

SAVA-X: Ego-to-Exo Imitation Error Detection via Scene-Adaptive View Alignment and Cross-Fusion

Supplementary Material

A. Implementation Details

A.1. Metric Computation

AUPRC computation. For a fixed temporal IoU threshold, we collect all predicted segments and their confidence scores $\{s_i\}_{i=1}^N$ together with binary labels $\{y_i\}_{i=1}^N$ indicating whether each prediction is a true positive ($y_i = 1$) or a false positive ($y_i = 0$). Let P denote the total number of positive ground-truth instances. Following the COCO protocol, we first construct a precision–recall curve from these predictions and then approximate the area under the curve (AUPRC) by averaging the interpolated precision values at 101 uniformly sampled recall points in $[0, 1]$. This yields an interpolated average precision, which we report as AUPRC. In our experiments, we apply this procedure separately for different “positive” definitions (e.g., correct vs. error events) by treating each target category as positive and all others as negative.

A.2. Base Encoder

Our base encoder follows the multi-scale temporal convolutional encoder in Wang et al. [5]. Given frame visual features $\mathbf{v} \in \mathbb{R}^{T \times d}$ from the backbone, we first reshape them to a 1D feature map of size $T \times d$ and feed them into a lightweight temporal pyramid. The first level uses a 1×1 Conv1D and GroupNorm to project the input feature dimension to the transformer hidden dimension. Subsequent levels apply 3×3 Conv1D with stride 2 and GroupNorm to progressively downsample the sequence, producing multi-scale features with decreasing temporal resolution and shared hidden dimension. For each level, we add a sine-based temporal positional encoding and pass the resulting feature maps and masks to the Deformable Transformer encoder. This design keeps the hidden sizes and number of feature levels consistent with [5], so that performance gains mainly stem from our cross-view modules rather than changes in the base encoder.

A.3. Deformable DETR

Our event detection head is a 1D multi-scale Deformable DETR that follows the design of Wang et al. [5]. Given the multi-scale temporal features and masks from the base encoder, we first flatten all levels into a single sequence and add level-specific embeddings and temporal positional encodings. The encoder then applies L_{enc} layers of multi-scale deformable self-attention and feed-

Algorithm 1 COCO-style AUPRC computation for a fixed tIoU

Require: Scores s_i , labels $y_i \in \{0, 1\}$ for $i = 1, \dots, N$; number of positives $P = \sum_i y_i$

Ensure: AUPRC value AP

- 1: Sort indices π such that $s_{\pi_1} \geq s_{\pi_2} \geq \dots \geq s_{\pi_N}$
 - 2: Initialize cumulative true/false positives: $TP[k] = \sum_{i=1}^k \mathbb{1}[y_{\pi_i} = 1]$, $FP[k] = \sum_{i=1}^k \mathbb{1}[y_{\pi_i} = 0]$
 - 3: **for** $k = 1$ to N **do**
 - 4: precision $[k] \leftarrow \frac{TP[k]}{\max(1, TP[k] + FP[k])}$
 - 5: recall $[k] \leftarrow \frac{TP[k]}{\max(1, P)}$
 - 6: **end for**
 - 7: **For** $k = N - 1$ down to 1 :
 - 8: precision $[k] \leftarrow \max(\text{precision}[k], \text{precision}[k + 1])$
 - 9: Initialize AP $\leftarrow 0$
 - 10: **for** $j = 0$ to 100 **do**
 - 11: $r_j \leftarrow j/100$ {101 uniformly sampled recall points}
 - 12: Find the smallest index k such that recall $[k] \geq r_j$
 - 13: **if** such k exists **then**
 - 14: $p(r_j) \leftarrow \text{precision}[k]$
 - 15: **else**
 - 16: $p(r_j) \leftarrow 0$
 - 17: **end if**
 - 18: AP $\leftarrow \text{AP} + p(r_j)$
 - 19: **end for**
 - 20: AP $\leftarrow \text{AP}/101$
 - 21: **return** AP
-

forward networks to produce a unified feature memory. The decoder takes a fixed set of learnable query embeddings and performs, at each of its L_{dec} layers, (i) self-attention among queries and (ii) multi-scale deformable cross-attention to the encoded memory using normalized reference points, followed by a feed-forward network. As in the original Deformable DETR, we use iterative reference-point refinement via a small MLP head attached to each decoder layer, and share the hidden dimension and attention configuration with [5]. This keeps the detection head identical to the PDVC Deformable DETR implementation, so that the performance gains in our experiments mainly arise from the proposed cross-view modules rather than changes in the underlying transformer architecture.

