

# Imagine Before Concentration: Diffusion-Guided Registers Enhance Partially Relevant Video Retrieval

## Supplementary Material

### A. More Details on Method

#### A.1. Derivation of ELBO in Eq. (2)

Given the following predictive function:

$$p_{\theta,\phi}(Q|V) = \int \underbrace{p_{\theta}(Q|V, r)}_{\text{concentration}} \cdot \underbrace{p_{\phi}(r|V)}_{\text{imagination}} dr, \quad (18)$$

where  $p_{\phi}(r|V)$  denotes the video branch that first generates global registers, and  $p_{\theta}(Q|V, r)$  models the register-augmented cross-modal matching. Our training objective is to maximize  $p_{\theta,\phi}(Q|V)$ . We introduce a variational posterior  $q_{\varphi}(r|Q_a)$  to approximate the true posterior  $p(r|Q_a)$ .

First, we rewrite Eq. (18) by introducing  $q_{\varphi}(r|Q_a)$ :

$$\begin{aligned} & \log p_{\theta,\phi}(Q|V) \\ &= \log \int p_{\theta}(Q|V, r) p_{\phi}(r|V) dr \\ &= \log \int p_{\theta}(Q|V, r) \frac{p_{\phi}(r|V)}{q_{\varphi}(r|Q_a)} q_{\varphi}(r|Q_a) dr. \end{aligned} \quad (19)$$

Next, we invoke Jensen’s inequality [5], leveraging the concavity of the log function, which satisfies:

$$\log \mathbb{E}[X] \geq \mathbb{E}[\log X]. \quad (20)$$

Thus, the logarithm can be moved outside the integral, yielding a tractable lower bound:

$$\begin{aligned} & \log p_{\theta,\phi}(Q|V) \\ & \geq \int q_{\varphi}(r|Q_a) \log \left( p_{\theta}(Q|V, r) \frac{p_{\phi}(r|V)}{q_{\varphi}(r|Q_a)} \right) dr. \end{aligned} \quad (21)$$

We subsequently decompose the logarithmic term within the integrand:

$$\begin{aligned} & \log \left( p_{\theta}(Q|V, r) \frac{p_{\phi}(r|V)}{q_{\varphi}(r|Q_a)} \right) \\ &= \log p_{\theta}(Q|V, r) + \log \frac{p_{\phi}(r|V)}{q_{\varphi}(r|Q_a)}. \end{aligned} \quad (22)$$

Hence, the lower bound becomes:

$$\begin{aligned} \log p_{\theta,\phi}(Q|V) & \geq \int q_{\varphi}(r|Q_a) \log p_{\theta}(Q|V, r) dr \\ & \quad - \int q_{\varphi}(r|Q_a) \log \frac{q_{\varphi}(r|Q_a)}{p_{\phi}(r|V)} dr. \end{aligned} \quad (23)$$

Therefore, we obtain the final form of the ELBO in Eq. (2):

$$\begin{aligned} \log p_{\theta,\phi}(Q|V) & \geq \mathbb{E}_{q_{\varphi}(r|Q_a)} [\log p_{\theta}(Q|V, r)] \\ & \quad - \mathbb{KL}[q_{\varphi}(r|Q_a) \| p_{\phi}(r|V)]. \end{aligned} \quad (24)$$

#### A.2. Relationship between Eq. (7) and $L_{\text{dre}}$

Following [7], Bayes’ rule gives:

$$q(\hat{\mathbf{q}}_t | \hat{\mathbf{q}}_{t-1}, \hat{\mathbf{q}}_0) = \frac{q(\hat{\mathbf{q}}_{t-1} | \hat{\mathbf{q}}_t, \hat{\mathbf{q}}_0) q(\hat{\mathbf{q}}_t | \hat{\mathbf{q}}_0)}{q(\hat{\mathbf{q}}_{t-1} | \hat{\mathbf{q}}_0)}. \quad (25)$$

