

Novelty Assessment Report

Paper: SAM 3: Segment Anything with Concepts

PDF URL: <https://openreview.net/pdf?id=r35clVtGzw>

Venue: ICLR 2026 Conference Submission

Year: 2026

Report Generated: 2025-12-29

Abstract

We present Segment Anything Model (SAM) 3, a unified model that detects, segments, and tracks objects in images and videos based on concept prompts, which we define as either short noun phrases (e.g., “yellow school bus”), image exemplars, or a combination of both. Promptable Concept Segmentation (PCS) takes such prompts and returns segmentation masks and unique identities for all matching object instances. To advance PCS, we build a scalable data engine that produces a high-quality dataset with 4M unique concept labels, including hard negatives, across images and videos. Our model consists of a vision backbone shared between an image-level detector and a memory-based video tracker. Recognition and localization are decoupled with a presence head, which significantly boosts detection accuracy. SAM 3 delivers a 2x gain over existing systems in both image and video PCS, and improves previous SAM capabilities in interactive visual segmentation tasks. We open source SAM 3 along with our new Segment Anything with Concepts (SA-Co) benchmark.

Disclaimer

This report is **AI-GENERATED** using Large Language Models and WisPaper (a scholar search engine). It analyzes academic papers' tasks and contributions against retrieved prior work. While this system identifies **POTENTIAL** overlaps and novel directions, **ITS COVERAGE IS NOT EXHAUSTIVE AND JUDGMENTS ARE APPROXIMATE**. These results are intended to assist human reviewers and **SHOULD NOT** be relied upon as a definitive verdict on novelty.

Note that some papers exist in multiple, slightly different versions (e.g., with different titles or URLs). The system may retrieve several versions of the same underlying work. The current automated pipeline does not reliably align or distinguish these cases, so human reviewers will need to disambiguate them manually.

If you have any questions, please contact: mingzhang23@m.fudan.edu.cn

Core Task Landscape

This paper addresses: **Promptable Concept Segmentation in Images and Videos**

A total of **50 papers** were analyzed and organized into a taxonomy with **16 categories**.

Taxonomy Overview

The research landscape has been organized into the following main categories:

- **Foundation Model Architectures and Core Mechanisms**
- **Prompt Design and Learning Strategies**
- **Domain Adaptation and Transfer Learning**
- **Task-Specific Applications and Extensions**
- **Few-Shot and Personalized Segmentation**
- **Generative Models for Segmentation**

Complete Taxonomy Tree

- Promptable Concept Segmentation in Images and Videos Survey Taxonomy
- Foundation Model Architectures and Core Mechanisms
 - Unified Promptable Segmentation Frameworks ★ (4 papers)
 - [0] SAM 3: Segment Anything with Concepts (Anon et al., 2026) [View paper](#)
 - [5] Image segmentation using text and image prompts (Timo LÅ¼addecke, 2022) [View paper](#)
 - [16] Segment everything everywhere all at once (Zou, 2023) [View paper](#)
 - [33] UniVS: Unified and Universal Video Segmentation with Prompts as Queries (Minghan Li, 2024) [View paper](#)
 - Video Object Segmentation with Memory Mechanisms (4 papers)
 - [7] VISA: Reasoning Video Object Segmentation via Large Language Models (Yan, 2024) [View paper](#)
 - [8] Medical SAM 2: Segment medical images as video via Segment Anything Model 2 (Zhu Jia-yuan, 2024) [View paper](#)
 - [15] SAM-I2V: Upgrading SAM to Support Promptable Video Segmentation with Less than 0.2% Training Cost (Haiyang Mei, 2025) [View paper](#)
 - [46] Biomedical SAM 2: Segment Anything in Biomedical Images and Videos (Yan, 2024) [View paper](#)
 - Attention and Diffusion-Based Segmentation Mechanisms (3 papers)
 - [14] Seg4diff: Unveiling open-vocabulary segmentation in text-to-image diffusion transformers (Kim Chaehyun, 2025) [View paper](#)
 - [23] Unleashing text-to-image diffusion models for visual perception (Wenliang Zhao, 2023) [View paper](#)
 - [42] From text to mask: Localizing entities using the attention of text-to-image diffusion models (Changming Xiao, 2024) [View paper](#)
- Prompt Design and Learning Strategies
 - Learnable and Adaptive Prompt Mechanisms (4 papers)
 - [12] Learnable prompt for few-shot semantic segmentation in remote sensing domain (Steve Andreas Immanuel, 2024) [View paper](#)
 - [13] Prompt-driven dynamic object-centric learning for single domain generalization (Deng Li, 2024) [View paper](#)
 - [20] Learning to prompt segment anything models (Huang Jiaying, 2024) [View paper](#)
 - [44] Promptonomyvit: Multi-task prompt learning improves video transformers using synthetic scene data (Roei Herzig, 2024) [View paper](#)
 - Cross-Modal and Multi-Modal Prompt Fusion (3 papers)
 - [3] Prompt-driven referring image segmentation with instance contrasting (Chao Shang, 2024) [View paper](#)
 - [38] Prompting segmentation with sound is generalizable audio-visual source localizer (Wang Yao-ting, 2024) [View paper](#)
 - [43] Multimodal Prompt-Guided Bidirectional Fusion for Referring Remote Sensing Image Segmentation (Yingjie Li, 2025) [View paper](#)
 - Training-Free and Zero-Shot Prompt Strategies (3 papers)
 - [9] Towards training-free open-world segmentation via image prompt foundation models (Lv Tang, 2025) [View paper](#)

