



# 第二十届中国图象图形学学会青年科学家会议

## Not All Pairs are Equal: Hierarchical Learning for Average-Precision-Oriented Video Retrieval

Yang Liu<sup>1</sup>, Qianqian Xu<sup>2,\*</sup>, Peisong Wen<sup>1,2</sup>, Siran Dai<sup>4</sup>, Qingming Huang<sup>1,2,3,\*</sup>

1 School of Computer Science and Technology, University of Chinese Academy of Sciences 3 BDKM, University of Chinese Academy of Sciences

2 Institute of Computing Technology, Chinese Academy of Sciences 4 Institute of Information Engineering, Chinese Academy of Sciences

This paper was presented with the **Honourable Mention Award** at the ACM MultiMedia 2024 Conference.



### Alignment of the Objective with the Metric

#### □ Evaluation Metric: Average Precision (AP)

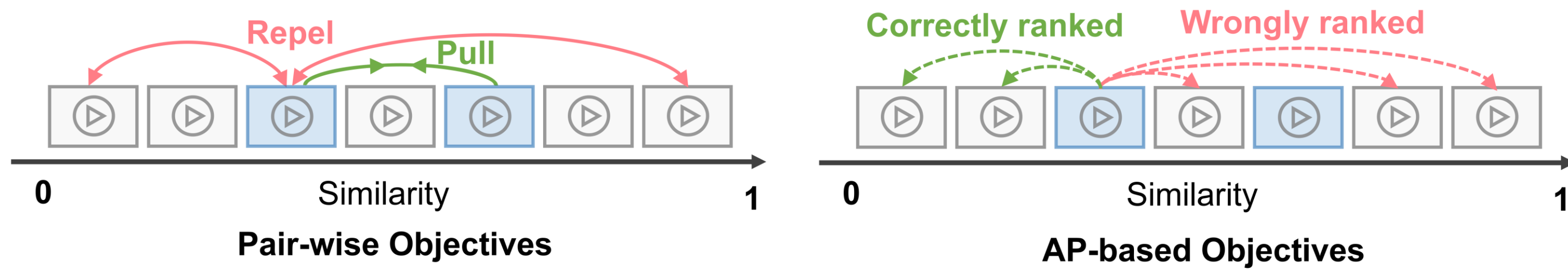
- Evaluates the overall rankings of relevant videos  $AP = \frac{1}{n} \sum_{i=1}^n \frac{i}{r_i}$
- Assigns larger weights on higher-ranked instances

#### □ Previous Training Objectives: Pair-wise Objectives

- Pull the positive instances closer and repel the negative ones
- ✗ Treat all mis-ranked pairs equally
- ✗ Mismatch with the evaluation metric

#### □ New Training Objectives: AP-based Objectives

- Rectify the wrongly ranked positive-negative pairs in the list
- ✓ Fill the gap between training objectives and evaluation metrics



