# Zero-Shot Object Segmentation via Attentive Graph Neural Networks

June 3, 2025

## 1 Introduction

Zero-shot video object segmentation (ZVOS) is the task of identifying the primary objects in a video sequence without having any prior knowledge about the object classes. The paper addresses this problem by proposing a novel Attentive Graph Neural Network (AGNN) model. Specifically, AGNN builds a fully connected graph where video frames are represented as nodes, and relations between frame pairs are represented as edges. Through an attention-based mechanism and iterative message passing, AGNN captures much richer high-order relations among frames compared to previous models. This enables a more comprehensive understanding of the video content and results in more accurate foreground estimation. At the time of its publication, AGNN set a new state-of-the-art performance for the ZVOS task.

## 2 Problem Definition

Zero-shot video object segmentation (ZVOS) refers to the task of segmenting the primary foreground objects in a video sequence without relying on any prior object category information or manual annotations at test time.

Formally, given a sequence of video frames $\mathcal{I} = \{I_i \in \mathbb{R}^{w \times h \times 3} \mid i = 1, ..., N\}$, where each frame $I_i \in \mathbb{R}^{H \times W \times 3}$ is a color image of height $H$ and width $W$ over $C$ channels, the goal is to predict a corresponding sequence of binary segmentation masks $\mathcal{S} = \{S_i \in \{0,1\}^{w \times h} \mid i = 1, ..., N\}$. A value of 1 indicates a foreground pixel belonging to a primary object, while 0 indicates background.

Unlike semi-supervised video object segmentation, where ground-truth annotations are provided for the first frame, ZVOS requires the model to autonomously infer the foreground regions across all frames based solely on the video content itself, making it a highly challenging and unsupervised task.
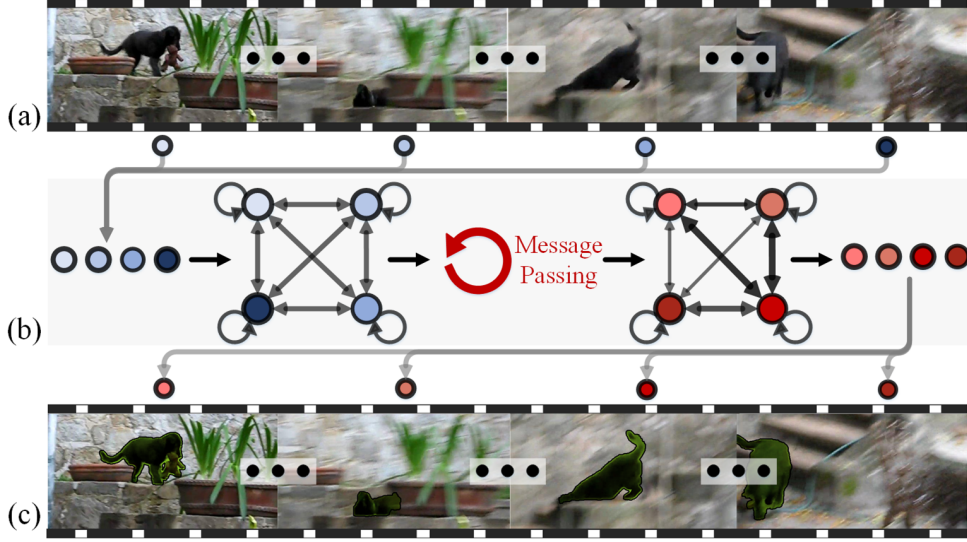
Figure 1: AGNN for ZVOS

## 3  Methodology

### 3.1  Background: Graph Neural Networks

Graph Neural Networks (GNN) are a powerful tool for collectively aggregate information form data that can be represented through graphs.

Given a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, to each node $v_i \in \mathcal{V}$ is associated a *node embedding* (or *node representation*) $\mathbf{v}_i$ and to each edge $e_{ij} = (v_i, v_j) \in \mathcal{E}$ is associated an *edge embedding* $\mathbf{e}_{ij}$.

For each node $v_i$, through aggregating information from the neighboring nodes $\mathcal{N}_i$, we learn updated node embedding $\mathbf{h}_i$. The network maps the initial graph representation $\mathcal{G}$ to the node outputs $\{\mathbf{o}_i\}_{i=1}^{|\mathcal{V}|}$, through a read-out phase based on the final node representations.

The learning is articulated in two steps:

- **Parametric message passing phase**: this first phase is run for $K$ steps and recursively propagates messages from neighboring nodes and updates nodes representations.

$$\mathbf{m}_i^k = \sum_{v_j \in \mathcal{N}_i} \mathbf{m}_{j,i}^k = \sum_{v_j \in \mathcal{N}_i} M(\mathbf{h}_j^{k-1}, \mathbf{e}_{i,j}^{k-1}) \qquad (message\ aggregation)$$

$$\mathbf{h}_i^k = U(\mathbf{h}_i^{k-1}, \mathbf{m}_i^{k-1}) \qquad (node\ representation\ update)$$

- **Read-out phase**: the final node representation $\mathbf{h}_i^K$ are mapped to the output labels $\mathbf{o}_i$ via a read-out function:

$$\mathbf{o}_i = R(\mathbf{h}_i^K)$$

Where M is a *message function*, U is a *state update function* and R is a *read-out function*, all three of which are learned differentiable functions. Notice that the message received from node $i$ at iteration $k$ depends on the neighboring nodes state and edges from previous iteration $k-1$. The initial node state for node $v_i$ is $\mathbf{h}_i^0 = \mathbf{v}_i$

**Benefits and limitations of Graph Neural Networks.** GNNs provide a natural and flexible framework for modeling relational data, enabling the learning of node representations that

capture both local and global graph structure. They are widely applicable across many fields such as computer vision, natural language processing, and bioinformatics. However, traditional GNNs can struggle with scalability when applied to large or fully connected graphs, as the cost of message passing grows with the number of edges. Additionally, vanilla GNNs often focus only on local neighborhoods, which may limit their ability to capture higher-order relationships unless specifically designed to propagate information globally, as is done in AGNN.