- [25] Training-free open-ended object detection and segmentation via attention as prompts (Lin Zhiwei, 2024) [View paper](#)
- [39] Prompting to Adapt Foundational Segmentation Models (Jie Hu, 2024) [View paper](#)
- Domain Adaptation and Transfer Learning
 - Medical Image Segmentation Adaptation (6 papers)
 - [1] Segment anything model for medical image segmentation: Current applications and future directions (Zhang, 2024) [View paper](#)
 - [2] Promise: Prompt-driven 3d medical image segmentation using pretrained image foundation models (Hao Li, 2024) [View paper](#)
 - [19] nninteractive: Redefining 3d promptable segmentation (Isensee, 2025) [View paper](#)
 - [22] Promptable cancer segmentation using minimal expert-curated data (Wang Yipei, 2025) [View paper](#)
 - [24] Interactive Segmentation Model for Placenta Segmentation from 3D Ultrasound images (Li Hao, 2024) [View paper](#)
 - [35] One-prompt to segment all medical images (Junde Wu, 2024) [View paper](#)
 - Remote Sensing and Aerial Imagery Adaptation (4 papers)
 - [6] Sopseg: Prompt-based small object instance segmentation in remote sensing imagery (Wang Chen-hao, 2025) [View paper](#)
 - [26] Attention Prompt-Driven Source-Free Adaptation for Remote Sensing Images Semantic Segmentation (Kuiliang Gao, 2024) [View paper](#)
 - [32] DiffRIS: Enhancing Referring Remote Sensing Image Segmentation with Pre-trained Text-to-Image Diffusion Models (Dong Zhe, 2025) [View paper](#)
 - [45] Progressive Self-Prompting Segment Anything Model for Salient Object Detection in Optical Remote Sensing Images (Xiao-ning Zhang, 2025) [View paper](#)
 - Source-Free and Prompt-Driven Domain Adaptation (2 papers)
 - [40] Poda: Prompt-driven zero-shot domain adaptation (Mohammad Fahes, 2023) [View paper](#)
 - [49] Enhancing Label-efficient Medical Image Segmentation with Text-guided Diffusion Models (Feng, 2024) [View paper](#)
- Task-Specific Applications and Extensions
 - Referring and Reasoning-Based Segmentation (3 papers)
 - [4] Reasoning Segmentation for Images and Videos: A Survey (Shen Yiqing, 2025) [View paper](#)
 - [21] Driving Referring Video Object Segmentation with Vision-Language Pre-trained Models (Zhou ZiKun, 2024) [View paper](#)
 - [28] Unified Multi-Modality Video Object Segmentation Using Reinforcement Learning (Mingjie Sun, 2024) [View paper](#)
 - Interactive and Real-Time Segmentation Systems (4 papers)
 - [27] Rethinking interactive image segmentation with low latency high quality and diverse prompts (Qin Liu, 2024) [View paper](#)
 - [37] Llava-interactive: An all-in-one demo for image chat, segmentation, generation and editing (Wei-Ge Chen, 2023) [View paper](#)
 - [47] Visual Prompt Selection Framework for Real-Time Object Detection and Interactive Segmentation in Augmented Reality Applications (Eunghyeol Song, 2024) [View paper](#)
 - [50] SkipClick: Combining Quick Responses and Low-Level Features for Interactive Segmentation in Winter Sports Contexts (Robin Schäfer, 2025) [View paper](#)
 - Open-World and Open-Vocabulary Segmentation (2 papers)
 - [30] Segprompt: Boosting open-world segmentation via category-level prompt learning (Muzhi Zhu, 2023) [View paper](#)
 - [36] Faster Interactive Segmentation of Identical-Class Objects With One Mask in High-Resolution Remotely Sensed Imagery (Zhang Zhili, 2025) [View paper](#)
 - Specialized Application Domains (5 papers)
 - [11] An Interactive Prompt Based Network for Urban Floods Area Segmentation Using UAV Images (Lingfei Shi, 2025) [View paper](#)
 - [17] Zero-shot surgical tool segmentation in monocular video using Segment Anything Model 2 (Ange Lou, 2024) [View paper](#)
 - [18] Depthwise-dilated convolutional adapters for medical object tracking and segmentation using the segment anything model 2 (Guoping Xu, 2025) [View paper](#)
 - [41] An unsupervised approach towards promptable defect segmentation in laser-based additive manufacturing by segment anything (Israt Zarin Era, 2023) [View paper](#)
 - [48] Performance and nonadversarial robustness of the segment anything model 2 in surgical video segmentation (Yiqing Shen, 2024) [View paper](#)
- Few-Shot and Personalized Segmentation
 - One-Shot and Few-Shot Adaptation (2 papers)
 - [10] MWVOS: Mask-free weakly supervised video object segmentation via promptable foundation model (Zhenghao Zhang, 2025) [View paper](#)
 - [31] Personalize Segment Anything Model with One Shot (Zhang, 2023) [View paper](#)
 - Personalized and Subject-Specific Segmentation (1 papers)
 - [34] Subject-Diffusion: Open Domain Personalized Text-to-Image Generation without Test-time Fine-tuning (Jian Ma, 2023) [View paper](#)
- Generative Models for Segmentation (1 papers)
 - [29] Adding Conditional Control to Text-to-Image Diffusion Models (Lvmin Zhang, 2023) [View paper](#)

