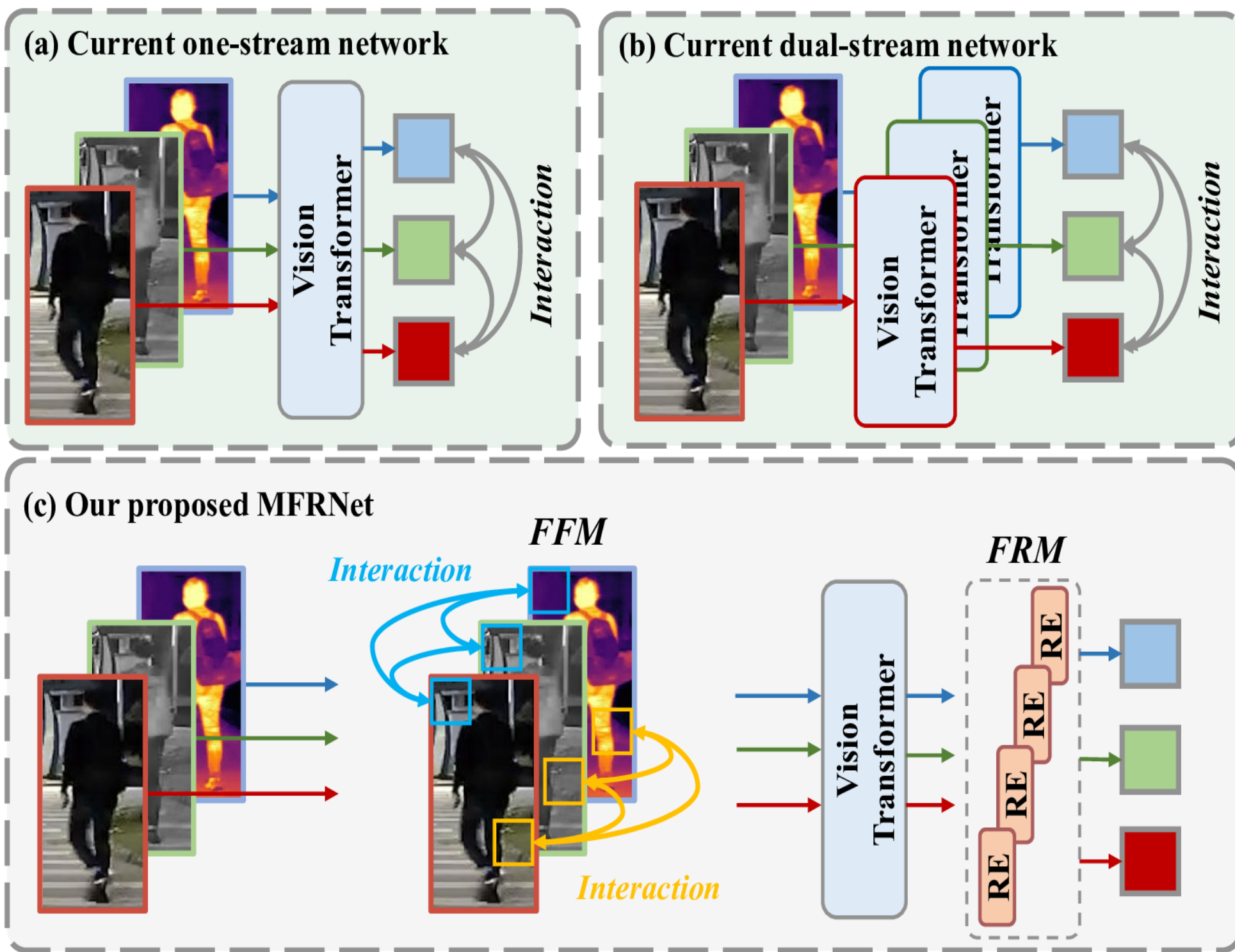


## Motivation



**Fig. 1: Motivation of proposed MFRNet.**

(a) and (b) show the current mainstream one-stream and two-stream networks, respectively. Besides the limitation of neglecting the pixel-level alignment characteristics of multi-modal images, current methods also face the challenge of balancing the modality-specific and modality-shared representation in a unified network.