Armed with this new equation, we can retry the derivation resuming from the ELBO in Eq. (7) by viewing  $\hat{\mathbf{q}}$  as  $\hat{\mathbf{q}}_0$ :

$$\begin{aligned} & \log p(\hat{\mathbf{q}}_0) \\ &= \log \int p(\hat{\mathbf{q}}_{0:T}) d\hat{\mathbf{q}}_{1:T} \\ &= \log \mathbb{E}_{q(\hat{\mathbf{q}}_{1:T} | \hat{\mathbf{q}}_0)} \left[ \frac{p(\hat{\mathbf{q}}_{0:T})}{q(\hat{\mathbf{q}}_{1:T} | \hat{\mathbf{q}}_0)} \right] \\ & \geq \underbrace{\mathbb{E}_{q(\hat{\mathbf{q}}_1 | \hat{\mathbf{q}}_0)} [\log p_{\phi}(\hat{\mathbf{q}}_0 | \hat{\mathbf{q}}_1)]}_{\text{(reconstruction term)}} - \underbrace{\mathbb{KL}(q(\hat{\mathbf{q}}_T | \hat{\mathbf{q}}_0) \| p(\hat{\mathbf{q}}_T))}_{\text{(prior matching term)}} \\ & \quad - \underbrace{\sum_{t=2}^T \mathbb{E}_{q(\hat{\mathbf{q}}_t | \hat{\mathbf{q}}_0)} [\mathbb{KL}(q(\hat{\mathbf{q}}_{t-1} | \hat{\mathbf{q}}_t, \hat{\mathbf{q}}_0) \| p_{\phi}(\hat{\mathbf{q}}_{t-1} | \hat{\mathbf{q}}_t))]}_{\text{(denoising matching term)}}, \end{aligned} \quad (26)$$

where (i) the reconstruction term corresponds to the negative reconstruction error over  $\hat{\mathbf{q}}_0$ ; (ii) the prior matching term is constant with no trainable parameters and can thus be ignored during optimization; and (iii) the denoising matching terms constrain  $p_{\phi}(\hat{\mathbf{q}}_{t-1} | \hat{\mathbf{q}}_t)$  to align with the tractable ground-truth transition  $q(\hat{\mathbf{q}}_{t-1} | \hat{\mathbf{q}}_t, \hat{\mathbf{q}}_0)$  [7]. Consequently,  $\phi$  is optimized to iteratively recover  $\hat{\mathbf{q}}_{t-1}$  from  $\hat{\mathbf{q}}_t$ . Following [4], the denoising matching terms can be simplified as

$$\sum_{t=2}^T \mathbb{E}_{t,\epsilon} \left[ \|\epsilon - \epsilon_{\phi}(\hat{\mathbf{q}}_t, t)\|_2^2 \right], \quad (27)$$

where  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and  $\epsilon_{\phi}(\hat{\mathbf{q}}_t, t)$  is parameterized by a neural network (e.g., U-Net [4]) to predict the noise  $\epsilon$  that generates  $\hat{\mathbf{q}}_t$  from  $\hat{\mathbf{q}}_0$  in the forward process [7]. A detailed derivation of Eq. (7) and Eq. (26) is provided in [7].

$L_{\text{dre}}$  extends Eq. (27) by incorporating a conditioning variable  $\mathbf{c}$ . Inspired by [4], the reconstruction term in Eq. (26) has a relatively minor effect; hence,  $L_{\text{dre}}$  is employed to optimize the denoising matching terms in Eq. (26), thereby providing an approximate estimation of Eq. (26) for training.

#### A.3. Further details of DreamPRVR architecture

**Diffusion Register Estimator (DRE)** As illustrated in Fig. 7(a), the DRE block follows an MLP-based architecture incorporating Layer Normalization [1], activation functions, and linear projection layers. We use  $N_{\text{dre}} = 2$  blocks.