Narrative

Core task: promptable concept segmentation in images and videos. The field has evolved around enabling flexible, user-guided segmentation through diverse prompt modalities—ranging from points and boxes to text descriptions and even audio cues. The taxonomy reveals six major branches that capture this landscape. Foundation Model Architectures and Core Mechanisms anchor the field with unified frameworks like SAM Concepts[0] and Segment Everything Everywhere[16], which provide general-purpose segmentation engines adaptable to multiple prompt types. Prompt Design and Learning Strategies explore how to optimize prompt representations, whether through learnable tokens as in Learning Prompt SAM[20] or dynamic adaptation schemes. Domain Adaptation and Transfer Learning address the challenge of moving these models into specialized contexts such as medical imaging (SAM Medical Survey[1], Medical SAM 2[8]) or remote sensing (Few-shot Remote Sensing[12]), often with minimal retraining. Task-Specific Applications and Extensions branch into niche problems like surgical tool segmentation (Zero-shot Surgical Tool[17]) and urban flood mapping (Urban Floods Interactive[11]), while Few-Shot and Personalized Segmentation focuses on tailoring models to individual users or rare concepts (Personalize SAM[31]). Finally, Generative Models for Segmentation leverage diffusion architectures (ControlNet[29], Unleashing Diffusion Perception[23]) to integrate segmentation with synthesis.

A particularly active tension lies between building universal, training-free frameworks versus domain-specific fine-tuning: works like Training-free Open-World[9] and PODA Zero-shot[40] pursue broad generalization without additional data, while medical and remote sensing branches emphasize adaptation to specialized distributions. SAM Concepts[0] sits squarely within the Foundation Model Architectures branch, specifically under Unified Promptable Segmentation Frameworks, alongside Text Image Prompts[5] and UniVS Prompts Queries[33]. Compared to these neighbors, SAM Concepts[0] emphasizes a holistic approach to handling diverse concept-level prompts within a single architecture, whereas Text Image Prompts[5] pioneered early multimodal prompt fusion and UniVS Prompts

Queries[33] focuses on unifying prompt and query mechanisms for video understanding. This positioning reflects a trend toward consolidating multiple prompt modalities into cohesive systems that balance generality with task-specific performance.

Related Works in Same Category

The following **3 sibling papers** share the same taxonomy leaf node with the original paper:

1. Image segmentation using text and image prompts

Authors: Timo Ldecke, Alexander Ecker, Alexander S. Ecker | **Year/Venue:** 2022 | **URL:** [View paper](#)

Abstract

Image segmentation is usually addressed by training a model for a fixed set of object classes. Incorporating additional classes or more complex queries later is expensive as it requires re-training the model on a dataset that encompasses these expressions. Here we propose a system that can generate image segmentations based on arbitrary prompts at test time. A prompt can be either a text or an image. This approach enables us to create a unified model (trained once) for three common segmentation ...

Relationship Analysis

Both papers belong to the unified promptable segmentation frameworks category, as they support multiple prompt modalities (text, image exemplars, and geometric prompts) within a single architecture for segmentation tasks. They overlap in addressing promptable concept segmentation with text and image-based queries, enabling zero-shot and few-shot segmentation across diverse visual concepts. However, SAM 3 differs by introducing a dual detector-tracker architecture with memory-based video tracking, a presence head for decoupled recognition-localization, and a large-scale data engine producing 4M unique concept labels, while CLIPSeg focuses on a lightweight transformer decoder built on frozen CLIP encoders with visual prompt engineering techniques for combining masks and images.

2. Segment everything everywhere all at once

Authors: Zou, Xueyan, Yang Jianwei, Xueyan Zou, Zhang Hao, et al. (19 authors total) | **Year/Venue:** 2023 | **URL:** [View paper](#)

Abstract

In this work, we present SEEM, a promptable and interactive model for segmenting everything everywhere all at once in an image, as shown in Fig.1. In SEEM, we propose a novel decoding mechanism that enables diverse prompting for all types of segmentation tasks, aiming at a universal segmentation interface that behaves like large language models (LLMs). More specifically, SEEM is designed with four desiderata: i) Versatility. We introduce a new visual prompt to unify different spatial queries inc...

Relationship Analysis

Both papers belong to the unified promptable segmentation frameworks category, aiming to handle multiple prompt types within a single architecture for image and video segmentation. They overlap in supporting diverse prompts (text, points, boxes, masks, image exemplars) and addressing both image and video segmentation tasks with interactive refinement capabilities. However, SAM 3 focuses on promptable concept segmentation with a decoupled detector-tracker architecture, a novel presence head for recognition-localization decoupling, and a large-scale data engine producing 4M unique concept labels, while SEEM emphasizes a unified decoder with memory prompts for interactive refinement, compositional prompting through joint visual-semantic space alignment, and semantic-aware segmentation without the specialized data engine or concept-centric focus of SAM 3.