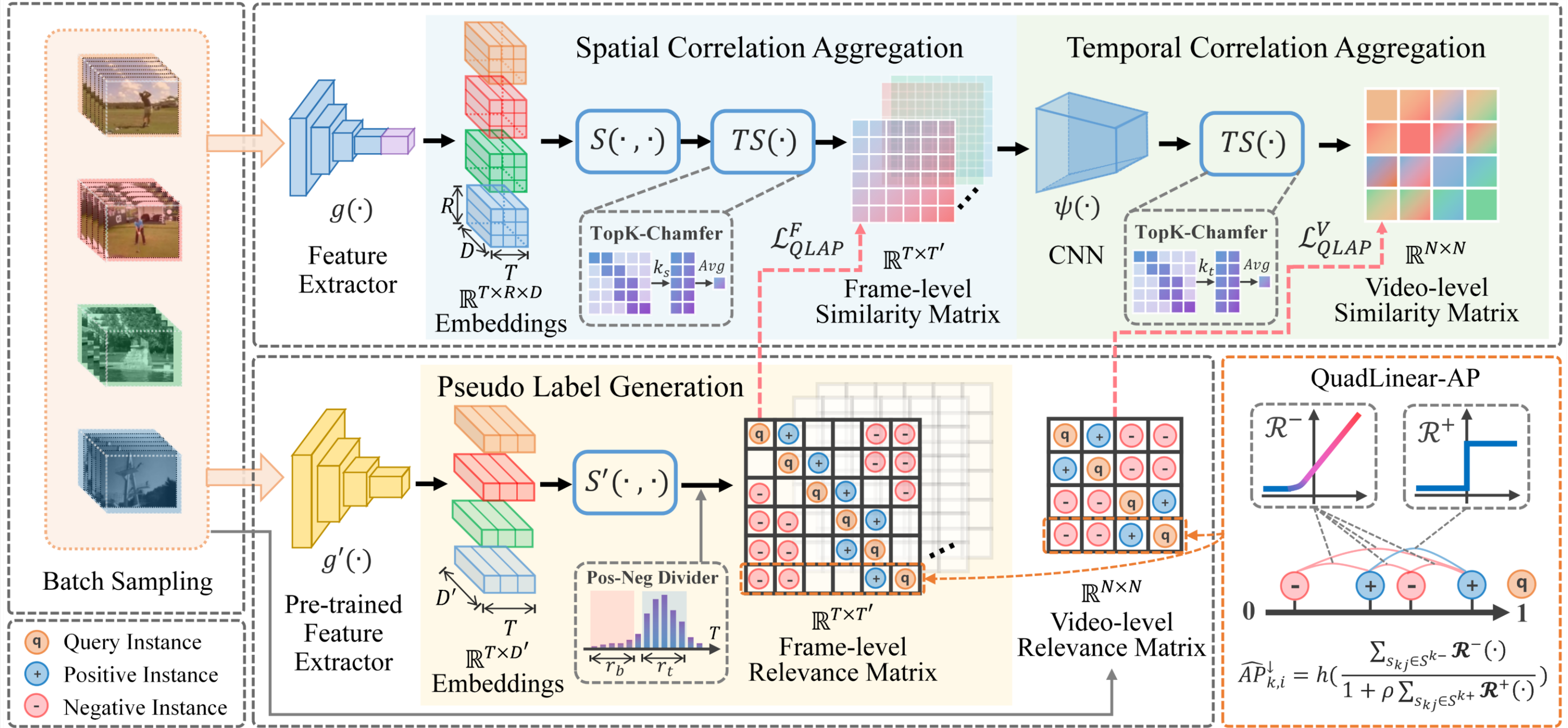
#### □ Two Challenges for Video-Oriented AP Optimization

- Current AP losses are suboptimal for **video-level** retrieval
- Noisy **frame-level** matching leads to a biased AP estimation

### Hierarchical Learning Framework

#### □ Overall Framework: HAP-VR

Hierarchical learning for Average-Precision-oriented Video Retrieval



#### □ Video-Oriented AP Optimization Algorithm

- **Step1:** Bottom-up video similarity measure
- **Step2:** Pseudo-label generation  $\mathcal{L} = \underbrace{\lambda_f \mathcal{L}_{QLAP}^F}_{\text{frame-level}} + \underbrace{\lambda_v \mathcal{L}_{QLAP}^V}_{\text{video-level}} + \mathcal{L}_{base}$
- **Step3:** Hierarchical AP optimization

