phoenix combines the original SAM loss function with our contrastive loss: $\mathcal{L} = \mathcal{L}_{SAM} + \lambda_{CMRL} \cdot \mathcal{L}_{CMRL}$, where \mathcal{L}_{SAM} is the linear combination of Focal Lin et al. (2017) and Dice Milletari et al. (2016) loss in a 20:1 ratio.

Figure 5d reveals how individual contrastive components (λ_{intra} for intra-class consistency, λ_{inter} for inter-class contrast, and λ_{self} for self-improvement regularization) affect refinement quality. We systematically evaluate each component by fixing our default configuration ($\lambda_{intra} = 0.4$, $\lambda_{inter} = 0.4$, $\lambda_{self} = 0.2$) and varying each weight individually from 0.0 to 1.0. All three components exhibit similar inverted U-shaped performance curves, with effectiveness peaking around 0.4 before gradually declining with higher values. This pattern suggests that while each contrastive component provides valuable learning signals, excessive emphasis on any single component can undermine the primary segmentation objective. The optimal balance places greater emphasis on intra-class consistency and inter-class contrast compared to self-improvement regularization, indicating that feature space organization (establishing clear boundaries between foreground and background while ensuring consistency within each class) is particularly crucial for effective refinement.

Meanwhile, as shown in Figure 5e, the overall CMRL scaling factor λ_{CMRL} reaches optimal performance at the relatively small value of 1×10^{-3} , achieving an AP¹ of 28.7%. Performance degrades significantly with higher values (dropping to 27.3% at $\lambda_{CMRL} = 0.1$), confirming that contrastive learning should complement rather than dominate the training process.

These findings provide important practical guidance: while our tri-directional contrastive approach significantly enhances mask refinement, it requires careful integration with conventional segmentation objectives. The sensitivity analysis demonstrates that a light touch with contrastive learning, applying it as a targeted enhancement to traditional segmentation loss, yields the most robust and effective refinement model.

1242 C.2 EFFECT OF ENCODER ARCHITECTURE

1243
1244 To provide guidance for practical deployment across different computational constraints, we investi-
1245 gate how various ViT [Dosovitskiy et al. \(2021\)](#) encoder architectures impact Phoenix’s performance.
1246 This analysis is important for understanding the trade-offs between model capacity and refinement
1247 quality, helping practitioners select the appropriate configuration for their specific requirements.

1248 During training, we freeze the image encoder and fine-tune only the lightweight decoder with
1249 mini-batch size 2 (total batch size 16 with 8 GPUs) using FP16 precision on V100 GPUs. In this
1250 setting, we measure the latency and VRAM of the image encoder using mini-batch size 2 inputs
1251 ($2 \times 3 \times 1024 \times 1024$). As shown in Table 8f and Figure 5f, we evaluate both standard ViT variants and
1252 the lightweight EfficientViT [Zhang et al. \(2024\)](#).

1253 The results demonstrate clear trade-offs between model capacity, performance, and computational
1254 efficiency. ViT-H achieves the best refinement quality (28.7% AP¹ and 46.9% AP²), leveraging its
1255 larger capacity to capture richer image features, but requires substantial computational resources
1256 (4.9GB VRAM, 3.5 FPS). ViT-L shows minimal performance degradation while offering improved
1257 efficiency. The smaller ViT-B encoder shows more substantial drops (22.4% AP¹ and 42.0% AP²)
1258 but significantly reduces computational requirements (2.8GB VRAM, 12.4 FPS). Most notably,
1259 EfficientViT-XL1 provides an excellent accuracy-speed trade-off with competitive performance
1260 (28.0% AP¹ and 45.9% AP²) while achieving remarkable efficiency (45.2 FPS and minimal 0.8GB
1261 VRAM usage), offering flexible deployment options for different computational constraints. The
1262 per-GPU computational requirements are highly efficient across all variants, which is much lower
1263 than typical memory-intensive vision models.

