

CoT-SEG: RETHINKING SEGMENTATION WITH CHAIN-OF-THOUGHT REASONING AND SELF-CORRECTION

Shiu-hong Kao[†]

HKUST, NUS

shkao@u.nus.edu

Chak Ho Huang[†]

HKUST

chhuangab@connect.ust.hk

Huaiqian Liu[†]

HKUST

hliudj@connect.ust.hk

Yu-Wing Tai

Dartmouth College

yu-wing.tai@dartmouth.edu

Chi-Keung Tang

HKUST

cktang@cs.ust.hk

ABSTRACT

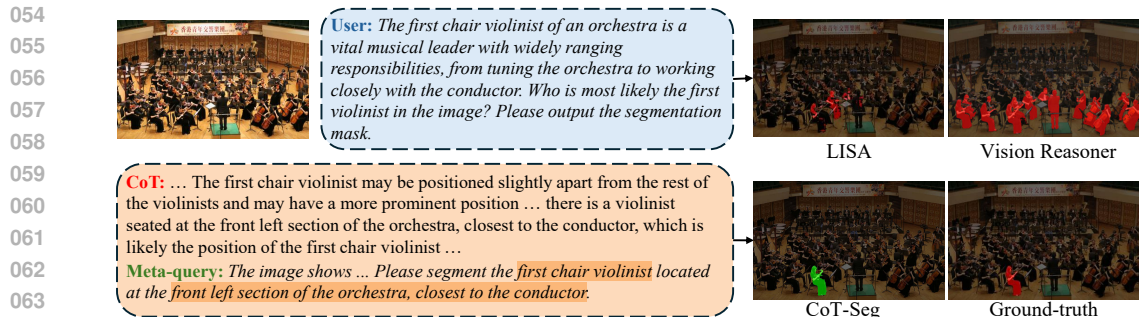
Existing works of reasoning segmentation often fall short in complex cases, particularly when addressing complicated queries and out-of-domain images. Inspired by the chain-of-thought reasoning, where harder problems require longer thinking steps/time, this paper aims to explore a system that can think step-by-step, look up information if needed, generate results, self-evaluate its own results, and recursively refine the results, in the same way humans approach harder questions. We introduce **CoT-Seg**, a framework that rethinks reasoning segmentation by combining **chain-of-thought reasoning** with recursive **self-correction**. Instead of fine-tuning, CoT-Seg leverages the inherent reasoning ability of pre-trained MLLMs (*e.g.*, GPT-4o) to decompose queries into meta-instructions, extract fine-grained semantics from images, and identify target objects even under implicit or complex prompts. Moreover, CoT-Seg incorporates a self-correction stage: the model evaluates its own segmentation against the original query and reasoning trace, identifies mismatches, and iteratively refines the mask. This tight integration of reasoning and correction significantly improves reliability and robustness, especially in ambiguous or error-prone cases. Furthermore, our CoT-Seg framework allows easy incorporation of retrieval-augmented reasoning, enabling the system to access external knowledge when the input lacks sufficient information. Our results highlight that combining chain-of-thought reasoning, self-correction, offers a powerful paradigm for vision language integration driven segmentation. Our project website is available at <https://danielshkao.github.io/cot-seg.html>.

1 INTRODUCTION

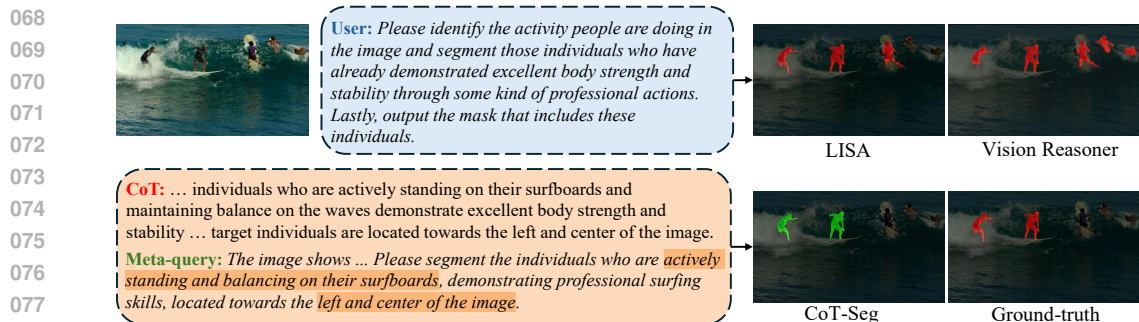
Reasoning segmentation represents a promising step toward vision-language integration, where a system generates a segmentation mask from complex and often implicit language queries. Recent progress has been driven by fine-tuning Multimodal Large Language Models (MLLMs), such as LISA (Lai et al., 2023), Seg-Zero (Liu et al., 2025a) and Vision Reasoner (Liu et al., 2025b), to produce segmentation outputs. Despite their success, these methods struggle with new cases that require nuanced reasoning, domain knowledge, or contextual inference which are the major challenges that humans naturally handle.

Consider the examples in Figures 1–2. Locating the first-chair violinist requires knowledge of orchestra seating arrangements, not just visual similarity. Differentiating surfers by posture demands reasoning about dynamic body positions. These examples highlight that stronger reasoning ability, together with mechanisms to evaluate and refine initial predictions, is essential for advancing reasoning segmentation.

[†]Equal contribution.



064 Figure 1: Finding the first violinist (concertmaster) is challenging among similar-looking musi-
065 cians. CoT-Seg reasons that they sit to the conductor’s left and generates a meta-query with relevant
066 *spatial* information, enabling more accurate segmentation than LISA and Vision Reasoner (No self-
067 correction was needed).



078 Figure 2: CoT-Seg reasons about the user’s query to segment surfers in the correct *pose*, capturing
079 only those who have popped up and are riding waves, unlike LISA and Vision Reasoner (No self-
080 correction was needed).

083 In this work, we present **CoT-Seg**, a training-free framework that rethinks reasoning segmentation
084 through the synergy of *chain-of-thought (CoT) reasoning* and *self-correction*. Rather than relying
085 on fine-tuning or additional training, CoT-Seg leverages the inherent reasoning ability of pre-trained
086 MLLMs (e.g., GPT-4o) to decompose queries into meta-instructions, extract detailed semantics,
087 and produce initial segmentation results. Crucially, CoT-Seg introduces a self-correction stage: the
088 model evaluates its own segmentation against the query and reasoning trace, identifies inconsistencies,
089 and refines the output through automatically generated meta-queries. This feedback loop allows
090 the system not only to think through the segmentation process but also to recognize and repair its
091 own mistakes.

092 Furthermore, we extend CoT-Seg with *retrieval-augmented reasoning*. When the query and image
093 lack sufficient information, CoT-Seg calls an external agent to retrieve relevant knowledge from the
094 web, integrating it into the reasoning process. This augmentation further strengthens its ability to
095 tackle ambiguous or knowledge-intensive cases.

096 Through extensive experiments on ReasonSeg and RefCOCO, we demonstrate that CoT-Seg subst-
097 antially outperforms existing methods. Our results show that integrating CoT reasoning, self-
098 correction provides a powerful paradigm for advancing reasoning-driven segmentation toward
099 human-level reliability.

102 2 RELATED WORK

104 **Image Segmentation and Reasoning Segmentation.** Image segmentation has evolved from early
105 graphical-model-based methods, such as Conditional Random Fields (CRFs) (Krähenbühl & Koltun,
106 2011; Chen et al., 2017) and region growing (Dias & Medeiros, 2019), to deep learning approaches
107 that utilize encoder-decoder architectures (Badrinarayanan et al., 2017), dilated convolutions (Yu
& Koltun, 2015), pyramid pooling (Zhao et al., 2017), and non-local operators (Liu et al., 2015).

Instance segmentation (He et al., 2017; Cheng et al., 2022) and panoptic segmentation (Kirillov et al., 2019; Cheng et al., 2020) further pushed the boundary to finer-grained understanding.

The emergence of foundation models for segmentation, especially the Segment Anything Model (SAM) (Kirillov et al., 2023), has revolutionized the field. By training on billions of masks and images, SAM enables promptable, zero-shot segmentation with multimodal inputs like points or bounding boxes. Leveraging SAM with Multimodal Large Language Models (MLLMs) has led to a new line of works on reasoning segmentation (Lai et al., 2024; Xia et al., 2024; Zhang et al., 2023a; He et al., 2024; Yao et al., 2025). These approaches generate segmentation masks conditioned on implicit or complex textual queries. However, combining MLLMs with SAM directly often fails in challenging scenarios, such as queries requiring domain knowledge, occluded objects, or intricate structures. In contrast, our work shows that *integrating chain-of-thought reasoning and self-correction* can improve robustness and accuracy in these difficult cases.