## 3.2 Background: Graph Attention Networks

Graph Attention Networks (GATs) are a class of graph neural networks that leverage attention mechanisms to improve the aggregation of information across a graph.

In traditional GNNs, node representations are updated by uniformly aggregating information from neighboring nodes. In contrast, GATs assign different importance weights to different neighbors, allowing the model to focus more on the most relevant nodes during message passing.

Formally, given a node $v_i$ and its neighbors $\mathcal{N}_i$, the attention coefficient $\alpha_{ij}$ between node $v_i$ and neighbor $v_j$ is computed as:

$$\alpha_{ij} = \frac{\exp\left(\text{LeakyReLU}\left(\mathbf{a}^\top [\mathbf{W}\mathbf{h}_i \,\|\, \mathbf{W}\mathbf{h}_j]\right)\right)}{\sum_{v_k \in \mathcal{N}_i} \exp\left(\text{LeakyReLU}\left(\mathbf{a}^\top [\mathbf{W}\mathbf{h}_i \,\|\, \mathbf{W}\mathbf{h}_k]\right)\right)}$$

where:

- $\mathbf{h}_i$ and $\mathbf{h}_j$ are the input features of nodes $v_i$ and $v_j$,

- $\mathbf{W}$ is a learnable weight matrix,

- $\mathbf{a}$ is a learnable attention weight vector,

- $\|$ denotes concatenation.

The new embedding of node $v_i$ is then computed as the weighted sum of transformed neighbor features:

$$\mathbf{h}'_i = \sigma\left(\sum_{v_j \in \mathcal{N}_i} \alpha_{ij} \mathbf{W}\mathbf{h}_j\right)$$

where $\sigma(\cdot)$ is a non-linear activation function, typically ReLU or ELU.

By incorporating attention, GATs can dynamically adjust the influence of each neighbor during aggregation, leading to more flexible and expressive models. They have been successfully applied to tasks such as node classification, graph classification, and link prediction.

In the AGNN model, a similar attention mechanism is adapted and extended to handle not just node-to-node relations but also pixel-wise relations within and between video frames, enabling the model to capture high-order dependencies for the video segmentation task.

# 4 Implementation

## 4.1 Network Architecture

Lets now see in detail the components of our GAN.

**FCN-Based Node Embedding**:

To extract node embeddings, the authors employ DeepLabV3, a classical Fully Convolutional Network (FCN) designed for semantic segmentation tasks.

$$\mathbf{h}_i^0 = \mathbf{v}_i = F_{DeepLab}(I_i) \in \mathbb{R}^{W \times H \times C}$$

Here, $\mathbf{h}_i^0$ represents the initial node embedding associated with frame $I_i$. It preserves both spatial information and high-level semantic features, which are crucial for accurate pixel-wise segmentation.

DeepLabV3 operates by applying atrous (dilated) convolutions at multiple scales to enlarge the receptive field without losing spatial resolution. In atrous convolution, filters are applied with defined gaps (controlled by a dilation rate) between kernel elements, allowing the network to capture larger context without increasing the number of parameters or reducing the feature map size.

To further enhance multi-scale feature extraction, DeepLabV3 incorporates an Atrous Spatial Pyramid Pooling (ASPP) module. ASPP applies multiple parallel atrous convolutions with different dilation rates, effectively capturing contextual information at multiple scales. This multi-branch design enables the model to simultaneously aggregate fine-grained local details and broader global structures.

As a result, DeepLabV3 can efficiently extract rich semantic information at various scales, making it well-suited for tasks requiring fine-grained spatial understanding such as video object segmentation.
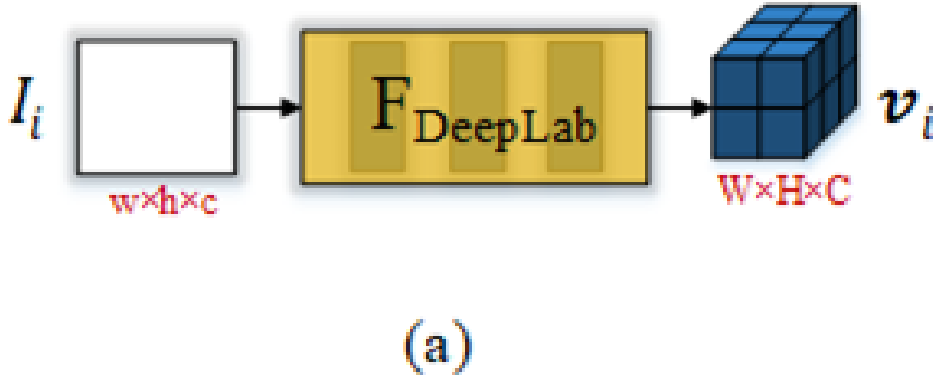


Figure 2: Node Embedding

**Intra-Attention Based Loop-Edge Embedding**:

A loop-edge $e_{i,i} \in \mathcal{E}$ connects a node to itself. In AGNN, the embedding $e_{i,i}^k$ associated with the loop-edge is formulated using an intra-attention mechanism. This mechanism enhances the node's feature representation by modeling long-range, multi-level dependencies across the spatial dimensions of the feature map.