3. UniVS: Unified and Universal Video Segmentation with Prompts as Queries

Authors: Minghan Li, Shuai Li, Ming-hui Li, Xindong Zhang, Lei Zhang | **Year/Venue:** 2024 • Computer Vision and Pattern Recognition | **URL:** [View paper](#)

Abstract

Despite the recent advances in unified image segmentation (IS), developing a unified video segmentation (VS) model remains a challenge. This is mainly because generic category-specified VS tasks need to detect all objects and track them across consecutive frames, while prompt-guided VS tasks require re-identifying the target with visual/text prompts throughout the entire video, making it hard to handle the different tasks with the same architecture. We make an attempt to address these issues and...

Relationship Analysis

Both papers belong to the unified promptable segmentation frameworks category, aiming to handle multiple prompt types within a single architecture for image and video segmentation. SAM 3 focuses on concept-based prompting using noun phrases and image exemplars to detect and segment all instances of a concept across images and videos, introducing a decoupled detector-tracker architecture with a novel presence head for open-vocabulary recognition. UniVS, in contrast, unifies different video segmentation tasks (VIS, VSS, VPS, VOS, RefVOS, PVOS) by using prompts as queries and converting predicted masks into pseudo visual prompts for tracking, eliminating the need for explicit inter-frame matching but focusing more on task unification rather than open-vocabulary concept segmentation with hard negatives.

Contributions Analysis

Overall novelty summary. The paper introduces Promptable Concept Segmentation (PCS), a task that unifies detection, segmentation, and tracking across images and videos using concept prompts (noun phrases, image exemplars, or both). It resides in the 'Unified Promptable Segmentation Frameworks' leaf alongside three sibling papers: Text Image Prompts, ControlNet-based methods, and UniVS Prompts Queries. This leaf represents a moderately populated research direction within the broader Foundation Model Architectures branch, which anchors the field's core segmentation engines. The taxonomy reveals that unified frameworks constitute one of sixteen leaf nodes across fifty papers, indicating a well-established but not overcrowded area.

The taxonomy tree shows that neighboring leaves include Video Object Segmentation with Memory Mechanisms (four papers on temporal tracking) and Attention and Diffusion-Based Segmentation Mechanisms (three papers on cross-modal alignment). The paper's integration of memory-based video tracking connects it to the video segmentation branch, while its decoupled recognition-localization architecture relates to attention mechanism research. The scope notes clarify that unified frameworks handle multiple prompt types within single architectures, distinguishing them from single-modality methods or domain-specific adaptations found in medical and remote sensing branches. This positioning suggests the work bridges foundational architecture design with video-specific temporal reasoning.

Among thirty candidates examined, the architecture contribution (decoupled recognition and localization) shows overlap with three prior works, while the PCS task formulation and data engine contributions appear more distinctive with zero refutable candidates each from ten examined papers. The limited search scope means these statistics reflect top-K semantic matches rather than exhaustive coverage. The architecture's presence head for decoupling recognition and localization appears to have precedent in the examined literature, whereas the concept-level prompt formulation and the human-AI data engine methodology show less direct overlap within the candidate set. The SA-Co benchmark contribution also lacks clear refutation among examined candidates.

Based on the limited thirty-candidate search, the work appears to make incremental architectural contributions while introducing a novel task formulation and benchmark. The taxonomy context reveals a moderately active research area with established sibling frameworks,

suggesting the field has matured beyond initial exploration but remains open to refinement. The analysis does not cover exhaustive prior work in video segmentation or concept-based retrieval systems outside the top semantic matches, leaving open questions about broader novelty claims.

This paper presents **3 main contributions**, each analyzed against relevant prior work:

Contribution 1: Promptable Concept Segmentation (PCS) task and SA-Co benchmark

Description: The authors introduce a new task called Promptable Concept Segmentation where users provide concept prompts (short noun phrases, image exemplars, or both) to detect, segment, and track all matching object instances in images and videos. They also create the SA-Co benchmark containing 214K unique concepts with exhaustive masks in 124K images and 1.7K videos.

This contribution was assessed against **10 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

1. Extending CLIP's Image-Text Alignment to Referring Image Segmentation

URL: [View paper](#)

Brief Assessment

Extending CLIP Referring[69] focuses on referring image segmentation (RIS) using natural language expressions to segment specific instances, not the broader PCS task of detecting and segmenting all matching instances of concepts across images and videos with noun phrases and exemplars.

2. Consensus-aware visual-semantic embedding for image-text matching

URL: [View paper](#)

Brief Assessment

Consensus Visual-Semantic[64] focuses on image-text matching using consensus knowledge from captioning corpora, not on concept segmentation with prompts or instance tracking in images/videos.

3. Referring Image Segmentation via Text Guided Multi-Level Interaction

URL: [View paper](#)

Brief Assessment

Text Multi-Level Interaction[70] focuses on referring image segmentation using natural language expressions with multi-level feature fusion, not on promptable concept segmentation with short noun phrases, image exemplars, or exhaustive instance tracking across images and videos.

4. Citetracker: Correlating image and text for visual tracking

URL: [View paper](#)

Brief Assessment

CiteTracker[61] focuses on visual object tracking using text descriptions generated from image patches, not on concept segmentation with exhaustive instance detection and tracking across images/videos as in the original paper's PCS task.

5. Tracking-forced Referring Video Object Segmentation

URL: [View paper](#)

Brief Assessment

Tracking-forced Referring[68] focuses on referring video object segmentation (RVOS) with language expressions for single target objects, not the open-vocabulary concept segmentation task with exhaustive instance detection and tracking across all matching objects.