1264 C.3 QUALITATIVE NOISE PATTERN ANALYSIS

1265
1266 A critical component of Phoenix is our adversarial mask perturbation (AMP) mechanism, which
1267 generates realistic training noise that closely approximates real-world segmentation failures. Figure 8
1268 provides additional noisy mask samples of our adversarial masks compared to traditional morpholog-
1269 ical noise masks. Unlike morphological operations that apply predictable geometric transformations
1270 directly to mask pixels, our approach generates semantically-aware perturbations by embedding
1271 space adversarial attacks. This enables the generation of diverse failure patterns, including contextual
1272 confusion, over- or under-segmentation errors, and boundary imprecision.

1273 The visual comparison in Figure 8 reveals that our adversarial masks closely resemble real seg-
1274 mentation failures from production models (1% PointWSSIS for instance segmentation, DIS-UNet
1275 for fine-grained segmentation), while morphological masks show limited, uniform patterns. This
1276 improved realism directly contributes to Phoenix’s superior mask refinement performance.
1277

1278 D ADDITIONAL QUALITATIVE RESULTS

1279
1280 We present comprehensive qualitative results to demonstrate Phoenix’s effectiveness across diverse
1281 segmentation tasks and challenging scenarios. These visual comparisons complement our quantitative
1282 analyses and provide insights into the specific types of improvements Phoenix achieves.
1283

1284 Figure 9 provides extensive qualitative comparisons on the LVIS dataset for instance segmentation
1285 tasks, showing Phoenix’s performance relative to SegRefiner and SAMRefiner across diverse object
1286 categories. The examples span various scales, from small objects to large objects, demonstrating
1287 Phoenix’s robust refinement capabilities. Particularly notable are the improvements in handling
1288 complex poses and occlusions, as seen in the human figures and animal examples. Phoenix consis-
1289 tently produces more accurate boundary delineations and better preserves object details compared to
1290 existing methods.

1291 Figure 10 shows additional qualitative results on the DIS benchmark, comparing Phoenix with
1292 SegRefiner and SAMRefiner on fine-grained segmentation tasks. The examples highlight Phoenix’s
1293 superior ability to handle complex structures and intricate boundaries. Namely, Phoenix successfully
1294 refines the delicate feather structures of various objects, preserving fine details that other methods
1295 struggle to capture. The bicycle examples demonstrate Phoenix’s effectiveness in handling complex
mechanical structures with thin components like spokes and cables. For architectural elements

1296 like windmills and playground equipment, Phoenix maintains structural integrity while refining
1297 boundary precision. The dinosaur skeleton example illustrates Phoenix’s capability with highly
1298 detailed, branching structures that require precise boundary delineation.

1299 These qualitative results consistently demonstrate Phoenix’s superior refinement quality across
1300 diverse scenarios, object categories, and segmentation tasks. The visual improvements align with
1301 our quantitative findings, confirming Phoenix’s effectiveness as a general-purpose mask refinement
1302 approach that can enhance segmentation quality across various application domains.
1303

1304 1305 1306 E FAILURE CASE ANALYSIS 1307

1308
1309 Despite Phoenix’s strong overall performance, we conduct a comprehensive analysis of challenging
1310 scenarios where our method faces limitations. Figure 11 presents representative failure cases orga-
1311 nized into three categories that highlight the current boundaries of our approach and provide insights
1312 for future improvements.

1313 **Ambiguous Target Objects:** This failure mode occurs when the target object itself is inherently
1314 difficult to distinguish or when multiple similar objects create boundary confusion. In Figure 11a
1315 first row, the unclear target object mask between two occluded giraffes presents an ambiguous
1316 segmentation scenario where even determining the correct target is challenging. The spatial overlap
1317 and visual similarity between the two giraffes make it difficult to establish which object should be
1318 segmented. Similarly, in Figure 11b second row, the noisy mask contains multiple objects, creating
1319 confusion about whether to remove or maintain the pot in the final segmentation. Such ambiguous
1320 scenarios arise when the input contains insufficient contextual information to resolve target object
1321 identity, making the refinement task inherently ill-defined.