Intra-attention has been shown to be complementary to convolutional operations, allowing the network to capture non-local interactions within the same frame, which is crucial for tasks such as semantic segmentation where distant but semantically related regions should influence each other.

$$\mathbf{e}_{i,i}^k = F_{IntraAtt}(\mathbf{h}_i^k) = \alpha \cdot SoftMax((W_f * \mathbf{h}_i^k)(W_h * \mathbf{h}_i^k)^T)(W_l * \mathbf{h}_i^k) + \mathbf{h}_i^k \quad \in \mathbb{R}^{W \times H \times C}$$

Where $\alpha$ is a learnable scalar parameter, $W$'s are learnable convolution kernels, $*$ is the convolution operation.
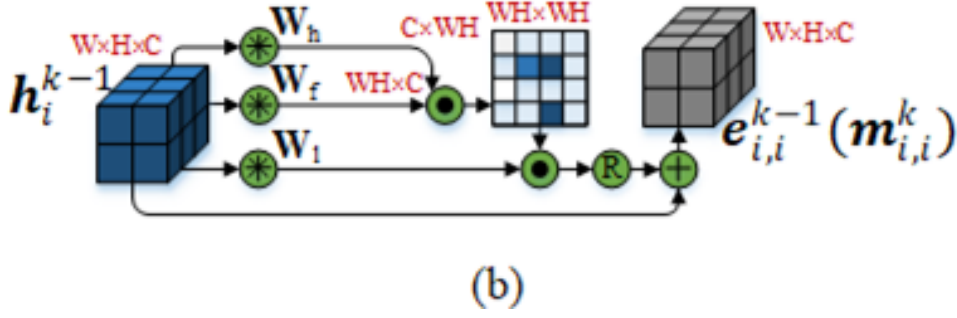
Figure 3: Intra-Attention

**Inter-Attention Based Line-Edge Embedding**:

A line-edge $e_{i,j} \in \mathcal{E}$ connects node $v_i$ to node $v_j$. Here, we compute an inter-attention mechanism to capture the bi-directional relations between node $v_i$ and node $v_j$.

$$\mathbf{e}_{i,j}^k = F_{InterAtt}(\mathbf{h}_i^k, \mathbf{h}_j^k) = \mathbf{h}_i^k \mathbf{W}_c \mathbf{h}_j^{kT} \quad \in \mathbb{R}^{(WH) \times (WH)}$$

$$\mathbf{e}_{j,i}^k = F_{InterAtt}(\mathbf{h}_j^k, \mathbf{h}_i^k) = \mathbf{h}_j^k \mathbf{W}_c^T \mathbf{h}_i^{kT} \quad \in \mathbb{R}^{(WH) \times (WH)}$$

Where:

- $\mathbf{e}_{i,j}^k = \mathbf{e}_{j,i}^{kT}$

- $\mathbf{W}_c \in \mathbb{R}^{C \times C}$ is a learnable weight matrix

- $\mathbf{h}_i^k \in \mathbb{R}^{(WH) \times C}$ are the 3D tensors flattened into matrix representation

The edge embedding $e_{i,j}^k$ obtained after applying inter-attention can be interpreted as the importance of node $v_i$ to node $v_j$ at iteration $k$. Specifically, it quantifies how much the feature representation of node $v_i$ should influence the update of node $v_j$ during the message passing phase. A higher attention score indicates that $v_i$ provides more relevant information to $v_j$, thereby guiding the propagation of meaningful features across the graph.
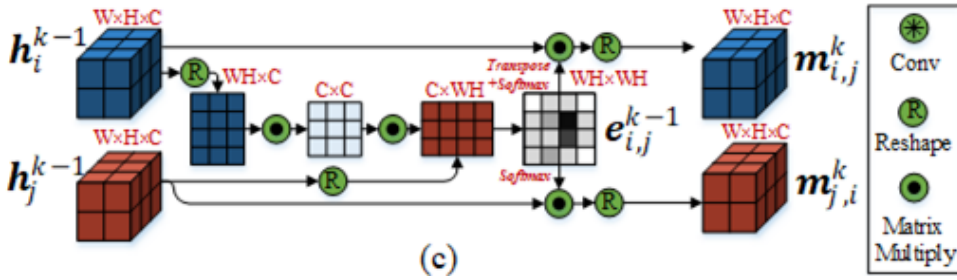


Figure 4: Inter-Attention

**Gated Message Aggregation**:

For the self-loop $e_{i,i}$ we view the loop-edge embedding $\mathbf{e}_{i,i}^k$ itself as the message

$$\mathbf{m}_{i,i}^k = \mathbf{e}_{i,i}^{k-1} \in \mathbb{R}^{W \times H \times C}$$

while for line-edges

$$\mathbf{m}_{j,i}^k = M(\mathbf{h}_i^{k-1}, \mathbf{e}_{j,i}^{k-1}) = SoftMax(\mathbf{e}_{j,i}^{k-1})\mathbf{h}_i^{k-1} \in \mathbb{R}^{(WH) \times C}$$

5

The SoftMax$(\cdot)$ function normalizes each row of its input. As a result, each row (corresponding to a spatial position) of the message $\mathbf{m}_{j,i}^k$ becomes a weighted combination of the rows (positions) of the feature map $\mathbf{h}_j^{k-1}$, where the weights are given by the corresponding column of $\mathbf{e}_{i,j}^{k-1}$.