### Gradient-Enhanced AP Optimization

#### □ Optimization Problem:

- Maximizing the AP score

$$\max_f AP(f) = \frac{1}{N} \sum_{k=1}^N AP_k(f)$$

#### □ Objective Reformulation:

- Minimizing the AP risk

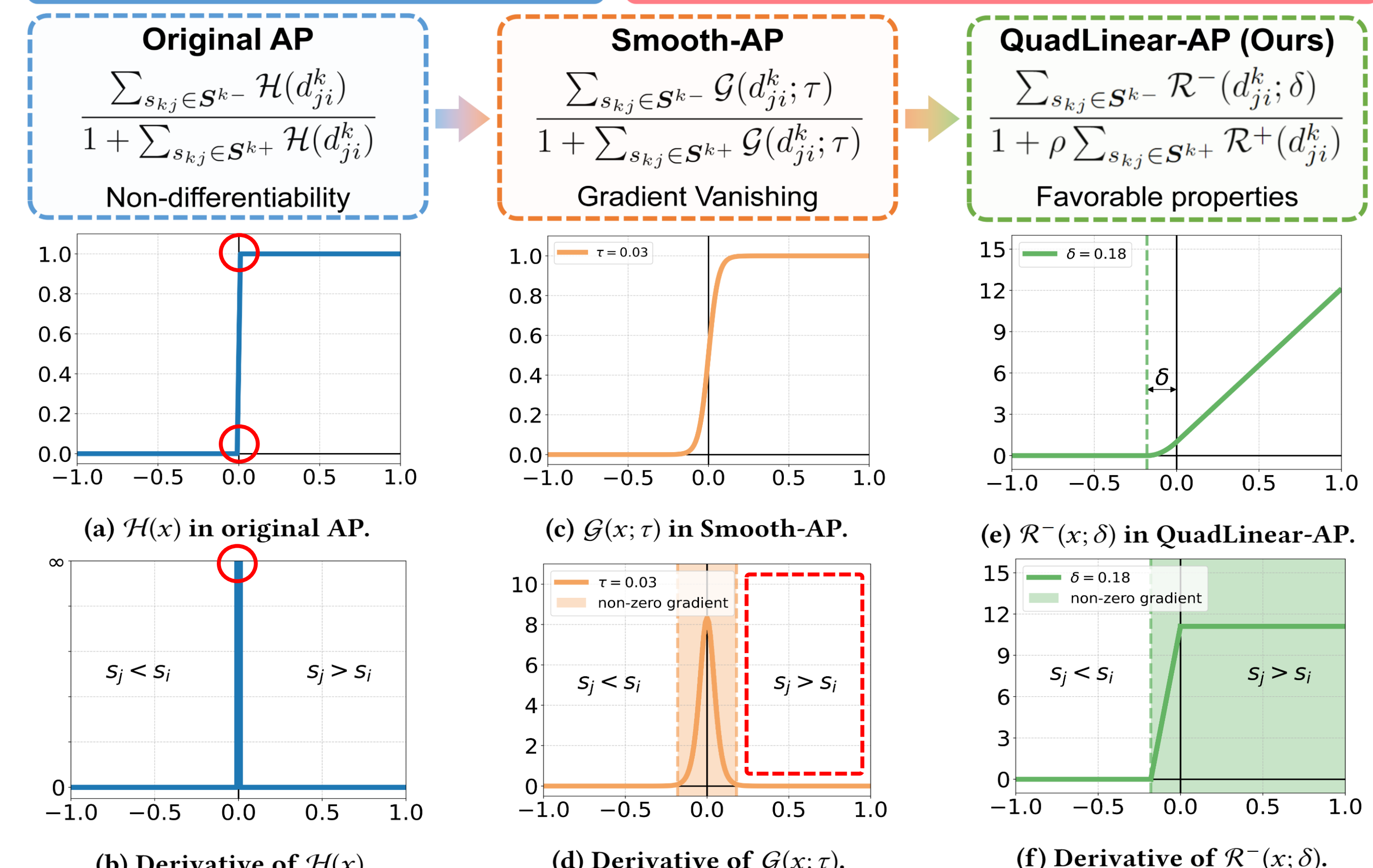
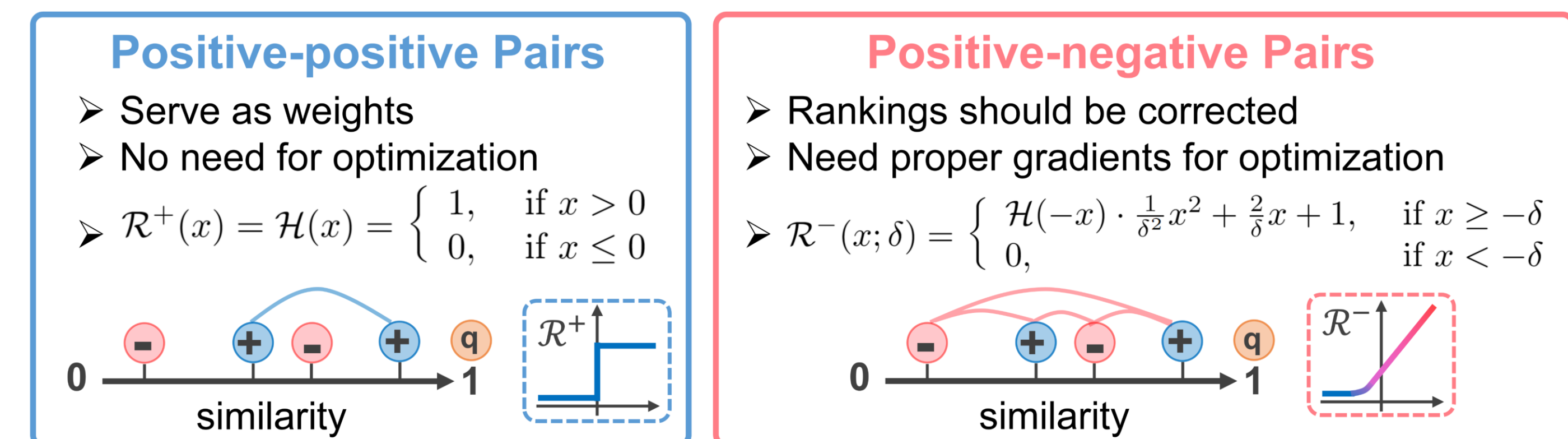
$$\min_f AP^\downarrow(f) = \frac{1}{N} \sum_{k=1}^N AP_k^\downarrow(f)$$

