

SAMWISE: Infusing Wisdom in SAM2 for Text-Driven Video Segmentation

Original

SAMWISE: Infusing Wisdom in SAM2 for Text-Driven Video Segmentation / Cuttano, Claudia; Trivigno, Gabriele; Rosi, Gabriele; Masone, Carlo; Averta, Giuseppe. - ELETTRONICO. - (2025), pp. 3395-3405. (2025 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Nashville (USA) 10-17 June 2025) [10.1109/CVPR52734.2025.00322].

Availability:

This version is available at: 11583/3004253 since: 2026-03-03T10:33:07Z

Publisher:

IEEE

Published

DOI:10.1109/CVPR52734.2025.00322

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

IEEE postprint/Author's Accepted Manuscript

©2025 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collecting works, for resale or lists, or reuse of any copyrighted component of this work in other works.

(Article begins on next page)

SAMWISE: Infusing Wisdom in SAM2 for Text-Driven Video Segmentation

Claudia Cattano¹ Gabriele Trivigno¹ Gabriele Rosi^{1,2} Carlo Masone^{1,2} Giuseppe Averta^{1,2}
¹ Politecnico di Torino ² Focoos AI
 {name.surname}@polito.it {name.surname}@focoos.ai

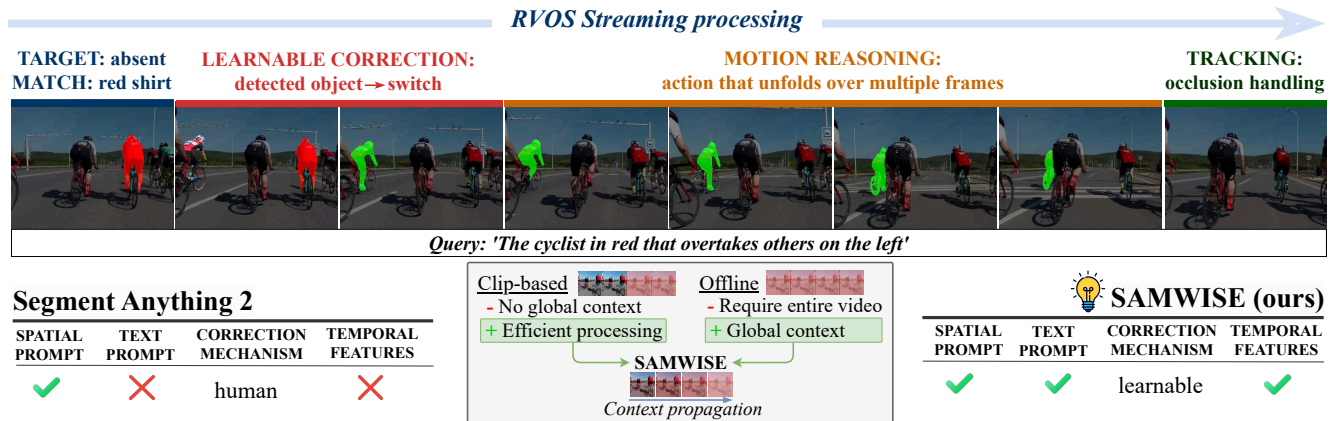


Figure 1. **SAMWISE**. Our approach infuses knowledge about natural language in the Segment-Anything 2 model, adding explicit temporal cues in the feature extraction for the task of streaming-based Referring Video Segmentation (RVOS). We use a learnable mechanism to mitigate the so-called *tracking bias*, *i.e.* SAM2 tendency to overlook a correct object once it becomes identifiable, due to its ongoing tracking of a different object. Our design enables effective streaming processing for RVOS, exploiting the memory from previous frames to propagate past context. The figure shows an example where the target object is not present in the first frame, leading SAM2 to start tracking the wrong one. Afterwards, when the correct object appears, our learnable correction mechanisms guides SAM2 to switch its tracking focus. By adding in its features the notion of temporal evolution, the model is able to recognize that the new object is more aligned with the provided textual query. Finally, we exploit SAM2 tracking skills and robustness to occlusions to keep following the object.

Abstract

Referring Video Object Segmentation (RVOS) relies on natural language expressions to segment an object in a video clip. Existing methods restrict reasoning either to independent short clips, losing global context, or process the entire video offline, impairing their application in a streaming fashion. In this work, we aim to surpass these limitations and design an RVOS method capable of effectively operating in streaming-like scenarios while retaining contextual information from past frames. We build upon the Segment-Anything 2 (SAM2) model, that provides robust segmentation and tracking capabilities and is naturally suited for streaming processing. We make SAM2 wiser, by empowering it with natural language understanding and explicit temporal modeling at the feature extraction stage, without fine-tuning its weights, and without outsourcing modality interaction to external models. To this end, we introduce a novel adapter module that injects temporal information

and multi-modal cues in the feature extraction process. We further reveal the phenomenon of tracking bias in SAM2 and propose a learnable module to adjust its tracking focus when the current frame features suggest a new object more aligned with the caption. Our proposed method, SAMWISE, achieves state-of-the-art across various benchmarks, by adding a negligible overhead of less than 5 M parameters. Code is available at <https://github.com/ClaudiaCattano/SAMWISE>.

1. Introduction

Referring video segmentation (RVOS) [9, 17, 24, 38, 43, 48] aims at segmenting and tracking specific objects of interest within video content, guided by natural language expressions [3, 10, 28]. Existing RVOS methods are mostly based on a *divide and conquer* paradigm, where the video is divided into shorter clips that are processed independently [3, 39, 43]. However, as demonstrated by MeViS [7], this

solution fails in examples that require taking into account long-term motion and global context. As a workaround to handle this challenge, the state-of-the-art method [11] processes the entire video in an *offline* fashion, first modeling trajectories of all instances throughout the entire clip and then selecting the most appropriate one. Albeit effective, this approach is not applicable when the model has access only to a portion of the video, for example when the data at inference time are presented in a streaming fashion or due to limitations in the computational resources. The trade-off of these two paradigms is schematized in Fig. 1. To this end, OnlineRefer [42] introduced a context propagation scheme for *online* RVOS but relies solely on past context from a single frame, limiting its ability to capture long-term dependencies. In this work, we investigate how to exploit the memory from past frames to design an RVOS method capable of retaining global context while operating within a streaming paradigm, *i.e.*, without requiring access to the whole video at once. This idea is inspired by the recent release of Segment-Anything 2 (SAM2) [35], a foundational model that has shown impressive capabilities in various Video Segmentation tasks thanks to a memory bank that allows to leverage long-range past information. Since SAM2 operates in a streaming fashion, extending this method to enable context-aware streaming processing in RVOS would appear a natural step. However, this entails some non-trivial challenges:

i) Text understanding. SAM2 original design accounts only for *spatial* prompts (e.g. points) and lacks mechanisms to interpret *semantic* prompts like text, which require reasoning over visual and textual modalities. While we are the first to address the challenge of adding textual prompts to SAM2, previous methods have explored this problem for SAM-1 at image-level. These solutions [19, 50] delegate visual-textual interaction to an off-the-shelf large VLM (like BEIT-3 [41], LLaVa [21]), which generates a multi-modal embedding that is used to prompt SAM-1.

ii) Temporal modeling. To segment the referred object throughout the video, it must be first *recognized* and then *tracked*. While the latter requires matching objects visual appearance across adjacent frame, the recognition problem entails modeling temporal evolution to reason over actions that unfold over multiple frames. However, SAM2 extracts frame features independently, lacking such reasoning.

iii) Tracking bias. In RVOS, the target object might be unrecognizable during certain time intervals, due to occlusions, presence of multiple instances or forthcoming actions, as in the first frames of Fig. 1. In such cases, SAM2 may start tracking an incorrect object that partially matches the textual prompt, and persist in following it, leading to what we denote as *tracking bias*. While SAM2 original design allows for a user to manually correct the prediction by providing a new prompt, such a strategy is not applicable in

tasks without a human-in-the-loop like RVOS.