In this way, the message function $M(\cdot)$ computes an edge-weighted message to be passed to neighboring nodes, effectively encoding the influence of node $v_j$ on node $v_i$. After message computation, $\mathbf{m}_{j,i}^k$ is reshaped back into a 3D tensor of size $W \times H \times C$.

However, not all neighboring nodes provide equally useful information. Some nodes may be noisy due to factors like camera shifts, occlusions, or out-of-view regions, and their messages could degrade the quality of the target node's representation. To address this, a learnable gating function $G(\cdot)$ is introduced to measure the confidence of each message $\mathbf{m}_{j,i}^k$, allowing the model to selectively emphasize reliable information during aggregation.

$$\mathbf{g}_{j,i}^k = G(\mathbf{m}_{j,i}^k) = \sigma(F_{GAP}(\mathbf{W}_g * \mathbf{m}_{j,i}^k + b_g)) \in [0,1]^C$$

where

- $F_{GAP}(\cdot)$ indicates the use of global-average pooling to generate channel-wise responses

- $\sigma(x) = \dfrac{1}{1 + e^{-x}}$

- $\mathbf{W}_g$ trainable convolution kernel

- $b_g$ trainable convolution bias

The function $F_{GAP}(\cdot)$ applies global average pooling across the spatial dimensions of the message tensor, compressing each feature map into a single scalar value. This generates channel-wise summaries that capture the overall importance of each feature channel across the entire spatial field.

The sigmoid activation $\sigma(\cdot)$ then squashes these channel responses into the range $[0,1]$, effectively acting as a soft gating mechanism. This allows the model to selectively control the influence of each channel in the message, where higher values indicate higher confidence and greater relevance of the information being propagated. So the model learns to filter or amplify different parts of the incoming messages based on reliability.

Finally, we collect messages from the neighbors and self-loop via gated summarization

$$\mathbf{m}_i^k = \sum_{v_j \in \mathcal{V}} \mathbf{g}_{j,i}^k \star \mathbf{m}_{j,i}^k \in \mathbb{R}^{W \times H \times C}$$

where

- "$\star$" is the channel-wise Hadamard product

The channel-wise Hadamard product refers to element-wise multiplication performed independently for each channel of a feature tensor.

Given a gating vector $\mathbf{g} \in [0,1]^C$ and a feature tensor $\mathbf{m} \in \mathbb{R}^{W \times H \times C}$, the operation is defined as:

$$(\mathbf{g} \star \mathbf{m})(x,y,c) = g_c \times m(x,y,c)$$

for each spatial position $(x,y)$ and channel $c$.

In this way, each channel of the feature tensor is modulated by a corresponding gating value, allowing the network to selectively emphasize or suppress information on a per-channel basis.

This mechanism improves robustness by filtering out noisy or unreliable message components before aggregation.

**ConvGRU Based Node-State Update**:

To update the node state of $v_i$ at iteration $k$ and preserve the spatial information conveyed by its previous state $\mathbf{h}_i^{k-1}$ and received message $\mathbf{m}_i^k$, we leverage convolutional GRU (ConvGRU).

$$\mathbf{h}_i^k = U_{ConvGRU}(\mathbf{h}_i^{k-1}, \mathbf{m}_i^k) \in \mathbb{R}^{W \times H \times C}$$

The ConvGRU is a convolutional extension of the standard Gated Recurrent Unit (GRU), designed to handle spatial data like images and feature maps. Instead of using fully connected operations, ConvGRU replaces all matrix multiplications with convolutional layers, preserving the spatial structure of the input across time or iterations.

A Gated Recurrent Unit (GRU) — whether standard or convolutional — always maintains a hidden state across time steps or iterations.

It uses two gating mechanisms:

- **Update gate**: controls how much new information to incorporate from the current input.

- **Reset gate**: controls how much of the previous memory to discard when computing the candidate hidden state.

In the context of AGNN:

- Each node maintains an embedding (a feature map) representing its current state.

- After aggregating new messages from neighboring nodes, the node uses the ConvGRU gates to decide how much of its past embedding to retain and how much new information to integrate.

Thus, ConvGRU introduces an explicit memory mechanism across message passing iterations, allowing each node to iteratively refine its state by selectively combining historical and incoming information.

In ConvGRU, all gating operations (reset gate, update gate, and candidate hidden state computation) are performed using convolutional layers instead of fully connected layers.

A convolutional layer applies the same set of learnable filters across all spatial locations of the feature map, meaning that the weights are shared spatially. This contrasts with a standard fully connected GRU, where different spatial positions would be updated independently without spatial consistency.

As a result:

- Every pixel or spatial position within a node's feature map is updated by the same convolutional operations.