(c) shows our proposed MFRNet and 'RE' refers to the representation expert. MFRNet inherits the idea of MoE and extends it with multi-modal fusion (FFM) and representation (FRM), allowing it to achieve fine-grained interaction and efficient representation for multi-modal data.

## Contribution

- (1) We propose a Modality Fusion and Representation Network (MFRNet) for multi-modal object re-identification, which inherits the idea of a sparse mixture of experts and extends it with multi-modal fusion and representation.
- (2) We introduce a Feature Fusion Module (FFM) and a Feature Representation Module (FRM). The former aims to achieve fine-grained interaction between multi-modal inputs, while the latter aims to achieve efficient and balanced feature extraction between modality-shared and modality-specific representations.
- (3) Extensive experiments on three public multi-spectral object ReID datasets including RGBNT201, RGBNT100, and MSVR310, verifying the superior performance of MFRNet.

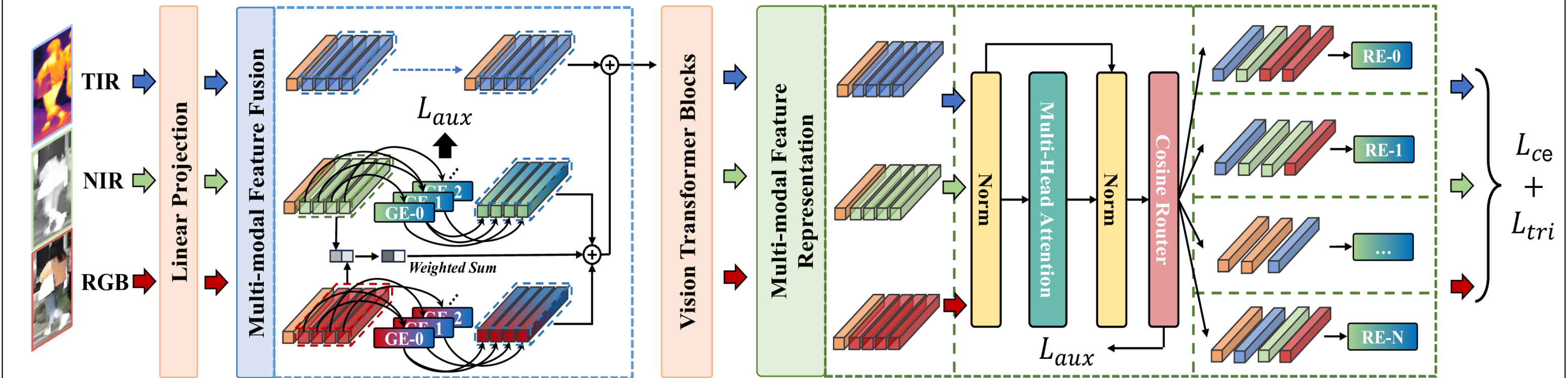
## Experiments

	Methods	mAP	R-1	R-5	R-10
Single	MUDeep (Qian et al., 2017)	23.8	19.7	33.1	44.3
	HACNN (Li et al., 2018)	21.3	19.0	34.1	42.8
	MLFN (Chang et al., 2018)	26.1	24.2	35.9	44.1
	PCB (Sun et al., 2018)	32.8	28.1	37.4	46.9
	OSNet (Zhou et al., 2019)	25.4	22.3	35.1	44.7
	CAL (Rao et al., 2021)	27.6	24.3	36.5	45.7
	HAMNet (Li et al., 2020)	27.7	26.3	41.5	51.7
Multi	PFNet (Zheng et al., 2021)	38.5	38.9	52.0	58.4
	IEEE (Wang et al., 2022)	47.5	44.4	57.1	63.6
	DENet (Zheng et al., 2023)	42.4	42.2	55.3	64.5
	UniCat (Crawford et al., 2023)	57.0	55.7	-	-
	HTT (Wang et al., 2024c)	71.1	73.4	83.1	87.3
	EDITOR (Zhang et al., 2024a)	66.5	68.3	81.1	88.2
	RSCNet (Yu et al., 2024)	68.2	72.5	-	-
	TOP-ReID (Wang et al., 2024a)	<u>72.3</u>	<u>76.6</u>	<u>84.7</u>	<u>89.4</u>
	Ours	<b>80.7</b>	<b>83.6</b>	<b>91.9</b>	<b>94.1</b>
	Methods	Params(M)	Flops(G)		
	HTT (Wang et al., 2024c)	85.6	33.1		
	EDITOR (Zhang et al., 2024a)	117.5	38.6		
	TOP-ReID (Wang et al., 2024a)	278.2	34.5		
	Ours	<b>57.1</b>	<b>22.1</b>		