Chain-of-Thought Reasoning in LLMs and MLLMs. Chain-of-Thought (CoT) reasoning improves reasoning performance in large language models by decomposing complex tasks into intermediate steps (Wei et al., 2022; Wang et al., 2022a; Zhang et al., 2023b; Lyu et al., 2023; Kojima et al., 2022). While CoT has been extensively explored in text-only LLMs, its integration into Multimodal LLMs (MLLMs) is more challenging. Existing approaches often rely on fine-tuning MLLMs with multimodal CoT datasets (Mondal et al., 2024; Zhang et al., 2023c; Lu et al., 2022) or introducing intermediate representations like graphs (Mitra et al., 2024) or code (Surís et al., 2023), which limit accessibility and scalability.

Recent works highlight the potential of *test-time CoT reasoning* in pre-trained LLMs (Snell et al., 2025) and its applications in visual reasoning (Guo et al., 2022; Lian et al., 2023), robotics (Hu et al., 2023), and multimodal planning (Yao et al., 2025). Inspired by these trends, our framework leverages carefully designed CoT prompts in a *training-free manner*, enabling MLLMs to reason over images and textual queries, evaluate initial segmentation outputs, and self-correct without additional training.

Self-Correction and Retrieval-Augmented Reasoning. While CoT provides step-by-step reasoning, errors in initial predictions can propagate if unchecked. Recent studies in reasoning with feedback (Zhao et al., 2025; He et al., 2025) demonstrate that self-evaluation and iterative refinement improve accuracy. Our method explicitly incorporates a *self-correction loop* for reasoning segmentation, allowing the model to detect inconsistencies and refine segmentation masks.

Furthermore, retrieval-augmented reasoning (Lewis et al., 2021; Komeili et al., 2021) has shown that external knowledge can enhance reasoning when input information is incomplete. CoT-Seg integrates retrieval mechanisms to access relevant knowledge at test time, enabling more robust segmentation under ambiguous or knowledge-intensive queries.

Overall, our work is positioned at the intersection of *reasoning segmentation, CoT-enabled MLLMs, self-correction, and retrieval augmentation*, combining these advances into a unified, training-free framework that achieves state-of-the-art performance in complex vision-language tasks.

3 METHOD

Given an image $I \in \mathbb{R}^{3 \times H \times W}$ and a textual query q , reasoning segmentation aims to predict a binary mask M corresponding to the object(s) referred by q . CoT-Seg achieves this by combining chain-of-thought reasoning, self-correction, and optional retrieval-augmented reasoning in a multi-agent framework. The system consists of three collaborating agents: the MLLM *Reasoner*, the *Segmentation Agent*, and the *Evaluator*.

Fig. 3 gives an overview of CoT-Seg, where the Reasoner analyzes the image and query using a chain-of-thought (CoT) process, generating an explicit meta-query that guides the Segmentation Agent. The Segmentation Agent produces an initial mask using the meta-query and its supported input types, such as text, points, bounding boxes, or scribbles. The Evaluator then analyzes the predicted mask in combination with the original query and image, identifying errors and synthesizing refinement meta-queries for self-correction. This iterative loop allows CoT-Seg to achieve robust zero-shot reasoning segmentation without additional training.

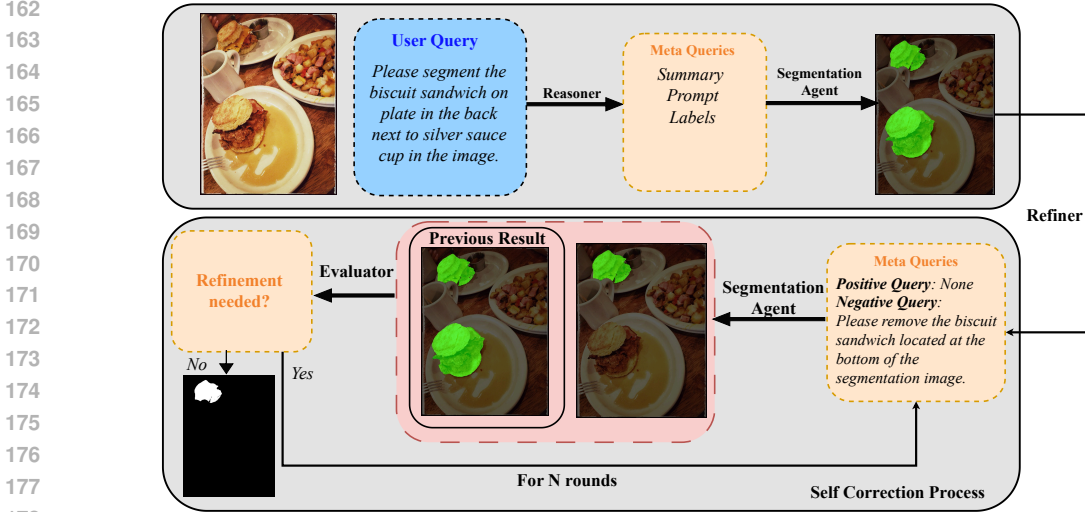


Figure 3: **Overview of CoT-Seg.** The pre-trained MLLM Reasoner generates a chain-of-thought (CoT) over the input image and query, producing an explicit meta-query that translates complex, implicit instructions into clear segmentation guidance. The Segmentation Agent predicts the initial mask, which is then optionally refined by the iterative refinement pipeline. The first-turn mask and original image are examined by the MLLM Evaluator which evaluates the mask and decide if any refinement is necessary. If it does require refinement, then it is passed onto the MLLM Refiner which produces two queries to correct for false positives and negatives. These queries are used inline with the segmentation agent to produce a refined mask for the next iteration of refinement.

3.1 MLLM REASONER

The Reasoner \mathcal{R} performs step-by-step chain-of-thought (CoT) reasoning to identify the target object(s) in the image. To achieve this, \mathcal{R} utilizes a series of *Question Proposers* that generate questions progressively from coarse to fine. Initially, coarse questions capture high-level scene context and object categories. Based on the answers, subsequent proposers generate finer-grained questions to localize the target objects, reasoning over attributes such as position, size, and relationships with other objects. This iterative process continues until sufficient information is collected to precisely identify the target or until it reaches max number of rounds.

Formally, each question-answer pair is generated autoregressively:

$$(Q_k, A_k) = \mathcal{R}(I, q, Q_{<k}, A_{<k}, \text{SegmentorCapabilities}), \quad k = 1, \dots, n, \quad (1)$$

where *SegmentorCapabilities* is defined as a textual description that informs the Reasoner of which input types the Segmentation Agent supports (e.g., text, points, bounding boxes, scribbles).

After completing all CoT steps, the Reasoner summarizes the collected information into a structured *meta-query* \tilde{q}_m , which is compatible with the Segmentation Agent aligning with *SegmentorCapabilities*. For non-textual inputs, such as points or scribbles, the meta-query is encoded in a JSON format specifying the input type, coordinates, and spatial attributes:

$$\tilde{q}_m = \mathcal{R}_{\text{summarize}}(\{Q_k, A_k\}_{k=1}^n, \text{SegmentorCapabilities}). \quad (2)$$

This structured meta-query is then passed to the Segmentation Agent to produce the initial mask, and subsequently to the Evaluator for self-correction if necessary. By combining coarse-to-fine question proposing with explicit summarization, the Reasoner ensures precise target localization and effective guidance for zero-shot segmentation.

3.2 REASONING SEGMENTATION AGENT

The Segmentation Agent \mathcal{A} predicts masks based on the meta-query \tilde{q}_m and its supported input types. It consists of a frozen vision encoder E , a mask decoder \mathcal{D} , and a vision-language model \mathcal{F}

for multimodal encoding e.g., (Lai et al., 2023; Zou et al., 2023b). The predicted mask is:

$$\hat{M} = \mathcal{A}(I, \tilde{q}_m) = \mathcal{D}(\mathcal{F}(I, \tilde{q}_m), E(I)). \quad (3)$$

By explicitly describing the segmentor’s input capabilities, both the Reasoner and Evaluator can adapt their CoT reasoning. If the segmentation agent cannot support a requested input type, the method may fail, highlighting the dependency on the segmentor’s flexibility. This design ensures that the meta-query generated by the Reasoner is always compatible with the segmentor.