- This parameter sharing drastically reduces the total number of learnable parameters compared to dense layers.

- It also preserves spatial structure, ensuring that neighboring pixels are treated consistently, which is crucial for pixel-level tasks like semantic segmentation.

**Read-Out Function**:

After $K$ message passing iterations, we obtain the final state $\mathbf{h}_i^K$ for node $v_i$. So, we can finally obtain our segmentation prediciton map $\hat{S}$ from a read-out funciton

$$\hat{S}_i = R_{FCN}([\mathbf{h}_i^K, \mathbf{h}_i^0]) \in [0,1]^{W \times H}$$

The read-out function consists of:

- Two 3x3 convolutional layers

- One final 1x1 convolutional layer

Again, to preserve spatial information, the readout function is implemented as a small FCN network, which has three convolution layers with a sigmoid function to normalize the prediction to $[0, 1]$.
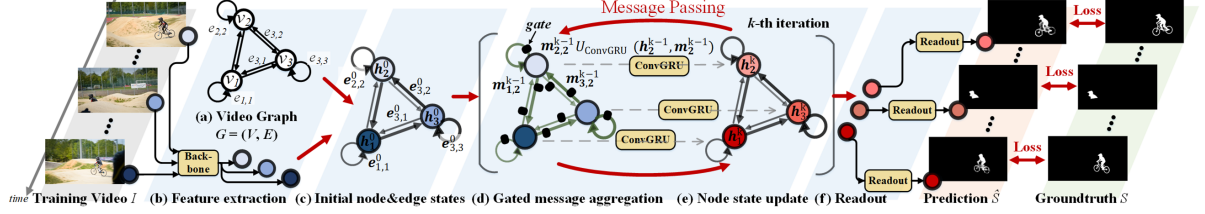


Figure 5: AGNN Model during training

**Loss Function**:

The training of AGNN uses a **weighted binary cross-entropy loss** defined at the pixel level. Given the ground-truth segmentation mask $S$ and the predicted foreground map $\hat{S}$, the loss is formulated as:

$$L(S, \hat{S}) = - \sum_{x=1}^{W \times H} \Big( (1 - \eta) S_x \log(\hat{S}_x) + \eta (1 - S_x) \log(1 - \hat{S}_x) \Big)$$

where:

- $S_x \in \{0, 1\}$ is the ground-truth label at pixel $x$.

- $\hat{S}_x \in [0, 1]$ is the predicted probability at pixel $x$.

- $\eta$ denotes the foreground-to-background pixel number ratio within $S$.

This weighting scheme helps to deal with the class imbalance problem, which is very common in segmentation tasks where the background usually occupies most of the image. By assigning more importance to the less frequent class (typically the foreground object), the model is encouraged to pay closer attention to correctly segmenting the object instead of being biased toward predicting background everywhere.

Since AGNN processes multiple frames at once, each with different object sizes and scene compositions, using a weighted loss also helps the model adapt better across frames. It forces the network to learn more robust and consistent object representations even when the foreground objects vary significantly in scale or visibility.

**Detailed Network Architecture**:

The whole model is end-to-end trainable, as all the functions in AGNN are parametrized by neural networks.

- The first 5 blocks of DeepLabV3 are used as backbone for features extraction

- For an input video $\mathcal{I}$, each frame $I_i$ has resolution $473 \times 473$

- The initial node representaion for node $v_i$ is $\mathbf{h}_i^0 \in \mathbb{R}^{60 \times 60 \times 256}$

- After $K$ message passing iterations, we obtain the segmentation prediction map $\hat{S} \in [0, 1]^{60 \times 60}$

## 4.2 Training Setup

To train our AGNN we leverage a random sampling strategy:

- Split each training video $\mathcal{I}$ of $N$ frames into $N'$ segments

- Randomly select one frame from each segment

Due to computational power limitation we randomly select 2 training videos and sample $N' = 3$ frames. The total number of iterations is set to $K = 3$.

## 4.3 Testing Setup

After training, we can use our AGNN model to perform per-pixel objects prediction in the following way:

1. Select an input video $\mathcal{I}$ of $N$ frames with $473 \times 473$ resolution

2. Split $\mathcal{I}$ into $T = \dfrac{N}{N'}$ subsets: $\{\mathcal{I}_1, ..., \mathcal{I}_T\}$

3. Each subset contains $N'$ frames: $\{I_t, I_{t+T}, ..., I_{N-T+t}\}$

4. Then feed each subset into the network and obtain segmentation maps for all the frames in the subset

5. Finally, apply Conditional Random Field as post-processing

## 4.4 DAVIS17 Validation.

During testing, we loaded the `checkpoint_epoch100.pth` model and ran the following steps for each video in DAVIS17's `val.txt` (20 videos):

1. Split each video of $N$ frames into subsets of size $N_{\text{test}} = 5$ (we used $\lceil N/5 \rceil$ subsets, repeating the last frame if needed to fill a subset).

2. For each subset of 5 frames:

   - Read each 480p image, resize to ($473 \times 473$), normalize (using ImageNet mean/std), and stack into a tensor of shape $(1, 5, 3, H, W)$.

   - Run a single forward pass through AGNN to obtain $\hat{P} \in [0, 1]^{(1,5,1,60,60)}$.