**Condition Generator** The condition  $\mathbf{c} \in \mathbb{R}^{N_r \times d}$  is obtained via a simple cross-attention mechanism between  $\mathbf{V}_v \in \mathbb{R}^{N_v \times d}$  and learnable parameters, as illustrated in Fig. 7 (b), and can be formulated as

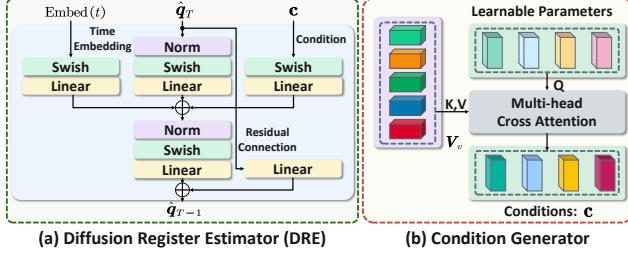


Figure 7. (a) Illustration of the proposed Diffusion Register Estimator Block (DRE).  $\text{Embed}(t)$ ,  $q_T$  and  $c$  denote the temporal embedding, the latent embedding corrupted by  $t$ -step noise and the guided condition, respectively. (b) Condition generator for DRE.

$$c = \text{CA}(LP, V_v, V_v) \in \mathbb{R}^{N_r \times d}, \quad (28)$$

where CA denotes cross-attention, and  $LP \in \mathbb{R}^{N_r \times d}$  represents learnable parameters.

**Asymmetric Attention Mask** We retain the Gaussian self-attention as in [6, 9] and instead define two cross-attention patterns through a designed masking strategy, as illustrated in Fig. 2(d). Given the global registers  $r_0$  and video embeddings  $V_o$ , the cross-attention configurations are defined as follows. For video embeddings:

$$\text{Query} = V_o, \quad \text{Key} = \text{Value} = \text{Concat}([V_o, r_0]). \quad (29)$$

For global registers:

$$\text{Query} = r_0, \quad \text{Key} = \text{Value} = V_o. \quad (30)$$

#### A.4. Learning Objectives

**Standard Similarity Retrieval Loss  $L_{\text{sim}}$**  Following prior works [3, 6, 10], we employ the widely adopted triplet loss [2]  $L^{\text{trip}}$  and InfoNCE loss [8, 11]  $L^{\text{ncc}}$  for PRVR. A text–video pair is treated as positive if the video contains a moment relevant to the query; otherwise, it is considered negative. Given a positive pair  $(Q, V)$ , the triplet ranking loss over a mini-batch  $\mathcal{B}$  is defined as follows:

$$L^{\text{trip}} = \frac{1}{n} \sum_{(Q, V) \in \mathcal{B}} \{ \max(0, m + S(Q^-, V) - S(Q, V)) + \max(0, m + S(Q, V^-) - S(Q, V)) \}, \quad (31)$$

where  $m$  denotes the margin,  $Q^-$  and  $V^-$  represent the negative text for  $V$  and the negative video for  $Q$ , respectively, and the similarity score  $S(\cdot, \cdot)$  is computed as in Eq. (20). The infoNCE loss is computed as:

$$L^{\text{ncc}} = -\frac{1}{n} \sum_{(Q, V) \in \mathcal{B}} \left\{ \log\left(\frac{S(Q, V)}{S(Q, V) + \sum_{Q_i^- \in \mathcal{N}_Q} S(Q_i^-, V)}\right) + \log\left(\frac{S(Q, V)}{S(Q, V) + \sum_{V_i^- \in \mathcal{N}_V} S(Q, V_i^-)}\right) \right\}, \quad (32)$$

where  $\mathcal{N}_Q$  and  $\mathcal{N}_V$  represent the negative texts and videos of  $V$  and  $Q$  within the mini-batch  $\mathcal{B}$ , respectively. Finally,  $L_{\text{sim}}$  is defined as:

#### Algorithm 1 Register Generation Process during Training

- 1: **Input:** Video features  $V_v$ , all textual features from the video  $q$ , condition features  $c$ , timesteps  $T$ , noise schedule  $\{\beta_t\}_{t=1}^T$ , diffusion register estimator (DRE)  $\epsilon_\phi$ , probabilistic variational sampler (PVS), textual perturbation sampler (TPS)
- 2: **Output:** Optimal Registers  $r_0$ , diffusion loss  $L_D$
- 3: Initialize loss accumulator  $L_D \leftarrow 0$
- 4: Precompute  $\alpha_t = 1 - \beta_t$  and  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$  for all  $t$
- 5: **Feature Encoding**
- 6:  $p(r_T | V_v) \leftarrow \text{PVS}(V_v)$  ▷ Using Eq. (8)
- 7: Sample  $r_T \leftarrow p(r_T | V_v) \sim \mathcal{N}(\mu_v, \sigma_v^2 I)$
- 8:  $p(\hat{q} | q) \leftarrow \text{TPS}(q)$  ▷ Using Eq. (6)
- 9: Sample  $\hat{q} \leftarrow p(\hat{q} | q) \sim \mathcal{N}(\mu_{\hat{q}}, \sigma_{\hat{q}}^2 I)$
- 10: **Forward Diffusion Process**
- 11:  $\hat{q}_0 \leftarrow \hat{q}$
- 12: **for**  $t = 1$  **to**  $T$  **do**
- 13: Sample  $\epsilon \sim \mathcal{N}(0, I)$
- 14:  $\hat{q}_t \leftarrow \sqrt{\alpha_t} \hat{q}_0 + \sqrt{1 - \alpha_t} \epsilon$  ▷ Add noise via Eq. (13)
- 15:  $\hat{\epsilon} \leftarrow \epsilon_\phi(\hat{q}_t, t, c)$  ▷ Predict noise
- 16:  $L_{\text{dre}} \leftarrow \|\epsilon - \hat{\epsilon}\|^2$  ▷ Calculate loss via Eq. (16)
- 17:  $L_D \leftarrow L_D + L_{\text{dre}}$  ▷ Accumulate loss
- 18: **end for**
- 19: **Reverse Generation Process**
- 20: Sample  $\hat{q}_T \leftarrow r_T$  ▷ Start generation
- 21: **for**  $t = T$  **to** 1 **do**
- 22: Predict noise  $\hat{\epsilon} \leftarrow \epsilon_\phi(\hat{q}_t, t, c)$
- 23: **if**  $t > 1$  **then**
- 24: Sample  $z \sim \mathcal{N}(0, I)$
- 25:  $\hat{q}_{t-1} \leftarrow \frac{1}{\sqrt{\alpha_t}} \left( \hat{q}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \hat{\epsilon} \right) + \sqrt{\beta_t} z$
- 26: **else**
- 27:  $\hat{q}_{t-1} \leftarrow \frac{1}{\sqrt{\alpha_t}} \left( \hat{q}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \hat{\epsilon} \right)$
- 28: **end if**
- 29: **end for**
- 30:  $r_0 \leftarrow \hat{q}_0$
- 31: **return**  $r_0, L_D$  ▷ Return features and loss

$$L_{\text{sim}} = L_c^{\text{trip}} + L_f^{\text{trip}} + \lambda_c L_c^{\text{ncc}} + \lambda_f L_f^{\text{ncc}}, \quad (33)$$

where  $f$  and  $c$  denote the objectives of the frame-scale and clip-scale branches, respectively and  $\lambda_f$  and  $\lambda_c$  are the corresponding hyper-parameters.