In this work, we aim at making SAM2 *wiser*, by addressing these limitations without fine-tuning SAM2 weights, thereby preserving its original capabilities, and without outsourcing modality interaction to external, heavy models. To overcome challenges *i)* and *ii)*, we design a learnable Adapter [12] module, named Cross-Modal Temporal Adapter (CMT), with two key principles in mind: a) enabling mutual interaction between visual and linguistic modalities; and b) encoding temporal cues into visual features. Then, to generate a prompt, we follow [21, 51] and employ a learnable MLP to project the sentence embedding for the SAM2 Mask Decoder, which then outputs the final segmentation mask. In this way, we can exploit SAM2 tracking capability to segment an object given a textual query across the video. Finally, to mitigate the *tracking bias* problem *iii)*, we introduce a lightweight Conditional Memory Encoder (CME) which detects when a candidate object, aligned with the text, appears in the frame, thus enabling SAM2 to dynamically refocus its tracking to the correct object as it becomes distinguishable.

Summarizing, this paper contributes with the following:

- We present SAMWISE, the first method that integrates natural language knowledge into SAM2 in an end-to-end solution tailored to address the challenges of RVOS. We introduce a novel adapter, namely Cross Modal Temporal (CMT) Adapter, which purposefully models temporal evolution and multi-modal interaction;
- We provide insight into the functioning of SAM2, highlighting the phenomenon of *tracking bias*, and introduce a learnable module (Conditional Memory Encoder) to adjust tracking based on new information;
- Our methods achieves state-of-the-art results both on traditional RVOS benchmarks (Ref-Youtube-VOS [38], Ref-DAVIS [17]), as well as the more challenging MeViS [7], without compromising SAM2 capabilities and adding less than 5M learnable parameters.

2. Related works

Referring Video Segmentation. In RVOS, the goal is to segment an object in a clip described with natural language [8, 17, 24]. Earlier works adapted image-based methods [2, 9, 17, 48], or used a spatio-temporal memory to attend to masks of previous frames [31, 38]. Subsequent works employ a DETR-like [4] structure to process multiple frames and text embeddings [3, 10, 29, 43]. All these methods process short clips independently, thus losing global context.

Recently, [7] showed how traditional RVOS benchmarks lack challenging captions that require to disambiguate between instances and their actions, as well as occlusions and dynamic queries, highlighting how they could be solved even with image-based methods. The MeViS dataset [7] tar-

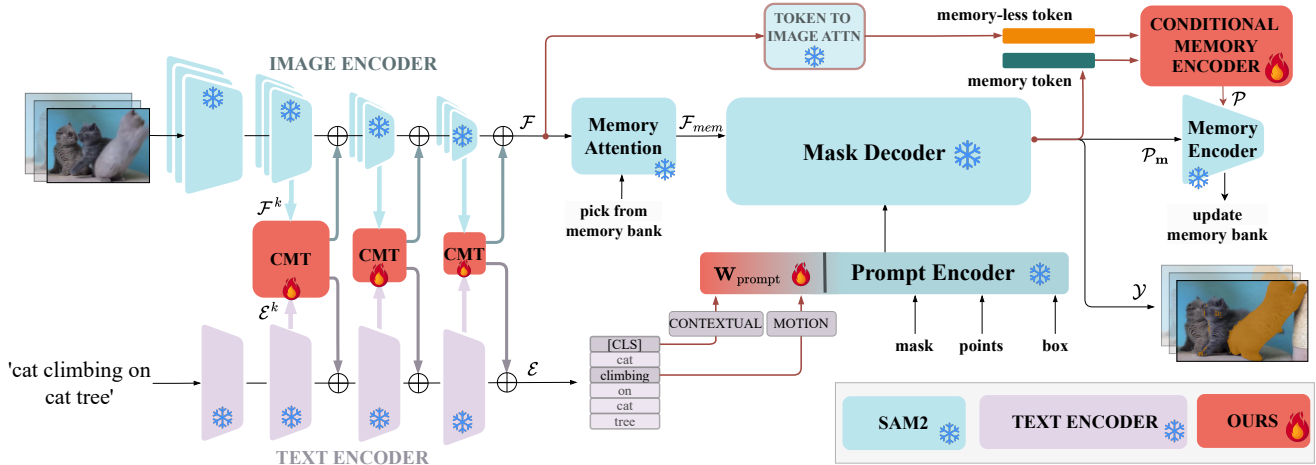


Figure 2. **Overview of SAMWISE.** We build on a frozen SAM2 and a frozen Text Encoder to segment images in video given a textual description. We incorporate the Cross-Modal Temporal Adapter (CMT) into the text and visual encoders at every intermediate layer k to model temporal dynamics within visual features while contaminating each modality with the other. Then, we extract the [CLS] and verb embeddings, namely Contextual and Motion prompts, from the adapted textual features and project them through a learnable MLP. The final embedding is used to prompt the Mask Decoder, which outputs the segmentation mask. Finally, the Conditional Memory Encoder detects when a new candidate object, aligned with the caption, appears in the frame, enabling SAM2 to dynamically refocus its tracking.

gets these scenarios, with challenging examples that previous image or clip-based methods fail to address. To this end, a few works proposed *offline* methods to explicitly model multiple object trajectories [11, 28], with the latter representing the state-of-the-art on MeViS. Concurrently, OnlineRefer [42] proposed a first attempt towards an *online* RVOS setting, with a query propagation scheme. However, its effectiveness is limited as predictions are based on a single frame. Our method builds on this paradigm by leveraging SAM2 memory bank to encode long-range past context.

Text-prompted Segment-Anything. Recent works have provided solutions to adapt SAM-1 for text-prompted segmentation. Grounded SAM [36] employs a two step pipeline where GroundingDINO [25] generates bounding boxes for SAM-1 to produce segmentation masks. Applying such pipeline in RVOS is problematic, as potential errors in the first frame are propagated throughout the whole video. To directly prompt SAM-1, RefSAM [20] exploits a projection layer to map the textual embedding into the prompt space, while [1, 19, 45] resort to large off-the-shelf VLM to generate a multi-modal embedding that is used to prompt SAM-1. Both solutions finetune the Mask Decoder, thereby compromising its capabilities on its original task. In contrast, our work is the first to propose an end-to-end model that incorporates textual knowledge within SAM2 without fine-tuning nor relying on external models.