$$\mathcal{R}(s, S) = 1 + \sum_{s' \in S} \mathcal{H}(s' - s)$$

$$d_{ji}^k = s_{kj} - s_{ki} \quad h(x) = \frac{x}{1+x}$$

#### □ Rethinking the Components of AP Risk: QuadLinear-AP

$$AP_k^\downarrow(f) = \frac{1}{|S^{k+}|} \sum_{s_{ki} \in S^{k+}} h\left(\frac{\sum_{s_{kj} \in S^{k-}} \mathcal{H}(d_{ji}^k)}{1 + \sum_{s_{kj} \in S^{k+}} \mathcal{H}(d_{ji}^k)}\right) \Rightarrow \widehat{AP}_k^\downarrow(f) = \frac{1}{|S^{k+}|} \sum_{s_{ki} \in S^{k+}} h\left(\frac{\sum_{s_{kj} \in S^{k-}} \mathcal{R}^-(d_{ji}^k; \delta)}{1 + \rho \sum_{s_{kj} \in S^{k+}} \mathcal{R}^+(d_{ji}^k; \delta)}\right)$$



#### □ Favorable Properties of QuadLinear-AP

- Differentiable AP optimization
- Suitable gradients for low AP area
- Continuous, Smooth, and Convex
- Monotonically increasing (non-strictly)
- Upper bound of Heaviside function

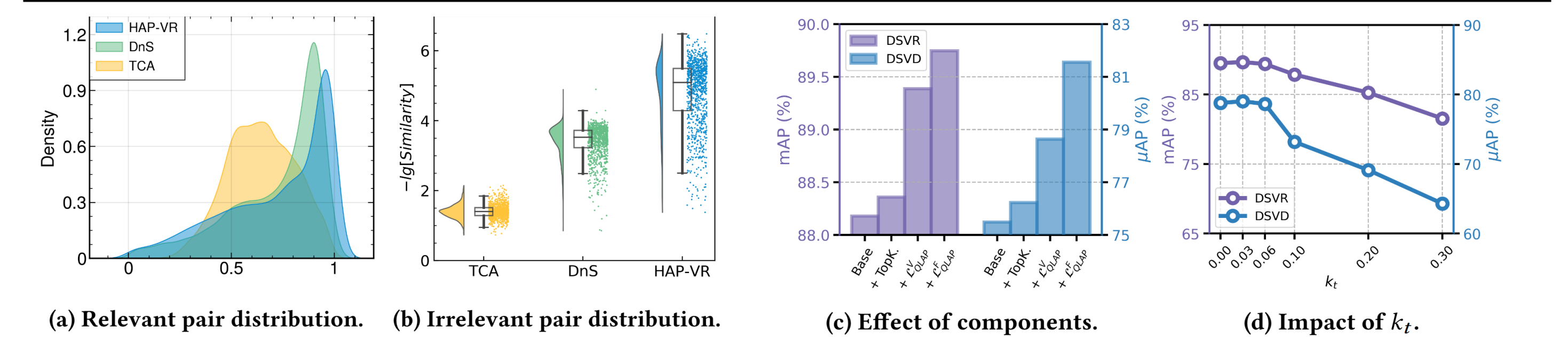
$$\mathcal{L}_{QLAP} = \frac{1}{N} \sum_{k=1}^N \widehat{AP}_k^\downarrow(f)$$