	Methods	M (RGB)	M (NIR)	M (TIR)	M (RGB+NIR)	M (RGB+TIR)	M (NIR+TIR)	Average
Single	MUDeep (Qian et al., 2017)	19.2	16.4	20.0	17.2	18.4	14.2	13.7
	HACNN (Li et al., 2018)	12.5	11.1	20.5	19.4	16.7	13.3	9.2
	MLFN (Chang et al., 2018)	20.2	18.9	21.1	19.7	17.6	11.1	13.2
	PCB (Sun et al., 2018)	23.6	24.2	24.4	25.1	19.9	14.7	20.6
	OSNet (Zhou et al., 2019)	19.8	17.3	21.0	19.0	18.7	14.6	12.3
	PFNet (Zheng et al., 2021)	-	-	31.9	29.8	25.5	25.8	-
	DENet (Zheng et al., 2023)	-	-	35.4	36.8	33.0	35.4	-
Multi	TOP-ReID (Wang et al., 2024a)	<u>54.4</u>	<u>57.5</u>	<u>64.3</u>	<u>67.6</u>	<u>51.9</u>	<u>54.5</u>	<u>35.3</u>
	Ours	<b>64.7</b>	<b>65.2</b>	<b>72.3</b>	<b>76.1</b>	<b>51.6</b>	<b>49.5</b>	<b>41.4</b>

	Methods	RGBNT100	MSVR310
Single	PCB (Sun et al., 2018)	57.2	83.5
	MGN (Wang et al., 2018)	58.1	83.1
	DMML (Chen et al., 2019)	58.5	82.0
	BoT (Luo et al., 2019)	78.0	95.1
	OSNet (Zhou et al., 2019)	75.0	95.6
	Circle Loss (Sun et al., 2020)	59.4	81.7
	HRCN (Zhao et al., 2021)	67.1	91.8
Multi	TransReID (He et al., 2021)	75.6	92.9
	AGW (Ye et al., 2022)	73.1	92.7
	HAMNet (Li et al., 2020)	74.5	93.3
	PFNet (Zheng et al., 2021)	68.1	94.1
	GAFNet (Guo et al., 2022)	74.4	93.4
	Graft (Yin et al., 2023)	76.6	94.3
	GPNet (He et al., 2023)	75.0	94.5
	PHT (Pan et al., 2023)	79.9	92.7
	UniCat (Crawford et al., 2023)	79.4	96.2
	CCNet (Zheng et al., 2022)	77.2	96.3
	HTT (Wang et al., 2024c)	75.7	92.6
	TOP-ReID (Wang et al., 2024a)	81.2	96.4
	EDITOR (Zhang et al., 2024a)	82.1	96.4
	RSCNet (Yu et al., 2024)	<u>82.3</u>	<u>96.6</u>
	Ours	<b>88.2</b>	<b>97.4</b>