Pre-Trained Knowledge Transfer. In recent years, the release of powerful pretrained models has sparked interest in the question of how to extend their skills to novel tasks, as full fine-tuning becomes increasingly impractical with growing model sizes [16, 33]. A powerful strat-

egy to address this problem relies on using Adapters [12], small trainable modules that enable efficient adaptation of pre-trained models. Following this paradigm, recent studies have explored adapting CLIP [34] for downstream tasks. At the image level, [44] inserts Transformer Decoder blocks within CLIP encoders, which entail costly Self-Attentions on all tokens. For video tasks, [40] places independent adapter modules within each encoder, whereas [15, 16, 27, 47] rely on a weight-sharing mechanism to project both modalities in a shared sub-space. Nevertheless, as features of each modality are independently extracted, none of these adapters allows explicit feature interaction, unlike our CMT, which also incorporates temporal modeling. Lastly, all these works start from a model that already includes a text encoder (CLIP), whereas ours is the first to propose an adapter for the Segment-Anything 2 model to add textual understanding, achieving robust performances while introducing less than 5 M parameters.

3. SAMWISE

Problem setting. Given an input video $\mathcal{V} = \{I_t\}_{t=1}^{T_V}$ with T_V frames and a referring expression, we aim to predict a set of binary masks $S = \{s_t\}_{t=1}^{T_V}$, $s_t \in \mathbb{R}^{H \times W}$ of the referred object. We tokenize the textual query in a set of L words, $E = \{e_l\}_{l=1}^L$, and add a global sentence representation token [CLS]. The tokens are then processed using a frozen text encoder to extract language features $\mathcal{E} \in \mathbb{R}^{L \times C_t}$. We process videos in a streaming fashion, collecting clips of T frames as they are available. Throughout the rest of the section, we use T to indicate clip length.

Overview. We first provide a brief discussion of the SAM2

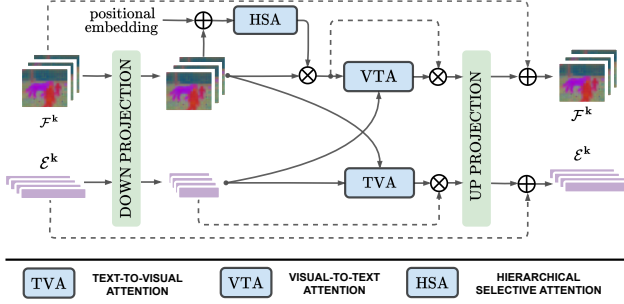


Figure 3. Architecture of our **Cross Modal Temporal (CMT)** Adapter, made up of Hierarchical Selective Attention (HSA) to model temporal cues, a Visual-to-Text Attention (VTA) and Text-to-Visual Attention (TVA) modules.

model (Sec. 3.1). We then outline the pipeline of our proposed SAMWISE, starting from the prompting strategy in Sec. 3.2. In Sec. 3.3, we detail our novel Cross-Modal Temporal Adapter. Lastly, in Sec. 3.5, we discuss our learnable correction strategy, named Conditional Memory Encoder, to address the issue of *tracking bias*.

3.1. Background: Segment-Anything

The Segment-Anything Model 2 (SAM2) builds upon SAM-1 [18] to tackle the task of Promptable Video Object Segmentation, *i.e.*, tracking an object in a video given a textual prompt. Following SAM-1, it consists of an *image encoder*, a *Prompt Encoder* and a *Mask Decoder*, which combines the image and prompt embeddings to predict segmentation masks. To enable video processing, SAM2 comes with a few modifications: *i*) the original ViT backbone is replaced by Hiera [37], roughly 3 times faster, which processes frames independently to provide hierarchical visual features. Hereinafter, we refer to them as *memory-less features* \mathcal{F} ; *ii*) frame embeddings are not directly fed to the Mask Decoder, but they are first *conditioned* on memories of past predictions from a *Memory Bank*. We refer to these conditioned features as *memory features* \mathcal{F}_{mem} . Lastly, *iii*) once the mask for the current frame is predicted, the *Memory Encoder* updates the Memory Bank. By design, SAM2 handles video frames as they become available, progressively encoding the past in its Memory Bank. We argue that this streaming approach is especially valuable in RVOS, enabling reasoning over a wide temporal horizon.

3.2. Prompting SAM2

To guide the SAM2 decoder, we use a Contextual Prompt, $\mathcal{E}_C \in \mathbb{R}^{1 \times C_t}$, which encodes the high-level semantic information for the given text query, emphasizing the essential aspects of the query while downplaying less relevant elements. To this end, we employ the [CLS] embedding of text features, \mathcal{E} . Furthermore, we also introduce a second prompt, the Motion Prompt $\mathcal{E}_M \in \mathbb{R}^{1 \times C_t}$, which cap-

tures action-related cues by using verb embeddings from \mathcal{E} . These prompts are concatenated and projected through a learnable three-layer MLP:

$$\rho = \mathbf{W}_{\text{prompt}}(\text{CAT}[\mathcal{E}_C, \mathcal{E}_M]). \quad (1)$$

In this way, the provided prompts encode both subject-related and motion-based information. Given that in our task the textual prompt is not referred a-priori to any particular frame, we prompt SAM2 at each frame, so that the model has to balance the influence of tracking while also considering the content of each frame. We discuss more in depth this aspect in Sec. 3.5.

3.3. Cross-Modal Temporal Adapter

An adapter consists of a linear down-projection (\mathbf{W}_{down}) to a bottleneck dimensionality, followed by an up-projection back (\mathbf{W}_{up}) in the original space, separated by a non-linear activation function σ . Formally, given an input feature $\mathbf{x} \in \mathbb{R}^{1 \times d}$, the adapter function is defined as:

$$\text{Adapter}(\mathbf{x}) = \mathbf{x} + \sigma(\mathbf{x}\mathbf{W}_{\text{down}})\mathbf{W}_{\text{up}} \quad (2)$$

We build on this popular Adapter framework [12] and propose a novel Cross-Modal Temporal Adapter (CMT) (see Fig. 3) which models temporal dynamics within visual features while enriching each modality with the other. Formally, given the visual features in a clip $\mathcal{F}^k \in \mathbb{R}^{T \times H_k \times W_k \times C_k}$ and the textual features $\mathcal{E}^k \in \mathbb{R}^{L \times C}$ extracted at layer k of the image and text encoders, respectively, the CMT can be formulated as:

$$\begin{aligned} \text{Adapter}(\mathcal{F}^k) &= \mathcal{F}^k + h(\mathcal{F}^k \mathbf{W}_{\text{down},v}, \mathcal{E}^k \mathbf{W}_{\text{down},t}) \mathbf{W}_{\text{up},v} \\ \text{Adapter}(\mathcal{E}^k) &= \mathcal{E}^k + h(\mathcal{E}^k \mathbf{W}_{\text{down},t}, \mathcal{F}^k \mathbf{W}_{\text{down},v}) \mathbf{W}_{\text{up},t} \end{aligned} \quad (3)$$

where $\mathbf{W}_{\text{down},v}$, $\mathbf{W}_{\text{down},t}$, $\mathbf{W}_{\text{up},v}$, $\mathbf{W}_{\text{up},t}$ are modality specific down- and up-projections weights and h is our proposed adapter function. The adapter output is summed with the original features, allowing the model to retain the original encoding while incorporating temporal and cross-modal reasoning. We integrate the Cross-Modal Temporal Adapter (CMT) into the frozen text and visual encoders at every intermediate layer k . In the following paragraphs we detail the temporal and cross-modal adaptation functions, which are tightly coupled in our Adapter module.