### Evaluation Results

#### □ Overall Performance of Our Proposed HAP-VR

- HAP-VR improves mAP and  $\mu$ AP effectively on multiple tasks

Method	Label	Trainset	Retrieval (mAP)					Detection ( $\mu$ AP)				
			EVVE	SVD	DSVR	CSVR	ISVR	EVVE	SVD	DSVD	CSVD	ISVD
DML <sup>†</sup> [32]	✓	VCDB (C&D)	61.10	85.00	52.80	51.40	44.00	75.50	/	39.00	36.50	30.00
TMK <sup>†</sup> [46]	✓	internal	61.80	86.30	52.40	50.70	42.50	/	/	/	/	/
TCA [53]	✓	VCDB (C&D)	63.08	<b>89.82</b>	86.81	82.31	69.61	76.90	56.93	69.09	62.28	49.24
ViSiL <sup>†</sup> [30]	✓	VCDB (C&D)	65.80	88.10	89.90	85.40	72.30	79.10	/	75.80	69.00	53.00
DnS (S <sub>a</sub> ) [34]	✓	DnS-100K	65.33	<b>90.20</b>	92.09	87.54	74.08	74.56	<b>72.24</b>	79.66	69.51	54.20
DnS (S <sub>b</sub> ) [34]	✓	DnS-100K	64.41	89.12	90.89	86.28	72.87	75.80	66.53	78.05	68.52	53.48
LAMV <sup>†</sup> [2]	✗	YFCC100M	62.00	88.00	61.90	58.70	47.90	80.60	/	55.40	50.00	38.80
VRL <sup>†</sup> [24]	✗	internal	/	/	90.00	85.80	70.90	/	/	/	/	/
ViSiL <sub>f</sub> <sup>†</sup> [30]	✗	ImageNet	62.70	/	89.00	84.80	72.10	74.60	/	66.90	59.50	45.90
S <sup>2</sup> VS [33]	✗	VCDB (D)	<b>67.17</b>	88.40	<b>92.53</b>	<b>87.73</b>	<b>74.51</b>	<b>80.72</b>	65.04	<b>86.12</b>	<b>77.41</b>	<b>63.26</b>
HAP-VR (Ours)	✗	VCDB (D)	<b>69.15</b>	89.00	<b>92.83</b>	<b>88.21</b>	<b>74.72</b>	<b>82.88</b>	<b>67.87</b>	<b>88.41</b>	<b>79.85</b>	<b>64.79</b>



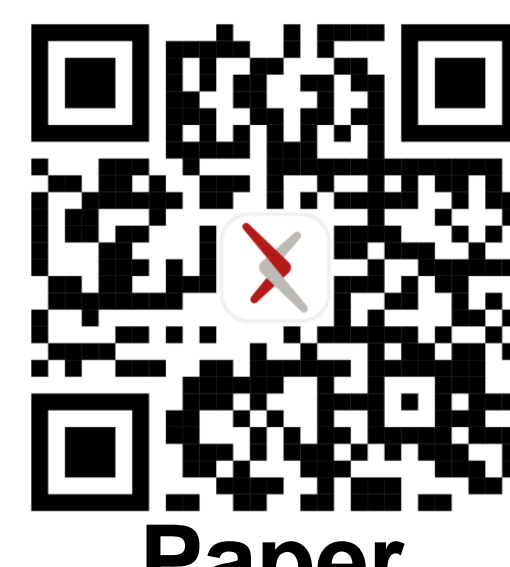
#### □ Effectiveness of Our Proposed QuadLinear-AP

- QuadLinear-AP outperforms previous pair-wise/AP-based losses

Losses	Retrieval (mAP)			Detection ( $\mu$ AP)		
	DSVR	CSVR	ISVR	DSVD	CSVD	ISVD
MAE	89.07	88.03	80.86	78.08	75.69	65.26
MSE	89.22	88.26	80.80	78.66	76.07	65.44
Contrastive [18]	88.67	88.09	80.97	75.12	74.23	67.41
Triplet [50]	88.11	87.77	<b>81.21</b>	72.94	73.18	69.23
Circle [55]	87.53	86.11	78.77	73.26	71.15	59.33
FastAP [6]	89.30	88.42	81.16	78.83	77.51	<b>69.95</b>
DIR [47]	<b>89.65</b>	<b>88.57</b>	80.64	78.50	76.22	65.42
BlackBox [45]	<b>89.70</b>	<b>88.55</b>	80.53	<b>80.07</b>	77.37	66.00
Smooth-AP [4]	89.36	88.33	80.73	79.85	<b>77.75</b>	68.42
QuadLinear-AP (Ours)	<b>90.80</b>	<b>89.68</b>	<b>81.31</b>	<b>82.92</b>	<b>80.03</b>	<b>71.45</b>

### Conclusions

- **Methodologically:** Propose a self-supervised framework (**HAP-VR**) for video retrieval to **bridge the gap** of the objective and metric.
- **Analytically:** Introduce a gradient-enhanced AP surrogate loss (**QuadLinear-AP**) and design a **hierarchical learning strategy** for **AP optimization** on both video and frame levels.
- **Empirically:** Experimental results demonstrate the effectiveness of our proposed framework on various video retrieval tasks.



Paper



Code



WeChat