6. Dual-level information interactive learning model for text-image person Re-identification

URL: [View paper](#)

Brief Assessment

Dual-level Person Re-ID[67] focuses on text-image person re-identification with instance-level and class-level learning, not on concept segmentation with text/image prompts for detecting and tracking object instances across images and videos.

7. Referring Video Object Segmentation With Cross-Modality Proxy Queries

URL: [View paper](#)

Brief Assessment

Cross-Modality Proxy[66] focuses on referring video object segmentation (RVOS) with textual expressions for single object tracking, not the open-vocabulary concept segmentation task with exhaustive instance detection across all matching objects using noun phrases and image exemplars.

8. Linguistic query-guided mask generation for referring image segmentation

URL: [View paper](#)

Brief Assessment

Linguistic Query Mask[62] focuses on referring image segmentation with single-instance text queries, not multi-instance concept detection and tracking across images/videos with noun phrases and exemplars.

9. Clip2: Contrastive language-image-point pretraining from real-world point cloud data

URL: [View paper](#)

Brief Assessment

CLIP2 Point Cloud[63] focuses on contrastive pretraining for 3D point cloud representation learning with language-image-point alignment, not on promptable concept segmentation with text/image exemplars for instance tracking in images and videos.

10. Lewis3d: Language-driven open-world instance-level 3d scene understanding

URL: [View paper](#)

Brief Assessment

Lewis3D[65] focuses on 3D scene understanding using point clouds and multi-view images, not 2D image/video concept segmentation with text and image exemplar prompts as in the original paper.

Contribution 2: SAM 3 architecture with decoupled recognition and localization

Description: The authors present SAM 3, a unified model consisting of a detector and tracker that share a vision encoder. A key architectural innovation is the presence head that decouples recognition (what) from localization (where), which significantly improves detection accuracy especially when training with challenging negative phrases.

This contribution was assessed against **10 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

1. Opensd: Unified open-vocabulary segmentation and detection

URL: [View paper](#)

Brief Assessment

[Final Audit Failure] The model insisted on a refutation claim but failed to provide verifiable evidence after multiple retries. Marked as cannot_refute for safety. Please manually verify the candidate text.

2. OV-DAR: Open-Vocabulary Object Detection and Attributes Recognition

URL: [View paper](#)

Brief Assessment

OV-DAR[56] focuses on decoupling detection and recognition for open-vocabulary object detection with attributes, while SAM 3's presence head decouples recognition (what) from localization (where) within a unified segmentation framework. The architectural approaches and task formulations differ substantially.

3. Generalized decoupled learning for enhancing open-vocabulary dense perception

URL: [View paper](#)

Brief Assessment

Generalized Decoupled Learning[52] focuses on enhancing CLIP for open-vocabulary dense perception by decoupling self-attention into content and context features, not on detector-tracker architectures with presence heads for concept segmentation.

4. Ultralytics YOLO Evolution: An Overview of YOLO26, YOLO11, YOLOv8 and YOLOv5 Object Detectors for Computer Vision and Pattern Recognition

URL: [View paper](#)

Brief Assessment

YOLO Evolution[59] focuses on object detection architectures in the YOLO family (YOLOv5-YOLO26), not on segmentation models with decoupled presence heads for open-vocabulary concept detection as in SAM 3.

5. LLaMA-Unidetector: An LLaMA-Based Universal Framework for Open-Vocabulary Object Detection in Remote Sensing Imagery

URL: [View paper](#)

Prior Art Analysis

LLaMA-Unidetector[58] demonstrates that the concept of decoupling recognition and localization in object detection architectures was already established prior to SAM 3. The candidate paper explicitly describes a 'decoupled learning strategy that separates localization and recognition' where a class-agnostic detector handles localization (distinguishing foreground/background) while a separate component handles recognition. This architectural approach directly parallels SAM 3's claimed innovation of using a presence head to decouple 'what' (recognition) from 'where' (localization), suggesting this design principle was not novel to SAM 3.

Evidence

Evidence 1 - **Rationale:** Both papers describe decoupling recognition and localization as a core architectural innovation. The candidate explicitly states this separation as a key feature of their framework. - **Original:** recognition and localization are decoupled with a presence head, which significantly boosts detection accuracy - **Candidate:** we introduce llama-unidetector, a universal framework that incorporates textual information into a closed-set detector, enabling the generalization to open-set scenarios. our llama-unidetector leverages a decoupled learning strategy that separates localization and recognition

Evidence 2 - **Rationale:** Both papers articulate the same fundamental problem and solution: separating the localization task (identifying where objects are) from the recognition task (identifying what they are). The candidate implements this through a two-stage pipeline, while the original uses a presence token, but the core architectural principle of decoupling these functions is identical. - **Original:** it can be difficult for each of the proposal queries to both recognize (what) and localize (where) an object in the image/frame. for the recognition component, contextual cues from the entire image are important. however, forcing proposal queries to understand the global context can be counterproduc... - **Candidate:** in the first stage, a class-agnostic detector identifies objects, distinguishing only between foreground and background. in the second stage, the detected foreground objects are passed through terraov-llm, a multimodal large language model (mllm), for recognition

6. What makes good open-vocabulary detector: A disassembling perspective

URL: [View paper](#)

Brief Assessment

[Final Audit Failure] The model insisted on a refutation claim but failed to provide verifiable evidence after multiple retries. Marked as cannot_refute for safety. Please manually verify the candidate text.