1322 **Totally Mislocalized Input Noisy Masks:** This represents the most severe failure mode where initial
1323 noisy masks are completely spatially displaced from the actual target objects. Examples include
1324 Figure 11a second and third rows, and Figure 11b first row, where the noisy masks bear no spatial
1325 correspondence to the ground-truth mask locations. In these cases, the refinement problem becomes
1326 fundamentally unsolvable because there is no meaningful overlap or spatial relationship between
1327 the input mask and the actual object boundaries. This failure mode reveals the inherent limitation
1328 of mask refinement approaches that depend critically on the initial quality of noisy masks. These
1329 failures underscore that refinement-based methods have fundamental prerequisites regarding input
1330 mask quality and cannot recover from arbitrary initialization errors.

1331 **Class-Agnostic Refinement Limitations:** Phoenix’s class-agnostic design, while enabling broad
1332 generalization, creates limitations when semantic understanding is required for proper refinement.
1333 Figure 11c demonstrates cases where Phoenix cannot identify or correct misclassified masks. In the
1334 first row, Phoenix cannot separate the person and bike from the merged mask because it lacks semantic
1335 understanding to distinguish between different object classes that have been incorrectly combined.
1336 The method treats the merged region as a single entity and refines its boundaries accordingly, without
1337 recognizing that it should be decomposed into separate semantic categories. In the second row,
1338 Phoenix cannot refine the misclassified region where the sail is incorrectly labeled as a person.
1339 While Phoenix successfully improves the mask’s spatial quality and boundary precision, it maintains
1340 the fundamental semantic error because our architecture focuses exclusively on visual boundary
1341 refinement without incorporating class-conditional reasoning.

1342 **Future Research Directions:** These failure modes directly connect to the research directions
1343 identified in our main paper analysis.

1344 (1) For ambiguous target objects, the integration of visual cues provides a promising solution
1345 by guiding the target object to be segmented. As demonstrated in Table 4f of our main paper,
1346 incorporating visual prompts (points or boxes) that provide explicit target object information achieves
1347 remarkable performance improvements. This approach directly addresses the target specification
1348 challenge by providing clear geometric guidance about which object should be the focus of refinement,
1349 effectively resolving ambiguity in multi-object scenarios like the occluded giraffes or mixed object
cases.

1350 (2) For totally mislocalized masks, visual prompts can also provide spatial anchoring, though the
1351 effectiveness depends on the degree of misalignment. When combined with additional spatial
1352 reasoning mechanisms, visual cues could help establish the correct target location even when initial
1353 masks are severely displaced.

1354 (3) For class-agnostic refinement limitations, our main paper identifies incorporating open vocabulary
1355 models [Ghiasi et al. \(2022\)](#); [Liang et al. \(2023\)](#) or text embeddings [Li et al. \(2023b\)](#); [Radford et al.](#)
1356 [\(2021\)](#) as a promising research direction to handle misclassified object masks. Since Phoenix operates
1357 in a class-agnostic manner and cannot correct semantic class errors inherent from real models,
1358 integrating vision-language understanding could enable the system to recognize and correct semantic
1359 misclassifications while maintaining spatial refinement capabilities. This multimodal extension
1360 would allow Phoenix to leverage both visual boundary information and semantic understanding,
1361 potentially resolving the sail-as-person misclassification by incorporating textual or semantic priors
1362 that distinguish between different object categories.

1363 These insights highlight both the current capabilities and future potential of mask refinement systems,
1364 pointing toward more robust, multimodal approaches that can handle increasingly diverse and
1365 challenging segmentation scenarios through the integration of visual prompts for target specification
1366 and semantic understanding for class-aware refinement.