3.3 EVALUATOR AND SELF-CORRECTION

The Evaluator \mathcal{J} assesses the quality of the mask generated by the Segmentation Agent and guides iterative refinement. It receives the original image I , the user query q , the predicted mask \hat{M} , and the SegmentorCapabilities as inputs. The Evaluator performs a chain-of-thought (CoT) reasoning process, similar to the Reasoner, to check whether the mask correctly covers the target objects and respects spatial and semantic constraints.

If refinement is needed, the Evaluator generates two types of meta-queries in a structured JSON format: \tilde{q}_P for false negatives and \tilde{q}_N for false positives. These queries specify the type of correction, spatial coordinates, and other relevant control signals compatible with the Segmentation Agent. Formally, the refinement process is:

$$S = \mathcal{J}_{\text{assess}}(I, \hat{M}, q, \text{SegmentorCapabilities}), \quad (4)$$

$$(\tilde{q}_P, \tilde{q}_N) = \mathcal{J}_{\text{refine}}(I, \hat{M}, q, S, \text{SegmentorCapabilities}), \quad (5)$$

$$s_P = \mathcal{A}(I, \tilde{q}_P), \quad s_N = \mathcal{A}(I, \tilde{q}_N), \quad (6)$$

$$s' = s + \gamma_P \cdot \text{ReLU}(s_P) - \gamma_N \cdot \text{ReLU}(s_N), \quad \hat{M}' = \{(i, j) \mid s'_{i,j} > 0\}, \quad (7)$$

where s implies the prediction score output by the segmentor, satisfying $\hat{M} = \{s_{i,j} \mid s_{i,j} > \text{threshold}\}$. This iterative self-correction loop continues until $S = 0$ (Correct Segmentation) or a maximum number of refinement rounds is reached. By using structured JSON communication, the Evaluator ensures compatibility with diverse Segmentation Agents and input modalities, enabling robust zero-shot segmentation with automated error correction. To ensure that \hat{M}' does not get worse than \hat{M} , which may also happen to humans after several refinement turns, \mathcal{J} will make a judgment whether to revert back to the previous segmentation \hat{M} as the chosen segmentation.

3.4 MULTIMODAL INPUT CONTROL

CoT-Seg supports diverse image-based controls in addition to textual queries, including points, bounding boxes, scribbles, and highlighted regions. The Reasoner \mathcal{R} is aware of the Segmentation Agent’s capabilities through the SegmentorCapabilities input. For non-textual inputs, it encodes the meta-query in JSON format specifying input type, coordinates, and spatial attributes. This allows both the Reasoner and Evaluator to generate compatible guidance and refinement instructions.

Given an image I and a control image I_{ann} , the Reasoner generates step-by-step CoT reasoning to interpret annotated regions and produce a meta-query \tilde{q}_m :

$$\tilde{q}_m = \mathcal{R}_{\text{summarize}}(\{Q_k, A_k\}_{k=1}^n, \text{SegmentorCapabilities}, I_{ann}), \quad (8)$$

which is then passed to the Segmentation Agent to produce the mask $\hat{M} = \mathcal{A}(I, \tilde{q}_m)$. The Evaluator can further refine the output via self-correction if necessary, using the same JSON format for multimodal control information.

3.5 RETRIEVAL-AUGMENTED REASONING

In cases where the input image and query do not provide sufficient information. CoT-Seg can optionally augment the Reasoner with an external retrieval step. Specifically, a Retrieval Agent is enabled to search for relevant information from the web or a knowledge database at the CoT step, which is then incorporated into the chain-of-thought reasoning.

270 For example, consider a query asking to segment a previously unknown object. The Reasoner might
271 lack sufficient internal knowledge to identify the individual. The Retrieval Agent searches for in-
272 formation about the person, such as reference images or textual descriptions, and provides these
273 as additional inputs to the Reasoner. The Reasoner then integrates the retrieved knowledge into its
274 CoT reasoning to generate a meta-query, e.g., specifying unique clothing, pose, or contextual cues,
275 which guides the Segmentation Agent to correctly segment the target. This mechanism allows CoT-
276 Seg to handle queries that require external or domain-specific knowledge, extending its reasoning
277 capabilities beyond the information present in the original input.

278 279 4 EXPERIMENTS

280 281 4.1 EXPERIMENTAL SETUP

282
283 We evaluate CoT-Seg on two recent and widely used benchmarks for reasoning segmentation: ReasonSeg (Lai et al., 2023) and RefCOCO (Kazemzadeh et al., 2014), which covers diverse objects
284 with compositional queries and fine-grained referring expressions. For fair comparison, we use publicly
285 released splits and follow the evaluation protocols of previous works (Lai et al., 2023; Xia et al.,
286 2024; Zhang et al., 2023a). We compare CoT-Seg against state-of-the-art reasoning segmentation
287 methods including LISA (Lai et al., 2023), GSVa (Xia et al., 2024), NextSeg (Zhang et al., 2023a),
288 and MultiSeg (He et al., 2024). As our method is training-free, we emphasize zero-shot evaluation to
289 highlight the effectiveness of inference-time reasoning and self-correction. Performance is measured
290 by Generalized Intersection-over-Union (gIoU) and Complete Intersection-over-Union (cIoU).
291

292 293 4.2 IMPLEMENTATION DETAILS

294 For the Reasoner and Evaluator modules, we use GPT-4o (Hurst et al., 2024) unless otherwise stated,
295 with system prompts tailored for CoT reasoning, summarization, and self-correction. The chain-of-
296 thought reasoning length is adaptively determined by the Reasoner, typically converging within
297 5–8 steps. The Segmentation Agent is instantiated with Vision Reasoner-7B (Liu et al., 2025b) by
298 default, though we also test compatibility with other SAM-based variants (Kirillov et al., 2023).
299 Structured communication between Reasoner, Evaluator, and Segmentation Agent is implemented
300 in JSON format to handle multimodal control inputs and to ensure capability alignment.

301 For retrieval-augmented reasoning, we employ a lightweight agent that queries the web using entity
302 names or context keywords extracted by the Reasoner. Retrieved data is passed back as either text
303 descriptions or reference images. To ensure reproducibility, all experiments are run on an NVIDIA
304 4090 GPU with 24GB memory, although the majority of reasoning computation occurs in the cloud-
305 hosted LLM.

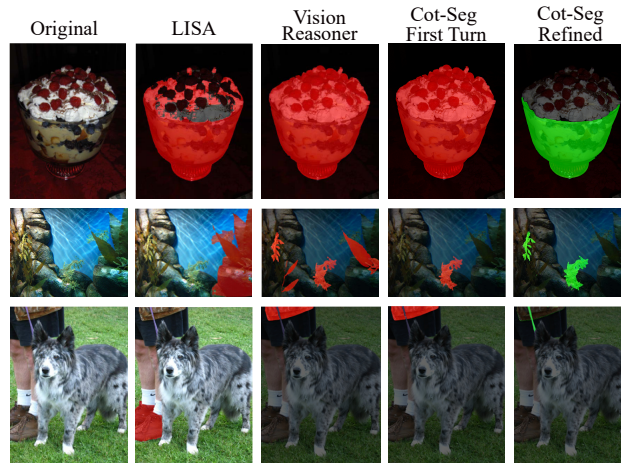
306 307 4.3 QUALITATIVE EVALUATION

308
309 We presented earlier qualitative comparisons in Figures 1–2. More results in Figure 4 demon-
310 strate how CoT-Seg progressively reasons about challenging queries and refines initial segmentation
311 masks, demonstrating CoT-Seg’s unique capabilities in: 1) resolving implicit queries with multi-step
312 reasoning; 2) correcting masks with fine-grained self-correction (e.g., removing false positives such
313 as ice cream and recovering missed objects like the leash in Figure 4); and 3) retrieval-augmented
314 reasoning for segmenting uncommon entities, such as identifying a new animal species (Figure 5)
315 by integrating retrieved textual and visual cues. These results show that CoT-Seg achieves higher
316 robustness in complex reasoning cases compared to prior methods that rely solely on direct prompt-
317 to-mask predictions.