Temporal Adaptation. Our approach aims to embed motion cues directly into the frame-level features of SAM2. Previous works based on Adapters either perform self-attention (SA) over all tokens in a clip [16], which is costly, or restrict the attention to the temporal axis for each pixel [22, 23]. We observe that, within a video, object motion across adjacent frames typically spans a localized region of the image [32]. Consequently, a given element of the feature map primarily benefits from interactions with its

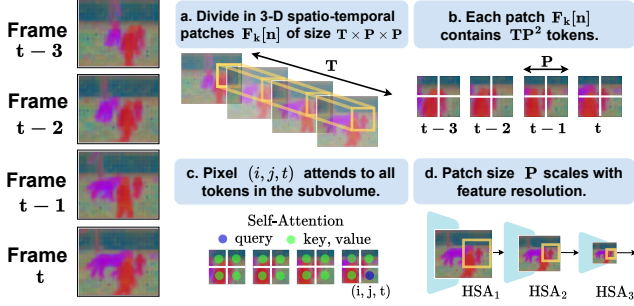


Figure 4. Scheme of our **Hierarchical Selective Attention** (HSA), modeling temporal evolution of features in our adapter.

spatial and temporal neighbors, rather than requiring long-range connections across the entire feature map. Building on this intuition, we introduce a Hierarchical Selective Attention (HSA) mechanism, illustrated in Fig. 4. By modeling interactions among spatially and temporally proximal regions, HSA reduces unnecessary computations while capturing motion-based context.

Formally, at layer k , given the set of feature maps for a T -frames clip: $\mathcal{F}^k \in \mathbb{R}^{T \times H_k \times W_k \times C_k}$, we decompose this feature volume into non-overlapping, 3-D spatio-temporal patches of size $T \times P \times P$, obtaining $N = H_k W_k / P^2$ sub-volumes. These sub-volumes, considered pixel-wise, can be represented as a set of tokens $F^k[n] = \{x_{i,j,t}^{k,n} \in \mathbb{R}^{C_k} : i \in 1, \dots, P, j \in 1, \dots, P, t \in 1, \dots, T\}$. To encode spatio-temporal positioning, to each vector we add a spatial ($e[i, j]$) and a temporal ($e[t]$) sinusoidal positional embeddings, in 2-D and 1-D formats, respectively. Specifically: $x_{i,j,t}^{k,n} = x_{i,j,t}^{k,n} + e[i, j] + e[t]$. Each sub-volume contains $M = P^2 * T$ tokens, on which we perform self-attention as follows:

$$\mathbf{x}_{i,j,t}^{k,n} := SA \left(\left\{ \mathbf{x}_{i',j',t'}^{k,n} \right\}_{\substack{i'=1..P \\ j'=1..P \\ t'=1..T}} \right). \quad (4)$$

At each layer k of the feature extraction process, the patch size P is progressively scaled, as depicted in Fig. 4-d. This scaling adapts the sub-volume to the hierarchy of feature resolution, encoding information at multiple scales.

Cross-Modal Adaptation. To unify text and visual representations, we encourage modality interaction from early stages of the feature extraction process through two symmetric operations: Visual-to-Text Attention (VTA) and Text-to-Visual Attention (TVA).

Within the former, each visual feature, already enriched with temporal information through the HSA, attends to the full textual expression, allowing the model to identify candidate regions within the image based on both categorical details (*e.g.*, the subject described in the text) and motion

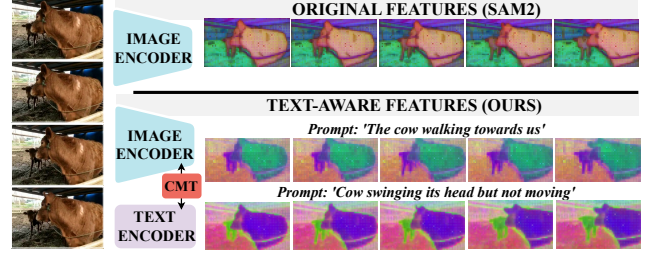


Figure 5. **Cross Modal Temporal Adapter**: we show via PCA that our CMT provides contextualized visual features based on the given textual prompt, compared to SAM2 original ones.

cues (*e.g.*, actions), facilitating early alignment with the prompt, as visible in Fig. 5.

Formally, at layer k , we consider the feature of each frame in the clip, *i.e.* $\mathcal{F}^k[t] \in \mathbb{R}^{H_k \times W_k \times C_k}$, $t = 1, \dots, T$, and the set of textual embeddings $\mathcal{E}^k \in \mathbb{R}^{L \times C}$ to compute:

$$\mathcal{F}^k[t] := \mathcal{F}^k[t] * CA(\mathcal{F}^k[t], \mathcal{E}^k). \quad (5)$$

In parallel, as the meaning of a caption can shift significantly depending on the visual content of the associated image [6], we aim at contextualizing the textual query with the semantics provided by the visual modality. To this end, the TVA progressively enriches the linguistic tokens $\mathcal{E}^k \in \mathbb{R}^{L \times C}$ with information from the visual feature maps, averaged over the video clip \mathcal{F}_{avg}^k :

$$\mathcal{E}^k := \mathcal{E}^k * CA(\mathcal{E}^k, \mathcal{F}_{avg}^k). \quad (6)$$

3.4. Mask prediction

At the end of the feature extraction process, we obtain the adapted visual and linguistic features, respectively \mathcal{E} and \mathcal{F} . To perform the final prediction, we extract the prompt ρ as in Eq. (1), while the Memory Attention module generates the *memory features* \mathcal{F}_{mem} by conditioning the visual features \mathcal{F} on past predictions from the Memory Bank. The prompt ρ is fed into the frozen Mask Decoder \mathcal{D}_{dec} , which generates the output mask $\mathcal{P}_M \in \mathbb{R}^{1 \times H \times W}$ and the mask token $\tau_m \in \mathbb{R}^{1 \times d}$, *i.e.* an embedding representing the segmented object. Formally:

$$\begin{aligned} \tau_m, \mathcal{P}_M &= \mathcal{D}_{dec}(\mathcal{F}_{mem}, \rho), \\ \mathcal{Y} &= \mathcal{P}_M > 0, \end{aligned} \quad (7)$$

where $\mathcal{Y} \in \mathbb{R}^{1 \times H \times W}$ denotes the output binary segmentation mask. Finally, the Memory Encoder updates the memory bank with \mathcal{P}_M .