## The Proposed MFRNet

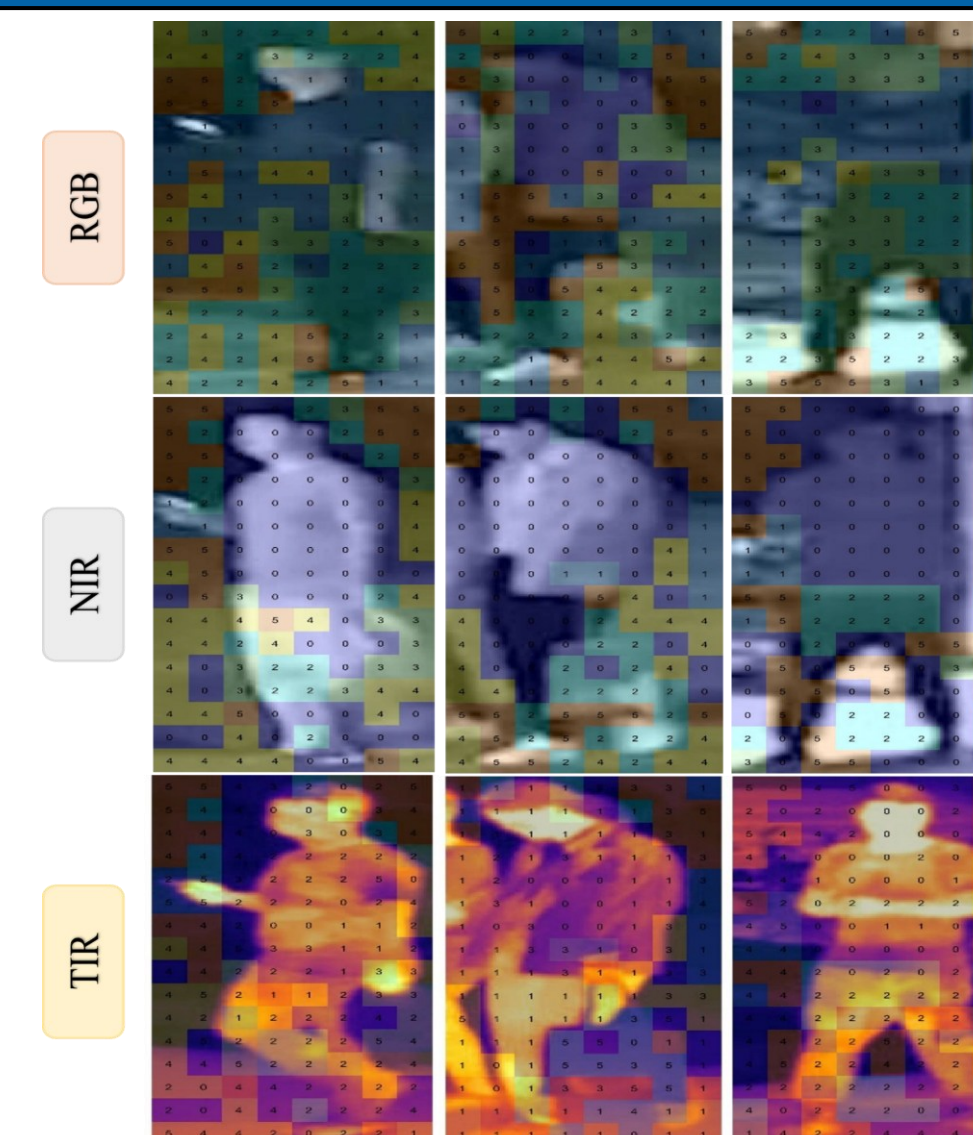


**Fig. 2: Overall architecture of MFRNet.**

MFRNet is built from the basic Transformer blocks inherited from vanilla ViT, with two new modules (feature fusion module and feature representation module) added to adapt to multi-modal object re-identification tasks.

In the blue box, the Feature Fusion Module (FFM) employs multiple simple generators to adaptively provide fine-grained interaction information, while in the green box, the Feature Representation Module (FRM) employs diverse representation experts to extract and combine modality-specific and modality-shared features. Here, GE refers to the generation experts in FFM, while RE represents the representation experts in FRM. Notably, the figure illustrates the NIR+RGB→TIR interaction, with the same interaction applying to NIR and RGB.

## Visualization Expert



**Fig. 3: Visualization Result.**

It can be observed that Expert 0 concentrates on extracting pedestrian features in the NIR modality, while Experts 1, 2, and 3 focus on pedestrian regions in the RGB and TIR modalities. Additionally, Experts 4 and 5 tend to extract background features across all modalities. As our initial hypothesis predicted, semantically similar content achieves knowledge sharing through selecting analogous experts, while modality-specific representations that resist fusion maintain their distinct characteristics via dedicated expert allocation.