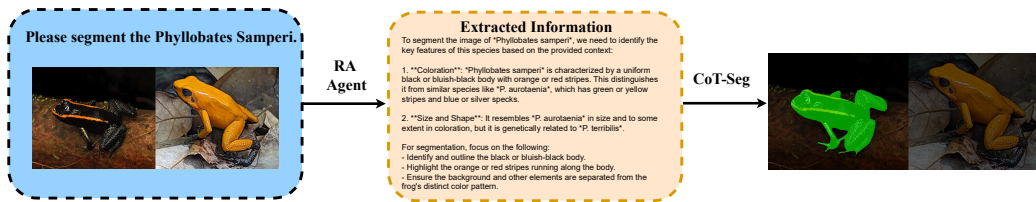
318 319 4.4 QUANTITATIVE EVALUATION

320 Tables 1 and 2 summarize quantitative comparisons across benchmarks. Overall, CoT-Seg achieves
321 SOTA or competitive results in both benchmarks, with the most improvements on ReasonSeg, where
322 high-level reasoning and domain knowledge are essential, while producing improved results after
323 self-correction. Notably, our training-free pipeline achieves higher cIoU than LISA, Seg-Zero and
Vision-Reasoner on RefCoco demonstrating the effectiveness of inference-time reasoning and self-

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338



339 Figure 4: Queries for each row are: 1. A fruit salad is a refreshing and delicious dessert that often
340 consists of a variety of fruits mixed together. What object in the picture could be used to hold and
341 serve such a dessert? 2. Please segment leafy sea dragons in this image.
342 3. What is the object that the person in the picture is holding onto while walking his dog?
343



344
345
346
347
348
349
350
351 Figure 5: A recently discovered species of frog unrecognizable to GPT-4o. With retrieval augmented
352 (RA) reasoning CoT-Seg was able to segment the frog based on its appearance descriptions from the
353 retrieval agent.

354
355 correction without requiring additional data or fine-tuning. CoT improves the general score while
356 auto-correction improves the result of hard tasks.

358 4.5 ABLATION STUDIES

360 **Impact of Self-Correction** Tables 1–2 compare performance with and without the refinement
361 module, showing how iterative refinement improves robustness in ambiguous or cluttered scenes.
362 Qualitative examples are shown in Figure 4. Through progressive refinement, CoT-Seg is capable of
363 correcting the missed object in the first turn results. Overall, CoT-Seg shows SOTA or competitive
364 results with minor deviations discussed in section 4.4.

365 **Effect of Chain-of-Thought Length** We vary the number of reasoning steps (e.g.,
366 2, 4, 8) to study the tradeoff between reasoning depth and segmentation quality.

368
369 Table 1: Referring expression segmentation results on RefCOCO (Kazemzadeh et al., 2014) dataset. The cIoU metrics of each split are reported. Baselines excerpted from (Lai et al., 2024).
370
371 Table 2: Quantitative evaluation on the test set of ReasonSeg (Lai et al., 2023). (ft) means fine-tuning on the train set. † is reproduced with the official released weights with 8-bit quantization.
372

Method	Val.	Test-A	Test-B
MCN (Luo et al., 2020)	62.4	64.2	59.7
VLT (Ding et al., 2021)	67.5	70.5	65.2
CRIS (Wang et al., 2022b)	70.5	73.2	66.1
LAVT (Yang et al., 2022)	72.7	75.8	68.8
LISA-7B (Fine tuned on ReferSeg) (Lai et al., 2024)	74.9	79.1	72.3
Seg-Zero 7B (Liu et al., 2025a)	-	80.3	-
Vision-Reasoner-7B (Liu et al., 2025b)	-	78.9	-
GSA-7B(ft) (Xia et al., 2024)	77.2	78.9	73.5
CoT-Seg	76.3	80.9	72.7
CoT-Seg with self-correction	77.2	80.9	73.8

Method	ReasonSeg (overall)	
	gIoU	cIoU
OVSeg (Liang et al., 2023)	26.1	20.1
GRES (Liu et al., 2023)	21.3	22.0
X-Decoder (Zou et al., 2023a)	21.7	16.3
SEEM (Zou et al., 2023b)	24.3	18.7
Seg-Zero-7B (Liu et al., 2025a)	57.5	52.0
Vision-Reasoner-7B (Liu et al., 2025b)	63.6	-
LISA-13B (Lai et al., 2023)	44.8	45.8
LISA-13B-Llama2 (ft) [†] (Lai et al., 2023)	50.0	51.9
LISA-13B-LLaVA1.5 (ft) (Lai et al., 2023)	61.3	62.2
CoT-Seg	66.0	58.8
CoT-Seg with self correction	66.7	60.4

Table 3 tabulates the results where all of the experiments use a maximum of two rounds of refinements for self-correction running on *RefCoCo Test-A*. The results show that the length of chain of thoughts is not critical to performance, with a length of 4 producing the best score among the tested fixed lengths. Notably, fixed CoT length is outperformed by variational length determined by the MLLM. The results indicate that two reasoning steps usually suffice while overthinking with too many steps may lower the accuracy, with varying lengths depending on the input results in the best accuracy.

Segmentor Compatibility In our quantitative experiments, the Segmentation Agent can be different SAM-based or open-vocabulary segmentation backbones. We analyze how their built-in capabilities (text-only prompts, multimodal prompts, or interactive point-based control) affect downstream performance. Table 4 tabulates the results, highlighting the importance of segmentor capability descriptions in guiding Reasoner and Segmentator collaboration.

MLLM Agent Variants We evaluate CoT-Seg with different MLLM backbones, such as GPT-4o, Gemma 3 12b, and Qwen2.5-VL-7B on *RefCoCo Test-A* with maximum of 2 rounds of refinement. Table 5 tabulates the results, which reveal how reasoning depth, hallucination tendency, and multimodal grounding influence segmentation quality and stability, showing the tradeoffs between proprietary and open-source models in reasoning-driven segmentation. For earlier VL models such as Qwen2.5, when given two segmentations they cannot determine which one is better so they can only fulfill the CoT part and not the auto-correction part of our framework.

Multimodal Input Control Our framework can be used for multiple kinds of input including but not limited to bounding box, point, and scribble annotations, demonstrating the flexibility of JSON-based multimodal reasoning and how CoT and auto correction works for all general reasoning strategies, Figure 6. CoT works especially well on improving segmentation based on rough human input, providing important text info for the segmentation agent.

5 CONCLUSION

We introduced **CoT-Seg**, a training-free framework that rethinks reasoning segmentation by integrating chain-of-thought reasoning and self-correction with off-the-shelf MLLMs and segmentation agents. Our method enables step-by-step reasoning to synthesize meta-queries, self recursive evaluation for refinement, and retrieval-augmented reasoning for knowledge gaps. Extensive experiments demonstrate that CoT-Seg substantially improves zero-shot segmentation performance across multiple benchmarks. This work highlights the untapped potential of inference-time reasoning and self-correction in bridging vision-language understanding with precise segmentation.

Table 3: CoT length experiment on Test-A of *RefCoco* (Kazemzadeh et al., 2014).

CoT Length	RefCoco gIoU	Test-A cIoU
CoT-Seg with self-correction using CoT length 2	79.9	79.4
CoT-Seg with self-correction using CoT length 4	80.3	79.5
CoT-Seg with self-correction using CoT length 8	80.1	79.4
CoT-Seg with self-correction using CoT variational length	80.1	80.9

Table 4: Segmentor experiment without self-correction on Test-A of *RefCoco* (Kazemzadeh et al., 2014).

Segmentor	RefCoco gIoU	Test-A cIoU
CoT-Seg with Vision-Reasoner-7B + SAM-HQ	80.1	80.9
CoT-Seg with LISA	77.8	79.2
CoT-Seg with GroundedSAM	51.4	61.9

Table 5: Different MLLM experiments on Test-A of *RefCoco* (Kazemzadeh et al., 2014).