3.5. Conditional Memory Encoder

We identify as *tracking bias* the phenomenon of SAM2 tracking the wrong object when the correct one is not yet identifiable in the video, and persist in following it. This

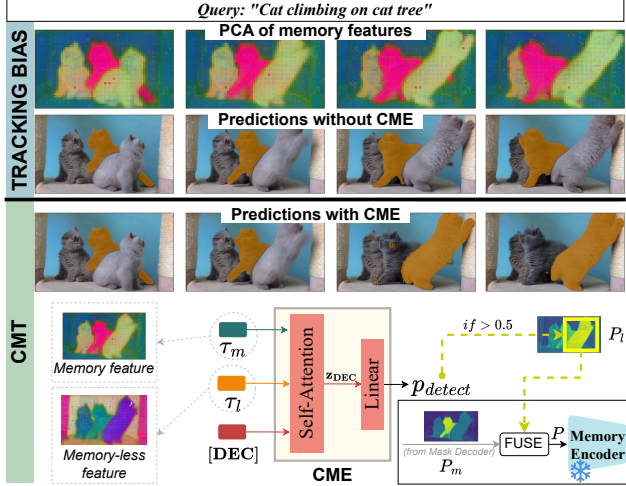


Figure 6. Effect of our **Conditional Memory Encoder**. The caption above requires disambiguating multiple instances of the same class (e.g., “cat”) by identifying a specific action (e.g., “climbing”). Since none of the instances perform this action initially, the model begins tracking the wrong instance and fails to correct itself once the target object performs the action. The top section visualizes the effect of this *tracking bias* in the *memory features* of SAM2. Our CME detects that the cat that starts *climbing* is more aligned with the caption, and encodes its presence in the memory bank, allowing SAM2 to switch its focus.

bias, as exemplified in Fig. 6, is encoded in the memory features, which are propagated to subsequent frames through the Memory Encoder. On the other hand, we observe that the memory-less features: i) contain an *unbiased* representation of the current frames, ii) are aligned with the textual prompt via our CMT (cf. Fig. 5), and iii) can thus be used to propose candidate instances that match the prompt without being biased by past predictions. Building on these intuitions, we derive a *memory-less token* τ_l from a cross-attention between the unbiased feature maps and the prompt. Such token represents a summary of the visual features that match the prompt. The idea is to compare it with the mask token τ_m generated by the Mask Decoder, to detect when they represent different objects, *i.e.*, to detect when SAM2 is tracking an object that is not the one *currently* most aligned with the caption. Formally:

$$\tau_l = CA(\mathcal{F}, \rho) \quad (8)$$

We note that we initialize (and keep frozen) the weights of the cross-attention with those from SAM2 Mask Decoder. We introduce a small learnable module, named Conditional Memory Encoder (CME), to detect such situations. When a new object is detected, a naive solution would be to compute its mask and use it to re-prompt the model at the given frame, just like a user would do, forcing SAM2 to switch its prediction. However, since the prediction computed on the *memory-less* features does not have access to past video

context, it might generate false positives. Thus, we propose a *soft assignment*, obtained by encoding the masks of both objects in the memory bank. Essentially, the CME allows SAM2 to ‘see’ other objects beyond the currently tracked one, and balance the influence of past context with new information, to select the one that fits the prompt the most. In detail, our CME, illustrated in Fig. 6-bottom, concatenates the two tokens τ_m, τ_l with a learnable *decision token* [DEC], and performs a self-attention followed by a Linear classifier:

$$[z_{DEC}, z_{MT}, z_{ML}] = SA\left(\left[[\text{DEC}], \tau_m[t], \tau_l[t]\right]\right) \quad (9)$$

$$p_{detect} = \phi(z_{DEC})$$

where ϕ is a linear function $\mathbb{R}^d \rightarrow \mathbb{R}^1$. When detecting a candidate text-aligned object, (*i.e.*, $p_{detect} > 0.5$), instead of directly feeding the predicted output mask \mathcal{P}_m to the Memory Encoder, our module computes the unbiased output mask, namely $\mathcal{P}_L \in \mathbb{R}^{1 \times H \times W}$, to fuse it with \mathcal{P}_m :

$$\begin{aligned} \mathcal{P}_L &= \mathcal{D}_{dec}(\mathcal{F}, \rho) \\ \mathcal{M}(h, w) &= \mathbb{1}(h, w) [h, w : \mathcal{P}_L > 0] \\ \mathcal{P} &= \lambda * \mathcal{P}_L \circ \mathcal{M} + \mathcal{P}_m \circ (1 - \mathcal{M}) \end{aligned} \quad (10)$$

where $\mathcal{M}(h, w)$ is a binary mask whose value is zero except for the pixels corresponding to the object, and λ is an hyperparameter weighing the influence of the memory-less prediction. The resulting mask \mathcal{P} is fed to the Memory Encoder. We train the CME via self-supervision with a standard Cross-Entropy loss, by providing examples where the *memory-less* features highlight different objects w.r.t the one currently tracked. We discuss in detail our training protocol in the Supp. Mat..

4. Experimental results

Dataset. We evaluate our method on MeVis [7], Ref-Youtube-VOS [2] and Ref-Davis [17]. MeVis includes 2,006 videos and features a total of 28K annotations that capture various aspects of motion. Ref-Youtube-VOS enhances the original YouTube-VOS benchmark by incorporating textual descriptions. It contains a total of 3,978 videos and approximately 15K language expressions. Ref-DAVIS17 builds upon DAVIS17 dataset, adding more than 1.5K linguistic annotations to 90 videos.

Evaluation Metrics. We utilize standard evaluation metrics, region similarity (\mathcal{J}), contour accuracy (\mathcal{F}), and their average ($\mathcal{J} \& \mathcal{F}$). For MeVis and Ref-Youtube-VOS we conduct the evaluation using the official challenge servers; for Ref-DAVIS17, we used the official evaluation code.

Implementation Details. We employ Hiera-B [37] as visual extractor. As text encoder, we experiment with

