

1 1. First Appendix

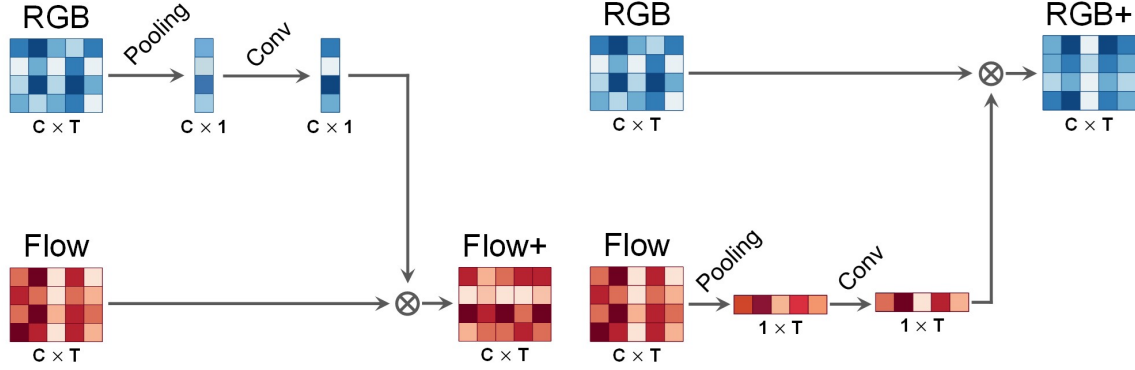


Figure 1: Diagram of the Dual-Modality Fusion (DMF), which consists of two parts: the left is RGB-generated channel attention for Flow enhancement, and the right is Flow-generated temporal attention for RGB enhancement.

2 1.1. Detail of Dual-Modality Fusion

3 When the DMF operation is stacked and used in subsequent fusion steps, it may lead to
 4 over-strengthening a small number of features and lead to over-fitting. Therefore, we set a
 5 hyper-parameter α to adjust the degree of DMF enhancement.

6 For the RGB and optical flow modalities of the video, we define two modality fusion
 7 strategies: *Concatenate Fusion* and *Mutual Fusion* (for simplicity, hereafter referred to as
 8 *concat* and *mutual*). *concat* can perform fusion between different modalities without in-
 9 formation loss, and it has been used in most of the mainstream WSTAL works at present
 10 Huang et al. (2022); Yang et al. (2022); Hong et al. (2021); He et al. (2022). While *mutual*
 11 enhances the implicit information shared between modalities before fusion, performing spe-
 12 cific tasks based on mutual fusion features can be regarded as double-checking between two
 13 modalities. This fusion strategy is referred to in a recent work Hong et al. (2021).

14 The *concat* fusion is denoted as:

$$concat = \mathcal{C}_{con}(cat(X_{rgb}, X_{flow})) \quad (1)$$

15 where \mathcal{C}_{con} is a 1×1 convolution for information interaction between channels, and *cat* rep-
 16 represents to concatenate X_{rgb} and X_{flow} along the channel dimension to achieve information
 17 fusion of the two modalities. And the *mutual* fusion is denoted as:

$$mutual = h_r(\mathcal{C}_r(X_{rgb})) + h_f(\mathcal{C}_f(X_{flow})) \quad (2)$$

18 where h_r and h_f denote to be the projectors for preserving the differences between modalities
 19 and avoiding their convergence to be completely the same.

20 For *mutual* fusion, we introduce contrastive learning for reinforcing the mutual infor-
 21 mation between modalities, i.e., to optimize the loss:

$$\mathcal{L}_{ml} = MSE(\mathcal{C}_r(X_{rgb}), \mathcal{C}_f(X_{flow})) \quad (3)$$

22 Combining the DMF operation and the fusion method, we produce various options for
 23 cross-modality fusion strategies. After experiments, our optimal fusion strategy is to use
 24 *mutual* fusion strategy with DMF feature enhancement as the input of T-RPN learning
 25 (Stage-1) and *concat* fusion without DMF enhancement as the input of action classification
 26 (Stage-2). We will make a detailed discussion of the reason why the heterogeneous feature
 27 fusion strategies in different stages of the Two-Stage Detection framework in Section 1.2.

Fusion Strategy	Stage-1				Stage-2				mAP@AVG
	RGB	Flow	Concat	Mutual	RGB	Flow	Concat	Mutual	
w/o Fusion	✓				✓				37.9
		✓				✓			41.4
w/ Fusion			✓				✓		44.4
			✓	✓			✓		45.5
				✓				✓	42.6
				✓				✓	44.3

Table 1: Performance comparison of using different fusion strategies on different stages.

Stage-1		Stage-2		mAP@AVG
Mutual	Mutual+	Concat	Concat+	
✓		✓		45.5
	✓	✓		47.0
✓			✓	45.9
	✓		✓	46.9

Table 2: Performance comparison with and without Cross-Modality Attention Network. + indicates that the input features are enhanced by DMF.

28 1.2. Fusion Strategy

29 Table 1 and Table 2 show the effect of the fusion mechanism of RGB and optical flow on
 30 the performance of WSTAL. According to the tabular results, first, we can observe that the
 31 model performance is much lower purely using single-modal information than the scheme
 32 using cross-modal fusion, indicating that both RGB and optical flow have irreplaceable
 33 value for WSTAL tasks. Besides, the motion information in the optical flow modality is
 34 relatively more important than the appearance information in the RGB modality.

35 Another important finding is that for Stage-1, the *mutual* fusion strategy is always
 36 better than the *concat* fusion strategy, while Stage-2 is just the opposite. We speculate

37 that the *mutual* fusion strategy can enhance the mutual information between modalities,
38 which can be regarded as the implicit expression of the action to be detected in the two
39 modalities. These mutual constraints between modalities make T-RPN screen out more
40 stringent foreground fragments, that is, higher-quality temporal region proposals. The
41 *concat* fusion strategy can be regarded as a union of two modalities, which provide sufficient
42 semantic information for action classification and benefit the recognition of Stage-2.

43 References

- 44 Bo He, Xitong Yang, Le Kang, Zhiyu Cheng, Xin Zhou, and Abhinav Shrivastava. Asm-loc:
45 action-aware segment modeling for weakly-supervised temporal action localization. In
46 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
47 pages 13925–13935, 2022.
- 48 Fa-Ting Hong, Jia-Chang Feng, Dan Xu, Ying Shan, and Wei-Shi Zheng. Cross-modal
49 consensus network for weakly supervised temporal action localization. In *Proceedings of*
50 *the 29th ACM International Conference on Multimedia*, pages 1591–1599, 2021.
- 51 Linjiang Huang, Liang Wang, and Hongsheng Li. Weakly supervised temporal action
52 localization via representative snippet knowledge propagation. In *Proceedings of the*
53 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3272–3281,
54 2022.
- 55 Zichen Yang, Jie Qin, and Di Huang. Acgnet: Action complement graph network for
56 weakly-supervised temporal action localization. In *Proceedings of the AAAI Conference*
57 *on Artificial Intelligence*, volume 36, pages 3090–3098, 2022.