   - Upsample each $60 \times 60$ probability map to the original frame resolution $480p$ via bilinear interpolation $\rightarrow \hat{p}_{h,w} \in [0, 1]$.

   - Binarize at threshold 0.5 to obtain a raw mask `raw_bin` $\in \{0, 1\}^{H \times W}$.

   - Apply dense CRF to refine the mask probabilities.

   - Threshold the CRF output at $0.5 \rightarrow \tilde{m} \in \{0, 1\}^{H \times W}$.

Once all 20 DAVIS17 validation videos were processed, we computed three scores:

- **Region similarity** $J$ (IoU): average pixel-IoU between each predicted mask and ground truth, over all frames.

- **Boundary accuracy** $F$: the contour F-measure between predicted vs. GT boundaries, with a 2-pixel tolerance.

- **Temporal stability** $T$: 1 minus normalized symmetric difference between consecutive predicted frames, normalized by GT union.

**1. Region similarity $J$.** For each frame, let $G \in \{0,1\}^{H \times W}$ be the ground-truth mask and $M \in \{0,1\}^{H \times W}$ the predicted mask. Then

$$
J = \begin{cases} 1, & \big|\{G=1\} \cup \{M=1\}\big| = 0, \\[2mm] \dfrac{\big|\{G=1\} \cap \{M=1\}\big|}{\big|\{G=1\} \cup \{M=1\}\big|}, & \text{otherwise.} \end{cases}
$$

Intuitively, $J$ measures the fraction of correctly segmented pixels relative to the union of ground-truth and prediction.

**2. Boundary accuracy $F$.** Define the sets of boundary pixels $\partial G$ and $\partial M$. Allowing a tolerance $d$ (e.g. $d = \lceil 0.008 \max(H,W) \rceil$ pixels), we compute

$$
P = \frac{\#\{\, p \in \partial M : \mathrm{dist}(p, \partial G) \leq d \,\}}{\#(\partial M)}, \quad R = \frac{\#\{\, g \in \partial G : \mathrm{dist}(g, \partial M) \leq d \,\}}{\#(\partial G)},
$$

and the boundary $F$-measure is their harmonic mean:

$$
F = \frac{2PR}{P + R}.
$$

This captures how well the predicted contour aligns, within $d$ px, to the true object edge.

**3. Temporal stability $T$.** To assess flicker or jitter, we compare consecutive predicted masks $M_t$ and $M_{t+1}$. Define the frame-to-frame IoU

$$
J_t^{\mathrm{tmp}} = \frac{|M_t \cap M_{t+1}|}{|M_t \cup M_{t+1}|} \quad (t = 1, \ldots, T-1).
$$

Then the video's temporal stability is

$$
T = \frac{1}{T-1} \sum_{t=1}^{T-1} J_t^{\mathrm{tmp}},
$$

so that higher $T$ means smoother, more coherent masks over time.

All three scores are averaged per-video and then over the dataset to yield the final $\overline{J}$, $\overline{F}$ and $\overline{T}$ reported in our experiments.
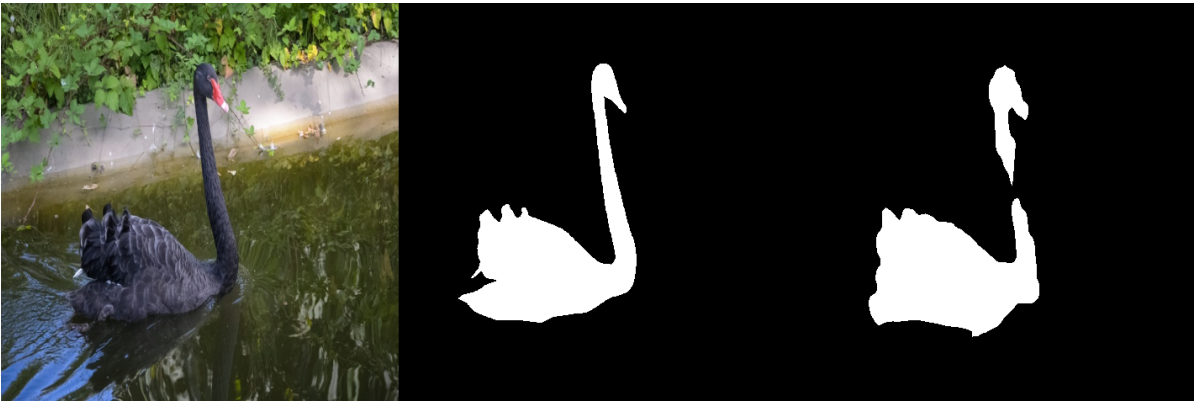
# 5 Results



Figure 6: Frame from Balckswan video: $\bar{J} = 0.886$; $\bar{F} = 0.701$

Frame from Breakdance $\bar{J} = 0.317$; $\bar{F} = 0.190$

## 5.1 Overall Performance

Table 1 reports the average values of the three key metrics on the DAVIS17 validation set:

$$\bar{J} = 0.635, \quad \bar{F} = 0.406, \quad \bar{T} = 0.723.$$

These numbers are comparable to those achieved by many zero-shot video object segmentation methods on DAVIS17 when trained from scratch (i.e., without any fine-tuning on DAVIS). An average Intersection-over-Union $\bar{J} \approx 0.63$–$0.65$ is respectable for a lightweight, single-model approach. The boundary F-measure $\bar{F} \approx 0.40$ indicates that predicted contours are often a few pixels away from the ground-truth edges. Finally, a temporal stability $\bar{T} \approx 0.72$ shows that consecutive masks overlap more than 70% of the time; top methods on DAVIS16 typically report $\bar{T}$ values in the range 0.75–0.85.