1368 F BROADER IMPACTS

1370 Phoenix offers significant potential to benefit various domains by improving segmentation quality
1371 across a wide range of applications. High-quality segmentation is foundational to many computer
1372 vision tasks, and our refinement approach provides substantial improvements without requiring
1373 architectural changes to base models or extensive retraining. The demonstrated ability to enhance
1374 both existing annotations and model outputs suggests Phoenix could reduce manual effort in dataset
1375 creation while improving the performance of downstream applications that rely on precise segmenta-
1376 tion. By focusing on a refinement paradigm that builds upon existing segmentation methods, Phoenix
1377 complements rather than replaces current approaches, allowing for integration into established
1378 workflows.

1379 While any advanced technology carries some responsibility for appropriate use, we have designed
1380 Phoenix to be broadly applicable to beneficial applications across diverse domains. We are committed
1381 to making this technology available to the research community upon publication to encourage further
1382 innovation and refinement of segmentation capabilities.

1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

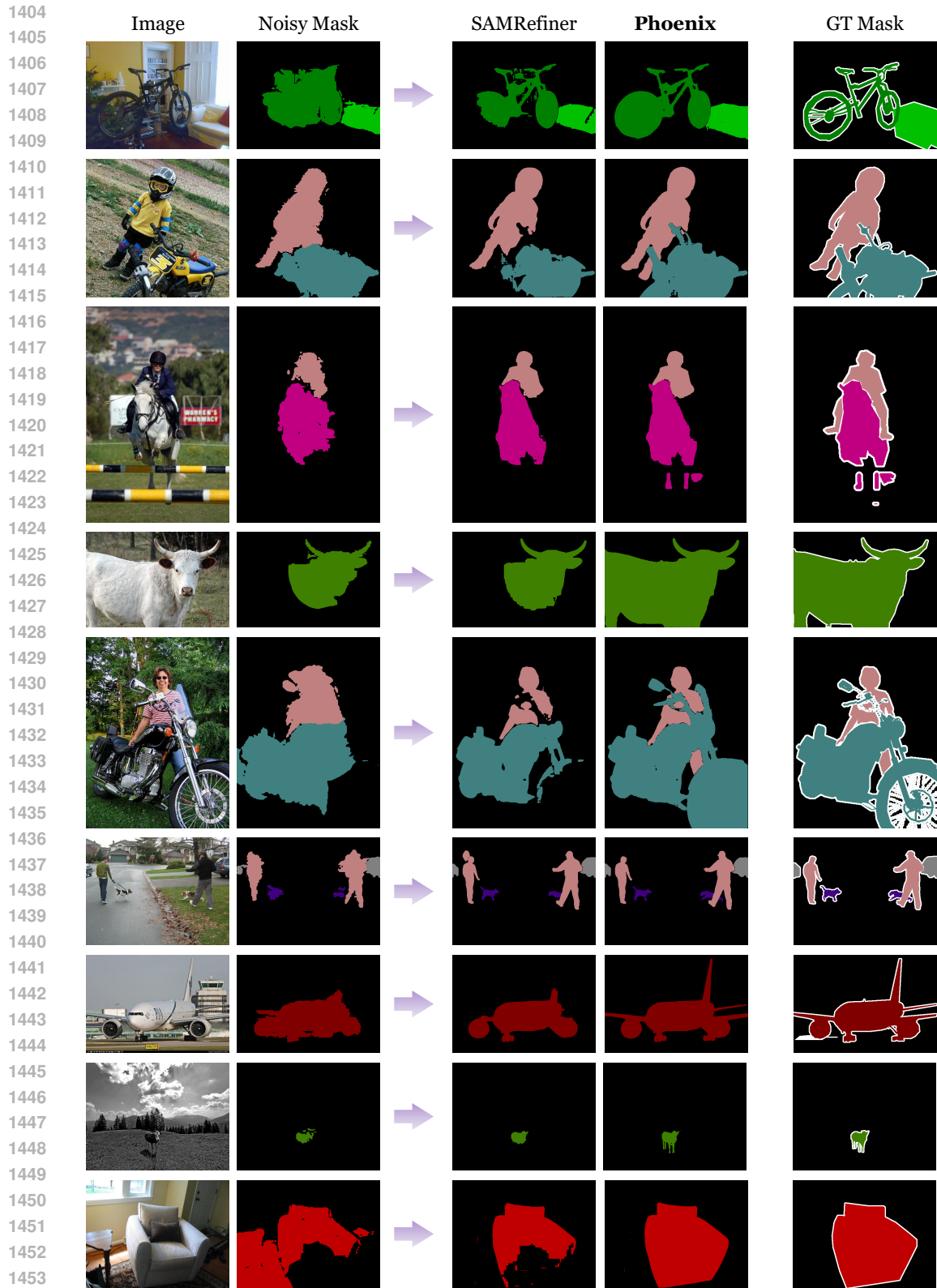
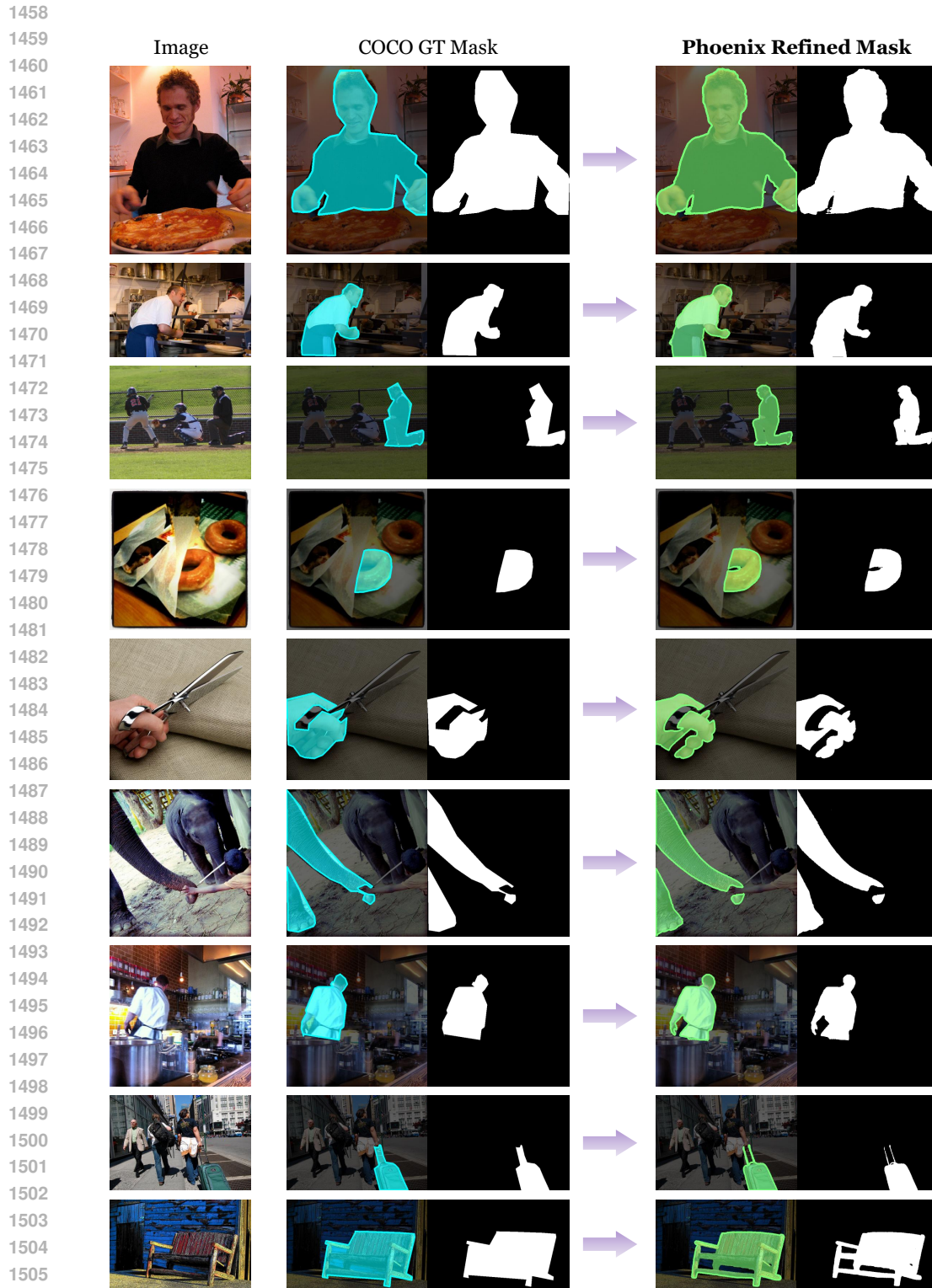
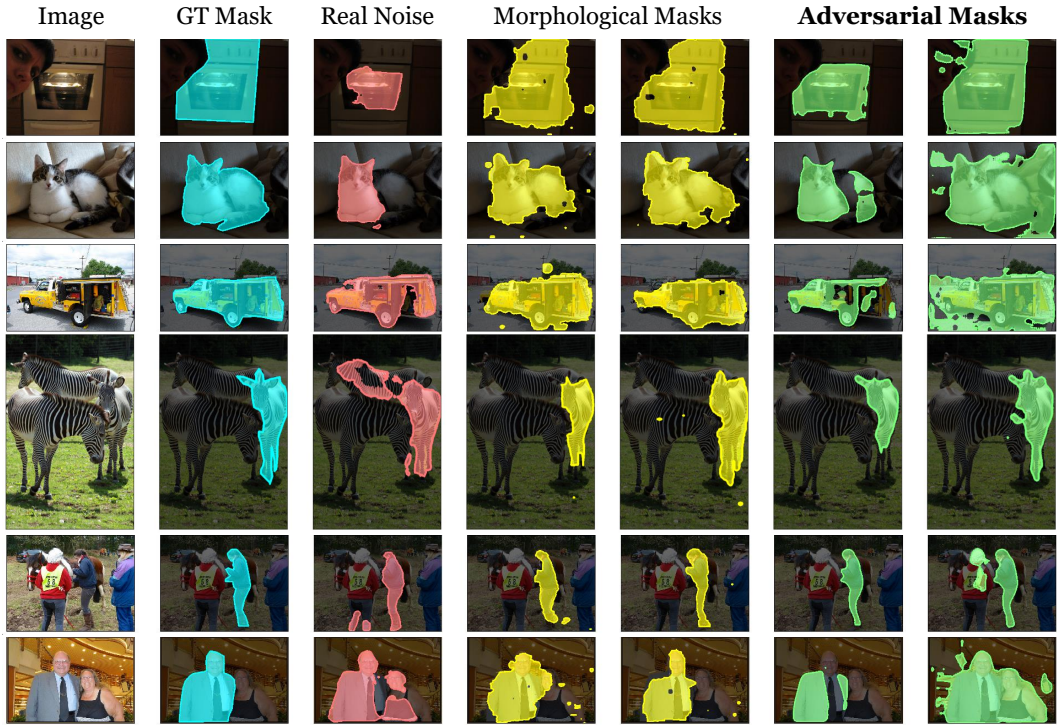