Method	Visual Encoder	Text Encoder	Total Params	MeViS			Ref-YouTube-VOS			Ref-DAVIS17		
				$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}	$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}	$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}
Large VLM based												
LISA [19] [CVPR'24]	ViT-H	LLaVa	7 B	37.2	35.1	39.4	53.9	53.4	54.3	64.8	62.2	67.3
VISA [45] [ECCV'24]	ViT-H	Chat-UniVi	7 B	43.5	40.7	46.3	61.5	59.8	63.2	69.4	66.3	72.5
One-Token-Seg-All [1] [NIPS'24]	ViT-H	Phi-3	3.8 B	42.3	39.4	45.2	61.7	60.2	63.3	67.7	63.8	71.5
MTTR [3] [CVPR'22]	V-Swin T	RoBERTa	-	30.0	28.8	31.2	55.3	54.0	56.6	-	-	-
TCE-RVOS [14] [WACV'24]	ResNet-50	RoBERTa	-	-	-	-	59.6	58.3	60.8	59.4	56.5	62.4
ReferFormer [43] [CVPR'22]	V-Swin B	RoBERTa	237 M	31.0	29.8	32.2	62.9	61.3	64.6	61.1	58.1	64.1
SOC [28] [NIPS'23]	V-Swin B	RoBERTa	220 M	-	-	-	66.0	64.1	67.9	64.2	61.0	67.4
OnlineRefer [42] [ICCV'23]	Swin L	RoBERTa	232 M	32.3	31.5	33.1	63.5	61.6	65.5	64.8	61.6	67.7
LMPM [7] [ICCV'23]	Swin T	RoBERTa	195 M	37.2	34.2	40.2	-	-	-	-	-	-
RefSAM [20] [arXiv]	ViT-B	T5	3 B	-	-	-	58.4	57.4	59.4	62.1	59.0	65.3
DsHmp [11] [CVPR'24]	V-Swin B	RoBERTa	339 M	-	-	-	67.1	65.0	69.1	64.9	61.7	68.1
DsHmp [11] [CVPR'24]	Swin T	RoBERTa	272 M	46.4	43.0	49.8	-	-	-	-	-	-
MUTR [46] [AAAI'24]	V-Swin B	RoBERTa	250 M	-	-	-	<u>67.5</u>	<u>65.4</u>	<u>69.6</u>	66.4	62.8	70.0
GroundingDINO [25]+SAM2	Hiera-B	BERT	240 M	37.7	34.9	40.5	<u>57.5</u>	<u>55.6</u>	<u>59.5</u>	66.4	62.8	69.9
SAMWISE (ours)	Hiera-B	CLIP-B	150 M	<u>48.3</u>	<u>45.4</u>	<u>51.2</u>	67.2	65.2	69.3	<u>68.5</u>	<u>65.6</u>	<u>71.5</u>
SAMWISE (ours)	Hiera-B	RoBERTa	202 M	49.5	46.6	52.4	69.2	67.8	70.6	70.6	67.4	74.5

Table 1. Comparison of SAMWISE against state-of-the-art RVOS methods on MeViS, Ref-YouTube-VOS and Ref-DAVIS datasets. We further include methods based on large VLMs for comparison. **Bold** and underline indicate the two top results.

CLIP [34] and RoBERTa [26]. We note that the text encoder and SAM2 weights are entirely frozen and we train only the Adapters and the CME module (4.2M parameters when using CLIP and 4.9M with RoBERTa). Following [10, 20, 28, 42, 43], we undergo pre-training for 6 epochs on RefCOCO+/g [30, 49] with a learning rate at $1e-4$ and fine-tune on Ref-YouTube-VOS [38] for 4 epochs with a learning rate of $1e-5$, using the Adam optimizer. The model trained on the Ref-YouTube-VOS is directly evaluated on DAVIS-17 [17]. On MeViS [7], we train for 1 epoch. We set $T = 8$.

4.1. Main Results

In this section, we compare against existing works in the literature, and ablate our contributions. In the Supp. Mat. we report additional qualitative results and ablations.

Baselines. To assess the validity of our approach, we divide the experimental comparison in the following categories:

- **Standard RVOS** methods: we compare against recent relevant works in RVOS. The main comparison is w.r.t. the previous state-of-the-art, namely DsHmp [11];
- Methods with **Context propagation**: OnlineRefer [42] was the first to propose this setting. RefSAM [20] relies on SAM1 to provide frame-level masks, and then propagates the mask token to subsequent frames. A baseline that we propose is GroundingDINO + SAM2, where we use the popular grounded detector to provide boxes for the first frame, and let SAM2 track the object;
- **Large VLM based**: Although these methods [1, 19, 45] are not comparable to ours, or previous ones, in terms of model size, we include them in the table to provide an interesting reference of performance.

Comparison with standard RVOS methods. Traditional

RVOS methods, such as ReferFormer [43] and MTTR [3], suffer a significant performance drop on the MeViS benchmark, as they are unable to solve queries which require to model long-term context. An exception is represented by LMPM and its follow up work DsHmp, which represents the state-of-the-art: these methods process the entire video clip at once, modeling multiple trajectories for all the instances in the video to select the one that fits the prompt the most. Despite this, SAMWISE outperforms DsHmp [11] on all three datasets, improving $\mathcal{J}\&\mathcal{F}$ of +3.1%, +2.1%, and +5.7%, while utilizing a smaller model in terms of total parameters. Notably, we achieve this by training only 4.9 M parameters out of 202 M. This result is particularly impressive, as offline methods exploit information from the entire video to handle challenges such as late-appearing objects or motion-dependent disambiguation, as opposed to our streaming approach. With respect to other methods, we outperform them significantly on MeViS, whereas the gap is smaller on Ref-YouTube-VOS and Ref-DAVIS, which contain more descriptive captions, and object-centric videos. Lastly, we also experiment with the text encoder of CLIP, which achieves state-of-the-art results on MeViS, and competitive performance on other benchmarks, while providing a more compact model with just 150 M params.

Comparison with Context Propagation methods. Our proposed baseline GroundingDINO [25]+SAM2, while obviously flawed, being forced to predict the desired instance based on the first frame only, achieves acceptable results on Ref-DAVIS, whereas on MeViS and Ref-YouTube-VOS its performance drops of 11.8% and 11.7%, respectively. Differently, SAMWISE, demonstrates excellent performance in both motion-dependent and static scenarios. Specifically,

MLP-only	Text-to-Visual	Visual-to-Text	HSA	CME	$\mathcal{J}\&\mathcal{F}$
✓					45.2
✓	✓				47.5
✓		✓			48.3
✓	✓				50.3
✓	✓	✓		✓	54.2
✓	✓	✓	✓	✓	55.5

Table 2. Ablation of our **Cross-Modal Temporal Adapter** (CMT). We show the effect of not using CMT (*i.e.* *MLP-only* to prompt SAM2), vs. adding one at a time its core components. Lastly, the **Conditional Memory Encoder** (CME) is added.

on MeViS, we outperform OnlineRefer [42] by +17.2%. On Ref-YouTube-VOS and Ref-DAVIS, the gap is of +5.7%, and +5.8%, respectively.

Comparison with Large-VLM based. While comparisons with Large-VLM based approaches are not standard in RVOS evaluations, we include them in this work to provide additional context. The VLM-based solutions [1, 19, 45] are designed to leverage the extensive reasoning capabilities of VLMs to address complex textual instructions and implicit descriptions that require world knowledge [19]. This leads to improved performance in tasks like MeViS, where reasoning over motion patterns is required. However, delegating cross-modal reasoning to these VLMs incurs in significant computational overhead, whereas SAMWISE incorporates visual-text interaction directly at the feature level. Notably, SAMWISE outperforms VISA, the best VLM-based competitor by a substantial margin, respectively +6%, +7.7%, +2.9% on the three benchmarks.

4.2. Ablation Studies

We conduct our ablations on MeViS, as it embodies the core challenges of *online* RVOS. We report results on the ‘valid_u’ set [7], employing CLIP-B as text encoder.

Making SAM2 Wisier. We start by showing, in Tab. 2, how each of the core components of our CMT Adapter progressively injects *wisdom* (*i.e.* knowledge about language and temporal context) into SAM2. The first line reports the result using the ‘naive’ solution of aligning the textual prompt to the visual features using a single learnable MLP [51]. While effective to some extent, the results show that allowing early interaction of the two modalities grants a substantial boost (+5.1% with both adapters). Adding explicit temporal feature modeling provides an additional improvements of +3.9%. These results sustain our intuition that adding frame-level alignment through a MLP is not enough to obtain robust performances, and that it is essential to couple the text and visual semantics, as well as modeling temporal context. Lastly, the table shows how adding our CME is effective in mitigating SAM2 *tracking bias*.