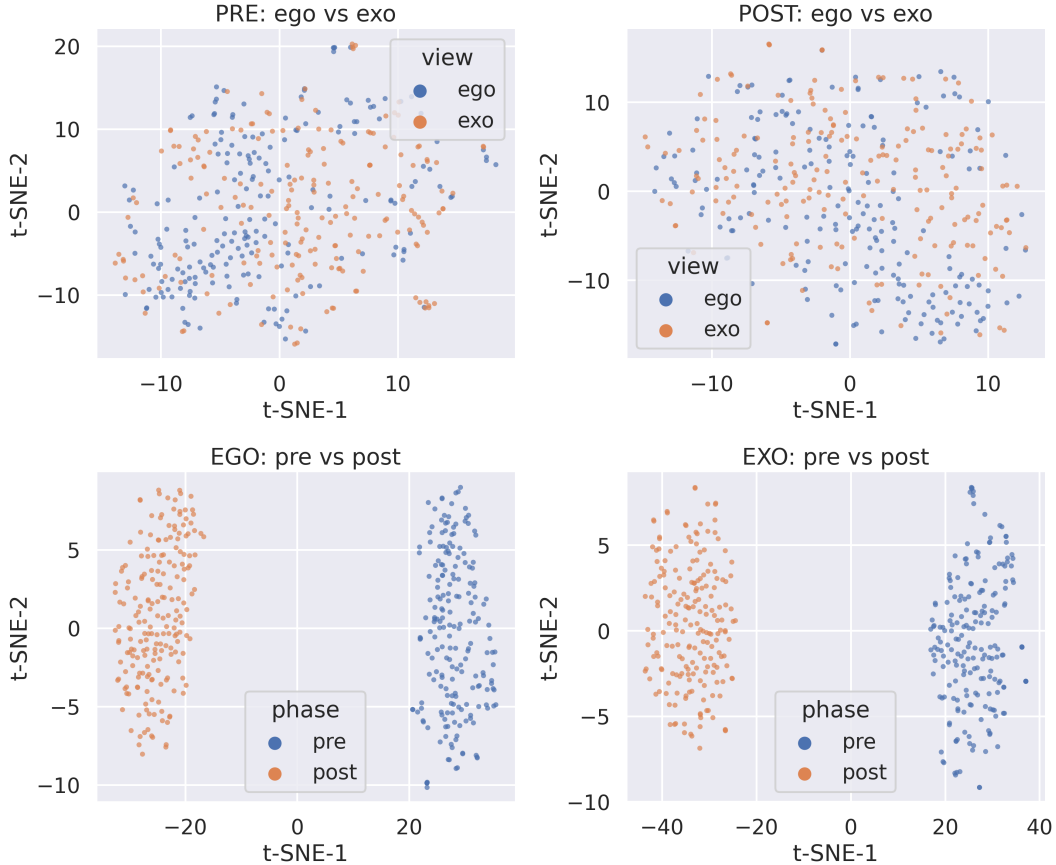


Figure A.1. t-SNE visualization of video-level features before and after SVE. **Top**: ego and exo features colored by view; SVE reduces the cross-view gap and yields more overlapped distributions. **Bottom**: ego (left) and exo (right) features colored by phase (pre vs. post), showing that SVE applies a structured, non-trivial transformation while avoiding representational collapse.

A.4. Task-specific heads

Except for the imitation-error prediction heads introduced in our work, all task heads follow Wang et al. [5]. Given the decoder features for each query, we apply a linear classification head to predict the event category and a regression MLP to predict the normalized start/end coordinates of the temporal segment.