**Query Diversity Loss  $L_{\text{div}}$**  Following Wang et al. [9], given a collection of text queries in the mini-batch  $\mathcal{B}$ , the query diversity loss is defined as:

$$\ell(i, j) = (1 + \cos(\mathbf{q}_i, \mathbf{q}_j)) \log(1 + e^{\omega(\cos(\mathbf{q}_i, \mathbf{q}_j) + \delta)}), \quad (34)$$

$$L_{\text{div}} = \frac{2}{M_q(M_q - 1)} \sum_{1 \leq i, j \leq M_q, i \neq j} \ell(i, j),$$

where  $\delta > 0$  is a margin factor,  $\omega > 0$  is a scaling factor and  $M_q$  is the number of text queries relevant to a video.

## A.5. Relationship between $L_{\text{DreamPRVR}}$ and $L_{\text{total}}$

$L_{\text{DreamPRVR}}$  is the theoretical training objective defined in Eq. (3). It consists of two components: (i) a KL-divergence term that enforces the registers to generate global contextual semantics consistent with the textual queries, and (ii) a likelihood term that strengthens video representation learning with register guidance, thereby facilitating improved cross-modal alignment and retrieval performance.

$L_{\text{total}}$  is the practical training objective, comprising four components:  $L_{\text{tssl}}$ ,  $L_{\text{pvs}}$ ,  $L_{\text{dre}}$ , and  $L_{\text{sim}}$ . Among them,  $L_{\text{tssl}}$ ,  $L_{\text{pvs}}$ , and  $L_{\text{dre}}$  jointly regularize the registers to generate text-consistent representations and capture richer textual semantics. These terms promote more effective register generation and correspond to optimizing the KL-divergence term in  $L_{\text{DreamPRVR}}$ . In addition,  $L_{\text{sim}}$  serves as the retrieval-oriented similarity learning objective, aiming to improve retrieval performance. This term aligns with maximizing the likelihood component in  $L_{\text{DreamPRVR}}$ .

## A.6. Register Generation Process

**Training Stage** Please refer to Algorithm 1.

**Inference Stage** The procedure follows Algorithm 1, with the forward diffusion process and TPS sampling omitted.

## References

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 1
- [2] Jianfeng Dong, Xirong Li, Chaoxi Xu, Xun Yang, Gang Yang, Xun Wang, and Meng Wang. Dual encoding for video retrieval by text. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(8):4065–4080, 2021. 2
- [3] Jianfeng Dong, Xianke Chen, Minsong Zhang, Xun Yang, Shujie Chen, Xirong Li, and Xun Wang. Partially relevant video retrieval. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 246–257, 2022. 2
- [4] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 1
- [5] Johan Ludwig William Valdemar Jensen. Sur les fonctions convexes et les inégalités entre les valeurs moyennes. *Acta mathematica*, 30(1):175–193, 1906. 1
- [6] Jun Li, Jinpeng Wang, Chaolei Tan, Niu Lian, Long Chen, Yaowei Wang, Min Zhang, Shu-Tao Xia, and Bin Chen. Hlformer: Enhancing partially relevant video retrieval with hyperbolic learning. *arXiv preprint arXiv:2507.17402*, 2025. 2
- [7] Calvin Luo. Understanding diffusion models: A unified perspective. *arXiv preprint arXiv:2208.11970*, 2022. 1
- [8] Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman. End-to-end learning of visual representations from uncurated instructional videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9879–9889, 2020. 2
- [9] Yuting Wang, Jinpeng Wang, Bin Chen, Tao Dai, Ruisheng Luo, and Shu-Tao Xia. Gmmformer v2: An uncertainty-aware framework for partially relevant video retrieval. *arXiv preprint arXiv:2405.13824*, 2024. 2
- [10] Yuting Wang, Jinpeng Wang, Bin Chen, Ziyun Zeng, and Shu-Tao Xia. Gmmformer: Gaussian-mixture-model based transformer for efficient partially relevant video retrieval. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2024. 2
- [11] Hao Zhang, Aixin Sun, Wei Jing, Guoshun Nan, Liangli Zhen, Joey Tianyi Zhou, and Rick Siow Mong Goh. Video corpus moment retrieval with contrastive learning. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 685–695, 2021. 2