7. Unbiased Region-Language Alignment for Open-Vocabulary Dense Prediction

URL: [View paper](#)

Prior Art Analysis

Unbiased Region-Language[60] demonstrates prior work on decoupling recognition and localization in open-vocabulary detection architectures. The candidate paper explicitly addresses the same architectural challenge of separating 'what' (recognition) from 'where' (localization) through a presence head mechanism. Both papers identify that forcing object queries to handle both recognition and localization creates conflicts, and both propose decoupling these tasks. The candidate's 'presence token' approach for global-level classification separate from local-level localization directly parallels SAM 3's 'presence head' that decouples recognition from localization.

Evidence

Evidence 1 - **Rationale:** Both papers propose decoupling mechanisms that separate global recognition (presence/classification) from local localization. The candidate's approach of decoupling foreground/background classification from region localization demonstrates the same architectural principle was applied in prior work. - **Original:** we decouple the recognition and localization steps by introducing a learned globalpresence token. this token is solely responsible for predicting whether the target concept in the form of a noun phrase

(np) is present in the image/frame, i.e. p(np is present in input). each proposal query q_i only ne... - **Candidate:** to address this bias, we propose aligning foreground and background regions separately, ensuring explicit semantic separation through distinct category sets. to achieve this, we introduce densevlm, an end-to-end framework designed for unbiased region-language alignment... a key feature of densevlm i...

8. A simple framework for open-vocabulary segmentation and detection

URL: [View paper](#)

Prior Art Analysis

Simple Open-Vocabulary[51] demonstrates prior work on decoupling recognition and localization in open-vocabulary detection architectures. The candidate paper explicitly describes a 'presence head' that decouples recognition (what) from localization (where) in their decoder architecture, addressing the same fundamental challenge that SAM 3 claims as novel. Both papers identify that forcing proposal queries to understand global context conflicts with the local nature of localization objectives, and both propose separating these concerns through dedicated architectural components.

Evidence

Evidence 1 - **Rationale:** Both papers provide empirical validation showing significant performance improvements from their decoupling approaches, demonstrating that the effectiveness of this architectural choice was already established in prior work. - **Original:** the presence head boosts cgf1 by +5.7 (9a), mainly improving image-level recognition ability measured by il mcc - **Candidate:** comparing the top two rows, we can find mask decoding conditioned on the gt concept and box significantly improves the quality (mask ap from 8.6 to 46.4), which even reaches a similar level to coco (46.4v.s. 53.2)

9. Multi-modal Prompts with Feature Decoupling for Open-Vocabulary Object Detection

URL: [View paper](#)

Brief Assessment

Multi-modal Feature Decoupling[54] focuses on open-vocabulary object detection with multi-modal prompts and feature decoupling, but the limited context provided does not contain sufficient technical details to assess whether it demonstrates prior work on decoupling recognition and localization in the same architectural manner as SAM 3's presence head approach.

10. Declip: Decoupled learning for open-vocabulary dense perception

URL: [View paper](#)

Brief Assessment

DeCLIP[53] focuses on decoupling self-attention in CLIP for open-vocabulary dense prediction tasks (detection/segmentation), not on decoupling recognition from localization in a unified detector-tracker architecture like SAM 3's presence head design.

Contribution 3: Scalable human- and AI-in-the-loop data engine

Description: The authors develop a data engine that iteratively generates annotated data through a feedback loop involving SAM 3, human annotators, and AI annotators (fine-tuned MLLMs). This engine produces high-quality training data with 4M unique phrases and 52M masks, plus synthetic data with 38M phrases and 1.4B masks, while doubling annotation throughput by delegating verification tasks to AI models.

This contribution was assessed against **10 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

1. An Integrated in Situ Image Acquisition and Annotation Scheme for Instance Segmentation Models in Open Scenes With a Human-Robot Interaction Approach

URL: [View paper](#)

Brief Assessment

Robot Interaction Annotation[78] focuses on eye-tracking-based annotation for robotic field applications, not on iterative data engine design with AI verifiers for segmentation at scale.

2. MedUHIP: Towards Human-In-the-Loop Medical Segmentation

URL: [View paper](#)

Brief Assessment

MedUHIP[73] focuses on medical image segmentation with human-in-the-loop interaction for refining individual segmentations, not on building a scalable data annotation engine for generating large-scale training datasets.

3. Understanding Novice's Annotation Process For 3D Semantic Segmentation Task With Human-In-The-Loop

URL: [View paper](#)

Brief Assessment

Novice Annotation 3D[75] focuses on human-in-the-loop annotation for 3D point cloud semantic segmentation with novice annotators, not on building a scalable data engine for 2D image/video segmentation with iterative model training and AI verifiers.

4. A unified microstructure segmentation approach via human-in-the-loop machine learning

URL: [View paper](#)

Brief Assessment

Unified Microstructure HITL[77] focuses on microstructure segmentation in materials science, not general-purpose visual segmentation. The domain and application are fundamentally different from SAM 3's open-vocabulary concept segmentation.