Hierarchical Selective Attention. In the top section of Tab. 3, we study the effect of the spatial patch size in our

	HSA Patch Size				
	Fixed Size			Hierarchical	
	1	4	8	8 / 4 / 2	16 / 8 / 4
$\mathcal{J}\&\mathcal{F}$	49.7	52.3	53.1	54.2	53.8
	CME vs Random choice				
	Never	Always	1 every 4	CME	
$\mathcal{J}\&\mathcal{F}$	54.2	50.7	52.4	55.5	

Table 3. Top: Ablation of the Patch size in our **HSA**, with the effect of a fixed size vs Hierarchical. Bottom: Effect of not predicting detections (*Never*), predicting them at every frame (*Always*), randomly (*1 in 4*) vs. using predictions of our CME.

HSA module, which models the temporal evolution of over a spatial patch of size P across the temporal axis. Using $P = 1$ is equivalent to processing each pixel independently across frames. The table shows that including spatial context, up to 8 pixels, is beneficial. Using a hierarchical patch size that scales with the feature map resolution yields a gain of +1.1% over the fixed sized alternative.

Conditional Memory Encoder. The bottom section of Tab. 3 provides insight into our CME module. The CME, essentially, detects whenever an object in the *unbiased* feature maps of the current frame displays higher alignment with the textual prompts w.r.t. the currently tracked one, but SAM2 fails in noticing it due to the *tracking bias* (Fig. 6). The table compares the effect of *Never* applying such strategy (*i.e.*, not using CME), doing it *Always* (*i.e.*, at every frame), or once every 4 frames. The results show that increasing the frequency of *artificial* detection worsens performances, adding noise to the tracking, whereas the predictions of our CME are beneficial, with a boost of +1.3%.

Adaptation strategy. In Tab. 4 we compare different adaptation strategies, including Full Fine-Tuning (FT), LoRa [13], AdaptFormer [5], and the proposed CMT.

	Full-FT	LoRa [13]	AdaptFormer [5]	CMT (ours)
MeViS $\mathcal{J}\&\mathcal{F}$	43.1	44.2	43.9	48.3

Table 4. Baselines: prior adapters and full-finetune, with CLIP-B.

5. Conclusion

In this work, we introduced SAMWISE, a novel approach for RVOS that builds upon the SAM2 model by incorporating i) natural language understanding, ii) temporal feature modeling, and iii) a learnable strategy to adjust tracking focus according to visual cues that emerge over time. SAMWISE achieves SOTA performance across benchmarks while adding less than 5M parameters, without modifying SAM2 weights or using external models for visual-text alignment. We obtain an effective pipeline for applications of streaming video segmentation, addressing limitations of existing RVOS approaches, which either lack long-term context or rely on single-frame context propagation.

Acknowledgements. Claudia Cuttano was supported by the Sustainable Mobility Center (CNMS) which received funding from the European Union Next Generation EU (Piano Nazionale di Ripresa e Resilienza (PNRR), Missione 4 Componente 2 Investimento 1.4 “Potenziamento strutture di ricerca e creazione di “campioni nazionali di R&S” su alcune Key Enabling Technologies”) with grant agreement no. CN_00000023.

Gabriele Trivigno, Giuseppe Averta and Carlo Masone were supported by FAIR - Future Artificial Intelligence Research which received funding from the European Union Next-GenerationEU (PIANO NAZIONALE DI RIPRESA E RESILIENZA (PNRR) – MISSIONE 4 COMPONENTE 2, INVESTIMENTO 1.3 – D.D. 1555 11/10/2022, PE00000013). This manuscript reflects only the authors’ views and opinions, neither the European Union nor the European Commission can be considered responsible for them.

We acknowledge the CINECA award under the ISCRA initiative, for the availability of high performance computing resources.

References

- [1] Zechen Bai, Tong He, Haiyang Mei, Pichao Wang, Ziteng Gao, Joya Chen, Lei Liu, Zheng Zhang, and Mike Zheng Shou. One token to seg them all: Language instructed reasoning segmentation in videos. *arXiv preprint arXiv:2409.19603*, 2024. 3, 7, 8
- [2] Miriam Bellver, Carles Ventura, Carina Silberer, Ioannis Kazakos, Jordi Torres, and Xavier Giro-i Nieto. A closer look at referring expressions for video object segmentation. *Multimedia Tools and Applications*, 82(3):4419–4438, 2023. 2, 6
- [3] Adam Botach, Evgenii Zheltonozhskii, and Chaim Baskin. End-to-end referring video object segmentation with multimodal transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4985–4995, 2022. 1, 2, 7
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020. 2
- [5] Shoufa Chen, Chongjian Ge, Zhan Tong, Jiangliu Wang, Yibing Song, Jue Wang, and Ping Luo. Adaptformer: Adapting vision transformers for scalable visual recognition. *Advances in Neural Information Processing Systems*, 35:16664–16678, 2022. 8
- [6] Henghui Ding, Chang Liu, Suchen Wang, and Xudong Jiang. Vlt: Vision-language transformer and query generation for referring segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(6):7900–7916, 2022. 5
- [7] Henghui Ding, Chang Liu, Shuting He, Xudong Jiang, and Chen Change Loy. MeViS: A large-scale benchmark for video segmentation with motion expressions. In *ICCV*, 2023. 1, 2, 6, 7, 8
- [8] Zihan Ding, Tianrui Hui, Shaofei Huang, Si Liu, Xuan Luo, Junshi Huang, and Xiaoming Wei. Progressive multimodal interaction network for referring video object segmentation. *The 3rd Large-scale Video Object Segmentation Challenge*, 8(10), 2021. 2
- [9] Kirill Gavriluk, Amir Ghodrati, Zhenyang Li, and Cees GM Snoek. Actor and action video segmentation from a sentence. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5958–5966, 2018. 1, 2
- [10] Mingfei Han, Yali Wang, Zhihui Li, Lina Yao, Xiaojun Chang, and Yu Qiao. Htm: Hybrid temporal-scale multimodal learning framework for referring video object segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13414–13423, 2023. 1, 2, 7
- [11] Shuting He and Henghui Ding. Decoupling static and hierarchical motion perception for referring video segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13332–13341, 2024. 2, 3, 7
- [12] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International conference on machine learning*, pages 2790–2799. PMLR, 2019. 2, 3, 4
- [13] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022. 8
- [14] Xiao Hu, Basavaraj Hampiholi, Heiko Neumann, and Jochen Lang. Temporal context enhanced referring video object segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 5574–5583, 2024. 7
- [15] Haojun Jiang, Jianke Zhang, Rui Huang, Chunjiang Ge, Zanlin Ni, Jiwen Lu, Jie Zhou, Shiji Song, and Gao Huang. Cross-modal adapter for text-video retrieval. *arXiv preprint arXiv:2211.09623*, 2022. 3
- [16] Xiaojie Jin, Bowen Zhang, Weibo Gong, Kai Xu, Xueqing Deng, Peng Wang, Zhao Zhang, Xiaohui Shen, and Jiashi Feng. Mv-adapter: Multimodal video transfer learning for video text retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27144–27153, 2024. 3, 4
- [17] Anna Khoreva, Anna Rohrbach, and Bernt Schiele. Video object segmentation with language referring expressions. In *Computer Vision—ACCV 2018: 14th Asian Conference on Computer Vision, Perth, Australia, December 2–6, 2018, Revised Selected Papers, Part IV 14*, pages 123–141. Springer, 2019. 1, 2, 6, 7
- [18] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023. 4
- [19] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9579–9589, 2024. 2, 3, 7, 8
- [20] Yonglin Li, Jing Zhang, Xiao Teng, Long Lan, and Xinwang Liu. Refsam: Efficiently adapting segmenting anything model for referring video object segmentation, 2024. 3, 7