CoT with different MLLMs	RefCoco gIoU	Test-A cIoU
CoT-Seg with self-correction using GPT-4o	80.0	80.9
CoT-Seg with self-correction using Gemma 3	80.1	80.3
CoT-Seg with self-correction using Qwen2.5-VL-7B	69.2	70.3



Figure 6: Multimodal inputs

REFERENCES

- 432
433
434 Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-
435 decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine*
436 *intelligence*, 39(12):2481–2495, 2017.
- 437 Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille.
438 Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and
439 fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):
440 834–848, 2017.
- 441 Bowen Cheng, Maxwell D Collins, Yukun Zhu, Ting Liu, Thomas S Huang, Hartwig Adam, and
442 Liang-Chieh Chen. Panoptic-deeplab: A simple, strong, and fast baseline for bottom-up panop-
443 tic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*
444 *recognition*, pp. 12475–12485, 2020.
- 445 Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-
446 attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF*
447 *conference on computer vision and pattern recognition*, pp. 1290–1299, 2022.
- 448 Philipe Ambrozio Dias and Henry Medeiros. Semantic segmentation refinement by monte carlo
449 region growing of high confidence detections. In *Computer Vision–ACCV 2018: 14th Asian*
450 *Conference on Computer Vision, Perth, Australia, December 2–6, 2018, Revised Selected Papers,*
451 *Part II 14*, pp. 131–146. Springer, 2019.
- 452 Henghui Ding, Chang Liu, Suchen Wang, and Xudong Jiang. Vision-language transformer and
453 query generation for referring segmentation. In *Proceedings of the IEEE/CVF international con-*
454 *ference on computer vision*, pp. 16321–16330, 2021.
- 455 Deng-Ping Fan, Ge-Peng Ji, Guolei Sun, Ming-Ming Cheng, Jianbing Shen, and Ling Shao. Cam-
456 ouflaged object detection. In *Proceedings of the IEEE/CVF conference on computer vision and*
457 *pattern recognition*, pp. 2777–2787, 2020.
- 458 Jiaxian Guo, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Boyang Li, Dacheng Tao, and
459 Steven CH Hoi. From images to textual prompts: Zero-shot vqa with frozen large language
460 models. *arXiv preprint arXiv:2212.10846*, 2022.
- 461 Jiayi He, Hehai Lin, Qingyun Wang, Yi Fung, and Heng Ji. Self-correction is more than refinement:
462 A learning framework for visual and language reasoning tasks. *arXiv preprint arXiv:2410.04055*,
463 2025.
- 464 Junwen He, Yifan Wang, Lijun Wang, Huchuan Lu, Jun-Yan He, Jin-Peng Lan, Bin Luo, and Xuan-
465 song Xie. Multi-modal instruction tuned llms with fine-grained visual perception. In *Proceedings*
466 *of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 13980–13990, 2024.
- 467 Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the*
468 *IEEE international conference on computer vision*, pp. 2961–2969, 2017.
- 469 Yingdong Hu, Fanqi Lin, Tong Zhang, Li Yi, and Yang Gao. Look before you leap: Unveiling the
470 power of gpt-4v in robotic vision-language planning. *arXiv preprint arXiv:2311.17842*, 2023.
- 471 Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Os-
472 trow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint*
473 *arXiv:2410.21276*, 2024.
- 474 Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to
475 objects in photographs of natural scenes. In *Proceedings of the 2014 conference on empirical*
476 *methods in natural language processing (EMNLP)*, pp. 787–798, 2014.
- 477 Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmen-
478 tation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
479 pp. 9404–9413, 2019.

- 486 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
487 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceed-*
488 *ings of the IEEE/CVF international conference on computer vision*, pp. 4015–4026, 2023.
- 489 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
490 language models are zero-shot reasoners. *arXiv preprint arXiv:2205.11916*, 2022.
- 492 Mojtaba Komeili, Kurt Shuster, and Jason Weston. Internet-augmented dialogue generation. *arXiv*
493 *preprint arXiv:2107.07566*, 2021.
- 494 Philipp Krähenbühl and Vladlen Koltun. Efficient inference in fully connected crfs with gaussian
495 edge potentials. *Advances in neural information processing systems*, 24, 2011.
- 497 Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Rea-
498 soning segmentation via large language model. *arXiv preprint arXiv:2308.00692*, 2023.
- 500 Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Rea-
501 soning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on*
502 *Computer Vision and Pattern Recognition*, pp. 9579–9589, 2024.
- 503 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,
504 Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe
505 Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks. *arXiv preprint*
506 *arXiv:2005.11401*, 2021.
- 507 Long Lian, Boyi Li, Adam Yala, and Trevor Darrell. Llm-grounded diffusion: Enhancing prompt
508 understanding of text-to-image diffusion models with large language models. *arXiv preprint*
509 *arXiv:2305.13655*, 2023.
- 511 Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yinan Zhao, Hang Zhang, Peizhao Zhang,
512 Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted
513 clip. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
514 pp. 7061–7070, 2023.
- 515 Chang Liu, Henghui Ding, and Xudong Jiang. Gres: Generalized referring expression segmentation.
516 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp.
517 23592–23601, 2023.
- 518 Wei Liu, Andrew Rabinovich, and Alexander C Berg. Parsenet: Looking wider to see better. *arXiv*
519 *preprint arXiv:1506.04579*, 2015.
- 521 Yuqi Liu, Bohao Peng, Zhisheng Zhong, Zihao Yue, Fanbin Lu, Bei Yu, and Jiaya Jia. Seg-
522 zero: Reasoning-chain guided segmentation via cognitive reinforcement. *arXiv preprint*
523 *arXiv:2503.06520*, 2025a.
- 524 Yuqi Liu, Tianyuan Qu, Zhisheng Zhong, Bohao Peng, Shu Liu, Bei Yu, and Jiaya Jia. Vision-
525 reasoner: Unified visual perception and reasoning via reinforcement learning. *arXiv preprint*
526 *arXiv:2505.12081*, 2025b.
- 528 Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord,
529 Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for
530 science question answering. *Advances in Neural Information Processing Systems*, 35:2507–2521,
531 2022.
- 532 Gen Luo, Yiyi Zhou, Xiaoshuai Sun, Liujuan Cao, Chenglin Wu, Cheng Deng, and Rongrong Ji.
533 Multi-task collaborative network for joint referring expression comprehension and segmentation.
534 In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, pp.
535 10034–10043, 2020.
- 537 Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki,
538 and Chris Callison-Burch. Faithful chain-of-thought reasoning. In *The 13th International Joint*
539 *Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter*
of the Association for Computational Linguistics (IJCNLP-AAACL 2023), 2023.

- 540 Chancharik Mitra, Brandon Huang, Trevor Darrell, and Roei Herzig. Compositional chain-of-
541 thought prompting for large multimodal models. In *Proceedings of the IEEE/CVF Conference*
542 *on Computer Vision and Pattern Recognition*, pp. 14420–14431, 2024.
- 543
- 544 Debjyoti Mondal, Suraj Modi, Subhadarshi Panda, Rituraj Singh, and Godawari Sudhakar Rao.
545 Kam-cot: Knowledge augmented multimodal chain-of-thoughts reasoning. *arXiv preprint*
546 *arXiv:2401.12863*, 2024.
- 547 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally
548 can be more effective than scaling model parameters. *ICLR*, 2025.
- 549
- 550 Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for
551 reasoning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp.
552 11888–11898, 2023.
- 553 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-
554 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models.
555 *arXiv preprint arXiv:2203.11171*, 2022a.
- 556
- 557 Zhaoqing Wang, Yu Lu, Qiang Li, Xunqiang Tao, Yandong Guo, Mingming Gong, and Tongliang
558 Liu. Cris: Clip-driven referring image segmentation. In *Proceedings of the IEEE/CVF conference*
559 *on computer vision and pattern recognition*, pp. 11686–11695, 2022b.
- 560 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
561 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. In *Advances*
562 *in neural information processing systems*, volume 35, pp. 24824–24837, 2022.
- 563 Zhuofan Xia, Dongchen Han, Yizeng Han, Xuran Pan, Shiji Song, and Gao Huang. Gsva: Gen-
564 eralized segmentation via multimodal large language models. In *Proceedings of the IEEE/CVF*
565 *Conference on Computer Vision and Pattern Recognition*, pp. 3858–3869, 2024.
- 566
- 567 Zhao Yang, Jiaqi Wang, Yansong Tang, Kai Chen, Hengshuang Zhao, and Philip HS Torr. Lavt:
568 Language-aware vision transformer for referring image segmentation. In *Proceedings of the*
569 *IEEE/CVF conference on computer vision and pattern recognition*, pp. 18155–18165, 2022.
- 570 Huanjin Yao, Jiaxing Huang, Yawen Qiu, Michael K. Chen, Wenzheng Liu, Wei Zhang, Wenjie
571 Zeng, Xikun Zhang, Jingyi Zhang, Yuxin Song, Wenhao Wu, and Dacheng Tao. Mmreason: An
572 open-ended multi-modal multi-step reasoning benchmark for mllms toward agi. *arXiv preprint*
573 *arXiv:2506.23563*, 2025.
- 574 Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv*
575 *preprint arXiv:1511.07122*, 2015.
- 576
- 577 Ao Zhang, Yuan Yao, Wei Ji, Zhiyuan Liu, and Tat-Seng Chua. Next-chat: An lmm for chat,
578 detection and segmentation. *arXiv preprint arXiv:2311.04498*, 2023a.
- 579
- 580 Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in
581 large language models. In *International Conference on Learning Representation*, 2023b.
- 582
- 583 Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. Multimodal
584 chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*, 2023c.
- 584
- 585 Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing
586 network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
587 2881–2890, 2017.
- 588
- 589 Xutong Zhao, Tengyu Xu, Xuewei Wang, Zhengxing Chen, Di Jin, Liang Tan, Yen-Ting, Zishun
590 Yu, Zhuokai Zhao, Yun He, Sinong Wang, Han Fang, Sarath Chandar, and Chen Zhu. Boosting
591 llm reasoning via spontaneous self-correction. *arXiv preprint arXiv:2506.06923*, 2025.
- 591
- 592 Xueyan Zou, Zi-Yi Dou, Jianwei Yang, Zhe Gan, Linjie Li, Chunyuan Li, Xiyang Dai, Harkirat
593 Behl, Jianfeng Wang, Lu Yuan, et al. Generalized decoding for pixel, image, and language. In
Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 15116–
15127, 2023a.