Table 1: Overall mean metrics on DAVIS17.

|  | $\bar{J}$ | $\bar{F}$ | $\bar{T}$ |
|---|---|---|---|
| Overall Mean | 0.635 | 0.406 | 0.723 |

## 5.2 Per-Video Analysis

Table 2 shows results for a subset of representative sequences. High-contrast, large-object videos such as *blackswan* and *lab-coat* yield excellent performance, with $\bar{J} > 0.88$, $\bar{F} > 0.63$, and $\bar{T} > 0.88$. For example, *blackswan* achieves $(J, F, T) = (0.886, 0.701, 0.906)$, while *lab-coat* reaches $(0.919, 0.637, 0.884)$. Similarly, *cows* obtains $(0.827, 0.507, 0.945)$, reflecting the smooth, coherent motion of the herd and a high temporal stability of 0.945. In contrast, sequences with rapid, erratic motion or low-contrast targets perform poorly. *breakdance* records $(0.317, 0.190, 0.120)$, and *india* obtains $(0.179, 0.200, 0.338)$, as the dancer's quick pose changes and the small, indistinct object, respectively, present significant challenges. Intermediate performance is observed on *camel, car-shadow*, with mean IoU around 0.77, boundary F between 0.36 and 0.55, and stability around 0.86–0.89, reflecting fairly easy region segmentation but less precise boundary alignment.

## 5.3 Discussion

The high scores on sequences such as *blackswan*, *lab-coat*, and *cows* can be attributed to large, high-contrast objects occupying most of the frame, minimal occlusions, and smooth, predictable motion. Under these conditions, our attention-based graph network and CRF refinement quickly "lock on" to the foreground, yielding strong region overlap ($J > 0.82$), accurate boundaries ($F > 0.50$), and stable temporal consistency ($T > 0.88$). Conversely, videos like *breakdance* and

Table 2: Selected per-video mean metrics on DAVIS17.

| Video | $\bar{J}$ | $\bar{F}$ | $\bar{T}$ |
|---|---|---|---|
| blackswan | 0.886 | 0.701 | 0.906 |
| lab-coat | 0.919 | 0.637 | 0.884 |
| cows | 0.827 | 0.507 | 0.945 |
| camel, car-shadow | 0.770 | 0.455 | 0.875 |
| breakdance | 0.317 | 0.190 | 0.120 |
| india | 0.179 | 0.200 | 0.338 |
| soapbox | 0.695 | 0.455 | 0.793 |

*india* present rapid, erratic motions, cluttered backgrounds, or very small, low-contrast objects that challenge a zero-shot model with a fixed receptive field. As a result, these sequences produce low region IoU ($J < 0.32$), poor boundary alignment ($F < 0.20$), and significant temporal flicker ($T < 0.35$).

Even in cases where $\bar{J} \approx 0.70$–$0.80$ (e.g., *plausible examples*), the corresponding $\bar{F}$ can be substantially lower (0.27–0.36), indicating that while the model correctly identifies coarse object regions, the predicted contours are often a few pixels offset from the ground-truth mask. Finally, an overall temporal stability of $\bar{T} = 0.723$ demonstrates that more than 70% of consecutive masks overlap well; sequences like *cows* (with $T = 0.945$) exhibit almost perfect frame-to-frame consistency, whereas jittery sequences like *breakdance* ($T = 0.120$) underscore the model's difficulty in maintaining stable predictions during rapid motion. A rule of thumb is that $T < 0.60$ indicates flickering or unstable predictions across frames.

# 6 Appendix A

## 6.1 Recap: Gated Recurrent Unit (GRU)

A Gated Recurrent Unit (GRU) is a recurrent cell that maintains a hidden state $\mathbf{h}_t$ over discrete time steps $t$, enabling the network to remember and update information dynamically.

1. **Reset gate $\mathbf{r}_t$:**
$$\mathbf{r}_t = \sigma\big(W_r\,[\mathbf{x}_t,\ \mathbf{h}_{t-1}]\ +\ \mathbf{b}_r\big) \quad \in [0,1]^{(\text{hidden\_dim})}$$

   Here:

   - $W_r \in \mathbb{R}^{H_{\text{hid}} \times (H_{\text{in}} + H_{\text{hid}})}$ is a learned weight matrix that mixes the current input $\mathbf{x}_t$ and the previous state $\mathbf{h}_{t-1}$.
   - $\mathbf{b}_r \in \mathbb{R}^{H_{\text{hid}}}$ is a learned bias vector.
   - $\sigma$ is the sigmoid, squashing each element to $[0,1]$.

   The reset gate controls how much of the previous memory to "forget" when computing the candidate state.

2. **Update gate $\mathbf{z}_t$:**
$$\mathbf{z}_t = \sigma\big(W_z\,[\mathbf{x}_t,\ \mathbf{h}_{t-1}]\ +\ \mathbf{b}_z\big)$$

   Analogous to the reset gate, $W_z$ and $\mathbf{b}_z$ determine how much of the new candidate state should replace the old state. A value near 1 means "keep the new," near 0 means "keep the old."