5. Reducing human annotation effort using self-supervised learning for image segmentation

URL: [View paper](#)

Brief Assessment

The candidate paper (Self-supervised Annotation Reduction[71]) focuses on reducing human annotation effort through self-supervised learning for image segmentation, not on building a comprehensive data engine with human and AI annotators for generating large-scale training data with diverse concepts and hard negatives as described in the original contribution.

6. Open-world point cloud semantic segmentation: A human-in-the-loop framework

URL: [View paper](#)

Brief Assessment

Open-world Human-in-Loop[72] focuses on interactive point cloud segmentation with human annotations during inference, not on building a scalable data annotation engine for training data generation. The candidate addresses a different problem domain (3D point cloud segmentation) with a different methodology (test-time human feedback for segmentation) rather than training data creation.

7. AI-human interactive pipeline with feedback to accelerate medical image annotation

URL: [View paper](#)

Brief Assessment

AI-human Feedback Pipeline[80] focuses on medical image annotation (prostate MRI segmentation) with a simple three-step iterative loop, whereas the original paper develops a comprehensive multi-phase data engine for open-vocabulary segmentation with sophisticated AI verifiers, hard negative mining, and ontology-based concept curation across diverse visual domains.

8. Towards Clinician-Preferred Segmentation: Leveraging Human-in-the-Loop for Test Time Adaptation in Medical Image Segmentation

URL: [View paper](#)

Brief Assessment

Clinician-Preferred Segmentation[76] focuses on test-time adaptation for medical image segmentation using clinician corrections, not on building a scalable data annotation engine for generating training datasets with human and AI annotators.

9. Transfer for medical image segmentation via iterative human-in-the-loop update: from labelled public to unlabelled clinical datasets for multi-organ segmentation in CT

URL: [View paper](#)

Brief Assessment

Transfer Human-in-Loop CT[79] focuses on rapid model transfer for medical image segmentation using human-in-the-loop for annotation refinement, not on building a scalable data engine for generating large-scale training datasets with diverse concepts and hard negatives.

10. LabelAny3d: Label any object 3d in the wild

URL: [View paper](#)

Brief Assessment

LabelAny3D[74] focuses on 3D bounding box annotation for monocular 3D detection using reconstruction-based methods, not on iterative data generation with human and AI annotators for segmentation tasks.

Appendix: Text Similarity Detection

No high-similarity text segments were detected across any compared papers.

References

- [0] SAM 3: Segment Anything with Concepts [View paper](#)
- [1] Segment anything model for medical image segmentation: Current applications and future directions [View paper](#)
- [2] Promise: Prompt-driven 3d medical image segmentation using pretrained image foundation models [View paper](#)
- [3] Prompt-driven referring image segmentation with instance contrasting [View paper](#)
- [4] Reasoning Segmentation for Images and Videos: A Survey [View paper](#)
- [5] Image segmentation using text and image prompts [View paper](#)
- [6] Soseg: Prompt-based small object instance segmentation in remote sensing imagery [View paper](#)
- [7] VISA: Reasoning Video Object Segmentation via Large Language Models [View paper](#)
- [8] Medical SAM 2: Segment medical images as video via Segment Anything Model 2 [View paper](#)
- [9] Towards training-free open-world segmentation via image prompt foundation models [View paper](#)
- [10] MWVOS: Mask-free weakly supervised video object segmentation via promptable foundation model [View paper](#)
- [11] An Interactive Prompt Based Network for Urban Floods Area Segmentation Using UAV Images [View paper](#)
- [12] Learnable prompt for few-shot semantic segmentation in remote sensing domain [View paper](#)
- [13] Prompt-driven dynamic object-centric learning for single domain generalization [View paper](#)
- [14] Seg4diff: Unveiling open-vocabulary segmentation in text-to-image diffusion transformers [View paper](#)
- [15] SAM-I2V: Upgrading SAM to Support Promptable Video Segmentation with Less than 0.2% Training Cost [View paper](#)
- [16] Segment everything everywhere all at once [View paper](#)
- [17] Zero-shot surgical tool segmentation in monocular video using Segment Anything Model 2 [View paper](#)
- [18] Depthwise-dilated convolutional adapters for medical object tracking and segmentation using the segment anything model 2 [View paper](#)
- [19] nninteractive: Redefining 3d promptable segmentation [View paper](#)
- [20] Learning to prompt segment anything models [View paper](#)
- [21] Driving Referring Video Object Segmentation with Vision-Language Pre-trained Models [View paper](#)
- [22] Promptable cancer segmentation using minimal expert-curated data [View paper](#)
- [23] Unleashing text-to-image diffusion models for visual perception [View paper](#)
- [24] Interactive Segmentation Model for Placenta Segmentation from 3D Ultrasound images [View paper](#)
- [25] Training-free open-ended object detection and segmentation via attention as prompts [View paper](#)
- [26] Attention Prompt-Driven Source-Free Adaptation for Remote Sensing Images Semantic Segmentation [View paper](#)
- [27] Rethinking interactive image segmentation with low latency high quality and diverse prompts [View paper](#)
- [28] Unified Multi-Modality Video Object Segmentation Using Reinforcement Learning [View paper](#)
- [29] Adding Conditional Control to Text-to-Image Diffusion Models [View paper](#)
- [30] Segprompt: Boosting open-world segmentation via category-level prompt learning [View paper](#)
- [31] Personalize Segment Anything Model with One Shot [View paper](#)
- [32] DiffRIS: Enhancing Referring Remote Sensing Image Segmentation with Pre-trained Text-to-Image Diffusion Models [View paper](#)
- [33] UniVS: Unified and Universal Video Segmentation with Prompts as Queries [View paper](#)
- [34] Subject-Diffusion: Open Domain Personalized Text-to-Image Generation without Test-time Fine-tuning [View paper](#)
- [35] One-prompt to segment all medical images [View paper](#)
- [36] Faster Interactive Segmentation of Identical-Class Objects With One Mask in High-Resolution Remotely Sensed Imagery [View paper](#)