594 Xueyan Zou, Jianwei Yang, Hao Zhang, Feng Li, Linjie Li, Jianfeng Wang, Lijuan Wang, Jian-
595 feng Gao, and Yong Jae Lee. Segment everything everywhere all at once. *Advances in neural*
596 *information processing systems*, 36:19769–19782, 2023b.
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

648 APPENDIX

649
650 A FURTHER DETAILS

651
652 **CoT First Turn Template**

653 We use this as a basic description for the LLM to propose questions answer pairs for CoT process,
654 we replace `<QUERY>` with the user query.

655 *You will serve as an agent for language-based image segmentation model. During each inference,*
656 *your task is to consider a query and describe a given image with chain of thoughts. You need to*
657 *provide details to help the segmentation model understand the image better. The target objects may*
658 *contain multiple layers, be blocked by other object, or be seamlessly embedded in their surround-*
659 *ings. Your description will be later sent to the segmentation as prompt. For example, if given an*
660 *image, you need to describe what can be seen in the image, the number of objects for each categories,*
661 *the position of the target object, the structure of the object, the number of layers of the object, etc.*
662 *The actual description depends on the given image. For the output, you need to follow the format:-*
663 *Question 1: Answer 1.- Question 2: Answer 2 ..., etc, where each pair of prompt and answer implies*
664 *the chain of thoughts, i.e., different levels or different part of the image understanding. For example,*
665 *the first prompt can be related to the overall style or background of the image. Finally, you need*
666 *to summarize the description based on your generated prompts and answers with strictly with the*
667 *format: Your summary here ... considering the prompt where the user is looking for ..., the object*
668 *of interest may be ... Then, based on the summary, you have to generate a pseudo-prompt to query*
669 *the segmentation model. This pseudo-prompt should contains the information about what is in the*
670 *image, what to segment, and where the target object is. It must strictly follow the format: - Prompt:*
671 *The image showsPlease segment the ... located at ... of the image. Lastly, please generate a list*
672 *of labels that would be passed to an object detector based on the summary and thinking process*
673 *that describes the object(s) that fits the user query and should be segmented and in strictly in the*
674 *format of cat. remote control. television. with a period separating every label and if there is multiple*
675 *word in a label then separate using space. Do not include other objects not given by the prompt.*
676 `<QUERY>`

677 **CoT Self-Correction Process**

678 - Reasoning process:

679 1. *Original image: The original image shows an underwater scene with various marine life*
680 *forms and corals. The Pagurian, or hermit crab, is located at the center-right of the image. It is*
681 *identifiable by its orange color and the shell it occupies.*

682 2. *Segmentation image: The segmentation image shows an isolated orange Pagurian on a white*
683 *background. The shell is not included in the segmentation.*

684 3. *Summary: The segmentation correctly isolates the Pagurian itself but fails to include the shell,*
685 *which is an integral part of the Pagurian's appearance in the original image. Therefore, the*
686 *segmentation is incomplete.*

687 - Correctness: False

688 - Meta-queries (Output if the correctness is false):

689 *Positive: Please also segment the shell of the Pagurian, located at the center-right of the original*
690 *image.*

691 *Negative: None*

692 - Labels:

693 *Positive: shell*

694 *Negative: None.*

695 **Full First Turn and Subsequent Self-Correction CoT**

696 **First Turn**

697 - Question 1: What is the overall setting of the image?

698 - Answer 1: The image shows a person standing on grass, with a dog in the foreground.

699 - Question 2: What are the main objects visible in the image?

700 - Answer 2: The main objects visible are a person, a dog, and a leash.

701 - Question 3: What is the person holding?

- Answer 3: The person is holding a leash.

702 - Question 4: Where is the object of interest located in the image?
 703 - Answer 4: The leash is located in the upper left part of the image, extending from the person’s
 704 hand to the dog.
 705 Summary: The image shows a person standing on grass with a dog. Considering the prompt where
 706 the user is looking for the object that the person is holding while walking their dog, the object of
 707 interest may be the leash.
 708 Prompt: The image shows a person standing on grass with a dog. Please segment the leash located
 709 at the upper left part of the image.
Self-Correction
 710 - Reasoning process:
 711 1. Original image: The original image shows a person standing on grass with a dog. There is a
 712 leash visible in the upper left part of the image, held by the person.
 713 2. Segmentation image: The segmentation image shows a portion of the person’s clothing and a
 714 small part of the leash in the upper right corner, isolated on a white background.
 715 3. Summary: The segmentation result does not correctly reflect the user query. The query asks
 716 for the leash to be segmented, but the segmentation image only includes a small part of it and part
 717 of the person’s clothing. The entire leash should be included, and the person’s clothing should be
 718 excluded.
 719 - Correctness: False
 720 - Meta-queries (Output if the correctness is false):
 721 1. Positive: Please also segment the entire leash, located at the upper left part of the original image.
 722 2. Negative: Please remove the person’s clothing, located at the upper right part of the segmentation
 723 image.
 724 - Labels:
 725 1. Positive: leash.
 726 2. Negative: person’s clothing.

727 **B ADDITIONAL RESULTS**

728 **Benchmarks on COD10K** Results in Table 6 show that CoT-Seg displays SOTA or competitive
 729 results regarding segmenting highly difficult camouflaged objects, where CoT reasoning will make
 730 the task easier.
 731
 732

733 Table 6: Quantitative evaluation on camouflaged images in *COD-10K* (Fan et al., 2020). Note that
 734 the models are blind to both dataset divisions. CoT-Seg with refinement improves on implicit query
 735 where with one round of segmentation may not be enough to fulfill the query.
 736

Method	Training free	Train split				Test split			
		Implicit query		Explicit query		Implicit query		Explicit query	
		gIoU	cIoU	gIoU	cIoU	gIoU	cIoU	gIoU	cIoU
LISA-13B-Llama2 (Lai et al., 2024)	No	65.0	60.3	66.5	59.7	63.4	55.8	66.5	58.0
Vision-Reasoner-7B (Liu et al., 2025b)	No	77.1	77.1	77.6	76.6	76.9	76.3	77.4	75.0
CoT-Seg	Yes	<u>77.9</u>	78.8	<u>77.5</u>	78.0	<u>77.7</u>	<u>77.9</u>	<u>77.5</u>	<u>74.9</u>
CoT-Seg with self-correction	Yes	78.0	<u>78.4</u>	77.4	<u>76.6</u>	78.2	78.6	77.6	<u>74.9</u>

743
 744 **Additional Visual Examples** We show additional self-correction examples in Figure 12 and some
 745 examples of RAG in Figure 10 and 11.
 746
 747

748 **C SIMILAR WORKS ANALYSIS AND COMPARISON**

749 **Vision Reasoner** We discuss the difference between our work and the concurrent work Vision-
 750 Reasoner (Liu et al., 2025b). To the best of our knowledge VisionReasoner uses reinforcement
 751 learning to generate the bounding boxes and segmentations. VisionReasoner has greatly improved
 752 on previous reasoning segmentation models as shown in Tables 1– 6 but still fails in some complicated
 753 cases where there are a large number of objects to be segmented Figures 1, 8 and 9 or when
 754 the prompt is very implicit. CoT-Seg in comparison, is zero-shot and can be easily plugged in to different
 755 models, offering high flexibility and achieves higher scores in all the test data in Tables 1– 6.