3. **Candidate hidden state $\tilde{\mathbf{h}}_t$:**
$$\tilde{\mathbf{h}}_t = \tanh\big(W_h\,[\mathbf{x}_t,\ \mathbf{r}_t \odot \mathbf{h}_{t-1}]\ +\ \mathbf{b}_h\big)$$

   where:

- $W_h \in \mathbb{R}^{H_{\text{hid}} \times (H_{\text{in}} + H_{\text{hid}})}$ and $\mathbf{b}_h$ are learned parameters.
- $\mathbf{r}_t \odot \mathbf{h}_{t-1}$ element-wise gates the old state, allowing the model to reset parts of the memory selectively.

4. **Final state update:**
$$\mathbf{h}_t = (1 - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \tilde{\mathbf{h}}_t.$$

The update gate $\mathbf{z}_t$ blends the old state and the candidate state, enabling the network to retain useful information over long sequences.

## 6.2 Convolutional GRU (ConvGRU)

To apply the same memory mechanism to spatial feature maps of shape $(C \times H \times W)$, we replace each linear transform $W[\cdot]$ with a *convolution*, yielding the ConvGRU.

1. **Concatenate inputs:** $\mathbf{u} = [\mathbf{x}_t, \ \mathbf{h}_{t-1}] \in \mathbb{R}^{B \times (C_{\text{in}} + C_{\text{hid}}) \times H \times W}$. This stacks the new input feature map and the previous hidden-state map along the channel dimension.

2. **Compute reset & update gates with one convolution:**
$$\begin{bmatrix} \mathbf{r}_t, \ \mathbf{z}_t \end{bmatrix} = \sigma\big(\text{Conv2d}_{(2C_{\text{hid}})}(\mathbf{u})\big),$$

where the conv layer has $(C_{\text{in}} + C_{\text{hid}})$ input channels and outputs $2\,C_{\text{hid}}$ channels, which are then split into $\mathbf{r}_t$ and $\mathbf{z}_t$.

- *Why convolution?* Spatially shared weights reduce parameters and preserve the $(H, W)$ grid.

3. **Candidate state with gated previous hidden:**
$$\tilde{\mathbf{u}} = [\,\mathbf{x}_t, \ \mathbf{r}_t \odot \mathbf{h}_{t-1}\,], \quad \tilde{\mathbf{h}}_t = \tanh\big(\text{Conv2d}_{C_{\text{hid}}}(\tilde{\mathbf{u}})\big).$$

Reset gate $\mathbf{r}_t$ selectively filters the old state before combining with the new input.

4. **Final update combining old and candidate:**
$$\mathbf{h}_t = (1 - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \tilde{\mathbf{h}}_t.$$

**Key advantages of ConvGRU:**

- *Spatial consistency:* Padding preserves spatial resolution $(H, W)$.

- *Parameter efficiency:* Convolutional filters are shared across all pixel locations.

- *Iterative memory:* Nodes maintain and refine their feature maps through multiple message-passing iterations.

# 7 Appendix B

## 7.1 Conditional Random Field

In semantic segmentation—even when the network predicts a smooth probability map—boundary details often remain fuzzy or bleed into the background. A Conditional Random Field (CRF) is a probabilistic graphical model that lets us "sharpen" those boundaries by fusing two sources of information:

- Unary potentials derived from network's softmax probabilities

- Pairwise potentials encoding preferences that neighboring pixels should agree, weighted by both spatial proximity and appearance similarity

After our network produces an *unary* foreground probability map $\hat{S} \in [0,1]^{H \times W}$, we define a CRF energy over the binary mask $X \in \{0,1\}^{H \times W}$:

$$E(X) = \sum_i \psi_u(x_i) + \sum_{i<j} \psi_p(x_i, x_j),$$

where

$$\psi_u(x_i) = \begin{cases} -\log \hat{S}_i, & x_i = 1 \\ -\log\big(1 - \hat{S}_i\big), & x_i = 0 \end{cases}$$

and

$$\psi_p(x_i, x_j) = \mu(x_i, x_j) \left( w_1 \exp\!\big(-\tfrac{\|p_i - p_j\|^2}{2\sigma_\alpha^2} - \tfrac{\|I_i - I_j\|^2}{2\sigma_\beta^2}\big) + w_2 \exp\!\big(-\tfrac{\|p_i - p_j\|^2}{2\sigma_\gamma^2}\big) \right).$$

Here:

- $i, j$ index pixels, with $p_i \in \mathbb{R}^2$ their image coordinates and $I_i \in \mathbb{R}^3$ their RGB values.

- $\mu(x_i, x_j) = 1$ if $x_i \neq x_j$ (Potts model), else 0.

- $w_1, w_2, \sigma_\alpha, \sigma_\beta, \sigma_\gamma$ are hyper-parameters controlling the strength and scale of the Gaussian kernels in position and color space.

Approximate minimum-energy inference is performed via mean-field iterations (as in "Dense-CRF" [Krähenbühl Koltun, 2011]), which refines $\hat{S}$ into a sharper mask respecting image edges. This post-processing typically costs ∼0.5s/frame.

# References

[1] Wenguan Wang, Xiankai Lu, Jianbing Shen, David J. Crandall, and Ling Shao. Zero-Shot Video Object Segmentation via Attentive Graph Neural Networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.