- [37] Llava-interactive: An all-in-one demo for image chat, segmentation, generation and editing [View paper](#)
- [38] Prompting segmentation with sound is generalizable audio-visual source localizer [View paper](#)
- [39] Prompting to Adapt Foundational Segmentation Models [View paper](#)
- [40] Poda: Prompt-driven zero-shot domain adaptation [View paper](#)
- [41] An unsupervised approach towards promptable defect segmentation in laser-based additive manufacturing by segment anything [View paper](#)
- [42] From text to mask: Localizing entities using the attention of text-to-image diffusion models [View paper](#)
- [43] Multimodal Prompt-Guided Bidirectional Fusion for Referring Remote Sensing Image Segmentation [View paper](#)
- [44] Promptonomyvit: Multi-task prompt learning improves video transformers using synthetic scene data [View paper](#)
- [45] Progressive Self-Prompting Segment Anything Model for Salient Object Detection in Optical Remote Sensing Images [View paper](#)
- [46] Biomedical SAM 2: Segment Anything in Biomedical Images and Videos [View paper](#)
- [47] Visual Prompt Selection Framework for Real-Time Object Detection and Interactive Segmentation in Augmented Reality Applications [View paper](#)
- [48] Performance and nonadversarial robustness of the segment anything model 2 in surgical video segmentation [View paper](#)
- [49] Enhancing Label-efficient Medical Image Segmentation with Text-guided Diffusion Models [View paper](#)
- [50] SkipClick: Combining Quick Responses and Low-Level Features for Interactive Segmentation in Winter Sports Contexts [View paper](#)
- [51] A simple framework for open-vocabulary segmentation and detection [View paper](#)
- [52] Generalized decoupled learning for enhancing open-vocabulary dense perception [View paper](#)
- [53] Declip: Decoupled learning for open-vocabulary dense perception [View paper](#)
- [54] Multi-modal Prompts with Feature Decoupling for Open-Vocabulary Object Detection [View paper](#)
- [55] Opensd: Unified open-vocabulary segmentation and detection [View paper](#)
- [56] OV-DAR: Open-Vocabulary Object Detection and Attributes Recognition [View paper](#)
- [57] What makes good open-vocabulary detector: A disassembling perspective [View paper](#)
- [58] LLaMA-Unidetector: An LLaMA-Based Universal Framework for Open-Vocabulary Object Detection in Remote Sensing Imagery [View paper](#)
- [59] Ultralytics YOLO Evolution: An Overview of YOLO26, YOLO11, YOLOv8 and YOLOv5 Object Detectors for Computer Vision and Pattern Recognition [View paper](#)
- [60] Unbiased Region-Language Alignment for Open-Vocabulary Dense Prediction [View paper](#)
- [61] Citetracker: Correlating image and text for visual tracking [View paper](#)
- [62] Linguistic query-guided mask generation for referring image segmentation [View paper](#)
- [63] Clip2: Contrastive language-image-point pretraining from real-world point cloud data [View paper](#)
- [64] Consensus-aware visual-semantic embedding for image-text matching [View paper](#)
- [65] Lowis3d: Language-driven open-world instance-level 3d scene understanding [View paper](#)
- [66] Referring Video Object Segmentation With Cross-Modality Proxy Queries [View paper](#)
- [67] Dual-level information interactive learning model for text-image person Re-identification [View paper](#)
- [68] Tracking-forced Referring Video Object Segmentation [View paper](#)
- [69] Extending CLIP's Image-Text Alignment to Referring Image Segmentation [View paper](#)
- [70] Referring Image Segmentation via Text Guided Multi-Level Interaction [View paper](#)
- [71] Reducing human annotation effort using self-supervised learning for image segmentation [View paper](#)
- [72] Open-world point cloud semantic segmentation: A human-in-the-loop framework [View paper](#)
- [73] MedUHIP: Towards Human-In-the-Loop Medical Segmentation [View paper](#)
- [74] Labelany3d: Label any object 3d in the wild [View paper](#)
- [75] Understanding Novice's Annotation Process For 3D Semantic Segmentation Task With Human-In-The-Loop [View paper](#)
- [76] Towards Clinician-Preferred Segmentation: Leveraging Human-in-the-Loop for Test Time Adaptation in Medical Image Segmentation [View paper](#)
- [77] A unified microstructure segmentation approach via human-in-the-loop machine learning [View paper](#)
- [78] An Integrated in Situ Image Acquisition and Annotation Scheme for Instance Segmentation Models in Open Scenes With a Human-Robot Interaction Approach [View paper](#)
- [79] Knowledge transfer for medical image segmentation via iterative human-in-the-loop update: from labelled public to unlabelled clinical datasets for multi-organ segmentation in CT [View paper](#)
- [80] AI-human interactive pipeline with feedback to accelerate medical image annotation [View paper](#)