- [21] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. 2
- [22] Ruyang Liu, Jingjia Huang, Ge Li, Jiashi Feng, Xinglong Wu, and Thomas H Li. Revisiting temporal modeling for clip-based image-to-video knowledge transferring. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6555–6564, 2023. 4
- [23] Ruyang Liu, Chen Li, Yixiao Ge, Thomas H. Li, Ying Shan, and Ge Li. Bt-adapter: Video conversation is feasible without video instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13658–13667, 2024. 4
- [24] Si Liu, Tianrui Hui, Shaofei Huang, Yunchao Wei, Bo Li, and Guanbin Li. Cross-modal progressive comprehension for referring segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):4761–4775, 2021. 1, 2
- [25] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 3, 7
- [26] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019. 7
- [27] Haoyu Lu, Yuqi Huo, Guoxing Yang, Zhiwu Lu, Wei Zhan, Masayoshi Tomizuka, and Mingyu Ding. Uniadapter: Unified parameter-efficient transfer learning for cross-modal modeling. *arXiv preprint arXiv:2302.06605*, 2023. 3
- [28] Zhuoyan Luo, Yicheng Xiao, Yong Liu, Shuyan Li, Yitong Wang, Yansong Tang, Xiu Li, and Yujiu Yang. Soc: Semantic-assisted object cluster for referring video object segmentation. *Advances in Neural Information Processing Systems*, 36, 2024. 1, 3, 7
- [29] Bo Miao, Mohammed Bennamoun, Yongsheng Gao, and Ajmal Mian. Spectrum-guided multi-granularity referring video object segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 920–930, 2023. 2
- [30] Varun K Nagaraja, Vlad I Morariu, and Larry S Davis. Modeling context between objects for referring expression understanding. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 792–807. Springer, 2016. 7
- [31] Seoung Wug Oh, Joon-Young Lee, Ning Xu, and Seon Joo Kim. Video object segmentation using space-time memory networks. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9226–9235, 2019. 2
- [32] Mandela Patrick, Dylan Campbell, Yuki Asano, Ishan Misra, Florian Metzke, Christoph Feichtenhofer, Andrea Vedaldi, and Joao F Henriques. Keeping your eye on the ball: Trajectory attention in video transformers. *Advances in neural information processing systems*, 34:12493–12506, 2021. 4
- [33] Matthew E Peters, Sebastian Ruder, and Noah A Smith. To tune or not to tune? adapting pretrained representations to diverse tasks. *arXiv preprint arXiv:1903.05987*, 2019. 3
- [34] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 3, 7
- [35] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 2
- [36] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, et al. Grounded sam: Assembling open-world models for diverse visual tasks. *arXiv preprint arXiv:2401.14159*, 2024. 3
- [37] Chaitanya Ryali, Yuan-Ting Hu, Daniel Bolya, Chen Wei, Haoqi Fan, Po-Yao Huang, Vaibhav Aggarwal, Arkabandhu Chowdhury, Omid Poursaeed, Judy Hoffman, et al. Hiera: A hierarchical vision transformer without the bells-and-whistles. In *International Conference on Machine Learning*, pages 29441–29454. PMLR, 2023. 4, 6
- [38] Seonguk Seo, Joon-Young Lee, and Bohyung Han. Urvos: Unified referring video object segmentation network with a large-scale benchmark. In *European Conference on Computer Vision*, 2020. 1, 2, 7
- [39] Jiajin Tang, Ge Zheng, and Sibeil Yang. Temporal collection and distribution for referring video object segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15466–15476, 2023. 1
- [40] Mengmeng Wang, Jiazheng Xing, Boyuan Jiang, Jun Chen, Jianbiao Mei, Xingxing Zuo, Guang Dai, Jingdong Wang, and Yong Liu. A multimodal, multi-task adapting framework for video action recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 5517–5525, 2024. 3
- [41] Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, et al. Image as a foreign language: Beit pretraining for all vision and vision-language tasks. *arXiv preprint arXiv:2208.10442*, 2022. 2
- [42] Dongming Wu, Tiancai Wang, Yuang Zhang, Xiangyu Zhang, and Jianbing Shen. OnlineRefer: A simple online baseline for referring video object segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2761–2770, 2023. 2, 3, 7, 8
- [43] Jiannan Wu, Yi Jiang, Peize Sun, Zehuan Yuan, and Ping Luo. Language as queries for referring video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4974–4984, 2022. 1, 2, 7
- [44] Zunnan Xu, Zhihong Chen, Yong Zhang, Yibing Song, Xiang Wan, and Guanbin Li. Bridging vision and language encoders: Parameter-efficient tuning for referring image segmentation. In *Proceedings of the IEEE/CVF International*

Conference on Computer Vision, pages 17503–17512, 2023.

[3](#)

- [45] Cilin Yan, Haochen Wang, Shilin Yan, Xiaolong Jiang, Yao Hu, Guoliang Kang, Weidi Xie, and Efstratios Gavves. Visa: Reasoning video object segmentation via large language models. *arXiv preprint arXiv:2407.11325*, 2024. [3](#), [7](#), [8](#)
- [46] Shilin Yan, Renrui Zhang, Ziyu Guo, Wenchao Chen, Wei Zhang, Hongyang Li, Yu Qiao, Hao Dong, Zhongjiang He, and Peng Gao. Referred by multi-modality: A unified temporal transformer for video object segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 6449–6457, 2024. [7](#)
- [47] Lingxiao Yang, Ru-Yuan Zhang, Yanchen Wang, and Xiaohua Xie. Mma: Multi-modal adapter for vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23826–23837, 2024. [3](#)
- [48] Linwei Ye, Mrigank Rochan, Zhi Liu, and Yang Wang. Cross-modal self-attention network for referring image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10502–10511, 2019. [1](#), [2](#)
- [49] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14*, pages 69–85. Springer, 2016. [7](#)
- [50] Yuxuan Zhang, Tianheng Cheng, Rui Hu, Heng Liu, Longjin Ran, Xiaoxin Chen, Wenyu Liu, Xinggang Wang, et al. Evfsam: Early vision-language fusion for text-prompted segment anything model. *arXiv preprint arXiv:2406.20076*, 2024. [2](#)
- [51] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. [2](#), [8](#)