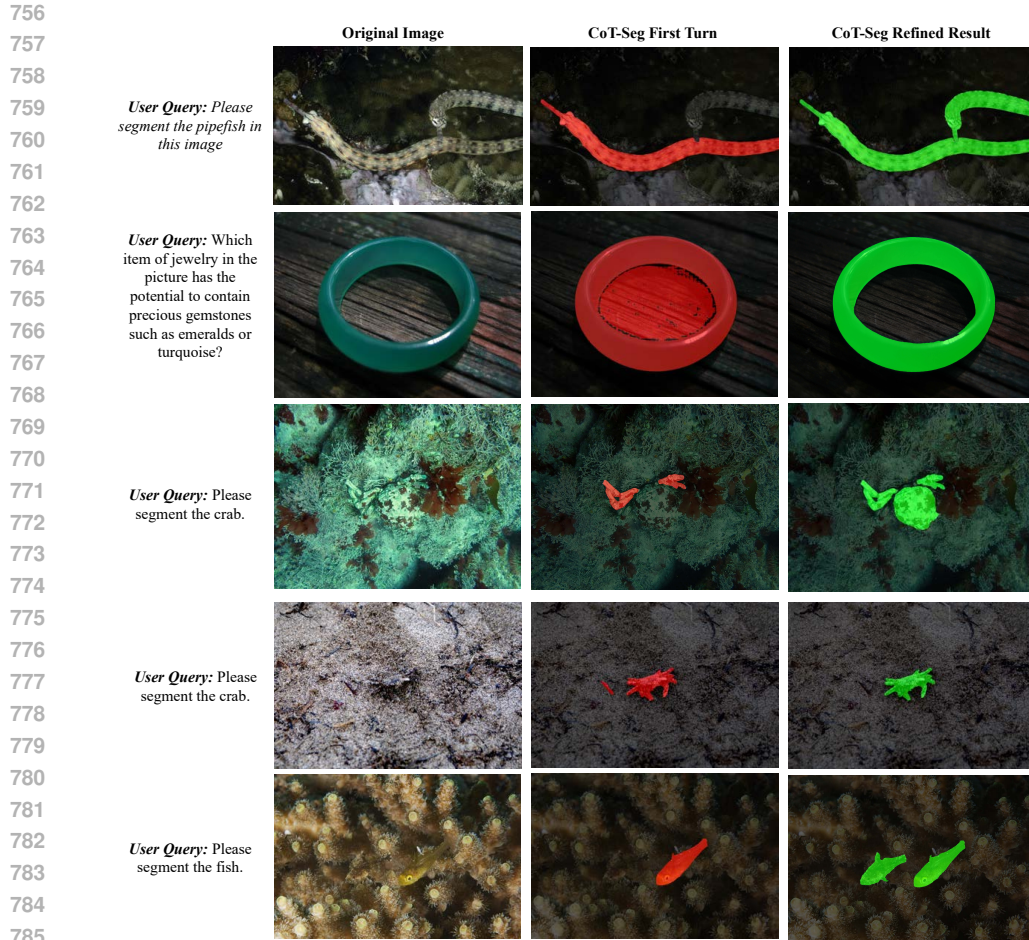


Figure 7: Additional self-correction results.

GSVA Table 7 shows CoT-Seg’s competitive performance on this standard referring segmentation benchmark with less emphasis on CoT deep reasoning for complex segmentation.

GSVA (Xia et al., 2024) also uses MLLM to guide segmentation. Specifically, GSVA uses MLLM to generate [SEG] tokens and prompt the segmentation model to support multiple object segmentation and a [NULL] token to reject absent object. In comparison, our approach uses chain of thought reasoning to assimilate and provide useful information to the segmentator agent, empowering our model to solve very implicit queries and achieve multiple-object segmentation in a training-free manner. Our auto-correction process further leverages MLLM to improve and obtain accurate segmentations that the segmentor agent cannot achieve on its own. In RefCOCO tests in Table 7, GSVA achieves slightly higher results, mainly because of training and finetuning on the RefCOCO training

Method	Val.	Test-A	Test-B
GSVA-Llama2-13B (Xia et al., 2024)	<u>77.7</u>	79.9	<u>74.2</u>
GSVA-Llama2-13B (ft) (Xia et al., 2024)	79.2	81.7	77.1
CoT-Seg	76.3	80.9	72.7
CoT-Seg with self-correction	<u>77.2</u>	<u>80.9</u>	73.8

Table 7: Quantitative comparison with GSVA (Xia et al., 2024) on RefCOCO

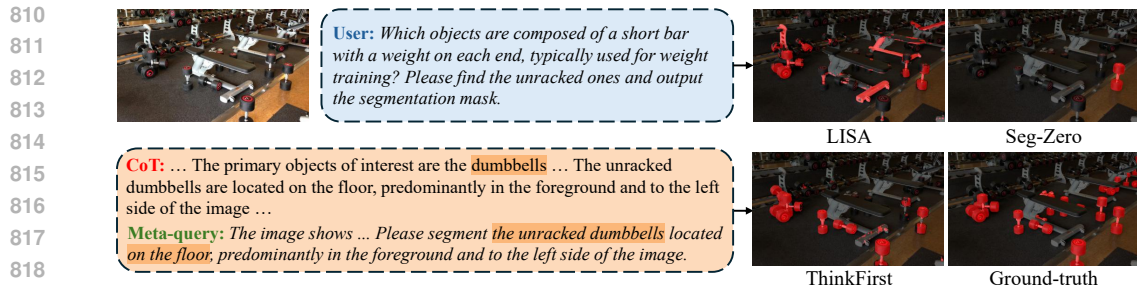


Figure 8: CoT-Seg reasons about the arrangement of dumbbells to segment those that are un racked, a more challenging task than simple detection (No self-correction was needed).

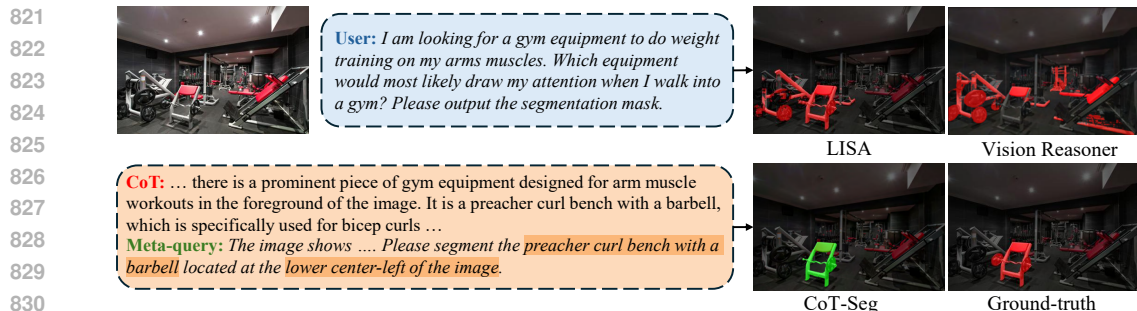


Figure 9: CoT-Seg identifies the gym equipment matching the user’s query for biceps, e.g., the preacher’s curl, reasoning about its function without any training (Self-correction was needed).



Figure 10: RAG Result

dataset getting higher accuracy in prompts containing numerical positional arguments, an example is shown in Figure 14. Further and similar training and finetuning of our model should be able to improve our results on these benchmarks, as well as incorporating our framework as a plugin to GSVA and other recent SOTA reasoning segmentation agents. On the other hand, GSVA does not focus on reasoning segmentation and as a cloU score of 43.4 for 7B model and 44.6 for 13B model on the Reasonseg dataset. Furthermore, inference on GSVA 7B(ft) shows that GSVA 7B(ft) was unable to get correct segmentation results when the prompt becomes more implicit such as Figures 1 and 2, the inference results are shown in Figure 13.

D REASONSEG-HARD: NEW EVALUATION DATASET FOR STRESS TESTING REASONING SEGMENTATION

Most queries/images in existing datasets are not too challenging, where auto-correction is not necessary. Specifically, in RefCoco and ReasonSeg, we found that only around 10% of the test cases require auto-correction, as the majority of the segmentation task is usually quite simple, where the segmentation agent alone can obtain the correct answer on the first turn.

