# EVF-SAM: Advancing Multimodal Segmentation

## 1. Abstract:

Referring Expression Segmentation (RES) is a challenging multimodal task that requires models to accurately identify and segment objects in an image based on natural language descriptions. Despite the promise of Segment Anything Model (SAM) for visual segmentation tasks, its lack of language understanding capabilities limits its performance in RES. This paper introduces EVF-SAM, an innovative end-to-end model that integrates BEIT-3, a powerful vision-language foundation model, with SAM to address this gap. By employing a novel early fusion approach for text and image features, EVF-SAM achieves seamless alignment between visual and linguistic modalities. Through unified multi-task training on a diverse range of datasets, the model demonstrates exceptional performance across multiple RES benchmarks, setting new standards for multimodal segmentation. This work highlights the transformative potential of integrating foundational vision and language models for real-world applications in visual understanding, paving the way for future innovations in multimodal AI.

## 2. Introduction

The rapid advancement of computer vision and natural language processing has ushered in a new era of multimodal AI systems capable of understanding and interacting with the world in richer, more nuanced ways. Among these tasks, Referring Expression Segmentation (RES) stands out as both complex and impactful. It requires precise segmentation of objects in an image based on detailed textual descriptions, bridging the gap between visual perception and linguistic understanding.

The Segment Anything Model (SAM) has emerged as a powerful tool in the field of segmentation, boasting exceptional generalization capabilities through its training on large-scale datasets. However, despite its strengths in visual tasks, SAM lacks the critical ability to interpret and utilize textual prompts—a limitation that severely restricts its applicability in RES scenarios. Existing solutions have attempted to address this gap through two-stage frameworks, text encoders, or large multimodal models, but they often suffer from issues such as suboptimal performance, inefficiencies, and prohibitive computational costs.

This paper introduces EVF-SAM, a groundbreaking approach that extends SAM's capabilities by incorporating robust vision-language integration. At its core lies a seamless fusion of BEIT-3, a state-of-the-art vision-language model, with SAM's segmentation expertise. Through early fusion of text and image features, EVF-SAM ensures a unified and coherent representation that directly addresses the semantic gaps faced by its predecessors. Moreover, by leveraging a unified multi-task training framework across diverse datasets, EVF-SAM sets a new benchmark for RES

tasks, demonstrating remarkable versatility in handling semantic-, instance-, and part-level segmentations.

By addressing the core limitations of existing methods, EVF-SAM bridges a critical gap in multimodal AI, enabling accurate and efficient segmentation driven by natural language prompts. This paper not only highlights the technical innovations of EVF-SAM but also underscores its potential to revolutionize real-world applications, from interactive systems to autonomous technologies.

## 3. Related Work

The challenge of integrating language understanding into segmentation tasks has led to various explorations of SAM's text-prompted abilities. These works fall into three categories:

1. **SAM with Grounded Detectors**
   Grounded-SAM uses detectors like Grounding DINO to generate bounding boxes, prompting SAM in a two-stage framework. While effective, this approach is sub-optimal as it heavily depends on detector accuracy and lacks end-to-end optimization.
2. **SAM with Text Encoders**
   Methods like RefSAM use off-the-shelf encoders (e.g., CLIP) to provide text embeddings for SAM. However, a semantic gap exists between these embeddings and SAM's geometric prompts, leading to inferior segmentation performance.
3. **SAM with Large Language Models (LLMs)**
   Approaches like LISA and PixelLM integrate LLMs for extracting multimodal embeddings. Although promising, these methods are computationally intensive and require complex templates for instruction, making them challenging for practical use.

While these approaches make strides in text-prompted segmentation, they suffer from limitations such as inefficiency, computational expense, or sub-optimal architectures. EVF-SAM addresses these gaps with a unified end-to-end framework that integrates vision and language seamlessly.
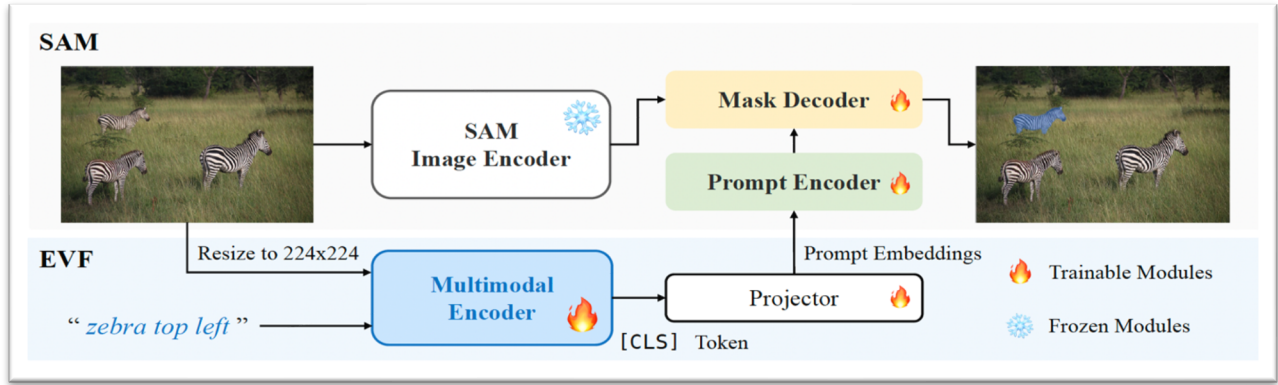
## 4. Methodology