Following PDVC, we additionally use a count head to estimate the number of events and a captioning head to generate a natural-language description for each detected segment, sharing the same hidden dimension and decoder outputs as in [5]. On top of these standard heads, we attach two lightweight binary imitation-error heads: (i) a fine-grained query-level error head that predicts whether each event segment corresponds to a correct or erroneous execution step, and (ii) a global video-level head that aggregates query features to predict the overall imitation quality of the sequence. These additions leave the original PDVC heads unchanged, ensuring that performance gains mainly stem from our cross-view alignment and error modeling rather than modifi-

cations to the baseline detection and captioning heads.

A.5. Losses

Dense video captioning loss. Following PDVC [5], we treat dense video captioning as a set prediction problem and supervise all decoder layers with a Hungarian-matching loss. Let $\{\hat{s}_i, \hat{c}_i, \hat{y}_i\}_{i=1}^N$ denote the predicted temporal segments (center-length parameterization), foreground scores, and caption word distributions for N event queries in one decoder layer, and let $\{s_j, y_j\}_{j=1}^{N_{gt}}$ be the ground-truth segments and captions. We first solve a bipartite matching between predictions and ground truths with cost

$$\mathcal{C}_{ij} = \alpha_{\text{giou}} \mathcal{L}_{\text{giou}}(\hat{s}_i, s_j) + \alpha_{\text{cls}} \mathcal{L}_{\text{cls}}(\hat{c}_i, \mathcal{K}[j \leq N_{gt}]) \quad (\text{A.1})$$

where $\mathcal{L}_{\text{giou}}$ is the temporal generalized IoU loss and \mathcal{L}_{cls} is the focal classification loss between foreground/background. Given the optimal assignment σ ,

the DVC loss for one decoder layer is

$$\begin{aligned}
\mathcal{L}_{\text{DVC}} &= \beta_{\text{giou}} \mathcal{L}_{\text{giou}} + \beta_{\text{cls}} \mathcal{L}_{\text{cls}} + \beta_{\text{ec}} \mathcal{L}_{\text{ec}} + \beta_{\text{cap}} \mathcal{L}_{\text{cap}} \\
\mathcal{L}_{\text{giou}} &= \frac{1}{N_{\text{gt}}} \sum_{j=1}^{N_{\text{gt}}} (1 - \text{GIoU}(\hat{\mathbf{s}}_{\sigma(j)}, \mathbf{s}_j)) \\
\mathcal{L}_{\text{cls}} &= \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{\text{focal}}(\hat{\mathbf{c}}_i, y_i^{\text{fg/bg}}) \\
\mathcal{L}_{\text{ec}} &= -\log p_{\text{ec}}(N_{\text{gt}}) \\
\mathcal{L}_{\text{cap}} &= \frac{1}{N_{\text{gt}}} \sum_{j=1}^{N_{\text{gt}}} \frac{1}{T_j} \sum_{t=1}^{T_j} -\log p(w_{j,t} | w_{j,<t}, \hat{\mathbf{z}}_{\sigma(j)})
\end{aligned} \tag{A.2}$$

where p_{ec} is the event-counter distribution over event numbers, $\hat{\mathbf{z}}_{\sigma(j)}$ is the matched query feature, and $w_{j,t}$ is the t -th word of the j -th ground-truth caption of length T_j . As in [5], we attach prediction heads to all decoder layers and sum the layer-wise losses.

Imitation-error classification losses. On top of the standard DVC loss, we introduce two lightweight error-prediction heads: a fine-grained query-level head and a global video-level head (Sec. A.4). For each matched query j , let $\hat{z}_{\sigma(j)}^{\text{fine}}$ be the scalar logit from the fine-grained error head and $e_j^{\text{fine}} \in \{0, 1\}$ be the corresponding binary error label (1 for erroneous execution, 0 for correct). We supervise this head with a binary cross-entropy (BCE) loss over the matched events:

$$\begin{aligned}
\mathcal{L}_{\text{err}}^{\text{fine}} &= \frac{1}{N_{\text{gt}}} \sum_{j=1}^{N_{\text{gt}}} \left[-e_j^{\text{fine}} \log \sigma(\hat{z}_{\sigma(j)}^{\text{fine}}) \right. \\
&\quad \left. - (1 - e_j^{\text{fine}}) \log (1 - \sigma(\hat{z}_{\sigma(j)}^{\text{fine}})) \right]
\end{aligned} \tag{A.3}$$