Given more recent and significant advancement in reasoning segmentation, a dataset update is due and necessary, which should contain more difficult cases such as closely connected objects and multiple objects similar to the object of interest, with implicit queries requiring complex reasoning to

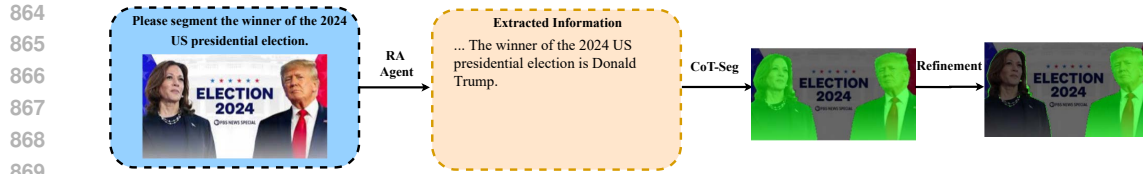


Figure 11: Retrieval-augmented CoT-Seg Result.

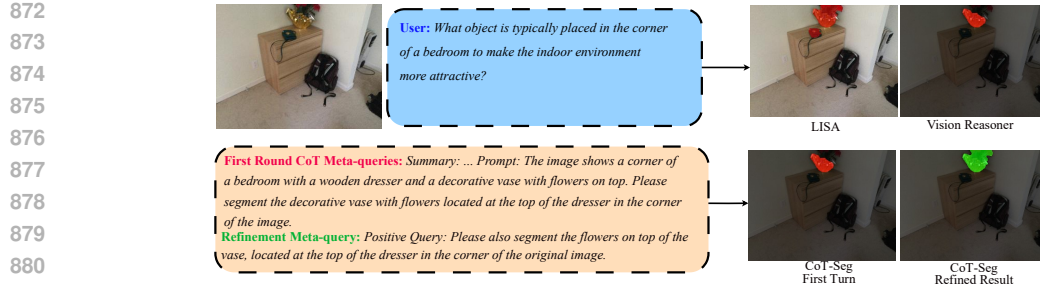


Figure 12: CoT-Seg can correct minor mistakes such as not segmenting the flowers with the vase.

understand for challenging segmentation task. For example, the new dataset should contain difficult scenarios such as those depicted in Figures 1 and 2, to validate stronger models to come in the future.

Thus, in this paper, we propose REASONSEG-HARD, a new evaluation dataset for stress testing reasoning segmentation. Specifically, we constructed a dataset with 213 image-query pairs consisting of 75 images and their respective queries sampled from ReasonSeg Test Split. Our results are shown in Table 8. We sample queries-images pairs that either require deeper and more thorough reasoning to identify object(s) of interest, specifically queries with context to be reasoned, or queries including complex objects inherently difficult to segment due to size, transparency and how well it is blended into the environment. For example, the dataset contains implicit queries such as “*When preparing for a festive event like Halloween, people often use certain objects to decorate their homes. What object in the picture would be suitable for this purpose?*” and excluding queries that may be too simple and obvious such as “*something that the person uses to fish*”. Figure 15 shows additional image/queries samples of our REASONSEG-HARD.

Method	gIoU	cIoU
LISA-13B-Llama2 [†] (Lai et al., 2023)	38.0	41.1
GSVA-7B (ft) (Xia et al., 2024)	40.9	37.8
Vision Reasoner (Liu et al., 2025b)	49.1	48.1
SegZero-7B (Liu et al., 2025a)	44.0	52.6
CoT-Seg	56.7	54.4
CoT-Seg with self-correction	58.6	57.4

Table 8: Reasoning segmentation evaluation with complex and implicit queries on our REASONSEG-HARD dataset. [†] is produced with the official released weights with 8-bit quantization.

E LLM USAGE

Large Language Models (LLMs) were only used to polish the English sparsely in the paper.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

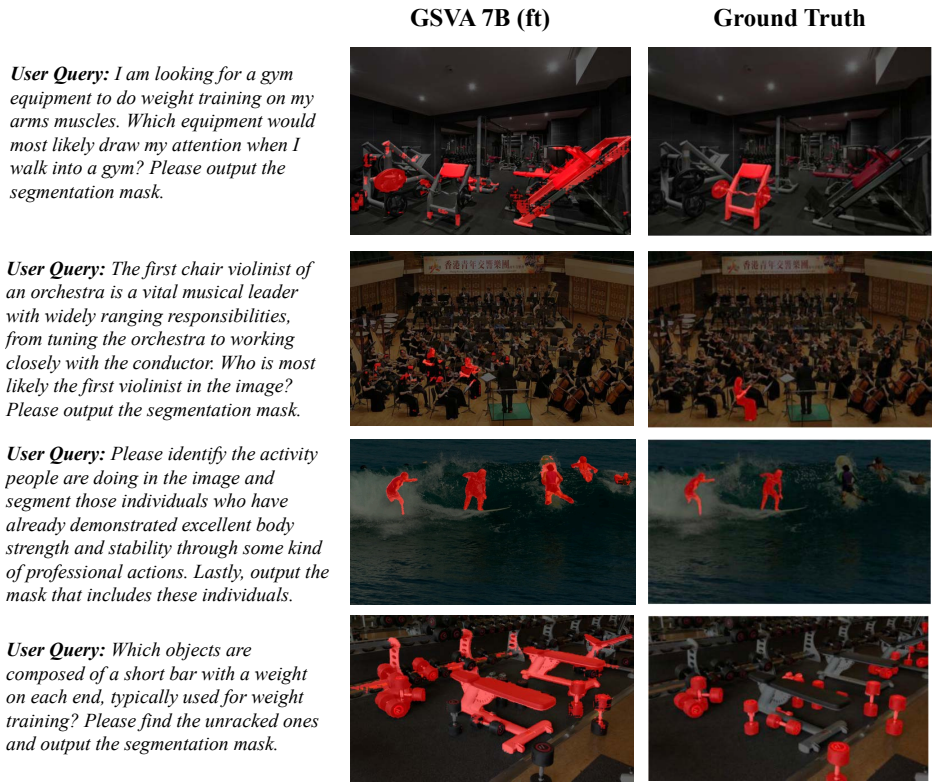


Figure 13: GSVA inference results.



Figure 14: GSVA - Prompt: Second from the right.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

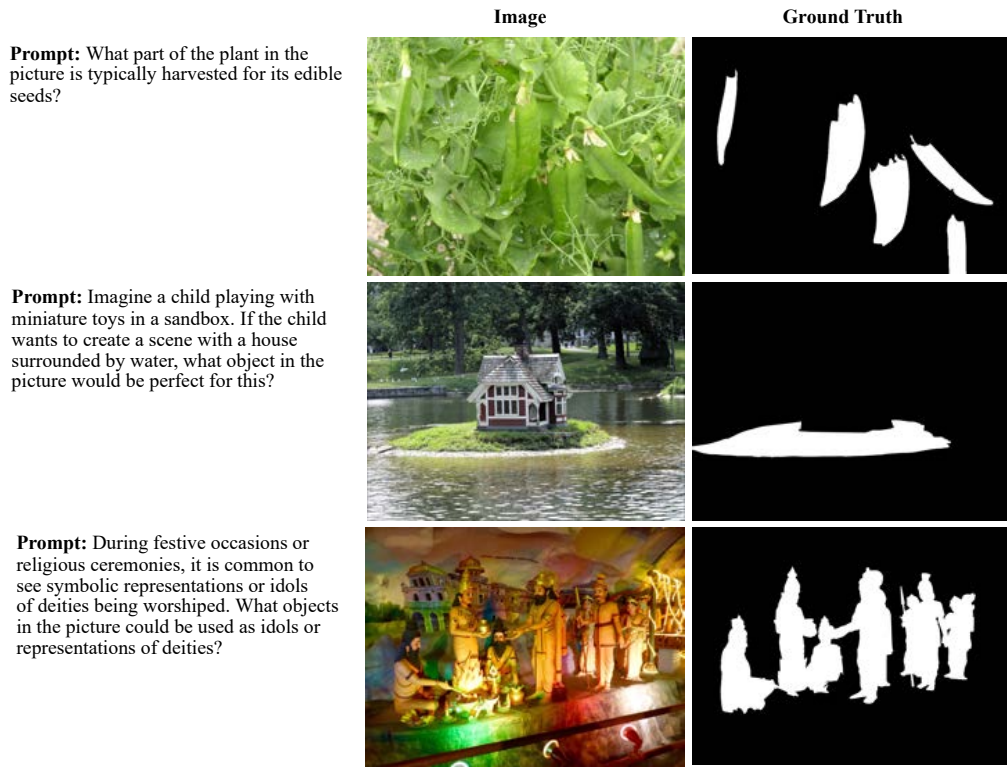


Figure 15: Samples of challenging examples in REASONSEG-HARD.