Figure 6: **Qualitative results for semantic segmentation refinement.** Examples show Phoenix’s effectiveness in refining coarse semantic masks across diverse scene categories. The method successfully corrects both over-segmentation and under-segmentation errors while maintaining semantic consistency, particularly excelling at complex boundary regions and multi-object scenarios.

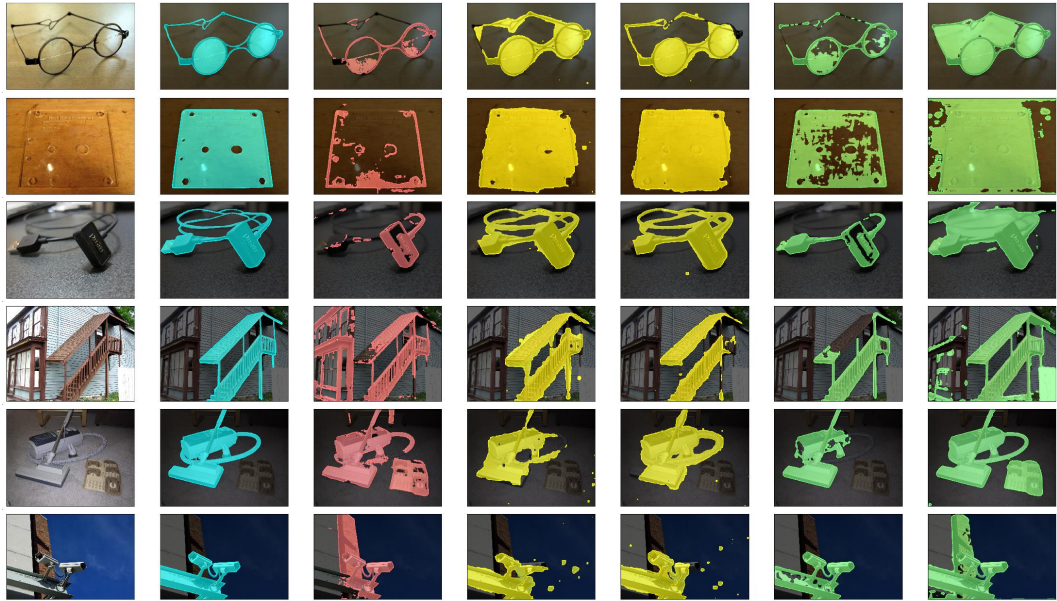


1507 **Figure 7: Qualitative results for human annotation refinement** for COCO 2017 *val* dataset. The
 1508 left columns show the original COCO ground truth masks. The right columns display Phoenix-refined
 1509 refined masks that achieve better alignment with high-quality mask annotations. Phoenix effectively corrects
 1510 annotation imperfections, particularly improving boundary precision and handling of fine details that
 1511 are often missed in standard polygon-based annotation workflows.

1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565



(a) Instance Segmentation on COCO train5K. Real Noisy Mask from 1% PointWSSIS



(b) Fine Segmentation on DIS5K. Real Noisy Mask from DIS-UNet

Figure 8: **Noise pattern analysis** demonstrating the superiority of adversarial mask perturbation over morphological methods. Our approach produces diverse, realistic failure patterns, including contextual errors, segmentation inconsistencies, and boundary imprecision that closely match real segmentation model failures, leading to improved refinement performance.

1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1600
 1601
 1602
 1603
 1604
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619



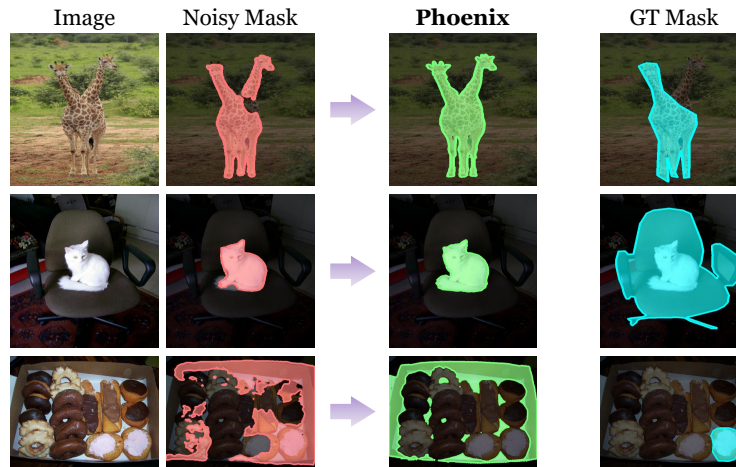
Figure 9: **Additional qualitative results for instance segmentation refinement.** Each row shows progression from noisy input masks through different refinement methods (SegRefiner Wang et al. (2023), SAMRefiner Lin et al. (2025)) to Phoenix’s output and ground truth. Phoenix consistently produces more accurate boundary delineation and better handling of complex object structures across diverse object categories.