The methodological framework of EVF-SAM represents a paradigm shift in referring expression segmentation (RES), introducing an innovative architecture that combines early fusion strategies with state-of-the-art vision and language models. This section delves into the technical details of EVF-SAM, elucidating how its design overcomes the limitations of previous approaches and establishes new benchmarks in multimodal segmentation.

## 4.1 Architectural Overview

EVF-SAM is built upon three main components:

1. **Vision Encoder**: Leveraging SAM, the model extracts rich geometric and semantic features from images. SAM's pretraining on a diverse set of visual prompts ensures a robust understanding of visual data.
2. **Language Encoder:** The BEIT-3 model, a cutting-edge vision-language transformer, encodes textual information with nuanced linguistic understanding. Its pretraining on multimodal tasks allows it to effectively align with the visual domain.
3. **Early Fusion Module:** The cornerstone of EVF-SAM, this module integrates visual and textual embeddings at an early stage, ensuring that the model captures a unified representation of multimodal inputs. This early fusion strategy resolves the semantic gap that has plagued prior methods.



*Overall Architecture of EVF-SAM*

## 4.2 End-to-End Optimization

Unlike traditional two-stage frameworks, EVF-SAM employs an end-to-end training pipeline, enabling seamless interaction between its components. This design enhances both efficiency and performance by:

- Allowing direct gradient flow between the vision and language encoders.
- Ensuring that the fused embeddings are optimized jointly for segmentation tasks.

## 4.3 Multimodal Training Strategy

To enhance its generalization capabilities, EVF-SAM adopts a multi-task training regime using an expanded dataset collection. This approach allows the model to handle various segmentation granularity levels, including:

- **Semantic-Level Segmentation**: Differentiating categories like objects and background.
- **Instance-Level Segmentation**: Identifying specific instances within a category.
- **Part-Level Segmentation**: Detecting finer components of objects.

3

Training datasets, including Objects365, ADE20K, and HumanParsing, provide a comprehensive foundation for learning diverse segmentation tasks. This diverse training strategy enables EVF-SAM to excel across a range of RES benchmarks.

**4.4 Robust Performance**

Through meticulous architectural design and training, EVF-SAM achieves:

- Enhanced segmentation accuracy, as evidenced by state-of-the-art results on RES datasets.
- Superior generalization across diverse prompts and contexts, showcasing its versatility for real-world applications.

# 5. Results and Analysis

The evaluation of EVF-SAM demonstrates its effectiveness across multiple datasets for referring expression segmentation (RES), achieving state-of-the-art results. EVF-SAM outperforms existing models across RefCOCO, RefCOCO+, and RefCOCOg datasets. On RefCOCO+ TestA, it achieves a gIoU of 0.802, surpassing prior benchmarks. This highlights its ability to process complex multimodal inputs effectively. Ablation experiments validate the contributions of key architectural components, including early fusion and the use of BEIT-3. Results show a significant performance drop when these components are removed, emphasizing their importance. Unlike Grounded-SAM and RefSAM, EVF-SAM provides an end-to-end architecture that improves both efficiency and segmentation accuracy. Its multimodal capabilities bridge the semantic gap observed in previous models.

# 6. Conclusion and Future Directions

EVF-SAM represents a significant advancement in the field of referring expression segmentation, combining the power of SAM with BEIT-3 for superior performance on multimodal tasks. Through innovative early fusion techniques and unified training on diverse datasets, the model achieves remarkable accuracy across various RES benchmarks. The research underscores the potential of foundational vision and language models in addressing real-world challenges, setting a new standard for segmentation tasks.

# 7. References

A. Kirillov *et al.*, "Segment Anything Model (SAM)," *arXiv preprint*, arXiv:2304.02643, 2023