where $\sigma(\cdot)$ is the sigmoid function. In addition, the global head aggregates all valid query features of a video into a single logit \hat{z}^{overall} that predicts the overall imitation quality, with binary label $e^{\text{overall}} \in \{0, 1\}$ (error-free vs. containing errors). We apply another BCE loss:

$$\begin{aligned}
\mathcal{L}_{\text{err}}^{\text{overall}} &= -e^{\text{overall}} \log \sigma(\hat{z}^{\text{overall}}) \\
&\quad - (1 - e^{\text{overall}}) \log (1 - \sigma(\hat{z}^{\text{overall}}))
\end{aligned} \tag{A.4}$$

Adaptive sampling regularization. We further regularize the adaptive sampler (AS) to avoid collapsed selections and redundant representations. Denote by $\{s_{x,t}\}_{t=1}^{T_x}$ and $\{s_{y,t}\}_{t=1}^{T_y}$ the normalized selection probabilities over Ego/Exo tokens. We add a selection-entropy regularizer [3] that encourages coverage instead

of concentrating mass on a few positions:

$$\begin{aligned}
\mathcal{L}_{\text{sel}} &= \frac{1}{\log T_x} \sum_{t=1}^{T_x} s_{x,t} \log(s_{x,t} + \varepsilon) \\
&\quad + \frac{1}{\log T_y} \sum_{t=1}^{T_y} s_{y,t} \log(s_{y,t} + \varepsilon)
\end{aligned} \tag{A.5}$$

where $\varepsilon > 0$ is a small constant for numerical stability. Let $u \in \{\text{exo}, \text{ego}\}$ index the view and $\hat{\mathbf{Z}}^u \in \mathbb{R}^{K_u \times d}$ be the matrix of gated active tokens (after selection and gating), with K_u tokens and feature dimension d . We further attach VICReg-style [1] variance and covariance penalties to suppress collapse and dimensional collinearity:

$$\begin{aligned}
\boldsymbol{\mu}^u &= \frac{1}{K_u} \sum_{i=1}^{K_u} \hat{\mathbf{Z}}_i^u, \quad \hat{\mathbf{Z}}_c^u = \hat{\mathbf{Z}}^u - \mathbf{1} \boldsymbol{\mu}^{u\top} \\
\mathcal{L}_{\text{var}}^u &= \frac{1}{d} \sum_{j=1}^d \left[\max(0, \gamma - \sqrt{\text{Var}(\hat{\mathbf{Z}}_{c,j}^u) + \varepsilon}) \right]^2 \\
\mathbf{C}^u &= \frac{1}{K_u - 1} \hat{\mathbf{Z}}_c^{u\top} \hat{\mathbf{Z}}_c^u, \quad \mathcal{L}_{\text{cov}}^u = \frac{1}{d} \sum_{\substack{i=1 \\ i \neq j}}^d \sum_{j=1}^d (\mathbf{C}_{ij}^u)^2 \\
\mathcal{L}_{\text{vic}} &= \mathcal{L}_{\text{var}}^{\text{exo}} + \mathcal{L}_{\text{var}}^{\text{ego}} + \mathcal{L}_{\text{cov}}^{\text{exo}} + \mathcal{L}_{\text{cov}}^{\text{ego}}
\end{aligned} \tag{A.6}$$

where $\gamma > 0$ is the variance lower bound and $\varepsilon > 0$ again ensures numerical stability.

B. Results

B.1. Correct Class Results

Table A.1 summarizes the AUPRC performance for the *correct* (non-error) class on EgoMe validation and test splits.

Among DVC-style baselines, PDVC [5] remains strong, while Exo2EgoDVC [2] does not consistently improve over PDVC despite explicitly transferring exocentric knowledge. TAL-only models (ActionFormer [6], TriDet [4]) perform clearly worse on AUPRC and tIoU, showing that off-the-shelf localization architectures are not sufficient for our fine-grained imitation setting.

Using only egocentric input further degrades PDVC, highlighting the benefit of multi-view information even when evaluating the correct class.