1620
 1621
 1622
 1623
 1624
 1625
 1626
 1627
 1628
 1629
 1630
 1631
 1632
 1633
 1634
 1635
 1636
 1637
 1638
 1639
 1640
 1641
 1642
 1643
 1644
 1645
 1646
 1647
 1648
 1649
 1650
 1651
 1652
 1653
 1654
 1655
 1656
 1657
 1658
 1659
 1660
 1661
 1662
 1663
 1664
 1665
 1666
 1667
 1668
 1669
 1670
 1671
 1672
 1673

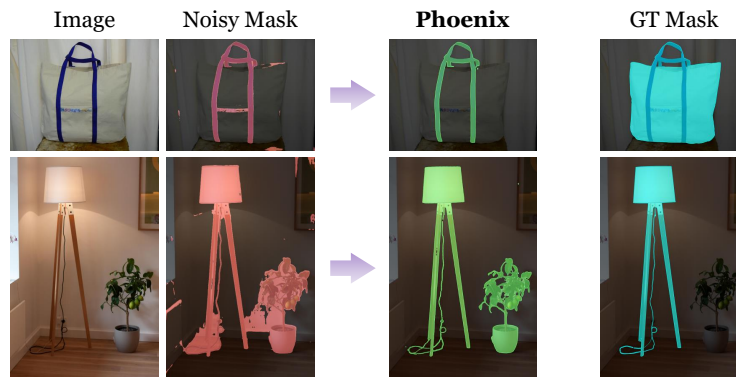


Figure 10: **Additional qualitative results for fine-grained segmentation refinement.** Examples demonstrate Phoenix’s capability to handle intricate object boundaries and thin structures across various categories. Phoenix shows superior boundary precision compared to baseline methods, particularly for objects with complex geometric features.

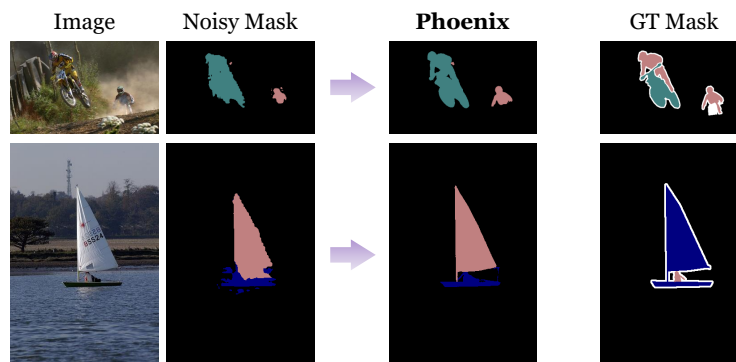
1674
 1675
 1676
 1677
 1678
 1679
 1680
 1681
 1682
 1683
 1684
 1685
 1686
 1687
 1688
 1689
 1690
 1691
 1692
 1693
 1694
 1695
 1696
 1697
 1698
 1699
 1700
 1701
 1702
 1703
 1704
 1705
 1706
 1707
 1708
 1709
 1710
 1711
 1712
 1713
 1714
 1715
 1716
 1717
 1718
 1719
 1720
 1721
 1722
 1723
 1724
 1725
 1726
 1727



(a) Instance Segmentation



(b) Fine Segmentation



(c) Semantic Segmentation

Figure 11: **Failure case analysis across different segmentation tasks.** (a) Instance segmentation failures include heavily occluded objects (giraffe), completely mislocalized masks (cat and chair), and ambiguous multi-object boundaries (food items). (b) Fine-grained segmentation failures involve complex geometric structures and spatial misalignment issues. (c) Semantic segmentation failures demonstrate class-agnostic limitations where Phoenix refines mask quality but cannot correct semantic misclassifications (sail classified as person). These cases highlight current method boundaries and inform future research directions.