Our SAVA-X achieves the best validation performance across all tIoU thresholds and mean AUPRC, and attains comparable or better test performance than PDVC, with particularly noticeable gains at high tIoU (e.g., +1.1 AUPRC@0.7 on validation). These results indicate that SAVA-X not only improves error detection, but also preserves or slightly enhances recognition and

Method	AUPRC on Validation					AUPRC on Test				
	0.3	0.5	0.7	Mean	tIoU	0.3	0.5	0.7	Mean	tIoU
<i>Dense Video Captioning (DVC) baselines</i>										
PDVC [5]	68.41	46.21	12.72	42.45	58.58	66.53	43.88	12.30	40.90	57.98
Exo2EgoDVC [2]	67.52	44.52	12.43	41.49	59.06	65.37	42.27	11.40	39.68	58.15
<i>Temporal Action Localization (TAL) baselines</i>										
ActionFormer [6]	65.87	31.68	4.54	34.03	48.89	63.25	29.20	4.17	32.20	48.25
TriDet [4]	65.05	32.25	4.35	33.88	49.05	62.45	30.29	4.30	32.35	49.02
<i>Only Egocentric Input</i>										
PDVC [5]	64.77	42.11	10.94	39.27	57.63	63.94	40.96	12.03	38.98	57.19
<i>Ours</i>										
SAVA-X	69.02	46.58	13.85	43.15	59.31	66.32	43.98	12.41	40.90	58.32

Table A.1. Comparison on EgoMe validation and test split. Left: results on *validation set*. Right: results on the *test set*. We report AUPRC for the correct class at multiple tIoU thresholds (0.3, 0.5, 0.7), their mean, and standalone temporal IoU (tIoU) for localization quality.

localization of correct executions, rather than trading off one class for the other.

B.2. TSNE

To better understand how SVE reshapes the feature space, we visualize pre- and post-SVE video-level features using t-SNE (Fig. A.1).

In the top row, we color points by view (ego vs. exo). Before SVE (top-left), ego and exo features form two partially separated clouds with a clear domain shift.

After SVE (top-right), the two distributions become much more interleaved, indicating that SVE effectively reduces the cross-view gap and encourages a more view-invariant embedding. In the bottom row, we fix the view and color points by phase (pre vs. post).

For both ego (bottom-left) and exo (bottom-right), pre- and post-SVE features form two well-separated clusters along the first t-SNE dimension, showing that SVE applies a non-trivial, structured transformation to the representations rather than a small perturbation.

Together, these plots suggest that SVE consistently aligns ego and exo distributions while preserving meaningful intra-view variability and avoiding collapse.

References

- [1] Adrien Bardes, Jean Ponce, and Yann LeCun. VICReg: Variance-invariance-covariance regularization for self-supervised learning. In *International Conference on Learning Representations*, 2022. 3
- [2] Takehiko Ohkawa, Takuma Yagi, Taichi Nishimura, Ryosuke Furuta, Atsushi Hashimoto, Yoshitaka Ushiku, and Yoichi Sato. Exo2egodvc: Dense video captioning of egocentric procedural activities using web instructional videos. In *2025 IEEE/CVF Winter Conference on Appli-*

cations of Computer Vision (WACV), pages 8324–8335. IEEE, 2025. 3, 4

- [3] Gabriel Pereyra, George Tucker, Jan Chorowski, Lukasz Kaiser, and Geoffrey Hinton. Regularizing neural networks by penalizing confident output distributions. In *International Conference on Learning Representations*, 2017. 3
- [4] Dingfeng Shi, Yujie Zhong, Qiong Cao, Lin Ma, Jia Li, and Dacheng Tao. Tridet: Temporal action detection with relative boundary modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18857–18866, 2023. 3, 4
- [5] Teng Wang, Ruimao Zhang, Zhichao Lu, Feng Zheng, Ran Cheng, and Ping Luo. End-to-End Dense Video Captioning with Parallel Decoding. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6827–6837, Montreal, QC, Canada, 2021. IEEE. 1, 2, 3, 4
- [6] Chenlin Zhang, Jianxin Wu, and Yin Li. Actionformer: Localizing moments of actions with transformers. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 492–510. Springer, 2022. 3, 4